

Multivariate Statistics and Machine Learning

Korbinian Strimmer

18 January 2024

Contents

Welcome	5
License	5
Preface	7
About the author	7
About the module	7
Prerequisites	8
Acknowledgements	8
1 Multivariate random variables	9
1.1 Essentials in multivariate statistics	9
1.2 Multivariate distributions	16
1.3 Multivariate normal distribution	17
2 Multivariate estimation in large sample and small sample settings	27
2.1 Overview	27
2.2 Empirical estimates	28
2.3 Maximum likelihood estimation	30
2.4 Sampling distribution of the empirical / maximum likelihood estimates	32
2.5 Small sample estimation	33
2.6 Full Bayesian multivariate modelling	37
2.7 Conclusion	39
3 Transformations and dimension reduction	41
3.1 Linear Transformations	41
3.2 Nonlinear transformations	45
3.3 General whitening transformations	48
3.4 Natural whitening procedures	56
3.5 Principal Component Analysis (PCA)	69
4 Unsupervised learning and clustering	75
4.1 Challenges in unsupervised learning	75

4.2	Hierarchical clustering	78
4.3	K-means clustering	84
4.4	Mixture models	90
4.5	Fitting mixture models to data and inferring the latent states . . .	95
4.6	Application of Gaussian mixture models	97
4.7	The EM algorithm	100
5	Supervised learning and classification	109
5.1	Aims of supervised learning	109
5.2	Bayesian discriminant rule or Bayes classifier	111
5.3	Normal Bayes classifier	112
5.4	The training step — learning QDA, LDA and DDA classifiers from data	115
5.5	Quantifying prediction error	119
5.6	Goodness of fit and variable ranking	120
5.7	Variable selection	123
6	Multivariate dependencies	127
6.1	Measuring the linear association between two sets of random variables	127
6.2	Canonical Correlation Analysis (CCA) aka CCA whitening	128
6.3	Vector correlation and RV coefficient	130
6.4	Limits of linear models and correlation	132
6.5	Mutual information as generalisation of correlation	134
6.6	Graphical models	140
7	Nonlinear and nonparametric models	147
7.1	Random forests	148
7.2	Gaussian processes	150
7.3	Neural networks	157
A	Further study	161
A.1	Recommended reading	161
A.2	Advanced reading	161
	Bibliography	163

Welcome

These are the lecture notes for MATH38161, a course in **Multivariate Statistics and Machine Learning** for third year mathematics students at the [Department of Mathematics of the University of Manchester](#).

The course text was written by [Korbinian Strimmer](#) from 2018–2024. This version is from 18 January 2024.

The notes will be updated from time to time. To view the current version visit the [online MATH38161 lecture notes](#).

You may also wish to download the MATH38161 lecture notes as [PDF in A4 format for printing](#) (double page layout) or as [6x9 inch PDF for use on tablets](#) (single page layout).

License

These notes are licensed to you under [Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License](#).

Preface

About the author

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

I have first taught this module in the winter semester 2018 at the University of Manchester, and subsequently all the following years until (including) the winter semester 2023.

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

If you are a University of Manchester student and enrolled in this module you will find on [Blackboard](#):

- a weekly learning plan for an 11 week study period,
- weekly worksheets with with examples (theory and application in R) and solutions in R Markdown, and
- exam papers of previous years.

Furthermore, there is also an [MATH38161 online reading list](#) hosted by the University of Manchester library.

Prerequisites

Multivariate statistics relies heavily on matrix algebra and vector and matrix calculus. For a refresher of the essentials please refer to the supplementary [Matrix and Calculus Refresher notes](#).

Furthermore, this module builds on earlier statistics modules, especially on likelihood estimation and Bayesian statistics as discussed, e.g., in the module [MATH27720 Statistics 2](#).

For an overview of essential probability distributions see the [Probability and Distributions Refresher notes](#).

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 Essentials in multivariate statistics

1.1.1 Why multivariate statistics?

In science we use experiments to learn about underlying mechanisms of interest, both deterministic and stochastic, to compare different models and to verify or reject hypotheses about the world. Statistics provides tools to quantify this procedure and offers methods to link data (experiments) with probabilistic models (hypotheses).

In univariate statistics with we use relatively simple approaches based on a single random variable or single parameter. However, in practise we often have to consider multiple random variables and multiple parameters, so we need more complex models and also be able to deal with more complex data. Hence, the need for multivariate statistical approaches and models.

Specifically, multivariate 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 modelling but nowadays most machine learning methods are fully based on statistical multivariate approaches, so the two fields are converging.

Multivariate models provide a means to learn dependencies and interactions among the components of the random variables which in turn allow us to draw conclusion about underlying mechanisms of interest (e.g. in biological or medical problems).

Two main tasks:

- unsupervised learning (finding structure, clustering)
- supervised learning (training from labelled 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

1.1.2 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 .

1.1.3 Multivariate data

Vector notation:

Samples from a multivariate distribution are *vectors* (not scalars as for univariate normal):

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \stackrel{\text{iid}}{\sim} F$$

Matrix and component notation:

All the data points are commonly collected into a matrix \mathbf{X} .

In statistics the convention is to store each data vector in the rows of the data matrix \mathbf{X} :

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^T = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & & & \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{pmatrix}$$

Therefore,

$$\mathbf{x}_1 = \begin{pmatrix} x_{11} \\ \vdots \\ x_{1d} \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} x_{21} \\ \vdots \\ x_{2d} \end{pmatrix}, \dots, \mathbf{x}_n = \begin{pmatrix} x_{n1} \\ \vdots \\ x_{nd} \end{pmatrix}$$

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 **statistics convention on data matrices** is *not* universal! In fact, in most of the machine learning literature in engineering and computer science the data samples are stored in the columns so that the variables appear in the rows (thus in the engineering convention the data matrix is transposed compared to the statistics convention).

In order to avoid confusion and ambiguity it is recommended to prefer vector notation to describe data over matrix or component notation (see also the section below on estimating covariance matrices for examples).

1.1.4 Mean of a random vector

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

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

1.1.5 Variance of a random vector

Recall the definition of mean and variance for a univariate random variable:

$$E(x) = \mu$$

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

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

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

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

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

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

1.1.6 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 \mathbf{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 separate entries does Σ have?

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

with $\sigma_{ij} = \sigma_{ji}$.

Number of separate entries: $\frac{d(d+1)}{2}$.

This number grows with the square of the 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 using standard algorithms such as [Gauss Jordan elimination](#).¹ Hence, computing Σ^{-1} is computationally expensive for large d !

1.1.7 Eigenvalue decomposition of Σ

Recall from linear matrix algebra that any real symmetric matrix has real eigenvalues and a complete set of orthogonal eigenvectors. These can be obtained by orthogonal eigendecomposition.

Applying eigenvalue decomposition to the covariance matrix yields

$$\Sigma = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$$

where \mathbf{U} is an orthogonal matrix containing the eigenvectors of the covariance matrix and

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

contains the corresponding eigenvalues λ_i .

Importantly, the eigenvalues of a covariance matrix are not only real-valued but are by construction further constrained to be non-negative. This can be seen by computing the quadratic form $\mathbf{z}^T \Sigma \mathbf{z}$ where \mathbf{z} is a non-random vector. For any non-zero \mathbf{z}

$$\begin{aligned} \mathbf{z}^T \Sigma \mathbf{z} &= \mathbf{z}^T \mathbb{E} \left((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \right) \mathbf{z} \\ &= \mathbb{E} \left(\mathbf{z}^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{z} \right) \\ &= \mathbb{E} \left(\left(\mathbf{z}^T (\mathbf{x} - \boldsymbol{\mu}) \right)^2 \right) \geq 0. \end{aligned}$$

Furthermore, with $\mathbf{y} = \mathbf{U}^T \mathbf{z}$ we get

$$\begin{aligned} \mathbf{z}^T \Sigma \mathbf{z} &= \mathbf{z}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \mathbf{z} \\ &= \mathbf{y}^T \mathbf{\Lambda} \mathbf{y} = \sum_{i=1}^d y_i^2 \lambda_i \end{aligned}$$

and hence all the $\lambda_i \geq 0$. Therefore the **covariance matrix Σ is always positive semi-definite**.

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

¹Specialised matrix algorithms improve this to about $O(d^{2.373})$. Matrices with special symmetries (e.g. diagonal and block diagonal matrices) or properties (e.g. orthogonal matrix) can also be inverted much easier.

1.1.8 Joint covariance matrix

Assume we have random vector \mathbf{z} with mean $E(\mathbf{z}) = \boldsymbol{\mu}_z$ and covariance matrix $\text{Var}(\mathbf{z}) = \boldsymbol{\Sigma}_z$.

Often it makes sense to partition the components of \mathbf{z} into two groups

$$\mathbf{z} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}$$

This induces a corresponding partition in the expectation

$$\boldsymbol{\mu}_z = \begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix}$$

where $E(\mathbf{x}) = \boldsymbol{\mu}_x$ and $E(\mathbf{y}) = \boldsymbol{\mu}_y$.

Furthermore, the joint covariance matrix for \mathbf{x} and \mathbf{y} can then be written as

$$\boldsymbol{\Sigma}_z = \begin{pmatrix} \boldsymbol{\Sigma}_x & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_y \end{pmatrix}$$

It contains the within-group covariance matrices $\boldsymbol{\Sigma}_x$ and $\boldsymbol{\Sigma}_y$ as diagonal elements and the cross-covariance matrix $\boldsymbol{\Sigma}_{xy} = \boldsymbol{\Sigma}_{yx}^T$ as off-diagonal element.

Note that the cross-covariance matrix $\boldsymbol{\Sigma}_{xy}$ is rectangular and not symmetric. We also write $\text{Cov}(\mathbf{x}, \mathbf{y}) = \boldsymbol{\Sigma}_{xy}$ and we can define cross-covariance directly by

$$\text{Cov}(\mathbf{x}, \mathbf{y}) = E((\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^T) = E(\mathbf{x}\mathbf{y}^T) - \boldsymbol{\mu}_x\boldsymbol{\mu}_y^T$$

1.1.9 Quantities related to the covariance matrix

1.1.9.1 Correlation matrix \mathbf{P}

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

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

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

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

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

Note the **variance-correlation decomposition**

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

where \mathbf{V} is a diagonal matrix containing the variances:

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

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

This is the definition of correlation written in matrix notation.

As with the covariance matrix, in many applications it makes sense to partition a joint correlation matrix

$$\mathbf{P}_z = \begin{pmatrix} \mathbf{P}_x & \mathbf{P}_{xy} \\ \mathbf{P}_{yx} & \mathbf{P}_y \end{pmatrix}$$

into within-group correlation matrices

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

and

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

and the cross-correlation matrix

$$\mathbf{P}_{xy} = \mathbf{V}_x^{-\frac{1}{2}} \boldsymbol{\Sigma}_{xy} \mathbf{V}_y^{-\frac{1}{2}}$$

with

$$\mathbf{P}_{yx} = \mathbf{P}_{xy}^T = \mathbf{V}_y^{-\frac{1}{2}} \boldsymbol{\Sigma}_{yx} \mathbf{V}_x^{-\frac{1}{2}}.$$

1.1.9.2 Precision matrix or concentration matrix

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

$\boldsymbol{\Omega}$ (“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 :

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

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} :

- 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).
- Σ^{-1} is a natural parameter from an exponential family perspective.

1.1.9.3 Partial correlation matrix

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

1.1.9.4 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:** $\det(\Sigma) = \prod_{i=1}^d \lambda_i$

The generalised variance $\det(\Sigma)$ is also known as the volume of Σ .

1.2 Multivariate distributions

1.2.1 Common distributions

In multivariate statistics we make use of multivariate distributions. These are typically generalisations of corresponding univariate distribution.

Among the most commonly used multivariate distributions are:

- The **multivariate normal distribution** $N_d(\mu, \Sigma)$ as a generalisation of univariate normal distribution $N(\mu, \sigma^2)$
- The **categorical distribution** $\text{Cat}(\pi)$ as a generalisation of the Bernoulli distribution $\text{Ber}(\theta)$
- The **multinomial distribution** $\text{Mult}(n, \pi)$ as a generalisation of binomial distribution $\text{Bin}(n, \theta)$

The above distribution have already been introduced earlier in [MATH27720 Statistics 2](#).

Conceptually, these multivariate generalisation work behave exactly the same as their univariate counterparts and are employed in the same settings.

1.2.2 Further multivariate distributions

For multivariate Bayesian analysis we also need to consider a number of further multivariate distributions:

- The **Dirichlet distribution** $\text{Dir}(\alpha)$ as the generalisation of the beta distribution $\text{Beta}(\alpha, \beta)$,
- The **Wishart distribution** as the generalisation of the gamma distribution $\text{Gam}(\alpha, \theta)$,
- The **inverse Wishart distribution** as the generalisation of the inverse gamma distribution $\text{Inv-Gam}(\alpha, \beta)$.

For technical details of the densities etc. of the multivariate distribution families we refer to the supplementary [Probability and Distribution refresher notes](#).

1.3 Multivariate normal distribution

The multivariate normal distribution is ubiquitous in multivariate statistics and hence it is important to discuss it in more detail.

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 $d = 1$

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

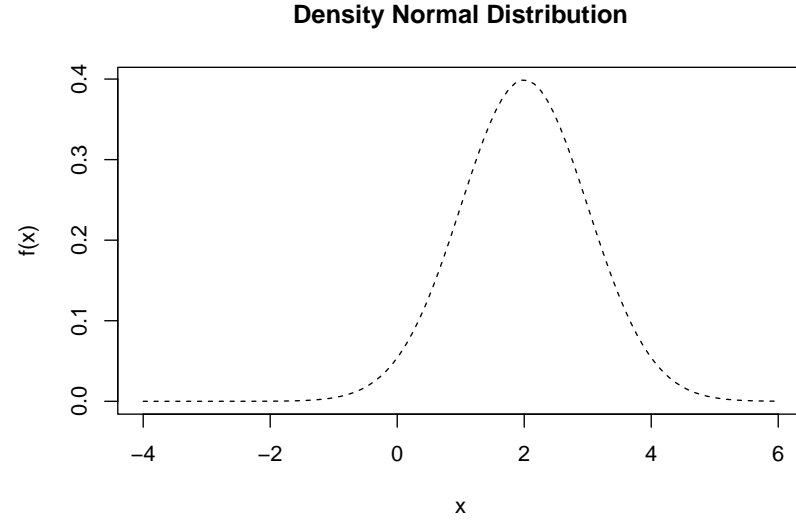
$$E(x) = \mu, \text{Var}(x) = \sigma^2$$

Probability Density Function:

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

Plot of univariate normal density :

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



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

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

Differential entropy:

$$H(F) = \frac{1}{2}(\log(2\pi\sigma^2) + 1)$$

Cross-entropy:

$$H(F_{\text{ref}}, F) = \frac{1}{2} \left(\frac{(\mu - \mu_{\text{ref}})^2}{\sigma^2} + \frac{\sigma_{\text{ref}}^2}{\sigma^2} + \log(2\pi\sigma^2) \right)$$

KL divergence:

$$D_{\text{KL}}(F_{\text{ref}}, F) = H(F_{\text{ref}}, F) - H(F_{\text{ref}}) = \frac{1}{2} \left(\frac{(\mu - \mu_{\text{ref}})^2}{\sigma^2} + \frac{\sigma_{\text{ref}}^2}{\sigma^2} - \log\left(\frac{\sigma_{\text{ref}}^2}{\sigma^2}\right) - 1 \right)$$

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 (und 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 model.

1.3.2 Multivariate normal model

Dimension d

$$\mathbf{x} \sim N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

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

Density:

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

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

- the density contains the precision matrix $\boldsymbol{\Sigma}^{-1}$
- to invert the covariance matrix $\boldsymbol{\Sigma}$ we need to invert its eigenvalues λ_i (hence we require that all $\lambda_i > 0$)
- the density also contains $\det(\boldsymbol{\Sigma}) = \prod_{i=1}^d \lambda_i \equiv$ product of the eigenvalues of $\boldsymbol{\Sigma}$
- note that $\det(2\pi\boldsymbol{\Sigma})^{-\frac{1}{2}} = \det(2\pi\mathbf{I}_d)^{-\frac{1}{2}} \det(\boldsymbol{\Sigma})^{-\frac{1}{2}} = (2\pi)^{-d/2} \det(\boldsymbol{\Sigma})^{-\frac{1}{2}}$

Special case: **standard multivariate normal** with

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

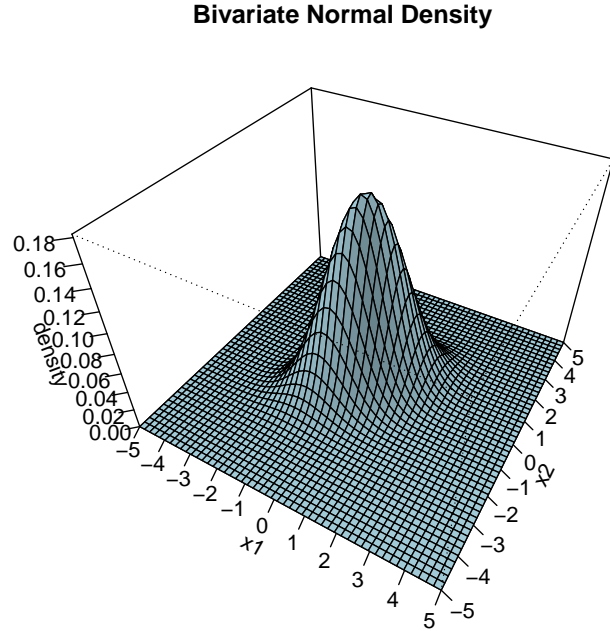
$$f(\mathbf{x}|\boldsymbol{\mu} = \mathbf{0}, \boldsymbol{\Sigma} = \mathbf{I}) = (2\pi)^{-d/2} \exp\left(-\frac{1}{2}\mathbf{x}^T\mathbf{x}\right) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x_i^2}{2}\right)$$

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

Misc:

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

Plot of the multivariate normal density:



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

An interactive [R Shiny web app](https://minerva.it.manchester.ac.uk/shiny/strimmer/bvn/) of the bivariate normal density plot is available online at <https://minerva.it.manchester.ac.uk/shiny/strimmer/bvn/>.

Differential entropy:

$$H = \frac{1}{2}(\log \det(2\pi\Sigma) + d)$$

Cross-entropy:

$$H(F_{\text{ref}}, F) = \frac{1}{2} \left\{ (\mu - \mu_{\text{ref}})^T \Sigma^{-1} (\mu - \mu_{\text{ref}}) + \text{Tr}(\Sigma^{-1} \Sigma_{\text{ref}}) + \log \det(2\pi\Sigma) \right\}$$

KL divergence:

$$\begin{aligned} D_{\text{KL}}(F_{\text{ref}}, F) &= H(F_{\text{ref}}, F) - H(F_{\text{ref}}) \\ &= \frac{1}{2} \left\{ (\mu - \mu_{\text{ref}})^T \Sigma^{-1} (\mu - \mu_{\text{ref}}) + \text{Tr}(\Sigma^{-1} \Sigma_{\text{ref}}) - \log \det(\Sigma^{-1} \Sigma_{\text{ref}}) - d \right\} \end{aligned}$$

1.3.3 Shape of the multivariate normal density

Now we 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 μ the condition for an ellipse is

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

where \mathbf{Q} is an orthogonal matrix whose column vectors indicate the direction of the axes. These are also called the **principal axes** of the ellipse.

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 $(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) = c$ where c is a constant. Using the eigenvalue decomposition of $\Sigma = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ we can rewrite this condition as

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

This implies that

- the contour lines of the multivariate normal density are indeed ellipses,
- the squared radii of the ellipse are proportional to the eigenvalues of Σ and
- the direction of the axes of the ellipse are given the principal axes and correspond to the column vectors in \mathbf{U} (i.e. the eigenvectors of Σ).

Equivalently, the positive square roots of the eigenvalues are proportional to the radii of the ellipse. Hence, for a singular covariance matrix with one or more $\lambda_i = 0$ the corresponding radii are zero.

An interactive [R Shiny web app](https://minerva.it.manchester.ac.uk/shiny/striimmer/bvn/) to play with the contour lines of the bivariate normal distribution is online at <https://minerva.it.manchester.ac.uk/shiny/striimmer/bvn/>.

1.3.4 Three types of covariances

Following the above we can parameterise a covariance matrix in terms of its i) volume, ii) shape and iii) orientation by writing

$$\Sigma = \kappa \mathbf{U} \mathbf{A} \mathbf{U}^T = \mathbf{U} (\kappa \mathbf{A}) \mathbf{U}^T$$

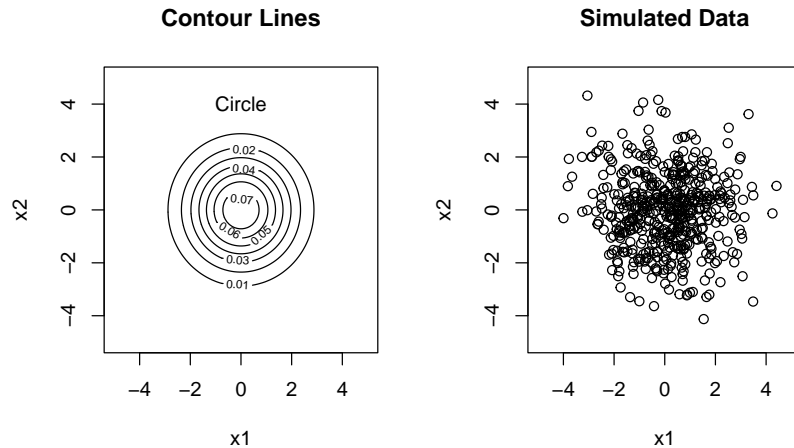
with $\mathbf{A} = \text{Diag}(a_1, \dots, a_d)$ and $\det(\mathbf{A}) = \prod_{i=1}^d a_i = 1$. Note that in this parameterisation the eigenvalues of Σ are $\lambda_i = \kappa a_i$.

- i) The **volume** is $\det(\Sigma) = \kappa^d$, determined by a single parameter κ . This parameter can be interpreted as the length of the side of a d -dimensional hypercube.
- ii) The **shape** is determined by the diagonal matrix \mathbf{A} with $d-1$ free parameters. Note that there are only $d-1$ and not d free parameters because of the constraint $\det(\mathbf{A}) = 1$.
- iii) The **orientation** is given by the orthogonal matrix \mathbf{U} , with $d(d-1)/2$ free parameters.

This leads to classification of covariances into three varieties:

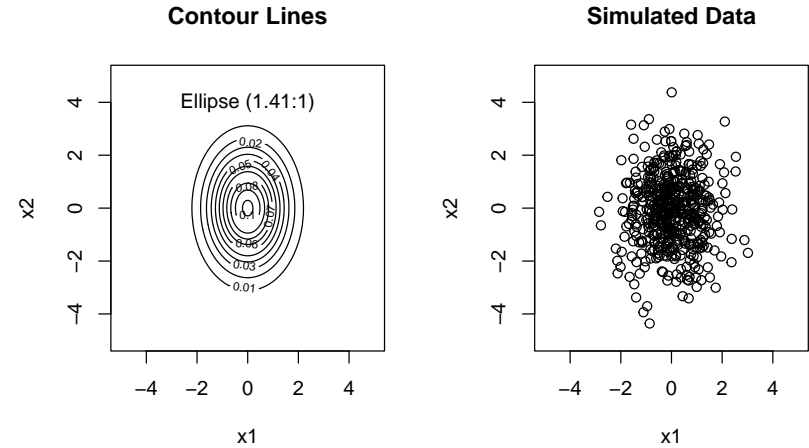
Type 1: spherical covariance $\Sigma = \kappa \mathbf{I}$, with spherical contour lines, 1 free parameter ($\mathbf{A} = \mathbf{I}$, $\mathbf{U} = \mathbf{I}$).

Example: $\Sigma = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$ with $\sqrt{\lambda_1/\lambda_2} = 1$:



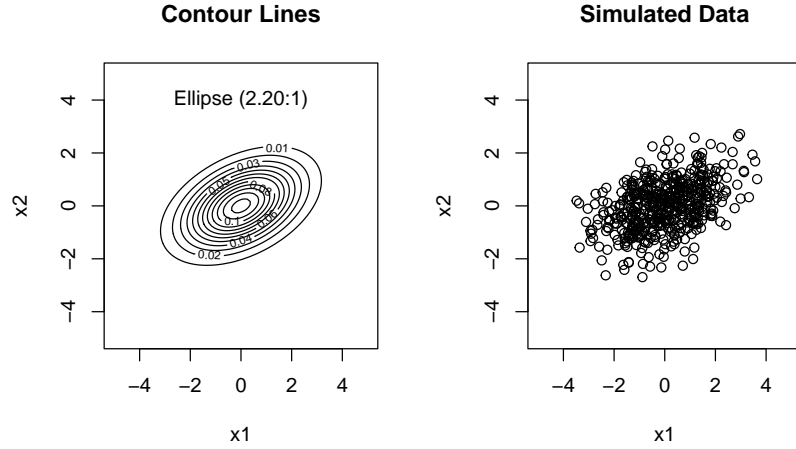
Type 2: diagonal covariance $\Sigma = \kappa \mathbf{A}$, with elliptical contour lines and the axes of the ellipse oriented parallel to the coordinates, d free parameters ($\mathbf{U} = \mathbf{I}$).

Example: $\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$ with $\sqrt{\lambda_1/\lambda_2} \approx 1.41$:



Type 3: general unrestricted covariance Σ , with elliptical contour lines, with the axes of the ellipse oriented according to the column vectors in \mathbf{U} , $d(d+1)/2$ free parameters.

Example: $\Sigma = \begin{pmatrix} 2 & 0.6 \\ 0.6 & 1 \end{pmatrix}$ with $\sqrt{\lambda_1/\lambda_2} \approx 2.20$:



1.3.5 Concentration of probability mass for small and large dimension

The density of the multivariate normal distribution has a bell shape with a single mode. Intuitively, we may assume that most of the probability mass is always concentrated around this mode, as it is in the univariate case ($d = 1$). While this is still true for small dimensions (small d) we now show that this intuition is incorrect for high dimensions (large d).

For simplicity we consider the standard multivariate normal distribution with dimension d

$$\mathbf{x} \sim N_d(\mathbf{0}, \mathbf{I}_d)$$

with a spherical covariance \mathbf{I}_d and sample \mathbf{x} . The squared Euclidean length of \mathbf{x} is $r^2 = \|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x} = \sum_{i=1}^d x_i^2$. The corresponding density of the d -dimensional standard multivariate normal distribution is

$$g_d(\mathbf{x}) = (2\pi)^{-d/2} e^{-\mathbf{x}^T \mathbf{x} / 2}$$

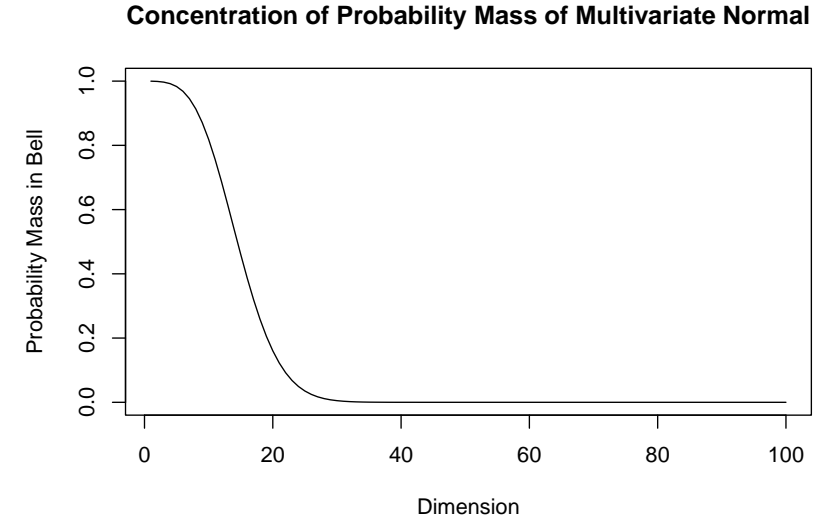
A natural way to define the main part of the “bell” of the standard multivariate normal as the set of all \mathbf{x} for which the density is larger than a specified fraction η (say 0.001) of the maximum value of the density $g_d(0)$ at the peak at zero. To formalise

$$B = \left\{ \mathbf{x} : \frac{g_d(\mathbf{x})}{g_d(0)} > \eta \right\}$$

which can be equivalently written as the set

$$B = \{ \mathbf{x} : \mathbf{x}^T \mathbf{x} = r^2 < -2 \log(\eta) = r_{\max}^2 \}$$

Each individual component in the sample \mathbf{x} is independently distributed as $x_i \sim N(0, 1)$, hence $r^2 \sim \chi_d^2$ is chi-squared distributed with degree of freedom d . The probability $\Pr(\mathbf{x} \in B)$ can thus be obtained as the value of the cumulative density function of a chi-squared distribution with d degrees of freedom at r_{\max}^2 . Computing this probability for fixed η as a function of the dimension d we obtain the following curve:



The above plot is for $\eta = 0.001$. You can see that for dimensions up to around $d = 10$ the probability mass is indeed concentrated in the center of the distribution but from $d = 30$ onwards it has moved completely to the tails.

Chapter 2

Multivariate estimation in large sample and small sample settings

2.1 Overview

In practical application of multivariate normal model we need to learn its parameters from observed data points:

$$x_1, x_2, \dots, x_n \stackrel{\text{iid}}{\sim} F_\theta$$

We first consider the case when there are many measurements available (n large), and then subsequently the case when the number of data points n is small compared to the dimensions and the number of parameters.

In a previous course in year 2 (see [MATH27720 Statistics 2](#)) the method of maximum likelihood as well as the essentials of Bayesian statistics were introduced. Below we apply these approaches to the problem of estimating the parameters of the multivariate normal distribution and also show how the main Bayesian modelling strategies extend to the multivariate case.

2.2 Empirical estimates

2.2.1 General principle

For large n we have thanks to the law of large numbers:

$$\underbrace{F}_{\text{true}} \approx \underbrace{\widehat{F}_n}_{\text{empirical}}$$

We now would like to estimate A which is a *functional* $A = m(F)$ of the distribution F — recall that a functional is a function that takes another function as argument. For example all standard distributional summaries such as the mean, the median etc. are derived from F and hence are functionals of F .

The *empirical estimate* is obtained by replacing the unknown true distribution F with the observed empirical distribution: $\widehat{A} = m(\widehat{F}_n)$.

For example, the expectation of a random variable is approximated/estimated as the average over the observations:

$$E_F(x) \approx E_{\widehat{F}_n}(x) = \frac{1}{n} \sum_{k=1}^n x_k$$

$$E_F(g(x)) \approx E_{\widehat{F}_n}(g(x)) = \frac{1}{n} \sum_{k=1}^n g(x_k)$$

Simple recipe to obtain an empirical estimator: simply replace the expectation operator by the sample average.

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

Note: the approximation of F by \widehat{F}_n is also the basis other approaches such as Efron's bootstrap method (1979)¹.

2.2.2 Empirical estimates of mean and covariance

Recall the definitions:

$$\mu = E(x)$$

and

$$\Sigma = E\left((x - \mu)(x - \mu)^T\right)$$

¹Efron, B. 1979. Bootstrap methods: Another look at the jackknife. The Annals of Statistics 7:1–26.
<https://doi.org/10.1214/aos/1176344552>

For the empirical estimate we replace the expectations by the corresponding sample averages.

These resulting estimators can be written in three different ways:

Vector notation:

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

$$\widehat{\Sigma} = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})(x_k - \hat{\mu})^T = \frac{1}{n} \sum_{k=1}^n x_k x_k^T - \hat{\mu} \hat{\mu}^T$$

Component notation:

The corresponding component notation with $\mathbf{X} = (x_{ki})$ and following the statistics convention with samples contained in rows of \mathbf{X} we get:

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

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

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

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 \mathbf{X} .

If the data matrix \mathbf{X} follows the **statistics convention** we can write

$$\hat{\mu} = \frac{1}{n} \mathbf{X}^T \mathbf{1}_n$$

$$\widehat{\Sigma} = \frac{1}{n} X^T X - \hat{\mu} \hat{\mu}^T$$

On the other hand, if X follows the **engineering convention** with samples in columns the estimators are written as:

$$\hat{\mu} = \frac{1}{n} X \mathbf{1}_n$$

$$\widehat{\Sigma} = \frac{1}{n} X X^T - \hat{\mu} \hat{\mu}^T$$

To avoid confusion when using matrix or component notation you need to always state which convention is used! In these notes we exclusively follow the statistics convention.

2.3 Maximum likelihood estimation

2.3.1 General principle

R.A. Fisher (1922)²: model-based estimators using the density or probability mass function

log-likelihood function:

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

likelihood = probability to observe data given the model parameters

Maximum likelihood estimate:

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

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

Recall from [MATH27720 Statistics 2](#) that maximum likelihood is closely linked to minimising the relative entropy (KL divergence) $D_{\text{KL}}(F, F_{\theta})$ between the unknown true model F to the specified model F_{θ} . Specifically, for large sample size n the model $F_{\hat{\theta}}$ fit by maximum likelihood is indeed the model that is closest to F .

²Fisher, R. A. 1922. On the mathematical foundations of theoretical statistics. Philosophical Transactions of the Royal Society A 222:309–368. <https://doi.org/10.1098/rsta.1922.0009>

Correspondingly, the great appeal of **maximum likelihood estimates** (MLEs) is that they **are optimal for large n** , i.e. so that **for large sample size no estimator can be constructed that outperforms the MLE** (note the emphasis on “for large n ”). 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 observed Fisher information which is related to the curvature of the log-likelihood function).

2.3.2 Maximum likelihood estimates of the parameters of the multivariate

normal distribution

We now derive the MLE of the parameters μ and Σ of the multivariate normal distribution. The corresponding log-likelihood function is

$$\begin{aligned} \log L(\mu, \Sigma) &= \sum_{k=1}^n \log f(x_k | \mu, \Sigma) \\ &= -\frac{nd}{2} \log(2\pi) - \frac{n}{2} \log \det(\Sigma) - \frac{1}{2} \sum_{k=1}^n (x_k - \mu)^T \Sigma^{-1} (x_k - \mu). \end{aligned}$$

Written in terms of the precision matrix $\Omega = \Sigma^{-1}$ this becomes

$$\log L(\mu, \Omega) = -\frac{nd}{2} \log(2\pi) + \frac{n}{2} \log \det(\Omega) - \frac{1}{2} \sum_{k=1}^n (x_k - \mu)^T \Omega (x_k - \mu).$$

First, to find the MLE for μ we compute

$$\nabla_{\mu} \log L(\mu, \Omega) = \frac{\partial \log L(\mu, \Omega)}{\partial \mu} = \sum_{k=1}^n \Omega (x_k - \mu)$$

noting that Ω is symmetric. Setting this equal to zero we get $\sum_{k=1}^n x_k = n \hat{\mu}_{ML}$ and thus

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

Next, to obtain the MLE for Ω we compute

$$\frac{\partial \log L(\mu, \Omega)}{\partial \Omega} = \frac{n}{2} \Omega^{-1} - \frac{1}{2} \sum_{k=1}^n (x_k - \mu)(x_k - \mu)^T$$

. Setting this equal to zero and substituting the MLE for μ we get

$$\widehat{\Omega}_{ML}^{-1} = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})(x_k - \hat{\mu})^T = \widehat{\Sigma}_{ML}.$$

See the supplementary [Matrix Refresher](#) notes for the relevant formulas in vector and matrix calculus.

Therefore, the MLEs are identical to the empirical estimates.

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

2.4 Sampling 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. The sample average is denoted by $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

1. Distribution of the estimate of the mean:

The empirical estimate of the mean is normally distributed:

$$\hat{\mu}_{ML} = \bar{x} \sim N_d\left(\mu, \frac{\Sigma}{n}\right)$$

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

2. Distribution of the covariance estimate:

The empirical and unbiased estimate of the covariance matrix is Wishart distributed:

Case of unknown mean μ (estimated by \bar{x}):

$$\hat{\Sigma}_{ML} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \sim W_d\left(\frac{\Sigma}{n}, n-1\right)$$

Since $E(\hat{\Sigma}_{ML}) = \frac{n-1}{n} \Sigma \implies \hat{\Sigma}_{ML}$ is biased, with $\text{Bias}(\hat{\Sigma}_{ML}) = \Sigma - E(\hat{\Sigma}_{ML}) = -\frac{\Sigma}{n}$.

Easy to make unbiased: $\hat{\Sigma}_{UB} = \frac{n}{n-1} \hat{\Sigma}_{ML} = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})(x_k - \bar{x})^T$ which is distributed as

$$\hat{\Sigma}_{UB} \sim W_d\left(\frac{\Sigma}{n-1}, n-1\right)$$

Hence $E(\hat{\Sigma}_{UB}) = \Sigma \implies \hat{\Sigma}_{UB}$ is unbiased.

But unbiasedness of an estimator is **not** a very relevant criterion in multivariate statistics, especially when the number of samples is small compared to the dimension (see further below).

Covariance estimator for known mean μ :

$$\frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T \sim W_d\left(\frac{\Sigma}{n}, n\right)$$

2.5 Small sample estimation

2.5.1 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)

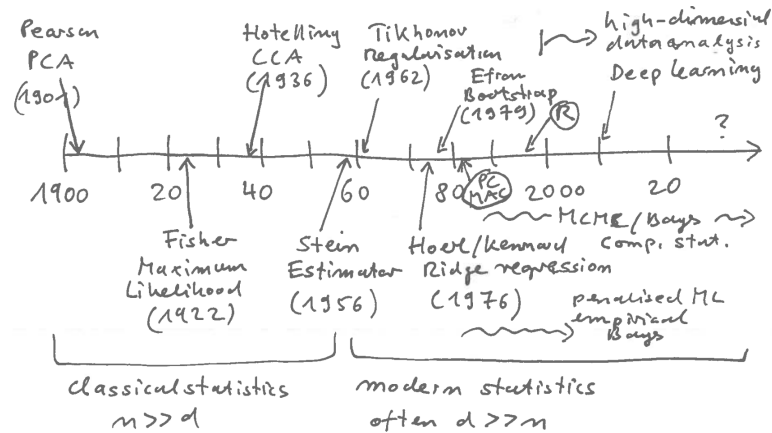
General problems of MLEs:

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
- 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).

2.5.2 Estimation of covariance matrix in small sample settings

Problems with ML estimate of Σ

1. Σ has $O(d^2)$ number of parameters! $\Rightarrow \hat{\Sigma}^{\text{MLE}}$ requires a lot of data!
 $n \gg d$ or d^2
2. if $n < d$ then $\hat{\Sigma}$ is positive **semi**-definite (even if the true Σ is positive

definite!)

$\Rightarrow \hat{\Sigma}$ will have **vanishing eigenvalues** (some $\lambda_i = 0$) and thus **cannot be inverted** and is singular!

Note that in many expression in multivariate statistics we actually need to use the inverse of the covariance matrix, e.g., in the density of the multivariate normal distribution, so it is essential that we obtain a non-singular invertible estimate of the covariance matrix.

Making the ML estimate of Σ invertible

There is a simple numerical trick credited to [A. N. Tikhonov](#) to make $\hat{\Sigma}$ invertible, by adding a small number (say $\varepsilon = 10^{-6}$ to the diagonal elements of $\hat{\Sigma}$:

$$S_{\text{Tik}} = \hat{\Sigma} + \varepsilon I$$

The resulting S_{Tik} is **positive definite** because the sum of a symmetric positive definite matrix (εI) and a symmetric positive semi-definite matrix ($\hat{\Sigma}$) is always positive definite.

However, while this simple regularisation results in an invertible matrix the estimator itself has not improved over the MLE, and the matrix S_{Tik} will also be poorly conditioned (i.e. large condition number).

Simple regularised estimate of Σ

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

Regularisation:

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

Idea: choose $\lambda \in [0, 1]$ such that S^* is better (e.g. in terms of MSE) than both S and I_d . Note that λ does not need to be small like ε .

This form of estimator corresponds to computing the mean of the Bayesian posterior by directly shrinking the MLE towards a prior mean (target):

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

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

- Instead of a diagonal target other options are possible, e.g. block-diagonal or banded covariances.
- If λ is not prespecified but learned from data (see below) then the resulting estimate is an empirical Bayes estimator.
- The resulting estimate will typically be biased as mixing in the target will increase the bias.

How to find optimal shrinkage / regularisation parameter λ ?

One way to do this is to minimise MSE (Mean Squared Error). This is also called L2 regularisation or Ridge regularisation.

Bias-variance trade-off: MSE is composed of squared bias and variance.

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

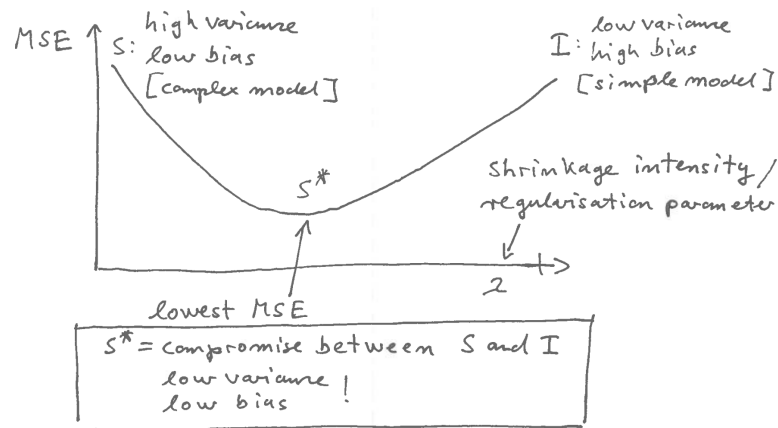
with $\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$

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

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

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

\Rightarrow overall, MSE is decreased



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)

In Worksheet 3 the empirical estimator of covariance is compared with the regularised covariance estimator implemented in the R package “corpcor”. This uses a regularisation similar as above (but for the correlation rather than the 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 [MATH27720 Statistics 2](#)).

2.6 Full Bayesian multivariate modelling

See also the [section about multivariate distributions in the Probability and Distribution refresher](#) for details about the distributions used below.

2.6.1 Three main scenarios

As discussed in [MATH27720 Statistics 2](#) there are three main Bayesian models in the univariate case that cover a large range of applications:

- 1) the beta-binomial model to estimate proportions
- 2) the normal-normal model to estimate means
- 3) the inverse gamma-normal model to estimate variances

Below we briefly sketch the extensions of these three elementary models to the multivariate case.

2.6.2 Dirichlet-multinomial model

This generalises the univariate beta-binomial model.

The Dirichlet distribution is useful as conjugate prior and posterior distribution for the parameters of a categorical distribution.

- Data: $D = \{x_1, \dots, x_n\}$ with $x_i \sim \text{Cat}(\pi)$
- MLE: $\hat{\pi}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$
- Prior parameters (Dirichlet in mean parameterisation): k_0, π_0 ,

$$\pi \sim \text{Dir}(\pi_0, k_0)$$

$$E(\pi) = \pi_0$$

- Posterior parameters: $k_1 = k_0 + n$, $\pi_1 = \lambda \pi_0 + (1 - \lambda) \hat{\pi}_{ML}$ with $\lambda = \frac{k_0}{k_1}$

$$\pi | D \sim \text{Dir}(\pi_1, k_1)$$

$$E(\pi | D) = \pi_1$$

- Equivalent update rule (for Dirichlet with α parameter): $\alpha_0 \rightarrow \alpha_1 = \alpha_0 + \sum_{i=1}^n x_i = \alpha_0 + n \hat{\pi}_{ML}$

2.6.3 Multivariate normal-normal model

This generalises the univariate normal-normal model.

The multivariate normal distribution is useful as conjugate prior and posterior distribution of the mean:

- Data: $D = \{x_1, \dots, x_n\}$ with $x_i \sim N_d(\mu, \Sigma)$ with known mean Σ
- MLE: $\hat{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$
- Prior parameters: k_0, μ_0

$$\mu \sim N_d \left(\mu_0, \frac{\Sigma}{k_0} \right)$$

$$E(\mu) = \mu_0$$

- Posterior parameters: $k_1 = k_0 + n, \mu_1 = \lambda \mu_0 + (1 - \lambda) \hat{\mu}_{ML}$ with $\lambda = \frac{k_0}{k_1}$

$$\mu | D \sim N_d \left(\mu_1, \frac{\Sigma}{k_1} \right)$$

$$E(\mu | D) = \mu_1$$

2.6.4 Inverse Wishart-normal model

This generalises the univariate inverse gamma-normal model for the variance.

The inverse Wishart distribution is useful as conjugate prior and posterior distribution of the covariance:

- Data: $D = \{x_1, \dots, x_n\}$ with $x_i \sim N_d(\mu, \Sigma)$ with known mean μ
- MLE: $\hat{\Sigma}_{ML} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$
- Prior parameters: κ_0, Σ_0

$$\Sigma \sim W_d^{-1}(\kappa_0 \Sigma_0, \kappa_0 + d + 1)$$

$$E(\Sigma) = \Sigma_0$$

- Posterior parameters: $\kappa_1 = \kappa_0 + n, \Sigma_1 = \lambda \Sigma_0 + (1 - \lambda) \hat{\Sigma}_{ML}$ with $\lambda = \frac{\kappa_0}{\kappa_1}$

$$\Sigma | D \sim W_d^{-1}(\kappa_1 \Sigma_1, \kappa_1 + d + 1)$$

$$E(\Sigma | D) = \Sigma_1$$

- Equivalent update rule: $v_0 \rightarrow v_1 = v_0 + n, \Psi_0 \rightarrow \Psi_1 = \Psi_0 + \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T = \Psi_0 + n \hat{\Sigma}_{ML}$

2.7 Conclusion

- Multivariate models are often high-dimensional with large number of parameters but often only a small number of samples are available. In this instance it is useful (and often necessary) to introduce additional information (via priors or by regularisation).
- Unbiased estimation, a highly valued property in classical univariate statistics when sample size is large and number of parameters is small, is typically not a good idea in multivariate settings and often leads to poor estimators.
- Regularisation introduces bias and reduces variance, minimising overall MSE. Likewise, Bayesian estimators also introduce bias and regularise (via the prior) and thus are useful in multivariate settings.

Chapter 3

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”.

3.1 Linear Transformations

3.1.1 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

Mean and variance of the original vector x :

$$E(x) = \mu_x$$

$$\text{Var}(x) = \Sigma_x$$

Mean and variance of the transformed random vector \mathbf{y} :

$$\mathbf{E}(\mathbf{y}) = \mathbf{a} + \mathbf{B}\boldsymbol{\mu}_x$$

$$\text{Var}(\mathbf{y}) = \mathbf{B}\boldsymbol{\Sigma}_x\mathbf{B}^T$$

Cross-covariance $\boldsymbol{\Phi} = \boldsymbol{\Sigma}_{xy} = \text{Cov}(\mathbf{x}, \mathbf{y})$ between \mathbf{x} and \mathbf{y} :

$$\boldsymbol{\Phi} = \text{Cov}(\mathbf{x}, \mathbf{B}\mathbf{x}) = \boldsymbol{\Sigma}_x\mathbf{B}^T$$

Note that $\boldsymbol{\Phi}$ is a matrix of dimensions $d \times m$ as the dimension of \mathbf{x} is d and the dimension of \mathbf{y} is m .

Cross-correlation $\boldsymbol{\Psi} = \mathbf{P}_{xy} = \text{Cor}(\mathbf{x}, \mathbf{y})$ between \mathbf{x} and \mathbf{y} :

$$\boldsymbol{\Psi} = \mathbf{V}_x^{-1/2}\boldsymbol{\Phi}\mathbf{V}_y^{-1/2}$$

where $\mathbf{V}_x = \text{Diag}(\boldsymbol{\Sigma}_x)$ and $\mathbf{V}_y = \text{Diag}(\mathbf{B}\boldsymbol{\Sigma}_x\mathbf{B}^T)$ are diagonal matrices containing the variances for the components of \mathbf{x} and \mathbf{y} . The dimensions of the matrix $\boldsymbol{\Psi}$ are also $d \times m$.

Special cases/examples:

Example 3.1. Univariate case ($d = 1, m = 1$): $y = a + bx$

- $\mathbf{E}(y) = a + b\mu$
- $\text{Var}(y) = b^2\sigma^2$
- $\text{Cov}(y, x) = b\sigma^2$
- $\text{Cor}(y, x) = \frac{b\sigma^2}{\sqrt{b^2\sigma^2}\sqrt{\sigma^2}} = 1$

Note that y can be predicted perfectly from x as $\text{Cor}(y, x) = 1$. This is because there is no error term in the transformation. See also the more general case with multiple correlation further below.

Example 3.2. Sum of two random univariate variables: $y = x_1 + x_2$, i.e. $a = 0$ and $\mathbf{B} = (1, 1)$

- $\mathbf{E}(y) = \mathbf{E}(x_1 + x_2) = \mu_1 + \mu_2$
- $\text{Var}(y) = \text{Var}(x_1 + x_2) = (1, 1) \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \sigma_1^2 + \sigma_2^2 + 2\sigma_{12} = \text{Var}(x_1) + \text{Var}(x_2) + 2\text{Cov}(x_1, x_2)$

Example 3.3. $y_1 = a_1 + b_1x_1$ and $y_2 = a_2 + b_2x_2$, i.e. $\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$ and

$$\mathbf{B} = \begin{pmatrix} b_1 & 0 \\ 0 & b_2 \end{pmatrix}$$

- $\mathbf{E}(\mathbf{y}) = \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} + \begin{pmatrix} b_1 & 0 \\ 0 & b_2 \end{pmatrix} \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \begin{pmatrix} a_1 + b_1\mu_1 \\ a_2 + b_2\mu_2 \end{pmatrix}$

- $\text{Var}(\mathbf{y}) = \begin{pmatrix} b_1 & 0 \\ 0 & b_2 \end{pmatrix} \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} \begin{pmatrix} b_1 & 0 \\ 0 & b_2 \end{pmatrix} = \begin{pmatrix} b_1^2\sigma_1^2 & b_1b_2\sigma_{12} \\ b_1b_2\sigma_{12} & b_2^2\sigma_2^2 \end{pmatrix}$
note that $\text{Cov}(y_1, y_2) = b_1b_2\text{Cov}(x_1, x_2)$

3.1.2 Squared multiple correlation

Squared multiple correlation $\text{MCor}(\mathbf{y}, \mathbf{x})^2$ is a scalar measure summarising the linear association between a scalar response variable y and a set of predictors $\mathbf{x} = (x_1, \dots, x_d)^T$. It is defined as

$$\begin{aligned} \text{MCor}(\mathbf{y}, \mathbf{x})^2 &= \boldsymbol{\Sigma}_{yx}\boldsymbol{\Sigma}_x^{-1}\boldsymbol{\Sigma}_{xy}/\sigma_y^2 \\ &= \mathbf{P}_{yx}\mathbf{P}_x^{-1}\mathbf{P}_{xy} \end{aligned}$$

If y can be perfectly linearly predicted by \mathbf{x} then $\text{MCor}(\mathbf{y}, \mathbf{x})^2 = 1$.

The empirical estimate of $\text{MCor}(\mathbf{y}, \mathbf{x})^2$ is the R^2 coefficient that you will find in any software for linear regression.

Example 3.4. Squared multiple correlation for an affine transformation.

Since we linearly transform \mathbf{x} into \mathbf{y} with no additional error involved we expect that for each component y_i in \mathbf{y} we have $\text{MCor}(y_i, \mathbf{x})^2 = 1$. This can be shown directly by computing

$$\begin{aligned} (\text{MCor}(y_1, \mathbf{x})^2, \dots, \text{MCor}(y_m, \mathbf{x})^2)^T &= \text{Diag}(\boldsymbol{\Sigma}_{yx}\boldsymbol{\Sigma}_x^{-1}\boldsymbol{\Sigma}_{xy}) / \text{Diag}(\boldsymbol{\Sigma}_y) \\ &= \text{Diag}(\mathbf{B}\boldsymbol{\Sigma}_x\boldsymbol{\Sigma}_x^{-1}\boldsymbol{\Sigma}_x\mathbf{B}^T) / \text{Diag}(\mathbf{B}\boldsymbol{\Sigma}_x\mathbf{B}^T) \\ &= \text{Diag}(\mathbf{B}\boldsymbol{\Sigma}_x\mathbf{B}^T) / \text{Diag}(\mathbf{B}\boldsymbol{\Sigma}_x\mathbf{B}^T) \\ &= (1, \dots, 1)^T \end{aligned}$$

3.1.3 Invertible location-scale transformation

If $m = d$ (square \mathbf{B}) and $\det(\mathbf{B}) \neq 0$ then the affine transformation is **invertible**.

Forward transformation:

$$\mathbf{y} = \mathbf{a} + \mathbf{B}\mathbf{x}$$

Back transformation:

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

Invertible transformations thus provide a one-to-one map between \mathbf{x} and \mathbf{y} .

Example 3.5. Mahalanobis transform

We assume $\mathbf{E}(\mathbf{x}) = \boldsymbol{\mu}_x$ and a positive definite covariance matrix $\text{Var}(\mathbf{x}) = \boldsymbol{\Sigma}_x$ with $\det(\boldsymbol{\Sigma}_x) > 0$.

The Mahalanobis transformation is given by

$$\mathbf{y} = \Sigma_x^{-1/2}(\mathbf{x} - \mu_x)$$

This corresponds to an affine transformation with $\mathbf{a} = -\Sigma_x^{-1/2}\mu_x$ and $\mathbf{B} = \Sigma_x^{-1/2}$.

Note that this transformation contains the **inverse principal matrix square root** $\Sigma_x^{-1/2}$. This can be obtained by eigendecomposition of $\Sigma_x = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ so that $\Sigma_x^{-1/2} = \mathbf{U}\mathbf{\Lambda}^{-1/2}\mathbf{U}^T$.

The mean and the variance of \mathbf{y} becomes

$$\mathbb{E}(\mathbf{y}) = \mathbf{0}$$

and

$$\text{Var}(\mathbf{y}) = \mathbf{I}_d$$

The Mahalanobis transforms performs three functions:

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

In the **univariate case** ($d = 1$) the coefficients reduce to $a = -\frac{\mu_x}{\sigma_x}$ and $B = \frac{1}{\sigma_x}$ and the Mahalanobis transform becomes

$$y = \frac{x - \mu_x}{\sigma_x}$$

i.e. it applies centering + standardisation.

The **Mahalanobis transformation** appears implicitly in many places in multivariate statistics, e.g. in the multivariate normal density. It is a particular example of a whitening transformation (of which there are infinitely many, see later in the course).

Example 3.6. Inverse Mahalanobis transformation

The inverse of the Mahalanobis transform is given by

$$\mathbf{y} = \mu_y + \Sigma_y^{1/2}\mathbf{x}$$

As the Mahalanobis transform is a whitening transform the inverse Mahalanobis transform is sometimes called the Mahalanobis colouring transformation. The coefficients in the affine transformation are $\mathbf{a} = \mu_y$ and $\mathbf{B} = \Sigma_y^{1/2}$.

Starting with $\mathbb{E}(\mathbf{x}) = \mathbf{0}$ and $\text{Var}(\mathbf{x}) = \mathbf{I}_d$ the mean and variance of the transformed variable are

$$\mathbb{E}(\mathbf{y}) = \mu_y$$

and

$$\text{Var}(\mathbf{y}) = \Sigma_y$$

3.1.4 Transformation of a density under an invertible location-scale transformation:

Assume $\mathbf{x} \sim F_x$ with density $f_x(\mathbf{x})$.

After linear transformation $\mathbf{y} = \mathbf{a} + \mathbf{B}\mathbf{x}$ we get $\mathbf{y} \sim F_y$ with density

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

Example 3.7. Transformation of standard normal with inverse Mahalanobis transform

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

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

Then the density after applying the inverse Mahalanobis transform

$\mathbf{y} = \mu_y + \Sigma_y^{1/2}\mathbf{x}$ is

$$\begin{aligned} f_y(\mathbf{y}) &= |\det(\Sigma_y^{1/2})|^{-1} (2\pi)^{-d/2} \exp\left(-\frac{1}{2}(\mathbf{y} - \mu_y)^T \Sigma_y^{-1/2} \Sigma_y^{-1/2} (\mathbf{y} - \mu_y)\right) \\ &= (2\pi)^{-d/2} \det(\Sigma_y)^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{y} - \mu_y)^T \Sigma_y^{-1} (\mathbf{y} - \mu_y)\right) \end{aligned}$$

$\Rightarrow \mathbf{y}$ has multivariate normal density $N_d(\mu_y, \Sigma_y)$

Application: e.g. random number generation: draw from $N_d(\mathbf{0}, \mathbf{I}_d)$ (easy!) then convert to multivariate normal by transformation (see Worksheet 4).

3.2 Nonlinear transformations

3.2.1 General transformation

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

with h an arbitrary vector-valued function

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

3.2.2 Delta method

Assume that we know the mean $\mathbb{E}(\mathbf{x}) = \mu_x$ and variance $\text{Var}(\mathbf{x}) = \Sigma_x$ of \mathbf{x} . Is it possible to say something about the mean and variance of the transformed random variable \mathbf{y} ?

$$\mathbb{E}(\mathbf{y}) = \mathbb{E}(h(\mathbf{x})) = ?$$

$$\text{Var}(\mathbf{y}) = \text{Var}(\mathbf{h}(\mathbf{x})) = ?$$

In general, for a transformation $\mathbf{h}(\mathbf{x})$ the exact mean and variance of the transformed variable cannot be obtained analytically.

However, we can find a **linear approximation** and then compute its mean and variance. This approximation is called the “Delta Method”, or the “law of propagation of errors”, and is credited to Gauss¹.

Linearisation of $\mathbf{h}(\mathbf{x})$ is achieved by a Taylor series approximation of first order of $\mathbf{h}(\mathbf{x})$ around \mathbf{x}_0 :

$$\mathbf{h}(\mathbf{x}) \approx \mathbf{h}(\mathbf{x}_0) + \underbrace{\mathbf{J}_h(\mathbf{x}_0)}_{\text{Jacobian matrix}} (\mathbf{x} - \mathbf{x}_0) = \underbrace{\mathbf{h}(\mathbf{x}_0) - \mathbf{J}_h(\mathbf{x}_0) \mathbf{x}_0}_{\mathbf{a}} + \underbrace{\mathbf{J}_h(\mathbf{x}_0)}_{\mathbf{B}} \mathbf{x}$$

If $\mathbf{h}(\mathbf{x})$ is scalar-valued then **gradient** $\nabla h(\mathbf{x})$ is given by the vector of partial correlations

$$\nabla h(\mathbf{x}) = \begin{pmatrix} \frac{\partial h(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial h(\mathbf{x})}{\partial x_d} \end{pmatrix}$$

where ∇ is the nabla operator.

The **Jacobian matrix** is the **generalisation of the gradient** if $\mathbf{h}(\mathbf{x})$ is vector-valued:

$$\mathbf{J}_h(\mathbf{x}) = \begin{pmatrix} \nabla h_1(\mathbf{x})^T \\ \nabla h_2(\mathbf{x})^T \\ \vdots \\ \nabla h_m(\mathbf{x})^T \end{pmatrix} = \begin{pmatrix} \frac{\partial h_1(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial h_1(\mathbf{x})}{\partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_m(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial h_m(\mathbf{x})}{\partial x_d} \end{pmatrix}$$

Note that in the Jacobian matrix by convention the gradient for each individual component of $\mathbf{h}(\mathbf{x})$ is contained in the *row* of the matrix so the number of rows corresponds to the dimension of $\mathbf{h}(\mathbf{x})$ and the number of columns to the dimension of \mathbf{x} .

First order approximation of $\mathbf{h}(\mathbf{x})$ around $\mathbf{x}_0 = \boldsymbol{\mu}_x$ yields $\mathbf{a} = \mathbf{h}(\boldsymbol{\mu}_x) - \mathbf{J}_h(\boldsymbol{\mu}_x) \boldsymbol{\mu}_x$ and $\mathbf{B} = \mathbf{J}_h(\boldsymbol{\mu}_x)$ and leads directly to the **multivariate Delta method**:

$$\begin{aligned} \mathbf{E}(\mathbf{y}) &\approx \mathbf{h}(\boldsymbol{\mu}_x) \\ \text{Var}(\mathbf{y}) &\approx \mathbf{J}_h(\boldsymbol{\mu}_x) \boldsymbol{\Sigma}_x \mathbf{J}_h(\boldsymbol{\mu}_x)^T \end{aligned}$$

The **univariate Delta method** is a special case:

$$\mathbf{E}(y) \approx h(\mu_x)$$

¹Gorroochurn, P. 2020. Who Invented the Delta Method, Really? The Mathematical Intelligencer 42:46–49. <https://doi.org/10.1007/s00283-020-09982-0>

$$\text{Var}(y) \approx \sigma_x^2 h'(\mu_x)^2$$

Note that the Delta approximation breaks down if $\text{Var}(\mathbf{y})$ is singular, for example if the first derivative (or gradient or Jacobian matrix) at $\boldsymbol{\mu}_x$ is zero.

Example 3.8. Variance of the odds ratio

The proportion $\hat{p} = \frac{n_1}{n}$ resulting from n repeats of a Bernoulli experiment has expectation $\mathbf{E}(\hat{p}) = p$ and variance $\text{Var}(\hat{p}) = \frac{p(1-p)}{n}$. What are the (approximate) mean and the variance of the corresponding odds ratio $\widehat{OR} = \frac{\hat{p}}{1-\hat{p}}$?

With $h(x) = \frac{x}{1-x}$, $\widehat{OR} = h(\hat{p})$ and $h'(x) = \frac{1}{(1-x)^2}$ we get using the Delta method $\mathbf{E}(\widehat{OR}) \approx h(p) = \frac{p}{1-p}$ and $\text{Var}(\widehat{OR}) \approx h'(p)^2 \text{Var}(\hat{p}) = \frac{p}{n(1-p)^3}$.

Example 3.9. Log-transform as variance stabilisation

Assume x has some mean $\mathbf{E}(x) = \mu$ and variance $\text{Var}(x) = \sigma^2 \mu^2$, i.e. the standard deviation $\text{SD}(x)$ is proportional to the mean μ . What are the (approximate) mean and the variance of the log-transformed variable $\log(x)$?

With $h(x) = \log(x)$ and $h'(x) = \frac{1}{x}$ we get using the Delta method $\mathbf{E}(\log(x)) \approx h(\mu) = \log(\mu)$ and $\text{Var}(\log(x)) \approx h'(\mu)^2 \text{Var}(x) = \left(\frac{1}{\mu}\right)^2 \sigma^2 \mu^2 = \sigma^2$. Thus, after applying the log-transform the variance does not depend any more on the mean!

3.2.3 Transformation of a probability density function under a general invertible transformation

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

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

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

$$f_y(\mathbf{y}) = |\det(\mathbf{J}_x(\mathbf{y}))| f_x(\mathbf{x}(\mathbf{y}))$$

where $\mathbf{J}_x(\mathbf{y})$ is the Jacobian matrix of the inverse transformation.

Special cases:

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

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

3.2.4 Normalising flows

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

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 for simplification and dimension reduction.

If you are interested in normalising flows then a good start to learn more are the review papers by Kobyzev et al (2021)² and Papamakarios et al. (2021)³.

3.3 General whitening transformations

3.3.1 Overview

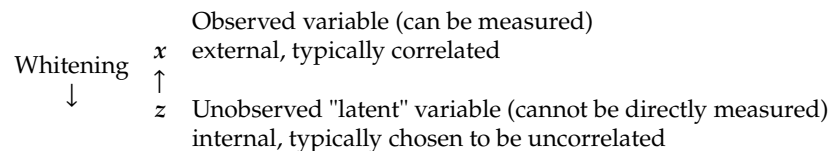
Whitening transformations are a special and widely used class of invertible location-scale transformations.

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

Whitening is **useful in preprocessing**, as they allow to **turn multivariate models into uncorrelated univariate models** (via decorrelation property). Some whitening transformations **reduce the dimension in an optimal way** (via compression property).

The *Mahalanobis* transform is a specific example of a **whitening transformation**. It is also known as “zero-phase component analysis” or short ZCA transform.

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



²Kobyzev et al. 2021. *Normalizing Flows: Introduction and Ideas*. *IEEE Trans. Pattern Anal. Mach. Intell.* **43**:3964-3979

³Papamakarios et al. 2021. *Normalizing Flows for Probabilistic Modeling and Inference*. *JMLR* **22**:1-64

3.3.2 Whitening transformation and whitening constraint

Starting point:

Random vector $x \sim F_x$ **not necessarily from multivariate normal**.

x has mean $E(x) = \mu$ and a positive definite (invertible) covariance matrix $\text{Var}(x) = \Sigma$.

Note that in the following we leave out the subscript x for the covariance of x unless it is needed for clarification.

The covariance can be split into positive variances V and a positive definite invertible correlation matrix P so that $\Sigma = V^{1/2} P V^{1/2}$.

Whitening transformation:

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

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

For Mahalanobis/ZCA whitening we already know from Example 3.5 that $W^{\text{ZCA}} = \Sigma^{-1/2}$.

In general, the whitening matrix W needs to satisfy a constraint:

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

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

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

3.3.3 Parameterisation of whitening matrix

Covariance-based parameterisation of whitening matrix:

A general way to specify a valid whitening matrix is

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

where Q_1 is an orthogonal matrix.

Recall that an orthogonal matrix Q has the property that $Q^{-1} = Q^T$ and as a consequence $Q^T Q = Q Q^T = I$.

As a result, the above W satisfies the whitening constraint:

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

Note the converse is also true: any whitening whitening matrix, i.e. any W satisfying the whitening constraint, can be written in the above form as $Q_1 = W \Sigma^{1/2}$ is orthogonal by construction.

⇒ instead of choosing W , we choose the orthogonal matrix Q_1 !

- recall that orthogonal matrices geometrically represent rotations (plus reflections).
- 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 $Q_1^{ZCA} = I$
- **whitening** can be interpreted as **Mahalanobis transformation** followed by further **rotation-reflection**

Correlation-based parameterisation of whitening matrix:

Instead of working with the covariance matrix Σ , we can express W also in terms of the corresponding correlation matrix $P = (\rho_{ij}) = V^{-1/2} \Sigma V^{-1/2}$ where V is the diagonal matrix containing the variances.

Specifically, we can specify the whitening matrix as

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

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

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

Conversely, any whitening matrix W can also be written in this form as $Q_2 = W V^{1/2} P^{1/2}$ is orthogonal by construction.

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

Both forms to write W using Q_1 and Q_2 are equally valid (and interchangeable).

Note that for the same W

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

and also

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

even though

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

3.3.4 Cross-covariance and cross-correlation for general whitening transformations

A useful criterion to characterise and to distinguish among whitening transformations is the cross-covariance and cross-correlation matrix between the original variable x and the whitened variable z :

a) **Cross-covariance** $\Phi = \Sigma_{xz}$ between x and z :

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

In component notation we write $\Phi = (\phi_{ij})$ where the row index i refers to x and the column index j to z .

Cross-covariance is linked with Q_1 ! Thus, choosing cross-covariance determines Q_1 (and vice versa).

Note that the above cross-covariance matrix Φ satisfies the condition $\Phi \Phi^T = \Sigma$.

The whitening matrix expressed in terms of cross-covariance is $W = \Phi^T \Sigma^{-1}$, so as required $W^T W = \Sigma^{-1} \Phi \Phi^T \Sigma^{-1} = \Sigma^{-1}$. Furthermore, Φ is the *inverse* of the whitening matrix, as $W^{-1} = (Q_1 \Sigma^{-1/2})^{-1} = \Sigma^{1/2} Q_1^{-1} = \Sigma^{1/2} Q_1^T = \Phi$.

b) **Cross-correlation** $\Psi = P_{xz}$ between x and z :

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

In component notation we write $\Psi = (\psi_{ij})$ where the row index i refers to x and the column index j to z .

Cross-correlation is linked with Q_2 ! Hence, choosing cross-correlation determines Q_2 (and vice versa). The whitening matrix expressed in terms of cross-correlation is $W = \Psi^T P^{-1} V^{-1/2}$.

Note that the factorisation of the cross-covariance $\Phi = \Sigma^{1/2} Q_1^T$ and the cross-correlation $\Psi = P^{1/2} Q_2^T$ into the product of a positive definite symmetric matrix and an orthogonal matrix are examples of a **polar decomposition**.

3.3.5 Inverse whitening transformation and loadings

Inverse transformation:

Recall that $z = Wx$. Therefore, the reverse transformation going from the whitened to the original variable is $x = W^{-1}z$. This can be expressed also in terms of cross-covariance and cross-correlation. With $W^{-1} = \Phi$ we get

$$x = \Phi z.$$

Furthermore, since $\Psi = V^{-1/2}\Phi$ we have $W^{-1} = V^{1/2}\Psi$ and hence

$$V^{-1/2}x = \Psi z.$$

The reverse whitening transformation is also known as colouring transformation (the previously discussed inverse Mahalanobis transform is one example).

Definition of loadings:

Loadings are the coefficients of the linear transformation from the latent variable back to the observed variable. If the variables are standardised to unit variance then the loadings are also called *correlation loadings*.

Hence, the cross-covariance matrix Φ plays the role of *loadings* linking the latent variable z with the original x . Similarly, the cross-correlation matrix Ψ contains the *correlation loadings* linking the (already standardised) latent variable z with the standardised x .

In the convention we use here the rows correspond to the original variables and the columns to the whitened latent variables.

Multiple correlation coefficients from z back to x :

We consider the backtransformation from the whitened variable z to the original variables x and note that the components of z are all uncorrelated with $P_z = I$. The squared multiple correlation coefficient $\text{MCor}(x_i, z)$ between each x_i and all

z is therefore just the sum of the corresponding squared correlations $\text{Cor}(x_i, z_j)^2$:

$$\begin{aligned} \text{MCor}(x_i, z)^2 &= P_{x_i z} P_z^{-1} P_{z x_i} = \\ &= \sum_{j=1}^d \text{Cor}(x_i, z_j)^2 \\ &= \sum_{j=1}^d \psi_{ij}^2 = 1 \end{aligned}$$

As shown earlier for a general linear one-to-one- transformation (which includes whitening as special case) the squared multiple correlation must be 1 because there is no error. We can confirm this by computing the **row sums of squares** of the cross-correlation matrix Ψ in matrix notation

$$\begin{aligned} \text{Diag}(\Psi \Psi^T) &= \text{Diag}(P^{1/2} Q_2^T Q_2 P^{1/2}) \\ &= \text{Diag}(P) \\ &= (1, \dots, 1)^T \end{aligned}$$

from which it is clear that the choice of Q_2 is not relevant.

Similarly, the **row sums of squares** of the cross-covariance matrix Φ equal the variances of the original variables, regardless of Q_1 :

$$\sum_{j=1}^d \phi_{ij}^2 = \text{Var}(x_i)$$

or in matrix notation

$$\begin{aligned} \text{Diag}(\Phi \Phi^T) &= \text{Diag}(\Sigma^{1/2} Q_1^T Q_1 \Sigma^{1/2}) \\ &= \text{Diag}(\Sigma) \\ &= (\text{Var}(x_1), \dots, \text{Var}(x_d))^T \end{aligned}$$

3.3.6 Summaries of cross-covariance Φ and cross-correlation Ψ resulting from whitening transformations

Matrix trace:

A simply summary of a matrix is its trace. For the cross-covariance matrix Φ the trace is the sum of all covariances between corresponding elements in x and z :

$$\text{Tr}(\Phi) = \sum_{i=1}^d \text{Cov}(x_i, z_i) = \sum_{i=1}^d \phi_{ii} = \text{Tr}(\Sigma^{1/2} Q_1^T)$$

Likewise, for the cross-correlation matrix Ψ the trace is the sum of all correlations between corresponding elements in \mathbf{x} and \mathbf{z} :

$$\text{Tr}(\Psi) = \sum_{i=1}^d \text{Cor}(x_i, z_i) = \sum_{i=1}^d \psi_{ii} = \text{Tr}(\mathbf{P}^{1/2} \mathbf{Q}_2^T)$$

In both cases the value of the trace depends on \mathbf{Q}_1 and \mathbf{Q}_2 . Interestingly, there is unique choice such that the trace is maximised.

Specifically, to maximise $\text{Tr}(\Phi)$ we conduct the following steps:

- i) Apply eigendecomposition to $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$. Note that $\mathbf{\Lambda}$ is diagonal with positive eigenvalues $\lambda_i > 0$ so Σ is positive definite and \mathbf{U} is an orthogonal matrix.
- ii) The objective function becomes

$$\begin{aligned} \text{Tr}(\Phi) &= \text{Tr}(\Sigma^{1/2} \mathbf{Q}_1^T) \\ &= \text{Tr}(\mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{U}^T \mathbf{Q}_1^T) \\ &= \text{Tr}(\mathbf{\Lambda}^{1/2} \mathbf{U}^T \mathbf{Q}_1^T \mathbf{U}) \\ &= \text{Tr}(\mathbf{\Lambda}^{1/2} \mathbf{B}) \\ &= \sum_{i=1}^d \lambda_i^{1/2} b_{ii}. \end{aligned}$$

Note that the product of two orthogonal matrices is itself an orthogonal matrix. Therefore, $\mathbf{B} = \mathbf{U}^T \mathbf{Q}_1^T \mathbf{U}$ is an orthogonal matrix and $\mathbf{Q}_1 = \mathbf{U}\mathbf{B}^T \mathbf{U}^T$.

- iii) As $\lambda_i > 0$ and all $b_{ii} \in [-1, 1]$ the objective function is maximised for $b_{ii} = 1$, i.e. for $\mathbf{B} = \mathbf{I}$.
- iv) In turn this implies that $\text{Tr}(\Phi)$ is maximised for $\mathbf{Q}_1 = \mathbf{I}$.

Similarly, to maximise $\text{Tr}(\Psi)$ we

- decompose $\mathbf{P} = \mathbf{G}\mathbf{\Theta}\mathbf{G}^T$ and then, following the above,
- find that $\text{Tr}(\Psi) = \text{Tr}(\mathbf{\Theta}^{1/2} \mathbf{G}^T \mathbf{Q}_2^T \mathbf{G})$ is maximised for $\mathbf{Q}_2 = \mathbf{I}$.

Squared Frobenius norm and total variation:

Another way to summarise and dissect the association between \mathbf{x} and the corresponding whitened \mathbf{z} is the squared Frobenius norm and the total variation based on Φ and Ψ .

The squared Frobenius norm (Euclidean) norm is the sum of squared elements of a matrix.

If we consider the squared Frobenius norm of the cross-covariance matrix, i.e. the sum of squared covariances between \mathbf{x} and \mathbf{z} ,

$$\|\Phi\|_F^2 = \sum_{i=1}^d \sum_{j=1}^d \phi_{ij}^2 = \text{Tr}(\Phi^T \Phi) = \text{Tr}(\Sigma)$$

we find that this equals the **total variation** of Σ and that it does not depend on \mathbf{Q}_1 . Likewise, computing the squared Frobenius norm of the cross-correlation matrix, i.e. the sum of squared correlations between \mathbf{x} and \mathbf{z} ,

$$\|\Psi\|_F^2 = \sum_{i=1}^d \sum_{j=1}^d \psi_{ij}^2 = \text{Tr}(\Psi^T \Psi) = \text{Tr}(\mathbf{P}) = d$$

yields the total variation of \mathbf{P} which also does not depend on \mathbf{Q}_2 . Note this is because the squared Frobenius norm is invariant against rotations and reflections.

Proportion of total variation:

We can now compute the contribution of each whitened component z_j to the total variation. The sum of squared covariances of each z_j with all x_1, \dots, x_d is

$$h_j = \sum_{i=1}^d \text{Cov}(x_i, z_j)^2 = \sum_{i=1}^d \phi_{ij}^2$$

with $\sum_{j=1}^d h_j = \text{Tr}(\Sigma)$ the total variation. In vector notation the contributions are written as the **column sums of squares** of Φ

$$\mathbf{h} = (h_1, \dots, h_d)^T = \text{Diag}(\Phi^T \Phi) = \text{Diag}(\mathbf{Q}_1 \Sigma \mathbf{Q}_1^T).$$

The relative contribution of z_j versus the total variation is

$$\frac{h_j}{\text{Tr}(\Sigma)}.$$

Crucially, in contrast to total variation, the contributions h_j depend on the choice of \mathbf{Q}_1 .

Similarly, the sum of squared correlations of each z_j with all x_1, \dots, x_d is

$$k_j = \sum_{i=1}^d \text{Cor}(x_i, z_j)^2 = \sum_{i=1}^d \psi_{ij}^2$$

with $\sum_{j=1}^d k_j = \text{Tr}(\mathbf{P}) = d$. In vector notation this corresponds to the **column sums of squares** of Ψ

$$\mathbf{k} = (k_1, \dots, k_d)^T = \text{Diag}(\Psi^T \Psi) = \text{Diag}(\mathbf{Q}_2 \mathbf{P} \mathbf{Q}_2^T).$$

The relative contribution of z_j with regard to the total variation of the correlation P is

$$\frac{k_j}{\text{Tr}(P)} = \frac{k_j}{d}.$$

As above, the contributions k_j depend on the choice of Q_2 .

Maximising the proportion of total variation:

Interestingly, it is possible to choose a unique whitening transformation such that the contributions are maximised, i.e. that the sum of the m largest contributions of h_j and k_j is as large as possible.

Specifically, we note that $\Phi^T \Phi$ and $\Psi^T \Psi$ are symmetric real matrices. For these type of matrices we know from Schur's theorem (1923) that the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ **majorise** the diagonal elements $p_1 \geq p_2 \geq \dots \geq p_d$. More precisely,

$$\sum_{i=1}^m \lambda_i \geq \sum_{i=1}^m p_i,$$

i.e. the sum of the largest m eigenvalues is larger than or equal to the sum of the m largest diagonal elements. The maximum (and equality) is only achieved if the matrix is diagonal, as in this case the diagonal elements are equal to the eigenvalues.

Therefore, the optimal solution to problem of maximising the relative contributions is obtained by computing the eigendecompositions $\Sigma = U \Lambda U^T$ and $P = G \Theta G^T$ and diagonalise $\Phi^T \Phi = Q_1 \Sigma Q_1^T$ and $\Psi^T \Psi = Q_2 P Q_2^T$ by setting $Q_1 = U^T$ and $Q_2 = G^T$, respectively. This yields for the maximised contributions

$$h = \text{Diag}(\Lambda) = (\lambda_1, \dots, \lambda_d)^T$$

and

$$k = \text{Diag}(\Theta) = (\theta_1, \dots, \theta_d)^T$$

with eigenvalues λ_i and θ_i arranged in decreasing order.

3.4 Natural whitening procedures

We now introduce several strategies (maximise correlation between individual components, maximise compression, structural constraints) to select an optimal whitening procedure.

Specifically, we discuss the following whitening transformations:

- **Mahalanobis** whitening, also known as **ZCA** (zero-phase component analysis) whitening in machine learning (based on covariance)
- **ZCA-cor** whitening (based on correlation)
- **PCA** whitening (based on covariance)

- **PCA-cor** whitening (based on correlation)
- **Cholesky** whitening

Thus, in the following we consider three main types (ZCA, PCA, Cholesky) of whitening.

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

3.4.1 ZCA whitening

Aim: remove correlations and standardise but otherwise make sure that the whitened vector z does not differ too much from the original vector x . Specifically, each latent component z_i should be as close as as possible to the corresponding original variable x_i :

$$\begin{aligned} z_1 &\leftrightarrow x_1 \\ z_2 &\leftrightarrow x_2 \\ z_3 &\leftrightarrow x_3 \\ &\dots \\ z_d &\leftrightarrow x_d \end{aligned}$$

One possible way to implement this is to compute the expected squared difference between the two centered random vectors z_c and x_c .

ZCA objective function: **minimise** $E(\|x_c - z_c\|_F^2)$ to find an optimal whitening procedure.

The ZCA objective function can be simplified as follows:

$$\begin{aligned} E(\|x_c - z_c\|_F^2) &= E(\|x_c\|_F^2) - 2E(\text{Tr}(x_c z_c^T)) + E(\|z_c\|_F^2) \\ &= \text{Tr}(E(x_c x_c^T)) - 2\text{Tr}(E(x_c z_c^T)) + \text{Tr}(E(z_c z_c^T)) \\ &= \text{Tr}(\text{Var}(x)) - 2\text{Tr}(\text{Cov}(x, z)) + \text{Tr}(\text{Var}(z)) \\ &= \text{Tr}(\Sigma) - 2\text{Tr}(\Phi) + d \end{aligned}$$

The same objective function can also be obtained by putting a diagonal constraint on the cross-covariance Φ . Specifically, we are looking for the Φ that is closest to the diagonal matrix I by **minimising**

$$\begin{aligned} \|\Phi - I\|_F^2 &= \|\Phi\|_F^2 - 2\text{Tr}(\Phi^T I) + \|I\|_F^2 \\ &= \text{Tr}(\Sigma) - 2\text{Tr}(\Phi) + d \end{aligned}$$

This will force the off-diagonal elements of Φ to be close to zero and thus leads to sparsity in the cross-covariance matrix.

The only term in the above that depends on the whitening transformation is $-2\text{Tr}(\Phi)$ as Φ is a function of Q_1 . Therefore we can use the following alternative objective:

ZCA equivalent objective: **maximise** $\text{Tr}(\Phi) = \text{Tr}(\Sigma^{1/2} Q_1^T)$ to find the optimal Q_1

Solution:

From the earlier discussion we know that the optimal matrix is

$$\mathbf{Q}_1^{\text{ZCA}} = \mathbf{I}$$

The corresponding whitening matrix for ZCA is therefore

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

and the cross-covariance matrix is

$$\mathbf{\Phi}^{\text{ZCA}} = \mathbf{\Sigma}^{1/2}$$

and the cross-correlation matrix

$$\mathbf{\Psi}^{\text{ZCA}} = \mathbf{V}^{-1/2} \mathbf{\Sigma}^{1/2}$$

Note that $\mathbf{\Sigma}^{1/2}$ is a symmetric positive definite matrix, hence its diagonal elements are all positive. As a result, the diagonals of $\mathbf{\Phi}^{\text{ZCA}}$ and $\mathbf{\Psi}^{\text{ZCA}}$ are positive, i.e. $\text{Cov}(x_i, z_i) > 0$ and $\text{Cor}(x_i, z_i) > 0$. Hence, for ZCA two corresponding components x_i and z_i are always positively correlated!

Proportion of total variation:

For ZCA with $\mathbf{Q}_1 = \mathbf{I}$ we find that $\mathbf{h} = \text{Diag}(\mathbf{\Sigma}) = \sum_{j=1}^d \text{Var}(x_j)$ with $h_i = \text{Var}(x_i)$ hence for ZCA the proportion of total variation contributed by the latent component z_i is the ratio $\frac{\text{Var}(x_i)}{\sum_{j=1}^d \text{Var}(x_j)}$.

Summary:

- ZCA/Mahalanobis transform is the unique transformation that minimises the expected total squared component-wise difference between \mathbf{x}_c and \mathbf{z}_c .
- In ZCA corresponding components in the whitened and original variables are always positively correlated. This facilitates the interpretation of the whitened variables.
- Use ZCA aka Mahalanobis whitening if you want to “just” remove correlations.

3.4.2 ZCA-Cor whitening

Aim: same as above but remove scale in \mathbf{x} first before comparing to \mathbf{z} .

ZCA-cor objective function: **minimise** $\mathbb{E}(\|\mathbf{V}^{-1/2} \mathbf{x}_c - \mathbf{z}_c\|_F^2)$ to find an optimal whitening procedure.

This can be simplified as follows:

$$\begin{aligned} \mathbb{E}(\|\mathbf{V}^{-1/2} \mathbf{x}_c - \mathbf{z}_c\|_F^2) &= \mathbb{E}(\|\mathbf{V}^{-1/2} \mathbf{x}_c\|_F^2) - 2\mathbb{E}(\text{Tr}(\mathbf{V}^{-1/2} \mathbf{x}_c \mathbf{z}_c^T)) + \mathbb{E}(\|\mathbf{z}_c\|_F^2) \\ &= \text{Tr}(\mathbb{E}(\mathbf{V}^{-1/2} \mathbf{x}_c \mathbf{x}_c^T \mathbf{V}^{-1/2})) - 2\text{Tr}(\mathbb{E}(\mathbf{V}^{-1/2} \mathbf{x}_c \mathbf{z}_c^T)) + \text{Tr}(\mathbb{E}(\mathbf{z}_c \mathbf{z}_c^T)) \\ &= \text{Tr}(\text{Cor}(\mathbf{x}, \mathbf{x})) - 2\text{Tr}(\text{Cor}(\mathbf{x}, \mathbf{z})) + \text{Tr}(\text{Var}(\mathbf{z})) \\ &= d - 2\text{Tr}(\mathbf{\Psi}) + d \\ &= 2d - 2\text{Tr}(\mathbf{\Psi}) \end{aligned}$$

The same objective function can also be obtained by putting a diagonal constraint on the cross-correlation $\mathbf{\Psi}$. Specifically, we are looking for the $\mathbf{\Psi}$ that is closest to the diagonal matrix \mathbf{I} by **minimising**

$$\begin{aligned} \|\mathbf{\Psi} - \mathbf{I}\|_F^2 &= \|\mathbf{\Psi}\|_F^2 - 2\text{Tr}(\mathbf{\Psi}^T \mathbf{I}) + \|\mathbf{I}\|_F^2 \\ &= d - 2\text{Tr}(\mathbf{\Psi}) + d \\ &= 2d - 2\text{Tr}(\mathbf{\Psi}) \end{aligned}$$

This will force the off-diagonal elements of $\mathbf{\Psi}$ to be close to zero and thus leads to sparsity in the cross-correlation matrix.

The only term in the above that depends on the whitening transformation is $-2\text{Tr}(\mathbf{\Psi})$ as $\mathbf{\Psi}$ is a function of \mathbf{Q}_2 . Thus we can use the following alternative objective instead:

ZCA-cor equivalent objective: **maximise** $\text{Tr}(\mathbf{\Psi}) = \text{Tr}(\mathbf{P}^{1/2} \mathbf{Q}_2^T)$ to find optimal \mathbf{Q}_2

Solution: same as above for ZCA but using correlation instead of covariance:

From the earlier discussion we know that the optimal matrix is

$$\mathbf{Q}_2^{\text{ZCA-Cor}} = \mathbf{I}$$

The corresponding whitening matrix for ZCA-cor is therefore

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

and the cross-covariance matrix is

$$\mathbf{\Phi}^{\text{ZCA-Cor}} = \mathbf{V}^{1/2} \mathbf{P}^{1/2}$$

and the cross-correlation matrix is

$$\mathbf{\Psi}^{\text{ZCA-Cor}} = \mathbf{P}^{1/2}$$

For the ZCA-cor transformation we also have $\text{Cov}(x_i, z_i) > 0$ and $\text{Cor}(x_i, z_i) > 0$ so that two corresponding components x_i and z_i are always positively correlated!

Proportion of total variation:

For ZCA-cor with $Q_2 = I$ we find that $k = \text{Diag}(\mathbf{P}) = d$ with all $k_i = 1$. Thus, in ZCA-cor each whitened component z_i contributes equally to the total variation $\text{Tr}(\mathbf{P}) = d$, with relative proportion $\frac{1}{d}$.

Summary:

- ZCA-cor whitening is the unique whitening transformation maximising the total correlation between corresponding elements in x and z .
- ZCA-cor leads to interpretable z because each individual element in z is (typically strongly) positively correlated with the corresponding element in the original x .
- As ZCA-cor is explicitly constructed to maximise the total pairwise correlations it achieves higher total correlation than ZCA.
- If x is standardised to $\text{Var}(x_i) = 1$ then ZCA and ZCA-cor are identical.

3.4.3 PCA whitening

Aim: remove correlations and at the same time compress information into a few latent variables. Specifically, we would like that the first latent component z_1 is maximally linked with all variables in x , followed by the second component z_2 and so on:

$$\begin{aligned} z_1 &\rightarrow x_1, x_2, \dots, x_d \\ z_2 &\rightarrow x_1, x_2, \dots, x_d \\ &\dots \\ z_d &\rightarrow x_1, x_2, \dots, x_d \end{aligned}$$

One way to measure the total association of the latent component z_j with all the original x_1, \dots, x_d is the sum of the corresponding squared covariances

$$h_j = \sum_{i=1}^d \text{Cov}(x_i, z_j)^2 = \sum_{i=1}^d \phi_{ij}^2$$

or equivalently the **column sum of squares** of Φ

$$\mathbf{h} = (h_1, \dots, h_d)^T = \text{Diag}(\Phi^T \Phi) = \text{Diag}(Q_1 \Sigma Q_1^T)$$

Each h_j is the contribution of z_j to $\text{Tr}(Q_1 \Sigma Q_1^T) = \text{Tr}(\Sigma)$ i.e. to the total variation based on Σ . As $\text{Tr}(\Sigma)$ is constant this implies that there are only $d - 1$ independent h_j .

In PCA-whitening we wish to concentrate most of the contributions to the total variation based on Σ in a small number of latent components.

PCA whitening objective function: find an optimal optimal Q_1 so that the resulting set $h_1 \geq h_2 \geq \dots \geq h_d$ in $\mathbf{h} = \text{Diag}(Q_1 \Sigma Q_1^T)$ majorizes any other set of relative contributions.

Solution:

Following the earlier discussion we apply Schur's theorem and find the optimal solution by diagonalising $\Phi^T \Phi$ through eigendecomposition of $\Sigma = \mathbf{U} \Lambda \mathbf{U}^T$. Hence, the optimal value for the Q_1 matrix is

$$Q_1^{\text{PCA}} = \mathbf{U}^T$$

However, recall that \mathbf{U} is not uniquely defined — you are free to change the columns signs. The corresponding whitening matrix is

$$\mathbf{W}^{\text{PCA}} = \mathbf{U}^T \Sigma^{-1/2} = \Lambda^{-1/2} \mathbf{U}^T$$

and the cross-covariance matrix is

$$\Phi^{\text{PCA}} = \mathbf{U} \Lambda^{1/2}$$

and the cross-correlation matrix is

$$\Psi^{\text{PCA}} = \mathbf{V}^{-1/2} \mathbf{U} \Lambda^{1/2}$$

Identifiability:

Note that all of the above (i.e. $Q_1^{\text{PCA}}, \mathbf{W}^{\text{PCA}}, \Phi^{\text{PCA}}, \Psi^{\text{PCA}}$) is not unique due to the sign ambiguity in the columns of \mathbf{U} .

Therefore, for identifiability reasons we may wish to impose a further constraint on Q_1^{PCA} or equivalently Φ^{PCA} . A useful condition is to require (for the given ordering of the original variables!) that Q_1^{PCA} has a positive diagonal or equivalently that Φ^{PCA} has a positive diagonal. This implies that $\text{Diag}(\mathbf{U}) > 0$ and $\text{Diag}(\Psi^{\text{PCA}}) > 0$, hence all pairs x_i and z_i are positively correlated.

It is particularly important to pay attention to the sign ambiguity when comparing different computer implementations of PCA whitening (and the related PCA approach).

Note that the actual objective of PCA whitening $\text{Diag}(\Phi^T \Phi)$ is not affected by the sign ambiguity since the column signs of Φ do not matter.

Proportion of total variation:

In PCA whitening the contribution h_i^{PCA} of each latent component z_i to the total variation based on the covariance $\text{Tr}(\Sigma) = \sum_{j=1}^d \lambda_j$ is $h_i^{\text{PCA}} = \lambda_i$. The fraction $\frac{\lambda_i}{\sum_{j=1}^d \lambda_j}$ is the relative contribution of each element in z to explain the total variation.

Thus, low ranking components z_i with small $h_i^{\text{PCA}} = \lambda_i$ may be discarded. In this way PCA whitening achieves both compression and dimension reduction.

Summary:

- PCA whitening is a whitening transformation that maximises compression with the sum of squared cross-covariances as underlying optimality criterion.
- There are sign ambiguities in the PCA whitened variables which are inherited from the sign ambiguities in eigenvectors.
- If a positive-diagonal condition on the orthogonal matrices is imposed then these sign ambiguities are fully resolved and corresponding components z_i and x_i are always positively correlated.

3.4.4 PCA-cor whitening

Aim: same as for PCA whitening but remove scale in x first. This means we use squared correlations rather than squared covariances to measure compression, i.e.

$$k_j = \sum_{i=1}^d \text{Cor}(x_i, z_j)^2 = \sum_{i=1}^d \psi_{ij}^2$$

or in vector notation the **column sum of squares** of Ψ

$$\mathbf{k} = (k_1, \dots, k_d)^T = \text{Diag}(\Psi^T \Psi) = \text{Diag}(Q_2 P Q_2^T)$$

Each k_j is the contribution of z_j to $\text{Tr}(Q_2 P Q_2^T) = \text{Tr}(P) = d$ i.e. the total variation based on P . As $\text{Tr}(P) = d$ is constant this implies that there are only $d - 1$ independent k_j .

In PCA-cor-whitening we wish to concentrate most of the contributions to the total variation based on P in a small number of latent components.

PCA-cor whitening objective function: find an optimal optimal Q_2 so that the resulting set $k_1 \geq k_2 \geq \dots \geq k_d$ in $\mathbf{k} = \text{Diag}(Q_2 P Q_2^T)$ majorizes any other set of relative contributions.

Solution:

Following the earlier discussion we apply Schur's theorem and find the optimal solution by diagonalising $\Psi^T \Psi$ through eigendecomposition of $P = G \Theta G^T$. Hence, the optimal value for the Q_2 matrix is

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

Again G is not uniquely defined — you are free to change signs of the columns. The corresponding whitening matrix is

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

The corresponding whitening matrix is

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

and the cross-covariance matrix is

$$\Phi^{\text{PCA-Cor}} = V^{1/2} G \Theta^{1/2}$$

and the cross-correlation matrix is

$$\Psi^{\text{PCA-Cor}} = G \Theta^{1/2}$$

Identifiability:

As with PCA whitening, there are sign ambiguities in the above because the column signs of G can be freely chosen. For identifiability we may wish to impose further constraints on $Q_2^{\text{PCA-Cor}}$ or equivalently on $\Psi^{\text{PCA-Cor}}$. A useful condition is to require (for the given ordering of the original variables!) that the diagonal elements of $Q_2^{\text{PCA-Cor}}$ are all positive or equivalently that $\Psi^{\text{PCA-Cor}}$ has a positive diagonal. This implies that $\text{Diag}(G) > 0$ and $\text{Diag}(\Phi^{\text{PCA-Cor}}) > 0$.

Note that the actual objective of PCA-cor whitening $\text{Diag}(\Psi^T \Psi)$ is not affected by the sign ambiguity since the column signs of Ψ do not matter.

Proportion of total variation:

In PCA-cor whitening the contribution $k_i^{\text{PCA-Cor}}$ of each latent component z_i to the total variation based on the correlation $\text{Tr}(P) = d$ is $k_i^{\text{PCA-Cor}} = \theta_i$. The fraction $\frac{\theta_i}{d}$ is the relative contribution of each element in z to explain the total variation.

Summary:

- PCA-cor whitening is a whitening transformation that maximises compression with the sum of squared cross-correlations as underlying optimality criterion.
- There are sign ambiguities in the PCA-cor whitened variables which are inherited from the sign ambiguities in the eigenvectors.
- If a positive-diagonal condition on the orthogonal matrices is imposed then these sign ambiguities are fully resolved and corresponding components z_i and x_i are always positively correlated.
- If x is standardised to $\text{Var}(x_i) = 1$, then PCA and PCA-cor whitening are identical.

3.4.5 Cholesky whitening

Cholesky matrix decomposition:

The Cholesky decomposition of a square matrix $A = LL^T$ requires a positive definite A and is unique. L is a lower triangular matrix with positive diagonal elements. Its inverse L^{-1} is also lower triangular with positive diagonal elements.

If D is a diagonal matrix with positive elements then DL is also a lower triangular matrix with a positive diagonal and the Cholesky factor for the matrix DAD .

Aim in Cholesky whitening:

Find a whitening transformation such that the cross-covariance Φ and cross-correlation Ψ have lower triangular structure. Specifically, we wish that the first latent variable z_1 is linked to all original variables x_1, \dots, x_d , the second latent variable z_2 is x_3, \dots, x_d , and so on, the last variable z_d is linked only with x_d .

$$\begin{aligned} z_1 &\rightarrow x_1, x_2, \dots, x_d \\ z_2 &\rightarrow x_2, \dots, x_d \\ &\dots \\ z_d &\rightarrow x_d \end{aligned}$$

We also impose that (as in the other whitening transformations discussed earlier) that the cross-correlation between each pair x_i and z_i is positive.

Thus, Cholesky whitening imposes a structural (sparsity!) constraint on the loadings, where the non-zero coefficients are all in the lower half whereas in the upper half the coefficients all vanish.

Therefore this procedure can be seen as a hybrid with ZCA and PCA whitening.

Solution: In order to find such a whitening transformation we use the Cholesky decomposition and apply it to the covariance matrix $\Sigma = LL^T$

The resulting whitening matrix is

$$W^{\text{Chol}} = L^{-1}$$

By construction, W^{Chol} is a lower triangular matrix with positive diagonal. The whitening constraint is satisfied as $(W^{\text{Chol}})^T W^{\text{Chol}} = (L^{-1})^T L^{-1} = (L^T)^{-1} L^{-1} = (LL^T)^{-1} = \Sigma^{-1}$.

The cross-covariance matrix is the inverse of the whitening matrix

$$\Phi^{\text{Chol}} = L$$

and the cross-correlation matrix is

$$\Psi^{\text{Chol}} = V^{-1/2} L$$

Both Φ^{Chol} and Ψ^{Chol} are lower triangular matrices with positive diagonal elements. Hence two corresponding components x_i and z_i are always positively correlated!

Finally, the corresponding orthogonal matrices are

$$Q_1^{\text{Chol}} = \Phi^T \Sigma^{-1/2} = L^T \Sigma^{-1/2}$$

and

$$Q_2^{\text{Chol}} = \Psi^T P^{-1/2} = L^T V^{-1/2} P^{-1/2}$$

Application to correlation instead of covariance:

We may also apply the Cholesky decomposition to the correlation rather than the covariance matrix. However, unlike for ZCA and PCA this does *not* lead to a different whitening transform:

Let's denote the Cholesky composition of the correlation matrix by $P = L_P L_P^T$. Then the corresponding whitening matrix is $W_p^{\text{Chol}} = L_p^{-1} V^{-1/2}$. As $P = V^{-1/2} \Sigma V^{-1/2} = V^{-1/2} L L^T V^{-1/2}$ we see that $L_P = V^{-1/2} L$ and hence $W_p^{\text{Chol}} = (V^{-1/2} L)^{-1} V^{-1/2} = L^{-1} = W^{\text{Chol}}$.

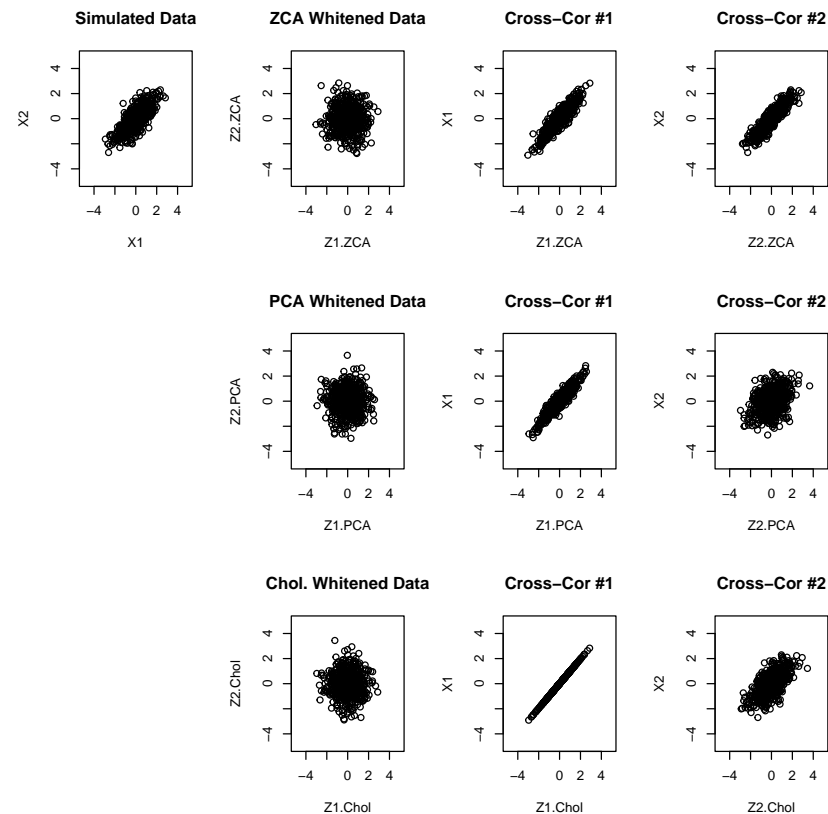
Dependence on the input order:

Cholesky whitening depends on the ordering of input variables. Each ordering of the original variables will yield a different triangular constraint and thus a different Cholesky whitening transform. For example, by inverting the ordering to x_d, x_{d-1}, \dots, x_1 we effectively enforce an upper triangular shape.

An alternative formulation of Cholesky whitening decomposes the precision matrix rather than the covariance matrix. This yields the upper triangular structure directly and is otherwise fully equivalent to Cholesky whitening based on decomposing the covariance.

3.4.6 Comparison of whitening procedures - simulated data

For comparison, here are the results of ZCA, PCA and Cholesky whitening applied to a simulated bivariate normal data set with correlation $\rho = 0.8$.



In column 1 you can see the simulated data as scatter plot.

Column 2 shows the scatter plots of the whitened data — as expected all three methods remove correlation and produce an isotropic covariance.

However, the three approaches differ in the cross-correlations. Columns 3 and 4 show the cross-correlations between the first two corresponding components (x_1 and z_1 , and x_2 and z_2) for ZCA, PCA and Cholesky whitening. As expected, in ZCA both pairs show strong correlation, but this is not the case for PCA and Cholesky whitening.

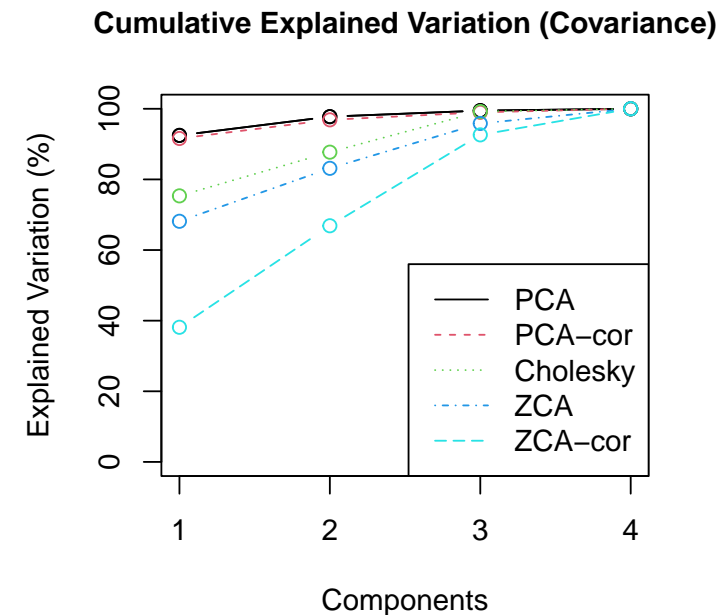
Note that for Cholesky whitening the first component z_1 is perfectly positively correlated with the original component x_1 because the whitening matrix is lower triangular with a positive diagonal and hence z_1 is just x_1 multiplied with a positive constant.

3.4.7 Comparison of whitening procedures - iris flowers

As an example we consider the well known [iris flower data set](#). It consists of botanical measures (sepal length, sepal width, petal length and petal width) for 150 iris flowers comprising three species ([Iris setosa](#), [Iris versicolor](#), [Iris virginica](#)). Hence this data set has dimension $d = 4$ and sample size $n = 150$.

We apply all discussed whitening transforms to this data, and then sort the whitened components by their relative contribution to the total variation. For Cholesky whitening we used the input order for the shape constraint.

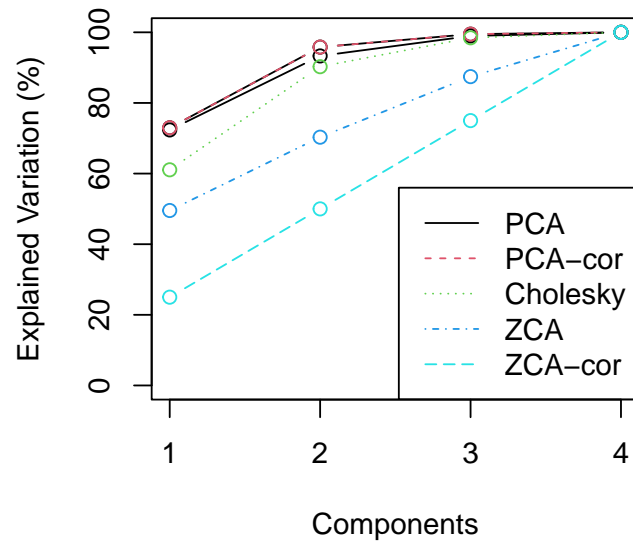
Here are the results for explained variation based on covariance loadings:



As expected, the two PCA whitening approaches compress the data most. On the other end of the spectrum, the ZCA whitening methods are the two least compressing approaches. Cholesky whitening is a compromise between ZCA and PCA in terms of compression.

Similar results are obtained based on correlation loadings - note how ZCA-cor provides equal weight for each latent variable.

Cumulative Explained Variation (Correlation)



3.4.8 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
Cholesky	triangular shaped \mathbf{W} , Φ and Ψ , sparsity

If data are standardised then Φ and Ψ will be the same and hence ZCA will become ZCA-cor and PCA becomes PCA-cor. The triangular shape constraint of Cholesky whitening depends on the ordering of the original variables.

Related methods not discussed in this course:

- Factor models: essentially whitening plus an additional error term, factors have rotational freedom just like in whitening
- Partial Least Squares (PLS): similar to Principal Components Analysis (PCA) but in a regression setting (with the choice of latent variables

depending on the response)

- Nonlinear dimension reduction methods such as SNE, tSNE, UMAP.

3.5 Principal Component Analysis (PCA)

3.5.1 PCA transformation

Principal component analysis was proposed in 1933 by Harald Hotelling⁴ and is very closely related to **PCA whitening**. The underlying mathematics was developed earlier in 1901 by Karl Pearson⁵ for the problem of orthogonal regression.

Assume random vector \mathbf{x} with $\text{Var}(\mathbf{x}) = \Sigma = \mathbf{U}\Lambda\mathbf{U}^T$. PCA is a particular orthogonal transformation of the original \mathbf{x} such that the resulting components are orthogonal:

$$\underbrace{\mathbf{t}^{\text{PCA}}}_{\text{Principal components}} = \underbrace{\mathbf{U}^T}_{\text{Orthogonal matrix}} \mathbf{x}$$

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

Note that while principal components are *orthogonal* they do *not* have unit variance. Instead, the variance of the principal components t_i is equal to the eigenvalues λ_i .

Thus PCA itself is *not* a whitening procedure but it is very closely linked to PCA whitening which is obtained by standardising the principal components to unit variance: $\mathbf{z}^{\text{PCA}} = \Lambda^{-1/2}\mathbf{t}^{\text{PCA}} = \Lambda^{-1/2}\mathbf{U}^T\mathbf{x} = \mathbf{U}^T\Sigma^{-1/2}\mathbf{x} = \mathbf{Q}_1^{\text{PCA}}\Sigma^{-1/2}\mathbf{x} = \mathbf{W}^{\text{PCA}}\mathbf{x}$

Compression properties:

The total variation is $\text{Tr}(\text{Var}(\mathbf{t}^{\text{PCA}})) = \text{Tr}(\Lambda) = \sum_{j=1}^d \lambda_j$. With principle components the fraction $\frac{\lambda_i}{\sum_{j=1}^d \lambda_j}$ can be interpreted as the proportion of variation contributed by each component in \mathbf{t}^{PCA} to the total variation. Thus, low ranking components in \mathbf{t}^{PCA} with low variation may be discarded, thus leading to a reduction in dimension.

⁴Hotelling, H. 1933. Analysis of a complex of statistical variables into principal components. Journal of Educational Psychology 24:417–441 (Part 1) and 24:498–520 (Part 2). <https://doi.org/10.1037/h0071325> and <https://doi.org/10.1037/h0070888>

⁵Pearson, K. 1901. On lines and planes of closest fit to systems of points in space. Philosophical Magazine 2:559–572. <https://doi.org/10.1080/14786440109462720>

3.5.2 Application to data

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

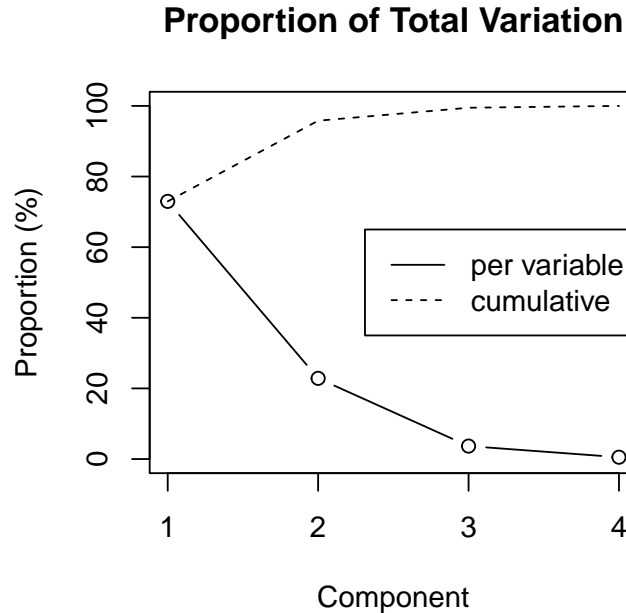
$$\underbrace{T}_{\text{Sample version of principal components}} = \underbrace{X}_{\text{Data matrix}} U$$

There are now two ways to obtain U :

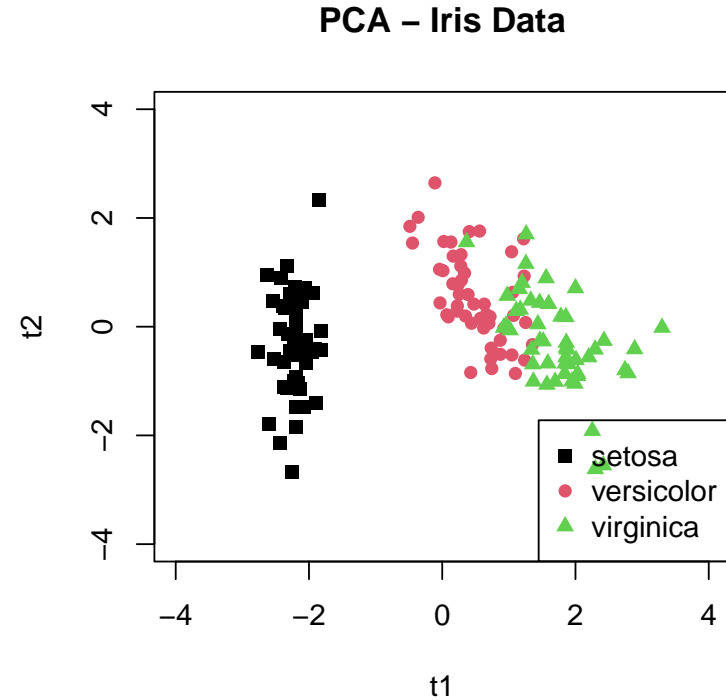
1. Estimate the covariance matrix, e.g. by $\hat{\Sigma} = \frac{1}{n} X_c^T X_c$ where X_c is the column-centred data matrix; then apply the eigenvalue decomposition on $\hat{\Sigma}$ to get U .
2. Compute the singular value decomposition of $X_c = VDU^T$. As $\hat{\Sigma} = \frac{1}{n} X_c^T X_c = U(\frac{1}{n} D^2)U^T$ you can just use U from the SVD of X_c and there is no need to compute the covariance.

3.5.3 Iris flower data example

We first standardise the data, then compute PCA components and plot the proportion of total variation contributed by each component. This shows that only two PCA components are needed to achieve 95% of the total variation:



A scatter plot of the the first two principal components is also informative:



This shows that there groupings among the 150 flowers, corresponding to the species, and that these groups can be characterised by the the principal components.

3.5.4 PCA correlation loadings

In an earlier section we have learned that for a general whitening transformation the cross-correlations $\Psi = \text{Cor}(x, z)$ play the role of correlation loadings in the inverse transformation:

$$V^{-1/2}x = \Psi z,$$

i.e. they are the coefficients linking the whitened variable z with the standardised original variable x . This relationship holds therefore also for PCA-whitening with $z^{\text{PCA}} = \Lambda^{-1/2}U^T x$ and $\Psi^{\text{PCA}} = V^{-1/2}U\Lambda^{1/2}$.

The classical PCA is not a whitening approach because $\text{Var}(t^{\text{PCA}}) \neq I$. However, we can still compute cross-correlations between x and the principal components

t^{PCA} , resulting in

$$\text{Cor}(x, t^{\text{PCA}}) = V^{-1/2} U \Lambda^{1/2} = \Psi^{\text{PCA}}$$

Note these are the same as the cross-correlations for PCA-whitening since t^{PCA} and z^{PCA} only differ in scale.

The inverse PCA transformation is

$$x = U t^{\text{PCA}}$$

In terms of standardised PCA components $z^{\text{PCA}} = \Lambda^{-1/2} t^{\text{PCA}}$ and standardised original components it becomes

$$V^{-1/2} x = \Psi \Lambda^{-1/2} t^{\text{PCA}}$$

Thus the cross-correlation matrix Ψ plays the role of *correlation loadings* also in classical PCA, i.e. they are the coefficients linking the standardised PCA components with the standardised original components.

3.5.5 PCA correlation loadings plot

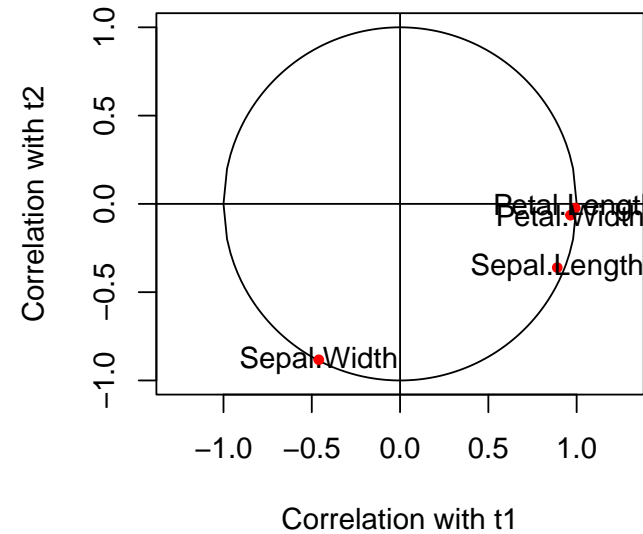
In PCA and PCA-cor whitening as well as in classical PCA the aim is compression, i.e. to find latent variables such that most of the total variation is contributed by a small number of components.

In order to be able to better interpret the top ranking PCA component we can use a visual device called *correlation loadings plot*. For this we compute the correlation between the PCA components 1 and 2 (t_1^{PCA} and t_2^{PCA}) with all original variables x_1, \dots, x_d .

For each original variable x_i we therefore have two numbers between -1 and 1, the correlation $\text{Cor}(x_i, t_1^{\text{PCA}}) = \psi_{i1}$ and $\text{Cor}(x_i, t_2^{\text{PCA}}) = \psi_{i2}$ that we use as coordinates to draw a point in a plane. Recall that the row sums of squares of the correlation loadings Ψ are all identical to 1. Hence, the sum of the squared loadings from just the first two components is also at most 1. Thus, 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.

As an example, here is the correlation loadings plot showing the cross-correlation between the first two PCA components and all four variables of the iris flower data set discussed earlier.

Correlation Loadings Plot Iris Data



The interpretation of this plot is discussed in Worksheet 5.

Chapter 4

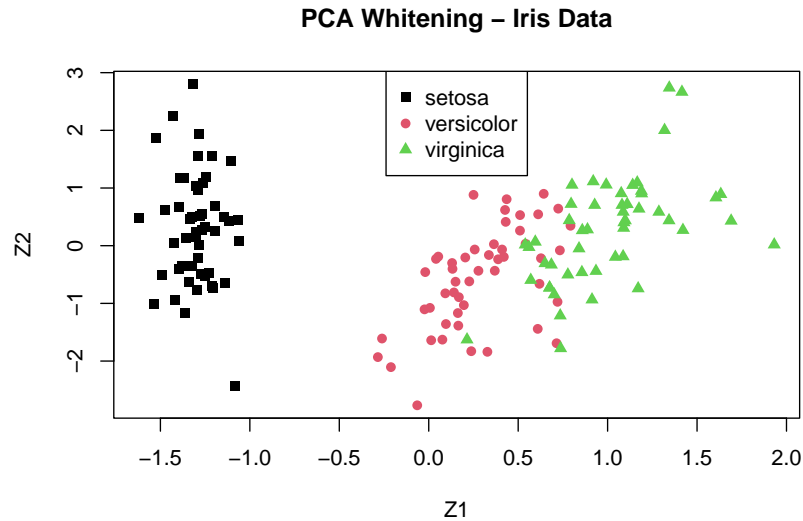
Unsupervised learning and clustering

4.1 Challenges in unsupervised learning

4.1.1 Objective

We observe data x_1, \dots, x_n for n objects (or subjects). Each sample x_i is a vector of dimension d . Thus, for each of the n objects / subjects we have measurements on d variables. The aim of unsupervised learning is to identify patterns relating the objects/subjects based on the information available in x_i . Note that in unsupervised learning we use *only* the information in x_i and nothing else.

For illustration consider the first two principal components of the Iris flower data (see e.g. Worksheet 5):



Clearly there is a group structure among the samples that is linked to particular patterns in the first two principal components.

Note that in this plot we have used additional information, the class labels (setosa, versicolor, virginica), to highlighting the true underlying structure (the three flower species).

In **unsupervised learning** the class labels are (assumed to be) unknown, and the aim is to infer the clustering and thus the classes labels.¹

There are many methods for clustering and unsupervised learning, both purely algorithmic as well as probabilistic. In this chapter we will study a few of the most commonly used approaches.

4.1.2 Questions and problems

In order to implement unsupervised learning we need to address a number of questions:

- how do we define clusters?
- how do we learn / infer clusters?
- how many clusters are there? (this is surprisingly difficult!)
- how can we assess the uncertainty of clusters?

Once we know the clusters we are also interested in:

¹In contrast, in **supervised learning** (to be discussed in a subsequent chapter) the class labels are known for a subset of the data (the training data set) and are required to learn a prediction function.

- which features define / separate each cluster?

(note this is a feature / variable selection problem, discussed in supervised learning).

Many of these problems and questions are highly specific to the data at hand. Correspondingly, there are many different types of methods and models for clustering and unsupervised learning.

In terms of representing the data, unsupervised learning tries to balance between the following two extremes:

- 1) all objects are grouped into a single cluster (low complexity model)
- 2) all objects are put into their own cluster (high complexity model)

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.

4.1.3 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!

\Rightarrow Clustering / partitioning / structure discovery is not easy!

\Rightarrow We cannot expect perfect answers or a single “true” clustering

In fact, as a model of the data many different clusterings may fit the data equally well.

⇒ We need to assess the uncertainty of the clustering

This can be done as part of probabilistic modelling or by resampling (e.g., bootstrap).

4.1.4 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.

4.2 Hierarchical clustering

4.2.1 Tree-like structures

Often, categorisations of objects are nested, i.e. there are sub-categories of categories etc. These can be naturally 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** explaining the biodiversity of life
- phylogenetic trees among species (e.g. vertebrates)
- 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 leaves upwards
- successively form partitions by first joining individual objects together, then recursively join groups of items together, until all is merged.

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.

4.2.2 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 one set. Create an internal node in the tree to represent this set.
- 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. The final node created is the root node.

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

- 1) a distance measure $d(a, b)$ between two individual elementary data points a and b .

This is typically one of the following:

- Euclidean distance $d(a, b) = \sqrt{\sum_{i=1}^d (a_i - b_i)^2} = \sqrt{(a - b)^T (a - b)}$
- Squared Euclidean distance $d(a, b) = (a - b)^T (a - b)$
- Manhattan distance $d(a, b) = \sum_{i=1}^d |a_i - b_i|$

- Maximum norm $d(a, b) = \max_{i \in \{1, \dots, d\}} |a_i - b_i|$

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

- 2) a distance measure $d(A, B)$ between two sets of objects $A = \{a_1, a_2, \dots, a_{n_A}\}$ and $B = \{b_1, b_2, \dots, b_{n_B}\}$ of size n_A and n_B , respectively.

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

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

4.2.3 Ward's clustering method

Another agglomerative hierarchical procedure is **Ward's minimum variance approach**² (see also³). In this approach in each iteration the two sets A and B are merged that lead to the **smallest increase in within-group variation**. The centroids of the two sets is given by $\mu_A = \frac{1}{n_A} \sum_{a_i \in A} a_i$ and $\mu_B = \frac{1}{n_B} \sum_{b_i \in B} b_i$.

The within-group sum of squares for group A is

$$w_A = \sum_{a_i \in A} (a_i - \mu_A)^T (a_i - \mu_A)$$

and is computed here on the basis of the difference of the observations a_i relative to their mean μ_A . However, since we typically only have pairwise distances available we don't know the group means so this formula can't be applied. Fortunately, it is also possible to compute w_A using only the pairwise differences using

$$w_A = \frac{1}{n_A} \sum_{a_i, a_j \in A, i < j} (a_i - a_j)^T (a_i - a_j)$$

This trick is employed in Ward's clustering method by constructing a distance measure between two sets A and B as

$$d(A, B) = w_{A \cup B} - w_A - w_B$$

and using as the distance between two elementary data points a and b the squared Euclidean distance

$$d(a, b) = (a - b)^T (a - b).$$

²Ward, J.H. 1963. Hierarchical grouping to optimize an objective function. JASA 58:236–244.

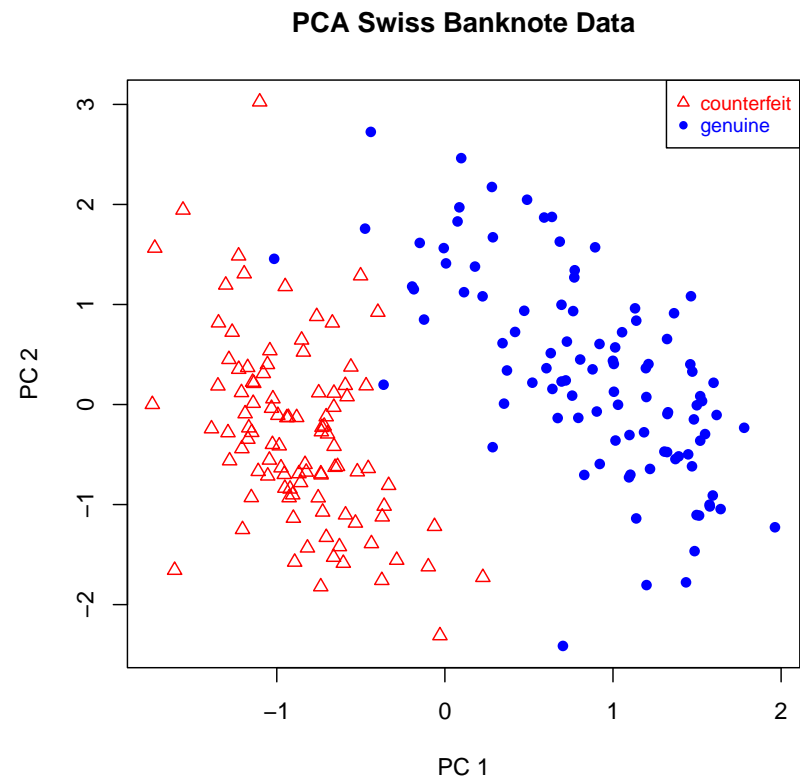
<https://doi.org/10.1080/01621459.1963.10500845>

³F. Murtagh and P Legendre. 2014. Ward's hierarchical agglomerative clustering method: which algorithms implement Ward's criterion? J. Classif. 31:274–295. <https://doi.org/10.1007/s00357-014-9161-z>

4.2.4 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 genuine and 100 are counterfeit. The measurements are: length, left width, right width, bottom margin, top margin, diagonal length of the bank notes.

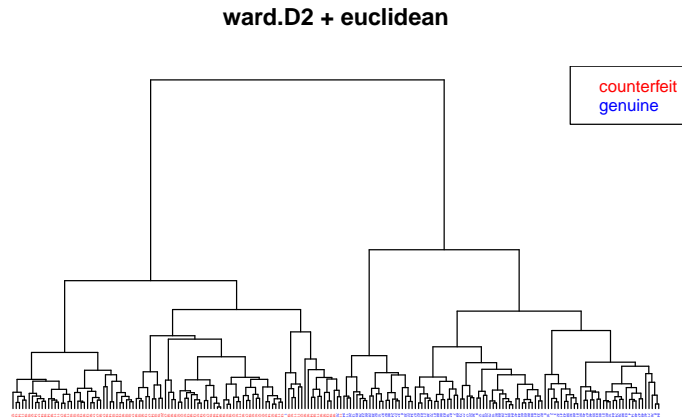
Plotting the first two PCAs of this data shows that there are indeed two well defined groups, and that these groups correspond precisely to the genuine and counterfeit banknotes:



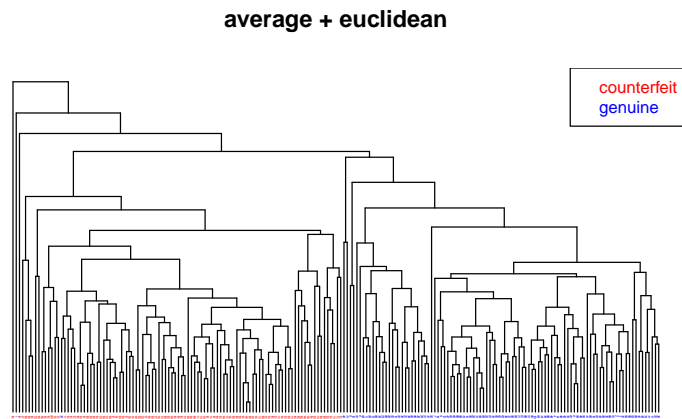
We now compare the hierarchical clusterings of the Swiss bank note data using four different methods using Euclidean distance.

An interactive [R Shiny web app](https://minerva.it.manchester.ac.uk/shiny/strimmer/hclust/) of this analysis (which also allows to explore further distance measures) is available online at <https://minerva.it.manchester.ac.uk/shiny/strimmer/hclust/>.

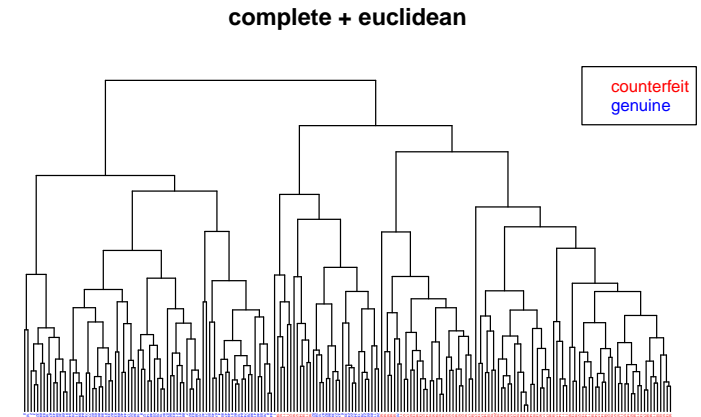
Ward.D2 (=Ward's method):



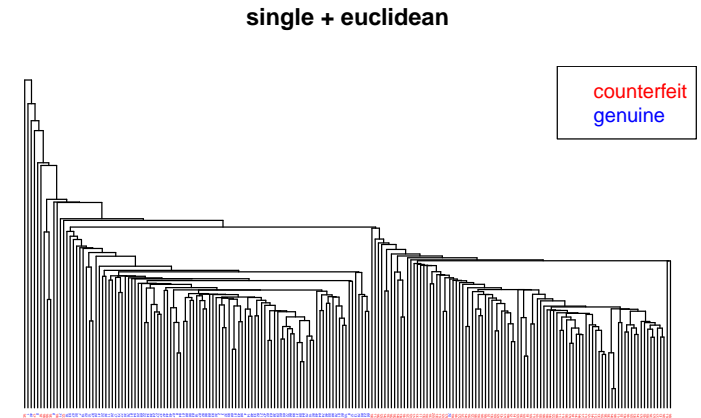
Average linkage:



Complete linkage:



Single linkage:



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).

4.2.5 Assessment of the uncertainty of hierarchical clusterings

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

- Sampling with replacement is used to generate a number of so-called bootstrap data sets (say $B = 200$) similar to the original one. Specifically, we create new data matrices by repeatedly randomly selecting columns (variables) from the original data matrix for inclusion in the bootstrap data matrix. Note that we sample columns as our aim is to cluster the samples.
- Subsequently, a hierarchical clustering is computed for each of the bootstrap data sets. As a result, we now have an “ensemble” of B bootstrap trees.
- Finally, analysis of the clusters (bipartitions) shown in all the bootstrap trees allows to count the clusters that appear frequently, and also those that appear less frequently. These counts provide a measure of the stability of the clusterings appearing in the original tree.
- Additionally, from the bootstrap tree we can also compute a consensus tree containing the most stable clusters. This can be viewed as an “ensemble average” of all the bootstrap trees.

A disadvantage of this procedure is that bootstrapping trees is computationally very expensive, as the original procedure is already time consuming but now needs to be repeated a large number of times.

4.3 K-means clustering

4.3.1 Set-up

- We assume that there are K groups (i.e. K is known in advance).
- For each group $k \in \{1, \dots, K\}$ we assume a group mean μ_k .
- Aim: partition the data points x_1, \dots, x_n into K non-overlapping groups.
- Each of the n data points 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 between the different groups (i.e. each group should differ from the other groups).

4.3.2 Algorithm

After running K -means we will get estimates of $\hat{\mu}_k$ of the group means, as well as allocations $y_i \in \{1, \dots, K\}$ of each data point x_i to one of the classes.

Initialisation:

At the start of the algorithm the n observations x_1, \dots, x_n are randomly allocated with equal probability to one of the K groups. The resulting assignment is y_1, \dots, y_n , with each $y_i \in \{1, \dots, K\}$. With $G_k = \{i | y_i = k\}$ 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 the group allocations y_i . Specifically, assign each data point x_i to the group k with the nearest $\hat{\mu}_k$. The distance is measured in terms of the Euclidean norm:

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

Steps 1 and 2 are repeated until the algorithm converges (i.e. until the group allocations don't change any more) or until a specified upper limit of iterations is reached.

4.3.3 Properties

K -means has been proposed in the 1950 to 1970s by various authors in diverse contexts⁴. Despite its simplicity K -means is, perhaps surprisingly, a very effective clustering algorithm. The main reason for this is the close connection of K -means with probabilistic clustering based on Gaussian mixture models (for details see later section).

Since the clustering depends on the initialisation it is often useful to run K -means several times with different starting group allocations of the data points. Furthermore, non-random or non-uniform initialisations can lead to improved and faster convergence, see the [K-means++](#) algorithm.

The clusters constructed in K -means have linear boundaries and thus form a [Voronoi tessellation](#) around the cluster means. Again, this can be explained by the close link of K -means with a particular Gaussian mixture model.

⁴H.-H. Bock. 2008. Origins and extensions of the k -means algorithm in cluster analysis. JEHPs 4, no. 2. <https://www.jehps.net/Decembre2008/Bock.pdf>

4.3.4 Choosing the number of clusters

Once the K -means algorithm has run we can assess the homogeneity and heterogeneity of the resulting clusters:

- a) the total within-group sum of squares SSW (in R: `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 (in R: `betweenss`), 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 the total sum of squares SST .

- 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).

Dividing 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 $T = B + W$.

In order to decide on the optimal number of clusters we run K -means for different settings for K and then choose the smallest K for which the explained variation B is not significantly worse compared to a clustering with a substantially larger K (see example below).

4.3.5 K-medoids aka PAM

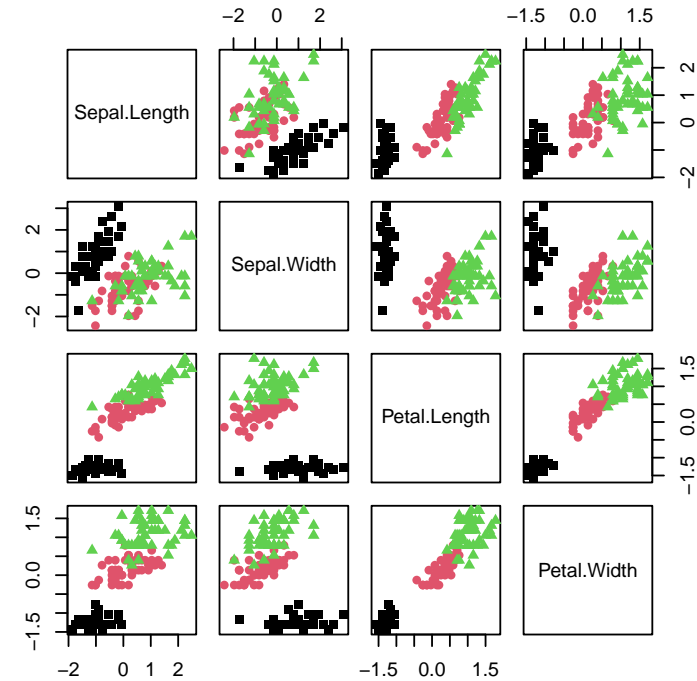
A closely related clustering method is K -medoids or PAM (“Partitioning Around Medoids”).

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.

4.3.6 Application of K -means to Iris data

Scatter plots of Iris data:



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

```
kmeans.out = kmeans(X.iris, 3)
kmeans.out
```

```
## K-means clustering with 3 clusters of sizes 47, 53, 50
##
## Cluster means:
##   Sepal.Length Sepal.Width Petal.Length Petal.Width
```

```
## 1  1.13217737  0.08812645  0.9928284  1.0141287
## 2 -0.05005221 -0.88042696  0.3465767  0.2805873
## 3 -1.01119138  0.85041372 -1.3006301 -1.2507035
##
## Clustering vector:
## [1] 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
## [38] 3 3 3 3 3 3 3 3 3 3 3 3 1 1 1 2 2 2 1 2 2 2 2 2 2 2 1 2 2 2 2 1
## [75] 2 1 1 1 2 2 2 2 2 2 2 1 1 2 2 2 2 2 2 2 2 2 2 2 2 2 1 2 1 1 1 1 2 1
## [112] 1 1 2 2 1 1 1 1 2 1 2 1 2 1 1 1 1 1 1 2 2 1 1 1 2 1 1 1 2 1 1
## [149] 1 2
##
## Within cluster sum of squares by cluster:
## [1] 47.45019 44.08754 47.35062
## (between_SS / total_SS =  76.7 %)
##
## Available components:
##
## [1] "cluster"      "centers"      "totss"        "withinss"     "tot.with
## [6] "betweenss"    "size"         "iter"         "ifault"       "tot.with
```

The corresponding total within-group sum of squares (SSW, tot.withinss) is

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

```
## [1] 138.8884
```

and the between-group sum of squares (SSB, betweenss) is

```
kmeans.out$betweenss
```

```
## [1] 457.1116
```

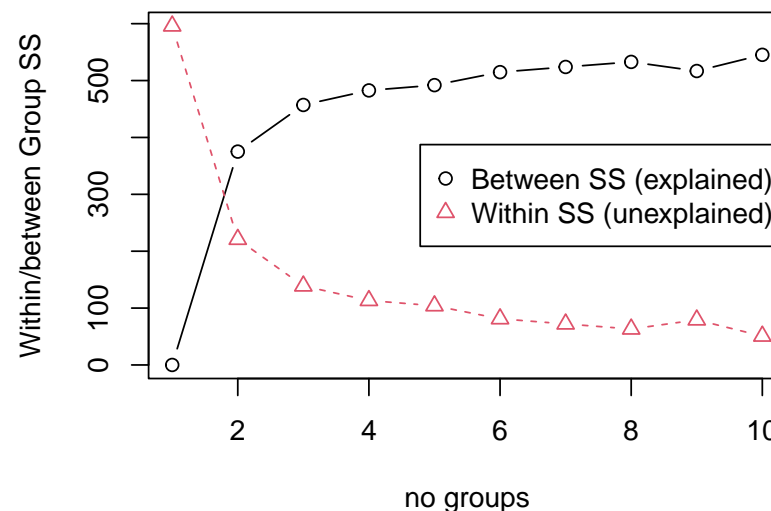
By comparing with the known class assignments we can assess the accuracy of K-means clustering:

```
table(L.iris, kmeans.out$cluster)
```

```
##
## L.iris      1  2  3
## setosa      0  0 50
## versicolor 11 39  0
## virginica   36 14  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.

4.3.7 Arbitrariness of cluster labels and label switching

It is important to realise that in unsupervised learning and clustering the labels of each group are assigned in an arbitrary fashion. Recall that for K groups there are $K!$ possibilities to attach the labels, corresponding to the number of permutations of K groups.

Thus, different runs of a clustering algorithm such as K-means may return the same clustering (groupings of samples) but with different labels. This phenomenon is called “label switching” and makes it difficult to automatise comparison of clusterings. In particular, one cannot simply rely on the automatically assigned group label, instead one needs to compare the actual members of the clusters.

A way to resolve the problem of label switching is to relabel the clusters using additional information, such as requiring that some samples are in specific groups (e.g.: sample 1 is always in group labelled “1”), and/or linking labels to orderings or constraints on the group characteristics (e.g.: the group with label “1” has always a smaller mean than group with label “2”).

4.4 Mixture models

4.4.1 Finite mixture model

- K groups / classes / categories, with finite K known in advance.
- Probability of class k : $\Pr(k) = \pi_k$ with $\sum_{k=1}^K \pi_k = 1$.
- Each class $k \in C = \{1, \dots, K\}$ is modelled by its own distribution F_k with own parameters θ_k .
- Density of class k : $f_k(x) = f(x|k)$.
- The conditional means and variances for each class $k \in C$ are $E(x|k) = \mu_k$ and $\text{Var}(x|k) = \Sigma_k$.
- The resulting mixture density for the observed variable x is

$$f_{\text{mix}}(x) = \sum_{k=1}^K \pi_k f_k(x)$$

Very often one uses **multivariate normal components** $f_k(x) = N(x|\mu_k, \Sigma_k)$
 \Rightarrow **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* (= terminology for a *random effects* regression model).

4.4.2 Total mean and variance of a multivariate mixture model

Using the **law of total expectation** we obtain the mean of the mixture density with multivariate x as follows:

$$\begin{aligned} E(x) &= E(E(x|k)) \\ &= E(\mu_k) \\ &= \sum_{k=1}^K \pi_k \mu_k \\ &= \mu_0 \end{aligned}$$

Note that we treat both x as well as k as random variables.

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

$$\begin{aligned} \underbrace{\text{Var}(x)}_{\text{total}} &= \underbrace{\text{Var}(E(x|k))}_{\text{explained / between-group}} + \underbrace{E(\text{Var}(x|k))}_{\text{unexplained / expected within-group / pooled}} \\ \Sigma_0 &= \text{Var}(\mu_k) + E(\Sigma_k) \\ &= \sum_{k=1}^K \pi_k (\mu_k - \mu_0)(\mu_k - \mu_0)^T + \sum_{k=1}^K \pi_k \Sigma_k \\ &= \Sigma_{\text{explained}} + \Sigma_{\text{unexplained}} \end{aligned}$$

Thus, the **total variance** decomposes into the **explained (between group) variance** and the **unexplained (expected within group, pooled) variance**.

4.4.3 Total variation

The **total variation** is given by the **trace of the covariance matrix**. The above decomposition for the total variation is

$$\begin{aligned} \text{Tr}(\Sigma_0) &= \text{Tr}(\Sigma_{\text{explained}}) + \text{Tr}(\Sigma_{\text{unexplained}}) \\ &= \sum_{k=1}^K \pi_k \text{Tr}((\mu_k - \mu_0)(\mu_k - \mu_0)^T) + \sum_{k=1}^K \pi_k \text{Tr}(\Sigma_k) \\ &= \sum_{k=1}^K \pi_k (\mu_k - \mu_0)^T (\mu_k - \mu_0) + \sum_{k=1}^K \pi_k \text{Tr}(\Sigma_k) \end{aligned}$$

If the covariances are replaced by their empirical estimates we obtain the $T = B + W$ decomposition of total variation familiar from K -means:

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

$$B = \text{Tr}(\hat{\Sigma}_{\text{explained}}) = \frac{1}{n} \sum_{k=1}^K n_k (\hat{\mu}_k - \hat{\mu}_0)^T (\hat{\mu}_k - \hat{\mu}_0)$$

$$W = \text{Tr}(\hat{\Sigma}_{\text{unexplained}}) = \frac{1}{n} \sum_{k=1}^K \sum_{i \in G_k} (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k)$$

4.4.4 Univariate mixture

For a univariate mixture ($d = 1$) with $K = 2$ components we get

$$\mu_0 = \pi_1 \mu_1 + \pi_2 \mu_2 ,$$

$$\sigma_{\text{within}}^2 = \pi_1 \sigma_1^2 + \pi_2 \sigma_2^2 = \sigma_{\text{pooled}}^2 ,$$

also known as pooled variance, and

$$\begin{aligned} \sigma_{\text{between}}^2 &= \pi_1 (\mu_1 - \mu_0)^2 + \pi_2 (\mu_2 - \mu_0)^2 \\ &= \pi_1 \pi_2^2 (\mu_1 - \mu_2)^2 + \pi_2 \pi_1^2 (\mu_1 - \mu_2)^2 \\ &= \pi_1 \pi_2 (\mu_1 - \mu_2)^2 \\ &= \left(\frac{1}{\pi_1} + \frac{1}{\pi_2} \right)^{-1} (\mu_1 - \mu_2)^2 \end{aligned}$$

The ratio of the between-group variance and the within-group variance is proportional (by factor of n) to the squared pooled-variance t -score:

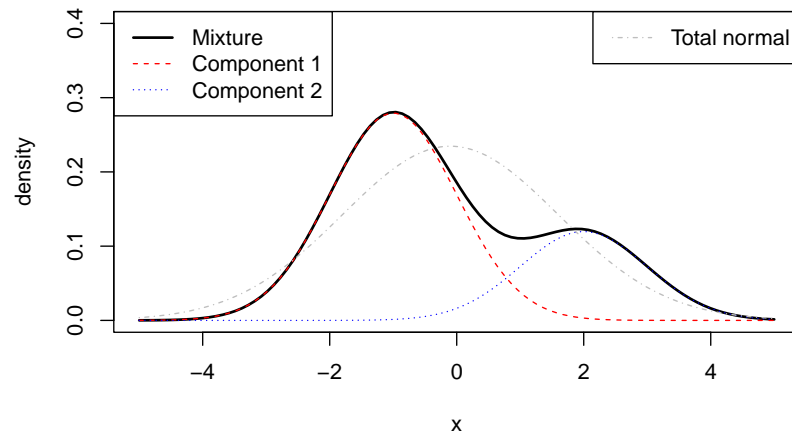
$$\frac{\sigma_{\text{between}}^2}{\sigma_{\text{within}}^2} = \frac{(\mu_1 - \mu_2)^2}{\left(\frac{1}{\pi_1} + \frac{1}{\pi_2}\right) \sigma_{\text{pooled}}^2} = \frac{t_{\text{pooled}}^2}{n}$$

If you are familiar with ANOVA (e.g. linear models course) you will recognise this ratio as the F -score.

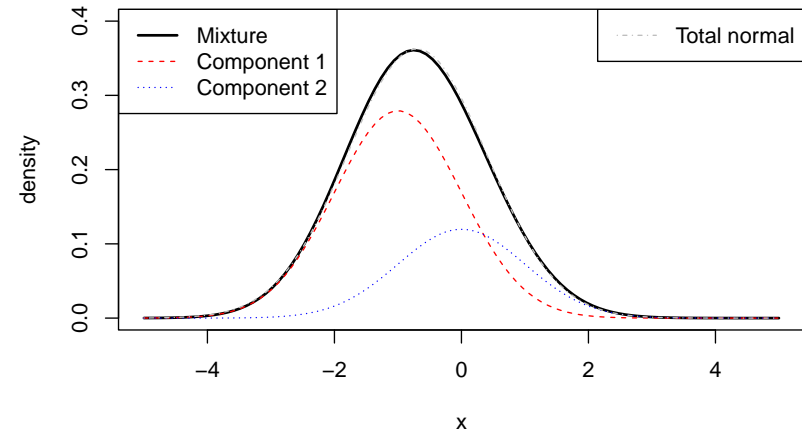
4.4.5 Example of mixtures

Mixtures can take on many different shapes and forms, so it is instructive to study a few examples. An interactive tool to visualise two component normal mixture is available online as [R Shiny web app](https://minerva.it.manchester.ac.uk/shiny/strimmer/mixture/) at <https://minerva.it.manchester.ac.uk/shiny/strimmer/mixture/>.

Bimodal Mixture of Two Normals



Unimodal Mixture of Two Normals



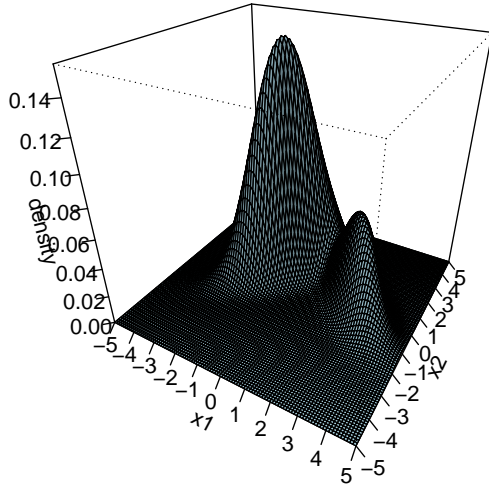
The first plot shows the bimodal density of a mixture distribution consisting of two normals with $\pi_1 = 0.7$, $\mu_1 = -1$, $\mu_2 = 2$ and the two variances equal to 1 ($\sigma_1^2 = 1$ and $\sigma_2^2 = 1$). Because the two components are well-separated there are two clear modes. The plot also shows the density of a normal distribution with the same total mean ($\mu_0 = -0.1$) and variance ($\sigma_0^2 = 2.89$) as the mixture distribution. Clearly the total normal and the mixture density are very different.

However, a two-component mixtures can also be unimodal. For example, if the mean of the second component is adjusted to $\mu_2 = 0$ then there is only a single mode and the total normal density with $\mu_0 = -0.7$ and $\sigma_0^2 = 1.21$ is now almost indistinguishable in form from the mixture density. Thus, in this case it will be very hard (or even impossible) to identify the two peaks from data.

Most mixtures we consider in this course are multivariate. For illustration, here is a plot of a mixture of two bivariate normals, with $\pi_1 = 0.7$, $\mu_1 = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$,

$$\Sigma_1 = \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix}, \mu_2 = \begin{pmatrix} 2.5 \\ 0.5 \end{pmatrix} \text{ and } \Sigma_2 = \begin{pmatrix} 1 & -0.7 \\ -0.7 & 1 \end{pmatrix}:$$

Mixture of Two Bivariate Normals



4.4.6 Sampling from a mixture model

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

This is done in a two-step process:

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

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

Here the vector $z = (z_1, \dots, z_K)^T$ indicates a hard group 0/1 allocation, with all components $z_{\neq k} = 0$ except for a single entry $z_k = 1$.

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

$$x \sim F_k$$

This two-stage sampling approach is also known as hierarchical generative model for a mixture distribution. This generative view is not only useful for simulating data from a mixture model but also highlights the role of the latent variable (the class allocation).

4.5 Fitting mixture models to data and inferring the latent states

In the following we denote by

- $X = (x_1, \dots, x_n)^T$ the data matrix containing the observations of n independent and identically distributed samples x_1, \dots, x_n , and
- $y = (y_1, \dots, y_n)^T$ the associated group memberships, as well as
- the parameters θ which for a Gaussian mixture model are $\theta = \{\pi, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K\}$.

4.5.1 Observed and latent variables

When we observe data from a mixture model we collect samples x_1, \dots, x_n . Associated with each observed x_i is the corresponding underlying class allocation y_1, \dots, y_n where each y_i takes on a value from $C = \{1, \dots, K\}$. Crucially, the class allocations y_i are unknown and cannot be directly observed, thus are latent.

- The **joint density** for observed and unobserved variables:

$$f(x, y) = f(x|y)\Pr(y) = f_y(x)\pi_y$$

The mixture density is therefore a **marginal density** as it arises from the joint density $f(x, y)$ by marginalising over the discrete variable y .

- Marginalisation: $f(x) = \sum_{y \in C} f(x, y)$

4.5.2 Complete data likelihood and observed data likelihood

If we know y in advance, i.e. if we know which sample belongs to a particular group, we can construct a *complete data log-likelihood* based on the joint distribution $f(x, y) = \pi_y f_y(x)$. The log-likelihood for θ given the both X and y is

$$\log L(\theta|X, y) = \sum_{i=1}^n \log f(x_i, y_i) = \sum_{i=1}^n \log (\pi_{y_i} f_{y_i}(x_i))$$

On the other hand, typically we do not know y and therefore use the marginal or mixture density $f(x)$ to construct the *observed data log-likelihood* (sometimes also called *incomplete data log-likelihood*) $f(x|\theta)$ as

$$\begin{aligned} \log L(\theta|X) &= \sum_{i=1}^n \log f(x_i|\theta) \\ &= \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k f_k(x_i) \right) \end{aligned}$$

The *observed data log-likelihood* can also be computed from the complete data likelihood function by marginalising over \mathbf{y}

$$\begin{aligned}\log L(\boldsymbol{\theta}|\mathbf{X}) &= \log \sum_{\mathbf{y}} L(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}) \\ &= \log \sum_{y_1, \dots, y_K} \prod_{i=1}^n f(\mathbf{x}_i, y_i) \\ &= \log \prod_{i=1}^n \sum_{k=1}^K f(\mathbf{x}_i, k) \\ &= \sum_{i=1}^n \log \left(\sum_{k=1}^K f(\mathbf{x}_i, k) \right)\end{aligned}$$

Clustering with a mixture model can be viewed as an *incomplete* or *missing* data problem (see also [MATH27720 Statistics 2](#)).

Specifically, we face the problem of

- fitting the model using only the observed data \mathbf{X} and
- simultaneously inferring the class allocations \mathbf{y} , i.e. states of the latent variable.

4.5.3 Fitting the mixture model to the observed data

For large sample size n the standard way to fit a mixture model is to employ maximum likelihood to find the MLEs of the parameters of the mixture model.

The direct way to fit a mixture model by maximum likelihood is to **maximise the observed data log-likelihood function** with regard to $\boldsymbol{\theta}$:

$$\hat{\boldsymbol{\theta}}^{ML} = \arg \max_{\boldsymbol{\theta}} \log L(\boldsymbol{\theta}|\mathbf{X})$$

Unfortunately, in practise evaluation and optimisation of the log-likelihood function can be difficult due to a number of reasons:

- The form of the observed data log-likelihood function prevents analytic simplifications (note the sum inside the logarithm) and thus can be difficult to compute.
- Because of the symmetries due to exchangeability of cluster labels the likelihood function is multimodal and thus hard to optimise. Note this is also linked to the general problem of label switching and non-identifiability of cluster labels — see the discussion for K -means clustering.
- Further identifiability issues can arise if (for instance) two neighboring components of the mixture model are largely overlapping and thus are too

close to each other to be discriminated as two different modes. In other words, it is difficult to determine the number of classes.

- Furthermore, the likelihood in Gaussian mixture models is singular if one of the fitted covariance matrices becomes singular. However, this can be easily addressed by using some form of regularisation (Bayes, penalised ML, etc.) or simply by requiring sufficient sample size per group.

4.5.4 Predicting the group allocation of a given sample

In probabilistic clustering the aim is to infer the latent states y_1, \dots, y_n for all observed samples $\mathbf{x}_1, \dots, \mathbf{x}_n$.

Assuming that the mixture model is known (either in advance or after fitting it) Bayes' theorem allows predict the probability that an observation \mathbf{x}_i falls in group $k \in \{1, \dots, K\}$:

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

Thus, for each of the n samples we get a probability mass function over the K classes with $\sum_{k=1}^K q_i(k) = 1$.

The posterior probabilities in $q_i(k)$ provide a so-called *soft assignment* of the sample \mathbf{x}_i to all classes rather than a 0/1 *hard assignment* to a specific class (as for example in the K -means algorithm).

To obtain at a hard clustering and to infer the most probable latent state we select the class with the highest probability

$$y_i = \arg \max_k q_i(k)$$

Thus, in probabilistic clustering we directly obtain an assessment of the uncertainty of the class assignment for a sample \mathbf{x}_i (which is not the case in simple algorithmic clustering such K -means). We can use this information to check whether there are several classes with equal or similar probability. This will be the case, e.g., if \mathbf{x}_i lies near the boundary between two neighbouring classes.

Using the interactive Shiny app for the univariate normal component mixture (online at <https://minerva.it.manchester.ac.uk/shiny/strimmer/mixture/>) you can explore the posterior probabilities of each class.

4.6 Application of Gaussian mixture models

4.6.1 Choosing the number of classes

In an application of a GMM we need to select a suitable value for K , i.e. 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:

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

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

In order to choose a suitable model we evaluate different models with different K and then choose the model that minimises AIC or 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.

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

4.6.2 Application of GMMs to Iris flower data

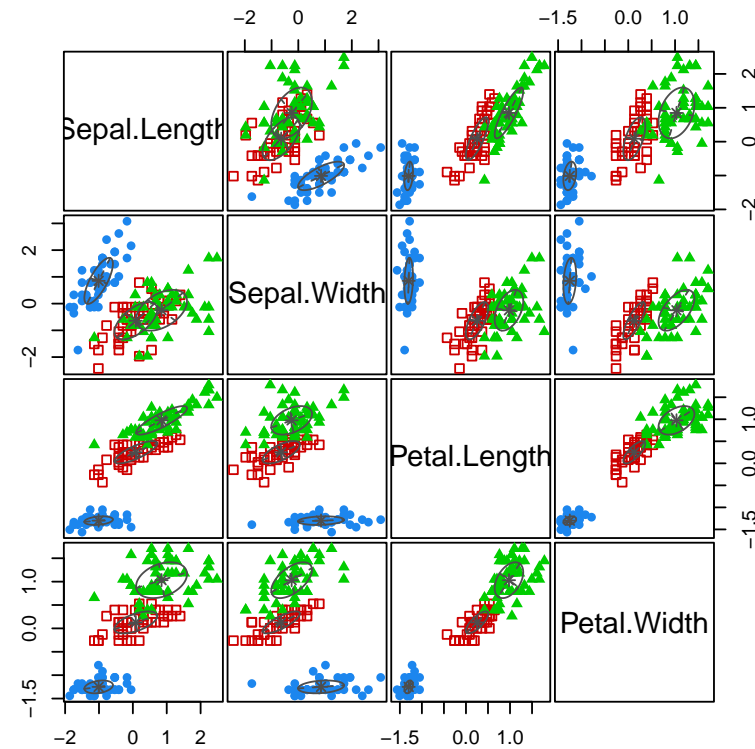
We now explore the application of Gaussian mixture models to the Iris flower data set we also investigated with PCA and K-means.

First, we fit a GMM with 3 clusters, using the R software “mclust” ⁵.

```
data(iris)
X.iris = scale((iris[, 1:4]), scale=TRUE) # center and standardise
L.iris = iris[, 5]

library("mclust")
gmm3 = Mclust(X.iris, G=3, verbose=FALSE)
plot(gmm3, what="classification")
```

⁵L. Scrucca L. et. al. 2016. mclust 5: Clustering, classification and density estimation using Gaussian finite mixture models. The R Journal 8:205–233. See <https://journal.r-project.org/archive/2016/RJ-2016-021/> and <https://mclust-org.github.io/mclust/>



The “mclust” software has used the following model when fitting the mixture:

```
gmm3$modelName
```

```
## [1] "VVV"
```

Here “VVV” is the name used by the “mclust” software for a model allowing for an individual unrestricted covariance matrix Σ_k for each class k .

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

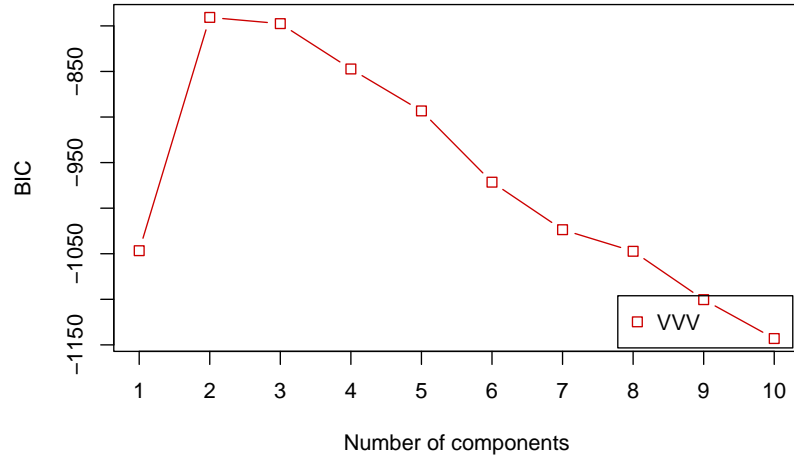
```
table(gmm3$classification, L.iris)
```

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

Note that in “mclust” 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 compute BIC for various numbers of groups we find that the model with the best $\text{BIC}_{\text{mclust}}$ is a model with 2 clusters but the model with 3 cluster has nearly as good a BIC:



4.7 The EM algorithm

4.7.1 Motivation

As discussed above, the observed data log-likelihood can be difficult to maximise directly due to its form as a log marginal likelihood. Intriguingly, it is possible to optimise it *indirectly* using the complete data log-likelihood, and what's more this also allows in many cases for an analytic expression of the maximisation step.

This method is called the EM algorithm and has been formally proposed and described by Arthur Dempster (1929–) and others in 1977⁶ but the algorithm was already known prior. It iteratively estimates both the parameters of the mixture model and the latent states. The key idea behind the EM algorithm is to capitalise on the simplicity of the complete data likelihood and to obtain

⁶Dempster, A. P, N. M. Laird and D. B. Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. JRSS B 39:1–38. <https://doi.org/10.1111/j.2517-6161.1977.tb01600.x>

estimates of θ by imputing the missing group allocations and then subsequently iteratively refining both the imputations and the estimates of θ .

More precisely, in the EM (=expectation-maximisation) algorithm we alternate between

- Step 1) updating the soft allocations of each sample using the current estimate of the parameters θ (obtained in step 2)
- Step 2) updating the parameter estimates by maximising the *expected* complete data log-likelihood. The expectation is taken with regard to the distribution over the latent states (obtained in step 1). Thus the complete data log-likelihood is averaged over the soft class assignments. For an exponential family (e.g. when the distributions for each group are normal) maximisation of the expected complete data log-likelihood can even be done analytically.

4.7.2 The EM algorithm

Specifically, the EM algorithm proceeds as follows:

1) Initialisation:

- Start with a guess of the parameters $\hat{\theta}^{(1)}$, then continue with “E” Step, Part A.
- Alternatively, start with a guess of the soft allocations for each sample $q_i(k)^{(1)}$, collected in the matrix $Q^{(1)}$, then continue with “E” Step, Part B. This may be derived from some prior information, e.g., from running *K*-means. Caveat: some particular initialisations correspond to invariant states and hence should be avoided (see further below).

2) E “expectation” step

- Part A: Use Bayes’ theorem to compute new probabilities of allocation to class k for all the samples x_i :

$$q_i(k)^{(b+1)} \leftarrow \frac{\hat{\pi}_k^{(b)} f_k(x_i | \hat{\theta}^{(b)})}{f(x_i | \hat{\theta}^{(b)})}$$

Note that to obtain $q_i(k)^{(b+1)}$ the current estimate $\hat{\theta}^{(b)}$ of the parameters of the mixture model is required.

- Part B: Construct the *expected* complete data log-likelihood function for θ using the soft allocations $q_i(k)^{(b+1)}$ collected in the matrix $Q^{(b+1)}$:

$$G(\theta | X, Q^{(b+1)}) = \sum_{i=1}^n \sum_{k=1}^K q_i(k)^{(b+1)} \log(\pi_k f_k(x_i))$$

Note that in the case that the soft allocations $Q^{(b+1)}$ turn into hard 0/1 allocations then $G(\theta|X, Q^{(b+1)})$ becomes equivalent to the complete data log-likelihood.

- 3) **M “maximisation” step** — Maximise the expected complete data log-likelihood to update the estimates of mixture model parameters:

$$\hat{\theta}^{(b+1)} \leftarrow \arg \max_{\theta} G(\theta|X, Q^{(b+1)})$$

- 4) Continue with 2) “E” Step until the series $\hat{\theta}^{(1)}, \hat{\theta}^{(2)}, \hat{\theta}^{(3)}, \dots$ has converged.

Crucially, maximisation of the expected complete data log-likelihood is typically much easier than maximisation of the observed data log-likelihood, and in many cases even analytically tractable, which is why in this case the EM algorithm is often preferred over direct maximisation of the observed data log-likelihood.

Note that to avoid singularities in the expected log-likelihood function we may need to adopt regularisation (i.e. penalised maximum likelihood or Bayesian learning) for estimating the parameters in the M-step.

4.7.3 EM algorithm for multivariate normal mixture model

We now consider the EM algorithm applied to the case when the conditional group distributions are normal, i.e. when applied to the Gaussian mixture model (GMM). In this case the two iterative steps in the EM algorithm can be expressed as follows:

E-step:

Update the soft allocations:

$$q_i(k)^{(b+1)} = \frac{\hat{\pi}_k^{(b)} N(x_i | \hat{\mu}_k^{(b)}, \hat{\Sigma}_k^{(b)})}{\hat{f}^{(b)}(x_i | \hat{\pi}_1^{(b)}, \dots, \hat{\pi}_K^{(b)}, \hat{\mu}_1^{(b)}, \dots, \hat{\mu}_K^{(b)}, \hat{\Sigma}_1^{(b)}, \dots, \hat{\Sigma}_K^{(b)})}$$

Correspondingly, the number of samples assigned to class k in the current step is

$$n_k^{(b+1)} = \sum_{i=1}^n q_i(k)^{(b+1)}$$

Note this is not necessarily an integer because of the soft allocations of samples to groups.

The expected complete data log-likelihood becomes:

$$G(\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K | X, Q^{(b+1)}) = \sum_{i=1}^n \sum_{k=1}^K q_i(k)^{(b+1)} \log(\pi_k f_k(x_i))$$

with

$$\log(\pi_k f_k(x_i)) = -\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1}(x - \mu_k) - \frac{1}{2} \log \det(\Sigma_k) + \log(\pi_k)$$

(Remark: we will encounter this expression again in the next chapter when discussing quadratic discriminant analysis).

M-step:

The maximisation of the expected complete data log-likelihood can be done analytically as it is a weighted version of the conventional single group multivariate normal log-likelihood. The resulting estimators are also weighted variants of the usual MLEs.

The updated estimates of the group probabilities are

$$\hat{\pi}_k^{(b+1)} = \frac{n_k^{(b+1)}}{n}$$

The updated estimates of the means are

$$\hat{\mu}_k^{(b+1)} = \frac{1}{n_k^{(b+1)}} \sum_{i=1}^n q_i(k)^{(b+1)} x_i$$

and the updated covariance estimates are

$$\hat{\Sigma}_k^{(b+1)} = \frac{1}{n_k^{(b+1)}} \sum_{i=1}^n q_i(k)^{(b+1)} (x_i - \mu_k^{(b+1)}) (x_i - \mu_k^{(b+1)})^T$$

Note that if $q_i(k)$ is a hard allocation (so that for any i only one class has weight 1 and all others weight 0) then all estimators above reduce to the usual MLEs.

In Worksheet 8 you can find a simple R implementation of the EM algorithm for univariate normal mixtures.

Similar analytical expressions can also be found in more general mixtures where the components are exponential families. As mentioned above this is one of the advantages of using the EM algorithm.

4.7.4 Convergence and invariant states

Under mild assumptions the EM algorithm is guaranteed to monotonically converge to local optima of the observed data log-likelihood⁷. Thus the series $\hat{\theta}^{(1)}, \hat{\theta}^{(2)}, \hat{\theta}^{(3)}, \dots$ converges to the estimate $\hat{\theta}$ found when maximising the observed data log-likelihood. However, the speed of convergence in the EM

⁷Wu, C.F. 1983. On the convergence properties of the EM algorithm. The Annals of Statistics 11:95–103. <https://doi.org/10.1214/aos/1176346060>

algorithm can sometimes be slow, and there are also situations in which there is no convergence at all to $\hat{\theta}$ because the EM algorithm remains in an invariant state.

An example of such an invariant state for a Gaussian mixture model is uniform initialisation of the latent variables $q_i(k) = \frac{1}{K}$, where K is the number of classes. With this we get in the M step $n_k = \frac{n}{K}$ and as parameter estimates

$$\begin{aligned}\hat{\pi}_k &= \frac{1}{K} \\ \hat{\mu}_k &= \frac{1}{n} \sum_{i=1}^n x_i = \bar{x} \\ \hat{\Sigma}_k &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T = \hat{\Sigma}\end{aligned}$$

Crucially, none of these actually depend on the group k ! Thus, in the E step when the next soft allocations are determined this leads to

$$q_i(k) = \frac{\frac{1}{K} N(x_i | \bar{x}, \hat{\Sigma})}{\sum_{j=1}^K \frac{1}{K} N(x_i | \bar{x}, \hat{\Sigma})} = \frac{1}{K}$$

After one cycle in the EM algorithm we arrive at the same soft allocation that we started with, and the algorithm is trapped in an invariant state! Therefore uniform initialisation should clearly be avoided!

You will explore this effect in practise in Worksheet 8.

4.7.5 Connection with the K -means clustering method

The K -means algorithm is very closely related to the EM algorithm and probabilistic clustering with a specific Gaussian mixture models.

Specifically, we assume a simplified model where the probabilities π_k of all classes are equal (i.e. $\pi_k = \frac{1}{K}$) and where the covariances Σ_k are all of the same spherical form $\sigma^2 I$. Thus, the covariance does not depend on the group, there is no correlation between the variables and the variance of all variables is the same.

First, we consider the “E” step. Using the mixture model above the soft assignment for the class allocation becomes

$$\log(q_i(k)) = -\frac{1}{2\sigma^2} (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k) + \text{const}$$

where const does not depend on k . This can be turned into a hard class allocation by

$$\begin{aligned}y_i &= \arg \max_k \log(q_i(k)) \\ &= \arg \min_k (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k)\end{aligned}$$

which is exactly the K -means rule to allocate samples to groups.

Second, in the “M” step we compute the parameters of the model. If the class allocations are hard the expected log-likelihood becomes the observed data likelihood and the MLE of the group mean is the average of samples in that group.

Thus, K -means can be viewed as an EM type algorithm to provide hard classification based on a simple restricted Gaussian mixture model.

4.7.6 Why the EM algorithm works — an entropy point of view

The iterative (soft) imputation of the latent states in the EM algorithm is intuitively clear.

However, in order to get a better understanding of EM we need to demonstrate

- i) why the expected observed log-likelihood needs to be maximised rather than, e.g., the observed log-likelihood with hard allocations, and
- ii) that applying the EM algorithm versus directly maximising the marginal likelihood both lead to the same fitted mixture model.

Intriguingly, both these aspects of the EM algorithm are easiest to understand from an entropy point of view, i.e. considering the entropy foundations of maximum likelihood and Bayesian learning — for details see [MATH27720 Statistics 2](#). Specifically, the reason for the need of using expectation is the link between likelihood and cross-entropy (which also is defined as expectation). Furthermore, the EM algorithm is an example of using an ELBO (“evidence lower bound”) to successively approximate the maximised marginal log-likelihood, with the bound only getting better in each step.

First, recall that the method of maximum likelihood results from minimising the KL divergence between an empirical distribution Q_x representing the observations x_1, \dots, x_n and the model family F_x^θ with parameters θ :

$$\hat{\theta}^{ML} = \arg \min_{\theta} D_{KL}(Q_x, F_x^\theta)$$

The KL divergence decomposes into a cross-entropy and an entropy part

$$D_{KL}(Q_x, F_\theta) = H(Q_x, F_\theta) - H(Q_x)$$

hence minimising the KL divergence with regard to θ is the same as maximising the function

$$\begin{aligned}-nH(Q_x, F_x^\theta) &= nE_{Q_x}(\log f(x|\theta)) \\ &= \sum_{i=1}^n \log f(x_i|\theta) \\ &= \log L(\theta|X)\end{aligned}$$

which is indeed the observed data log-likelihood for θ .

Second, we recall the chain rule for the KL divergence. Specifically, the KL divergence for the joint model forms an upper bound of the KL divergence for the marginal model:

$$\begin{aligned} D_{\text{KL}}(Q_{x,y}, F_{x,y}^{\theta}) &= D_{\text{KL}}(Q_x, F_x^{\theta}) + \underbrace{D_{\text{KL}}(Q_{y|x}, F_{y|x}^{\theta})}_{\geq 0} \\ &\geq D_{\text{KL}}(Q_x, F_x^{\theta}) \end{aligned}$$

Unlike for x we do not have observations about the latent state y . Nonetheless, we can model the joint distribution $Q_{x,y} = Q_x Q_{y|x}$ by assuming a distribution $Q_{y|x}$ over the latent variable.

The EM algorithm arises from iteratively decreasing the joint KL divergence $D_{\text{KL}}(Q_x Q_{y|x}, F_{x,y}^{\theta})$ with regard to both $Q_{y|x}$ and θ :

- 1) “E” Step: While keeping θ fixed we vary $Q_{y|x}$ to minimise the joint KL divergence. The minimum is reached at $D_{\text{KL}}(Q_{y|x}, F_{y|x}^{\theta}) = 0$. This is the case for $Q_{y|x} = F_{y|x}^{\theta}$, i.e. when the latent distribution $Q_{y|x}$ representing the soft allocations is computed by conditioning, i.e. using Bayes’ theorem.
- 2) “M” Step: While keeping $Q_{y|x}$ fixed the joint KL divergence is further minimised with regard to θ . This is equivalent to maximising the function $\sum_{k=1}^K \sum_{i=1}^n q(k|x_i) \log f(x_i, k|\theta)$ which is indeed the expected complete data log-likelihood.

Thus in the “E” step the first argument in the KL divergence is optimised (“I” projection) and in the “M” step it is the second argument that is optimised (“M” projection). In both steps the joint KL divergence always decreases and never increases. Furthermore, at the end of “E” step the joint KL divergence equals the marginal KL divergence. Thus, this procedure implicitly minimises the marginal KL divergence as well, and hence EM maximises the marginal log-likelihood.

An alternative way to look at the EM algorithm is in terms of cross-entropy. Using $H(Q_{x,y}) = H(Q_x) + H(Q_{y|x})$ we can rewrite the above upper bound for the joint KL divergence as an equivalent lower bound for n times the negative marginal cross-entropy:

$$\begin{aligned} -nH(Q_x, F_x^{\theta}) &= \underbrace{-nH(Q_x Q_{y|x}, F_{x,y}^{\theta})}_{\text{lower bound, ELBO}} + \underbrace{nH(Q_{y|x}) + nD_{\text{KL}}(Q_{y|x}, F_{y|x}^{\theta})}_{\geq 0} \\ &\geq \mathcal{F}(Q_x, Q_{y|x}, F_{x,y}^{\theta}) \end{aligned}$$

The lower bound is known as the “ELBO” (“evidence lower bound”). Then the EM algorithm arises by iteratively maximising the ELBO \mathcal{F} with regard to $Q_{y|x}$ (“E” step”) and θ (“M” step).

The entropy interpretation of the EM algorithm is due to Csiszàr and Tusnàdy (1984)⁸ and the ELBO interpretation was introduced by Neal and Hinton (1998)⁹.

⁸Csiszàr, I., and G. Tusnàdy. 1984. Information geometry and alternating minimization procedures. In Dudewicz, E. J. et al. (eds.) *Recent Results in Estimation Theory and Related Topics Statistics and Decisions*, Supplement Issue No. 1. pp. 205–237.

⁹Neal, R. M., and G. Hinton. 1998. A view of the EM algorithm that justifies incremental, sparse, and other variants. In Jordan, M.I. (eds.). *Learning in Graphical Models*. pp. 355–368. https://doi.org/10.1007/978-94-011-5014-9_12

Chapter 5

Supervised learning and classification

5.1 Aims of supervised learning

5.1.1 Supervised learning vs. unsupervised learning

Unsupervised learning:

Starting point:

- unlabelled data x_1, \dots, x_n .

Aim: find labels y_1, \dots, y_n to attach to each sample x_i .

For discrete labels y unsupervised learning is called *clustering*.

Supervised learning:

Starting point:

- labelled *training data*: $\{x_1^{\text{train}}, y_1^{\text{train}}\}, \dots, \{x_n^{\text{train}}, y_n^{\text{train}}\}$
- In addition, we have unlabelled *test data*: x^{test}

Aim: use training data to learn a function, say $h(x)$, to predict the label corresponding to the test data. The predictor function may provide a soft (probabilistic) assignment or a hard assignment of a class label to a test sample.

For y discrete supervised learning is called *classification*. For continuous y the label is called response and supervised learning becomes *regression*.

Thus, supervised learning is a two-step procedure:

- 1) Learn predictor function $h(x)$ using the training data x_i^{train} plus labels y_i^{train} .

- 2) Predict the label y^{test} for the test data \mathbf{x}^{test} using the estimated classifier function: $\hat{y}^{\text{test}} = \hat{h}(\mathbf{x}^{\text{test}})$.

5.1.2 Terminology

The function $h(\mathbf{x})$ that predicts the class y is called a *classifier*.

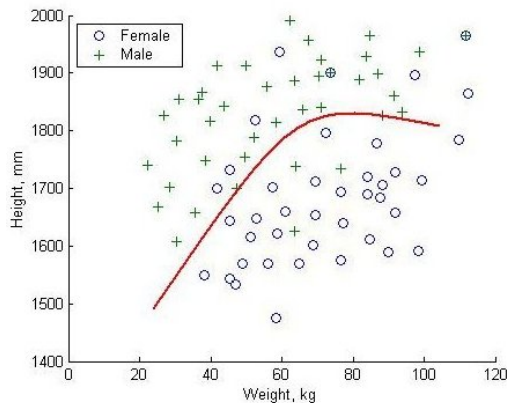
There are many types of classifiers, we focus here primarily on probabilistic classifiers (i.e. those that output probabilities for each possible class/label).

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 *decision boundary* between the classes is defined as the set of all \mathbf{x} for which the class assignment by the predictor $h(\mathbf{x})$ switches from one class to another.



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

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.

Common non-probabilistic methods for classification include:

- k -NN (Nearest Neighbors)
- SVM (support vector machine),
- random forest
- neural networks

Neural networks may in fact also be counted under probabilistic models as they essentially are high-dimensional complex nonlinear regression models. Just like a linear regression model can be fitted by “least squares” only without assuming an explicit probabilistic model neural networks are also most often trained by optimising a loss function.

5.2 Bayesian discriminant rule or Bayes classifier

Bayes classifiers are based on 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 $f_k(\mathbf{x}) = f(\mathbf{x}|\mathbf{k})$.
- prior probability of group k is $\Pr(k) = \pi_k$ with $\sum_{k=1}^K \pi_k = 1$
- marginal density is the mixture $f(\mathbf{x}) = \sum_{k=1}^K \pi_k f_k(\mathbf{x})$

For the moment we do not assume anything particular about the conditional densities $f_k(\mathbf{x})$ but later (see the following section) we will focus on normal densities and hence normal classifiers.

The posterior probability of group k is given by

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

This provides a “soft” classification

$$\mathbf{h}(\mathbf{x}^{\text{test}}) = (\Pr(k = 1|\mathbf{x}^{\text{test}}), \dots, \Pr(k = K|\mathbf{x}^{\text{test}}))^T$$

where each possible class $k \in \{1, \dots, K\}$ is assigned a probability to be the label for the test sample \mathbf{x} .

The *discriminant function* is defined as 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

this is the term $\log f(x)$. Hence we get for the Bayes discriminant function

$$d_k(x) = \log \pi_k + \log f_k(x).$$

For subsequent “hard” classification $h(x^{\text{test}})$ we then select the group/label for which the value of the discriminant function is maximised:

$$\hat{y}^{\text{test}} = h(x^{\text{test}}) = \arg \max_k d_k(x^{\text{test}}).$$

We have already encountered the Bayes classifier in the EM algorithm to predict the state of the latent variables (soft assignment) and in the K-means algorithm (hard assignment).

The discriminant functions $d_k(x)$ can be mapped back to the probabilistic class assignment by using the softargmax function (also known as softmax function):

$$\Pr(k|x) = \frac{\exp(d_k(x))}{\sum_{c=1}^K \exp(d_c(x))}.$$

In practise this will be calculated as

$$\Pr(k|x) = \frac{\exp(d_k(x) - d_{\max}(x))}{\sum_{c=1}^K \exp(d_c(x) - d_{\max}(x))}.$$

because subtracting $d_{\max}(x) = \max\{d_1(x), \dots, d_K(x)\}$ from all discriminant functions, and thus standardising the maximum of the discriminant functions to zero, avoids numerical overflow problems when computing the exponential function on a computer.

5.3 Normal Bayes classifier

5.3.1 Quadratic discriminant analysis (QDA) and Gaussian assumption

Quadratic discriminant analysis (QDA) is a special case of the Bayes classifier when all group-specific densities are multivariate normal with $f_k(x) = N(x|\mu_k, \Sigma_k)$. Note in particular that each group $k \in \{1, \dots, K\}$ has its own covariance matrix Σ_k .

Some calculation leads to the discriminant function for QDA:

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

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 squared **Mahalanobis distance** between x and μ_k : $d^{\text{Mahalanobis}}(x, \mu|\Sigma)^2 = (x - \mu)^T \Sigma^{-1}(x - \mu)$.
- 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).

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 find the minimum of the 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.

Decision boundaries for the QDA classifier can be either linear or nonlinear (quadratic curve). The decision boundary between any two classes i and j require that $d_i^{\text{QDA}}(x) = d_j^{\text{QDA}}(x)$, or equivalently $d_i^{\text{QDA}}(x) - d_j^{\text{QDA}}(x) = 0$, which is a quadratic equation.

5.3.2 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:

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

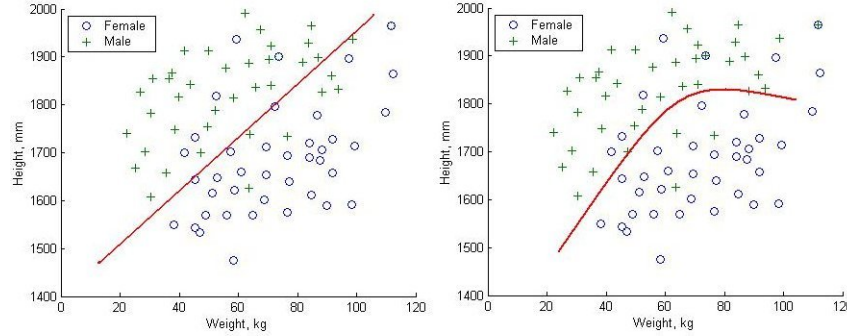
Note that term containing the log-determinant is now gone, and that LDA is essentially now a method that tries to minimize the Mahalanobis distance (while taking also into account the prior class probabilities).

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).

Comparison of linear decision boundaries 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.

5.3.3 Diagonal discriminant analysis (DDA)

For DDA we start with the same setting as for LDA, but now we simplify the model even further by additionally requiring a **diagonal covariance** containing only the variances (thus we assume that all correlations among the predictors x_1, \dots, x_d are zero):

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

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** and thus has linear decision boundaries.

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 that the underlying densities are normal.

However, don't let the label “naïve” mislead you as DDA and other “naïve” Bayes classifiers are often very effective methods for classification and prediction, especially in high-dimensional settings!

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

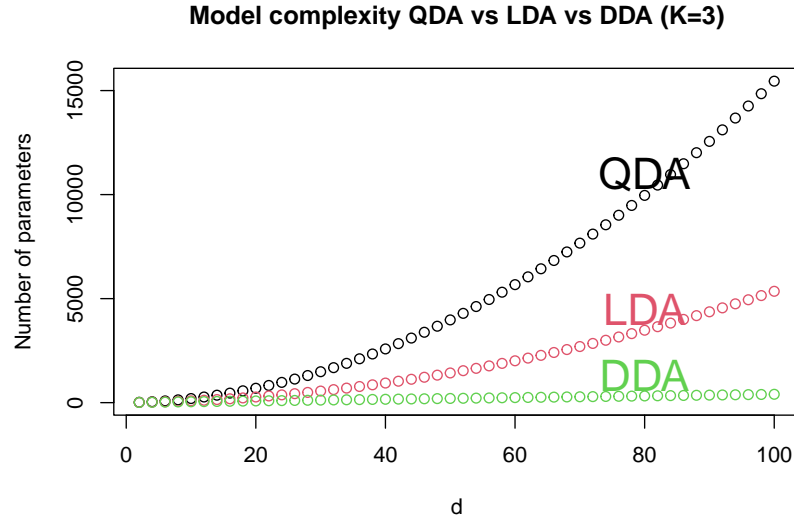
5.4.1 Number of model parameters

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 x_i^{train} and y_i^{train} :

- For QDA, LDA and DDA we need to learn π_1, \dots, π_K with $\sum_{k=1}^K \pi_k = 1$ and the mean vectors μ_1, \dots, μ_K
- For QDA we additionally require $\Sigma_1, \dots, \Sigma_K$
- For LDA we need Σ
- For DDA we estimate $\sigma_1^2, \dots, \sigma_d^2$.

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

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



5.4.2 Estimating the discriminant / predictor function

For QDA, LDA and DDA we learn the predictor by estimating the parameters of the discriminant function from the training data.

5.4.2.1 Large sample size

If the sample size of the training data set is sufficiently large compared to the model dimensions we can use maximum likelihood (ML) to estimate the model parameters. To be able to use ML we need a larger sample size for QDA and LDA (because full covariances need to be estimated) but for DDA a comparatively small sample size can be sufficient (which is one aspect why “naive” Bayes methods are very popular in practise).

To obtain the parameters estimates we use the known labels y_i^{train} to sort the samples x_i^{train} into the corresponding classes, and then apply the standard ML estimators. Let $g_k = \{i : y_i^{\text{train}} = k\}$ be the set of all indices of training sample belonging to group k , n_k the sample size in group k

The ML estimates of the class probabilities are the frequencies

$$\hat{\pi}_k = \frac{n_k}{n}$$

5.4. THE TRAINING STEP — LEARNING QDA, LDA AND DDA CLASSIFIERS FROM DA

and the ML estimate of the group means $k = 1, \dots, K$ are

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

The ML estimate of the global mean μ_0 (i.e. if we assume there is only a single class and ignore the group labels) is

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

Note the global mean is identical to the pooled mean (i.e. weighted average of the individual group means).

The ML estimates for the covariances Σ_k for QDA are

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

In order to get the ML estimate of the pooled variance Σ for use with LDA we compute

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

Note that the pooled variance Σ differs (substantially!) from the global variance Σ_0 that results from simply ignoring class labels and that is computed as

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

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/with-in group variance.

5.4.2.2 Small sample size

If the dimension d is large compared to the sample size then the number of parameters in the predictor function grows fast. Especially QDA but also LDA is data hungry and ML estimation becomes an ill-posed problem.

As discussed in Section 1.3 in this instance we need to use a regularised estimator for the covariance(s) such as estimators derived in the framework of penalised ML, Bayesian learning, shrinkage estimation etc. This also ensures that the estimated covariance matrices are positive definite (which is automatically guaranteed only for DDA if all variances are positive).

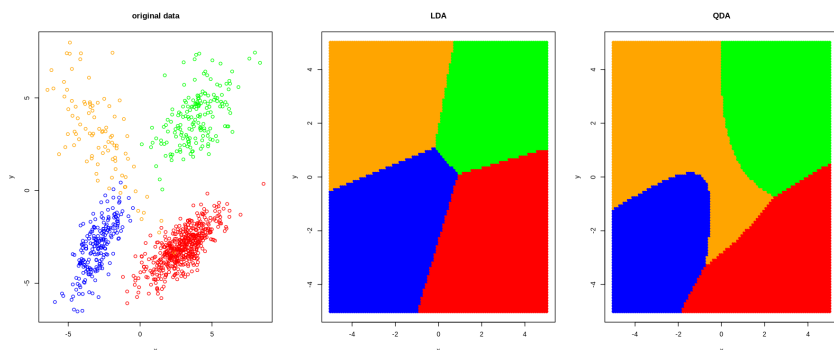
Furthermore, in small sample setting it is advised to reduce the number of parameters of the model. Thus using LDA or DDA is preferred over QDA. This can also prevent overfitting and lead to a predictor that generalises better.

To analyse high-dimensional data in the worksheets we will employ a regularised version of LDA and DDA using Stein-type shrinkage estimation as discussed in Section 1.3 and implemented in the R package “sda”.

5.4.3 Comparison of estimated decision boundaries: LDA vs. QDA

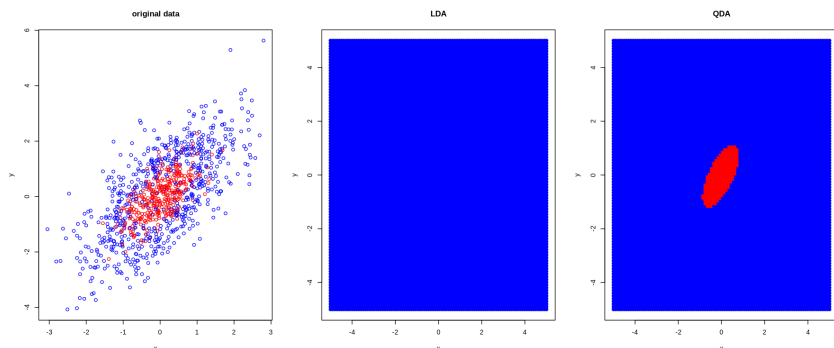
We compare two simple scenarios using simulated data.

Non-nested case ($K = 4$):



Both LDA and QDA clearly separate the 4 classes. Note the curved decision boundaries for QDA and the linear decision boundaries for LDA.

Nested case ($K = 2$):



In the nested case LDA fails to separate the two classes because there is no way to separate two nested classes with a simple linear boundary.

5.5 Quantifying prediction error

Once a classifier has been trained we are naturally interested in its performance to correctly classify previously unseen data points. This is useful for comparing different types of classifiers and also for comparing the same type of classifier using different sets of predictor variables.

5.5.1 Quantification of prediction error based on validation data

A measure of predictor error compares the predicted label \hat{y} with the true label y for validation data. A validation data set contains both the x_i and the associated label y_i but unlike the training data it has not been used for learning the predictor function.

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:

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

Alternatively, any other quantity derived from the confusion matrix (containing TP, TN, FP, FN) can be used.

The mean prediction error is 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 validation data set.

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. So this leads back to entropy and likelihood (see [MATH27720 Statistics 2](#)).

Once we have an estimate of the prediction error of a model we can use it to compare and choose among a set of candidate models, selecting those with a sufficiently low prediction error.

5.5.2 Estimation of prediction error using cross-validation

Unfortunately, quite often we do not have separate validation data available to evaluate a classifier.

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

Outline of cross-validation:

- 1) split the samples in the training data into a number (say K) parts (“folds”).
- 2) use each of the K folds as validation data and the other $K - 1$ folds as training data.
- 3) average over the resulting K individual estimates of prediction error, to get an overall aggregated predictor error, along with an error.

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

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 are able to train a reliable classifier from the remaining data). A typical value for K is 5 or 10, so that 80% respectively 90% of the samples are used for training and the other 20 % or 10% for validation.

If $K = n$ there are as many folds as there are samples and the validation data set consists only of a single data point. This is called “leave one out” cross-validation (LOOCV). There are analytic approximations for the prediction error obtained by LOOCV so that this approach is computationally inexpensive for some standard models (including regression).

Further reading:

To study the technical details of cross-validation: read **Section 5.1 Cross-Validation** in James et al. (2021) (R version) or James et al. (2023) (Python version).

5.6 Goodness of fit and variable ranking

As in linear regression we are interested in finding out whether the fitted mixture model is an appropriate model, and which particular predictor(s) x_j from $x = (x_1, \dots, x_d)^T$ are responsible prediction the outcome, i.e. for categorising 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.

5.6.1 LDA with $K = 2$ classes

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

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

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

The difference Δ_{12}^{LDA} is

$$\underbrace{\Delta_{12}^{LDA}}_{\text{log posterior odds}} = \underbrace{(\mu_1 - \mu_2)^T \Sigma^{-1} \left(x - \frac{\mu_1 + \mu_2}{2} \right)}_{\text{log Bayes factor } \log B_{12}} + \underbrace{\log \left(\frac{\pi_1}{\pi_2} \right)}_{\text{log prior odds}}$$

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 x . The *expected weight of evidence* assuming 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 :

$$E_{F_1}(\log B_{12}) = D_{KL}(F_1, F_2) = \frac{1}{2}(\mu_1 - \mu_2)^T \Sigma^{-1}(\mu_1 - \mu_2) = \frac{1}{2}\Omega^2$$

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

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

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 $SD(\hat{\mu}_1 - \hat{\mu}_2)$). T^2 is a measure of fit of the underlying two-component mixture.

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

$$\begin{aligned} \log B_{12} &= \left((\mu_1 - \mu_2)^T W^T \right) \left(W \left(x - \frac{\mu_1 + \mu_2}{2} \right) \right) \\ &= \omega^T \delta(x) \end{aligned}$$

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-cor whitening is to ensure that the whitened components are interpretable and stay maximally correlated to the original variables. However, you may also choose for example 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.

5.6.2 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).

The pooled centroid is given by $\mu_0 = \sum_{k=1}^K \pi_k \mu_k$ — this is the centroid if there would be only a single class. The corresponding probability (for a single class) 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.

5.7 Variable selection

In the previous we saw that in DDA the natural score for **ranking features** with regard to their relevance in separating the classes is the (squared) t -score, and for LDA a whitened version such as the squared correlation-adjusted t -score (based on ZCA-cor whitening) may be used. Once such a ranking has been established the question of a **suitable cutoff** arises, i.e. how many features need (or should) be retained in a model.

For large and high-dimensional models **feature selection can also be viewed as a form of regularisation and also dimension reduction**. Specifically, there

may be many variables/ features that do not contribute to the class prediction. Despite having in principle no effect on the outcome the presence of these “null variables” can nonetheless deteriorate (sometimes dramatically!) the overall predictive accuracy of a trained predictor, because they add noise and increase the model dimension. Therefore, variables that do not contribute to prediction should be filtered out in order to be able to construct good prediction models and classifiers.

5.7.1 Choosing a threshold by multiple testing using false discovery rates

The most simple way to determine a cutoff threshold is to use a standard technique for multiple testing.

For each predictor variable x_1, \dots, x_d we have a corresponding test statistic measuring the influence of this variable on the response, for example the t -scores and related statistics discussed in the previous section. In addition to providing an overall ranking the set of all these statistics can be used to determine a suitable cutoff by trying to separate two populations of predictor variables:

- “Null” variables that do not contribute to prediction
- “Alternative” variables that are linked to prediction

This can be done as follows:

- The distribution of the observed test statistics z_i is assumed to follow a two-component mixture where $F_0(z)$ and $F_A(z)$ are the distributions corresponding to the null and the alternative, $f_0(z)$ and $f_a(z)$ the densities, and π_0 and $\pi_A = 1 - \pi_0$ are the weights:

$$f(z) = \pi_0 f_0(z) + (1 - \pi_0) f_a(z)$$

- The null model is typically from a parametric family (e.g. normal around zero and with a free variance parameter) whereas the alternative is often modelled nonparametrically.
- After fitting the mixture model, often assuming some additional constraints to make the mixture identifiable, one can compute false discovery rates (FDR) as follows:

Local FDR:

$$\widehat{\text{fdr}}(z_i) = \widehat{\text{Pr}}(\text{null}|z_i) = \frac{\hat{\pi}_0 \hat{f}_0(z_i)}{\hat{f}(z_i)}$$

Tail-area-based FDR (=q-value):

$$\widehat{\text{Fdr}}(z_i) = \widehat{\text{Pr}}(\text{null}|Z > z_i) = \frac{\hat{\pi}_0 \hat{F}_0(z_i)}{\hat{F}(z_i)}$$

Note these are essentially p -values adjusted for multiple testing (by a variant of the Benjamini-Hochberg method).

By thresholding false discovery rates it is possible to identify those variables that clearly belong to each of the two groups but also those features that cannot easily be discriminated to fall into either group:

- “alternative” variables have low local FDR, e.g. $\widehat{\text{fdr}}(z_i) \leq 0.2$
- “null” variables have high local FDR, e.g. $\widehat{\text{fdr}}(z_i) \geq 0.8$
- features that cannot easily classified as null or alternative, e.g. $0.2 < \widehat{\text{fdr}}(z_i) < 0.8$

For feature selection in prediction settings we generally aim to remove only those variable that clearly belong to the null group, leaving all others in the model.

5.7.2 Variable selection using cross-validation

A conceptually simple but computationally more expensive approach to variable selection is to estimate the prediction error of the same type of predictor but with different sets of predictors using cross-validation, and then choosing a predictor that achieves good prediction accuracy while using only a small number of features.

This is a method that works very well in practise as is demonstrated in a number of problems in the worksheets.

Chapter 6

Multivariate dependencies

6.1 Measuring the linear association between two sets of random variables

6.1.1 Aim

The linear association between two scalar random variables x and y is measured by the correlation $\text{Cor}(x, y) = \rho$.

In this chapter we now would like to explore how to generalise correlation to the case of two random vectors. Specifically, we would like to find a scalar measure of association between two random vectors (or equivalently two sets of random variables) $\mathbf{x} = (x_1, \dots, x_p)^T$ and $\mathbf{y} = (y_1, \dots, y_q)^T$ that contains correlation and also multiple correlation as special case.

We assume a joint correlation matrix

$$\mathbf{P} = \begin{pmatrix} \mathbf{P}_x & \mathbf{P}_{xy} \\ \mathbf{P}_{yx} & \mathbf{P}_y \end{pmatrix}$$

with cross-correlation matrix $\mathbf{P}_{xy} = \mathbf{P}_{yx}^T$ and the within-group group correlations \mathbf{P}_x and \mathbf{P}_y . If the cross-correlations vanish, $\mathbf{P}_{xy} = 0$, then the two random vectors are uncorrelated, and the joint correlation matrix becomes a diagonal block matrix

$$\mathbf{P}_{\text{indep}} = \begin{pmatrix} \mathbf{P}_x & 0 \\ 0 & \mathbf{P}_y \end{pmatrix}.$$

To characterise the total association between \mathbf{x} and \mathbf{y} we are looking for a **scalar quantity** measuring the divergence of a distribution assuming the general joint correlation matrix \mathbf{P} from a distribution assuming the joint correlation matrix $\mathbf{P}_{\text{indep}}$ for uncorrelated \mathbf{x} and \mathbf{y} .

6.1.2 Known special cases

Ideally, in case of an univariate y but multivariate x this measure should reduce to the *squared multiple correlation* or *coefficient of determination*

$$\text{MCor}(x, y)^2 = \mathbf{P}_{yx} \mathbf{P}_x^{-1} \mathbf{P}_{xy}$$

This is well-known in linear regression to describe the strength of the total linear association between the predictors $\mathbf{x} = (x_1, \dots, x_p)^T$ and the response y .

To derive the squared multiple correlation we may proceed as follows. First we whiten the random vector x resulting in $\mathbf{z}_x = \mathbf{W}_x \mathbf{x} = \mathbf{Q}_x \mathbf{P}_x^{-1/2} \mathbf{V}_x^{-1/2} \mathbf{x}$ where \mathbf{Q}_x is an orthogonal matrix. The correlations between each component in \mathbf{z}_x and the response y are then

$$\boldsymbol{\omega} = \text{Cor}(\mathbf{z}_x, y) = \mathbf{Q}_x \mathbf{P}_x^{-1/2} \mathbf{P}_{xy}$$

As $\text{Var}(\mathbf{z}_x) = \mathbf{I}$ and thus the components in \mathbf{z}_x are uncorrelated we can simply add up the squared individual correlations to get as total association measure

$$\boldsymbol{\omega}^T \boldsymbol{\omega} = \sum_i \omega_i^2 = \mathbf{P}_{yx} \mathbf{P}_x^{-1} \mathbf{P}_{xy}$$

Note that the particular choice of the orthogonal matrix \mathbf{Q}_x for whitening x is not relevant for the squared multiple correlation.

Note that if the cross-correlations vanish ($\mathbf{P}_{xy} = 0$) then $\text{MCor}(x, y)^2 = 0$. If the correlation between the predictors vanishes ($\mathbf{P}_x = \mathbf{I}$) then $\text{MCor}(x, y)^2 = \sum_i \rho_{yx_i}^2$, i.e. it is the sum of the squared cross-correlations.

If there is only a single predictor x then $\mathbf{P}_{xy} = \rho$ and $\mathbf{P}_x = 1$ and the squared multiple correlation reduces to the squared Pearson correlation

$$\text{Cor}(x, y)^2 = \rho^2.$$

6.2 Canonical Correlation Analysis (CCA) aka CCA whitening

Canonical correlation analysis was invented by Harald Hotelling in 1936¹. CCA aims to characterise the linear dependence between two random vectors \mathbf{x} and \mathbf{y} by a set of **canonical correlations** λ_i .

CCA works by simultaneously whitening the two random vectors \mathbf{x} and \mathbf{y} where the whitening matrices are chosen in such a way that the cross-correlation matrix between the resulting whitened variables becomes diagonal, and the elements on the diagonal correspond to the **canonical correlations**.

¹Hotelling, H. 1936. Relations between two sets of variates. *Biometrika* 28:321–377. <https://doi.org/10.1093/biomet/28.3-4.321>

$$\begin{array}{ccc} \mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} & \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_q \end{pmatrix} & \begin{array}{l} \text{Var}(\mathbf{x}) = \boldsymbol{\Sigma}_x = \mathbf{V}_x^{1/2} \mathbf{P}_x \mathbf{V}_x^{1/2} \\ \text{Var}(\mathbf{y}) = \boldsymbol{\Sigma}_y = \mathbf{V}_y^{1/2} \mathbf{P}_y \mathbf{V}_y^{1/2} \end{array} \\ \text{Dimension } p & \text{Dimension } q & \end{array}$$

$$\text{Whitening of } \mathbf{x}: \mathbf{z}_x = \mathbf{W}_x \mathbf{x} = \mathbf{Q}_x \mathbf{P}_x^{-1/2} \mathbf{V}_x^{-1/2} \mathbf{x}$$

$$\text{Whitening of } \mathbf{y}: \mathbf{z}_y = \mathbf{W}_y \mathbf{y} = \mathbf{Q}_y \mathbf{P}_y^{-1/2} \mathbf{V}_y^{-1/2} \mathbf{y}$$

(note we use the correlation-based form of \mathbf{W})

Cross-correlation between \mathbf{z}_x and \mathbf{z}_y :

$$\text{Cor}(\mathbf{z}_x, \mathbf{z}_y) = \mathbf{Q}_x \mathbf{K} \mathbf{Q}_y^T$$

with $\mathbf{K} = \mathbf{P}_x^{-1/2} \mathbf{P}_{xy} \mathbf{P}_y^{-1/2}$.

Idea: we can choose suitable orthogonal matrices \mathbf{Q}_x and \mathbf{Q}_y by putting a structural constraint on the cross-correlation matrix.

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

Motivation: pairs of “modules” represented by components of \mathbf{z}_x and \mathbf{z}_y influencing each other (and not anyone other module).

$$\mathbf{z}_x = \begin{pmatrix} z_1^x \\ z_2^x \\ \vdots \\ z_p^x \end{pmatrix} \quad \mathbf{z}_y = \begin{pmatrix} z_1^y \\ z_2^y \\ \vdots \\ z_q^y \end{pmatrix}$$

$$\text{Cor}(\mathbf{z}_x, \mathbf{z}_y) = \begin{pmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_m \end{pmatrix}$$

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

6.2.1 How to make cross-correlation matrix $\text{Cor}(\mathbf{z}_x, \mathbf{z}_y)$ diagonal?

- Use Singular Value Decomposition (SVD) of matrix \mathbf{K} :

$$\mathbf{K} = (\mathbf{Q}_x^{\text{CCA}})^T \boldsymbol{\Lambda} \mathbf{Q}_y^{\text{CCA}}$$

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

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

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

Note that the signs of corresponding in columns in Q_x^{CCA} and Q_y^{CCA} are not identified. Traditionally, in an SVD the signs are chosen such that the singular values are positive. However, if we impose positive-diagonality on Q_x^{CCA} and Q_y^{CCA} , and thus positive-diagonality on the cross-correlations Ψ_x and Ψ_y , then the canonical correlations may take on both positive and negative values.

6.2.2 Related methods

- O2PLS: similar to CCA but using orthogonal projections rather than whitening.
- Vector correlation: aggregates the squared canonical correlations into a single overall measure (see below).

6.3 Vector correlation and RV coefficient

6.3.1 Vector alienation coefficient

In his 1936 paper introducing canonical correlation analysis Hotelling also proposed the *vector alienation coefficient* defined as

$$\begin{aligned} a(x, y) &= \frac{\det(P)}{\det(P_{\text{indep}})} \\ &= \frac{\det(P)}{\det(P_x) \det(P_y)} \end{aligned}$$

With $K = P_x^{-1/2} P_{xy} P_y^{-1/2}$ the vector alienation coefficient can be written (using the Weinstein-Aronszajn determinant identity and the formula for the determinant of block-structured matrices) as

$$\begin{aligned} a(x, y) &= \det(I_p - KK^T) \\ &= \det(I_q - K^T K) \\ &= \prod_{i=1}^m (1 - \lambda_i^2) \end{aligned}$$

where the λ_i are the singular values of K , i.e. the canonical correlations for the pair x and y . Therefore, the vector alienation coefficient is computed as a summary statistic of the canonical correlations.

If $P_{xy} = 0$ and thus x and y are uncorrelated then $P = P_{\text{indep}}$ and thus by construction the vector alienation coefficient $a(x, y) = 1$. Hence, the vector alienation coefficient is itself not a generalisation of the squared multiple correlation to the case of two random vectors as such a quantity should vanish in this case.

6.3.2 Rozeboom vector correlation

Instead, Rozeboom (1965)² proposed to use as squared *vector correlation* the complement of the vector alienation coefficient

$$\begin{aligned} \text{VCor}(x, y)^2 &= \rho_{xy}^2 = 1 - a(x, y) \\ &= 1 - \det(I_p - KK^T) \\ &= 1 - \det(I_q - K^T K) \\ &= 1 - \prod_{i=1}^m (1 - \lambda_i^2) \end{aligned}$$

If $P_{xy} = 0$ then $K = 0$ and hence $\text{VCor}(x, y)^2 = 0$.

Moreover, if either $p = 1$ or $q = 1$ the squared vector correlation reduces to the corresponding squared multiple correlation, which in turn for both $p = 1$ and $q = 1$ becomes the squared Pearson correlation.

You can find the derivation in Example Sheet 10.

Thus, Rozeboom's vector correlation indeed generalises both Pearson correlation and the multiple correlation coefficient.

6.3.3 RV coefficient

Another common approach to measure association between two random vectors is the *RV coefficient* introduced by Robert and Escoufier in 1976 as

$$RV(x, y) = \frac{\text{Tr}(\Sigma_{xy} \Sigma_{yx})}{\sqrt{\text{Tr}(\Sigma_x^2) \text{Tr}(\Sigma_y^2)}}$$

The main advantage of the RV coefficient is that it is easier to compute than the Rozeboom vector correlation as it uses the matrix trace rather than the matrix determinant.

²Rozeboom, W. W. 1965. Linear correlations between sets of variables. *Psychometrika* 30:57–71. <https://doi.org/10.1007/BF02289747>

For $q = p = 1$ the RV coefficient reduces to the squared correlation. However, the RV coefficient does *not* reduce to the multiple correlation coefficient for $q = 1$ and $p > 1$, and therefore the RV coefficient cannot be considered a coherent generalisation of Pearson and multiple correlation to the case when x and y are random vectors.

See also Worksheet 10.

6.4 Limits of linear models and correlation

6.4.1 Correlation measures only linear dependence

Linear models and measures of linear association (correlation) are very effective tools. However, it is important to recognise their limits especially when modelling complex nonlinear relationships.

A very simple demonstration of this is given by the following example. Assume x is a normally 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 running a small computer simulation:

```
x=rnorm(10000)
y = x^2
cor(x,y)
```

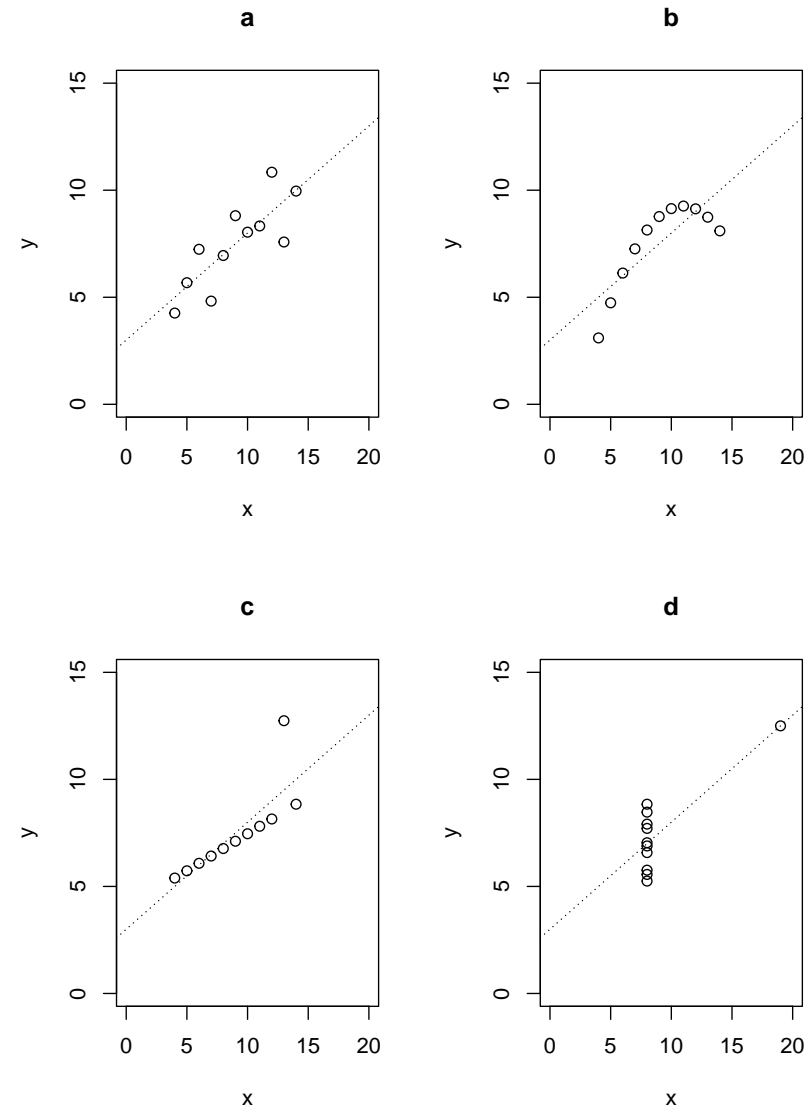
```
## [1] 0.009189134
```

Thus, correlation is (almost) zero even though x and y are dependent variables. This is because correlation only measures linear association, and the relationship between x and y is nonlinear.

6.4.2 Anscombe data sets

Using correlation (and more generally linear models) blindly can easily hide the underlying complexity of the analysed data. This is demonstrated by the classic “Anscombe quartet” of data sets presented in his 1973 paper ³ -

³Anscombe, F. J. 1973. Graphs in statistical analysis. The American Statistician 27:17-21, <http://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! However, 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 set “a” follows a linear model. Data set “b” represents a quadratic relationship. Data set “c” is linear but with an outlier that disturbs the linear relationship. Finally data set “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 Anscombe quartet will be analysed in the form of the “datasauRus” dozen - 13 highly nonlinear datasets that all share the same linear characteristics.

6.5 Mutual information as generalisation of correlation

6.5.1 Overview

A more general way than the vector correlation to measure multivariate association is mutual information (MI) which not only covers linear but also non-linear associations.

As we will see below the Rozeboom vector correlation arises naturally when computing the MI for the multivariate normal distribution, hence MI also recovers well-known measures of linear association (including multiple correlation and simple correlation), thus truly generalising correlation as measure of association.

6.5.2 Definition of mutual information

Recall the definition of Kullback-Leibler divergence, or relative entropy, between two distributions:

$$D_{\text{KL}}(F, G) := E_F \log \left(\frac{f(x)}{g(x)} \right)$$

Here F plays the role of the reference distribution and G is an approximating distribution, with f and g being the corresponding density functions (see [MATH27720 Statistics 2](#)).

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:

$$\text{MI}(x, y) = D_{\text{KL}}(F_{x,y}, F_x F_y) = E_{F_{x,y}} \log \left(\frac{f(x, y)}{f(x) f(y)} \right).$$

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.

6.5.3 Mutual information between two normal scalar variables

The KL divergence between two multivariate normal distributions F_{ref} and F is

$$D_{\text{KL}}(F_{\text{ref}}, F) = \frac{1}{2} \left\{ (\mu - \mu_{\text{ref}})^T \Sigma^{-1} (\mu - \mu_{\text{ref}}) + \text{Tr} \left(\Sigma^{-1} \Sigma_{\text{ref}} \right) - \log \det \left(\Sigma^{-1} \Sigma_{\text{ref}} \right) - d \right\}$$

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 $z = (x, y)^T$. The joint bivariate normal distribution is characterised by the mean $E(z) = \mu = (\mu_x, \mu_y)^T$ and the covariance matrix

$$\Sigma = \begin{pmatrix} \sigma_x^2 & \rho \sigma_x \sigma_y \\ \rho \sigma_x \sigma_y & \sigma_y^2 \end{pmatrix}$$

where $\text{Cor}(x, y) = \rho$. If x and y are independent then $\rho = 0$ and

$$\Sigma_{\text{indep}} = \begin{pmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{pmatrix}.$$

The product

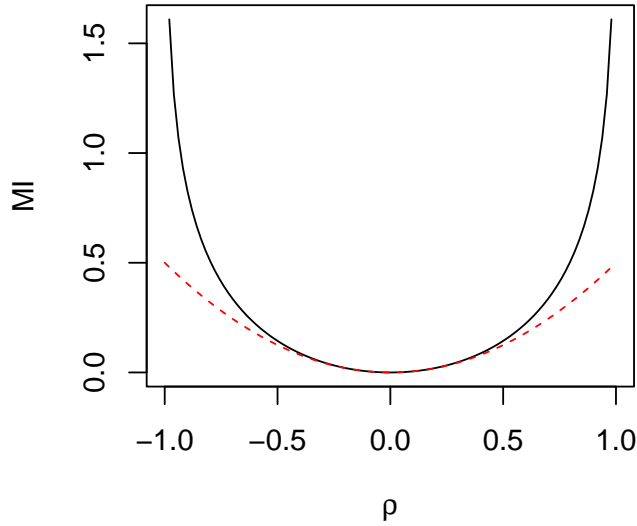
$$A = \Sigma_{\text{indep}}^{-1} \Sigma = \begin{pmatrix} 1 & \rho \frac{\sigma_y}{\sigma_x} \\ \rho \frac{\sigma_x}{\sigma_y} & 1 \end{pmatrix}$$

has trace $\text{Tr}(A) = 2$ and determinant $\det(A) = 1 - \rho^2$.

With this the mutual information between x and y can be computed as

$$\begin{aligned} \text{MI}_{\text{norm}}(x, y) &= D_{\text{KL}}(N(\mu, \Sigma), N(\mu, \Sigma_{\text{indep}})) \\ &= \frac{1}{2} \left\{ \text{Tr} \left(\Sigma_{\text{indep}}^{-1} \Sigma \right) - \log \det \left(\Sigma_{\text{indep}}^{-1} \Sigma \right) - 2 \right\} \\ &= \frac{1}{2} \left\{ \text{Tr}(A) - \log \det(A) - 2 \right\} \\ &= -\frac{1}{2} \log(1 - \rho^2) \\ &\approx \frac{\rho^2}{2} \end{aligned}$$

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 $2 \text{MI}_{\text{norm}}(x, y) \approx \rho^2$.

6.5.4 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

$$\Sigma = \begin{pmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{xy}^T & \Sigma_y \end{pmatrix}.$$

If x and y are independent then $\Sigma_{xy} = 0$ and

$$\Sigma_{\text{indep}} = \begin{pmatrix} \Sigma_x & 0 \\ 0 & \Sigma_y \end{pmatrix}.$$

The product

$$\begin{aligned} A &= \Sigma_{\text{indep}}^{-1} \Sigma = \begin{pmatrix} I_p & \Sigma_x^{-1} \Sigma_{xy} \\ \Sigma_y^{-1} \Sigma_{yx} & I_q \end{pmatrix} \\ &= \begin{pmatrix} I_p & V_x^{-1/2} P_x^{-1} P_{xy} V_y^{1/2} \\ V_y^{-1/2} P_y^{-1} P_{yx} V_x^{1/2} & I_q \end{pmatrix} \end{aligned}$$

has trace $\text{Tr}(A) = d$ and determinant

$$\begin{aligned} \det(A) &= \det(I_p - K K^T) \\ &= \det(I_q - K^T K) \end{aligned}$$

with $K = P_x^{-1/2} P_{xy} P_y^{-1/2}$. With $\lambda_1, \dots, \lambda_m$ the singular values of K (i.e. the canonical correlations between x and y) we get

$$\det(A) = \prod_{i=1}^m (1 - \lambda_i^2)$$

The mutual information between x and y is then

$$\begin{aligned} \text{MI}_{\text{norm}}(x, y) &= D_{\text{KL}}(N(\mu, \Sigma), N(\mu, \Sigma_{\text{indep}})) \\ &= \frac{1}{2} \left\{ \text{Tr} \left(\Sigma_{\text{indep}}^{-1} \Sigma \right) - \log \det \left(\Sigma_{\text{indep}}^{-1} \Sigma \right) - d \right\} \\ &= \frac{1}{2} \left\{ \text{Tr}(A) - \log \det(A) - d \right\} \\ &= -\frac{1}{2} \sum_{i=1}^m \log(1 - \lambda_i^2) \\ &= -\frac{1}{2} \log \left(\prod_{i=1}^m (1 - \lambda_i^2) \right) \\ &= -\frac{1}{2} \log(1 - \text{VCor}(x, y)^2) \end{aligned}$$

From the above we see that $\text{MI}_{\text{norm}}(x, y)$ is simply the sum of the MIs resulting from the individual canonical correlations λ_i with the same functional link between the MI and the squared correlation as in the bivariate normal case.

Furthermore we obtain that

$$\text{MI}_{\text{norm}}(x, y) = -\frac{1}{2} \log(1 - \text{VCor}(x, y)^2) \approx \frac{1}{2} \text{VCor}(x, y)^2$$

Thus, in the multivariate case $\text{MI}_{\text{norm}}(x, y)$ has again exactly the same functional relationship with the squared vector correlation $\text{VCor}(x, y)^2$ as the $\text{MI}_{\text{norm}}(x, y)$ for two univariate variables with squared Pearson correlation ρ^2 .

Thus, Rozeboom's vector correlation emerges as a special case of mutual information computed for jointly multivariate normally distributed variables.

6.5.5 Using MI for variable selection

A very general way to write down a model predicting y by x is as follows:

- $F_{y|x}$ is a conditional distribution of y given predictors x and
- F_y is the marginal distribution of y without predictors.

Typically $F_{y|x}$ is a complex model and F_y a simple model (no predictors). Note that the predictive model can assume any form (incl. nonlinear).

Intriguingly the expected KL divergence between the conditional and the marginal distribution

$$E_{F_x} D_{\text{KL}}(F_{y|x}, F_y) = \text{MI}(x, y)$$

is equal to mutual information between x and y ! Thus $\text{MI}(x, y)$ measures the impact of conditioning. If the MI is small (i.e. close to zero) then x is not useful in predicting y .

The above identity can be verified as follows. The KL divergence between $F_{y|x}$ and F_y is given by

$$D_{\text{KL}}(F_{y|x}, F_y) = E_{F_{y|x}} \log \left(\frac{f(y|x)}{f(y)} \right),$$

which is a random variable since it depends on x . Taking the expectation with regard to F_x (the distribution of x) we get

$$\begin{aligned} E_{F_x} D_{\text{KL}}(F_{y|x}, F_y) &= E_{F_x} E_{F_{y|x}} \log \left(\frac{f(y|x)f(x)}{f(y)f(x)} \right) \\ &= E_{F_{x,y}} \log \left(\frac{f(x, y)}{f(y)f(x)} \right) = \text{MI}(x, y). \end{aligned}$$

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.

6.5.6 Other measures of general dependence

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.

Besides mutual information there are others measures of general dependence between multivariate random variables.

Two important measures to capture nonlinear association that have been proposed in recent literature are

- i) [distance correlation](#) and
- ii) the [maximal information coefficient](#) (MIC and MIC_e).

6.6 Graphical models

6.6.1 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. For an overview see, e.g., the reference book by Murphy (2023).

6.6.2 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.

6.6.3 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 $p(x_1, x_2, \dots, x_d)$. 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.

6.6.4 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:

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

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 node and propagates recursively across 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.

6.6.5 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**.

Suppose x_i and x_j are two random variables/nodes from $\{x_1, \dots, x_d\}$, and the set $\{x_k\}$ represents all other variables/nodes with $k \neq i$ and $k \neq j$. Then the 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 for all variables $\{x_1, \dots, x_d\}$ factorises as

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

or equivalently

$$\frac{p(x_i, x_j, \dots, x_d)}{p(\{x_k\})} = p(x_i, x_j | \{x_k\}) = p(x_i | \{x_k\}) p(x_j | \{x_k\}).$$

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

6.6.5.1 Gaussian graphical model

Assuming that x_1, \dots, x_d are jointly normally 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 short **GGM**. A further alternative name commonly used is **covariance selection model**.

6.6.5.2 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 $\{x_k\}$.

6.6.6 Null distribution of the empirical correlation coefficient

Suppose we have two uncorrelated random variables x and y with $\rho = \text{Cor}(x, y) = 0$. After observing data x_1, \dots, x_n and y_1, \dots, y_n we compute the empirical covariance matrix $\hat{\Sigma}_{xy}$ and from it the empirical correlation coefficient $r = \widehat{\text{Cor}}(x, y)$.

The distribution of the empirical correlation assuming $\rho = 0$ is useful as null-model for testing whether the underlying correlation is in fact zero having observed empirical correlation r . If x and y are normally distributed with $\rho = 0$ the distribution of the empirical correlation r has mean $E(r) = 0$ and variance $\text{Var}(r) = \frac{1}{\kappa}$. Here κ is the degree of freedom of the null distribution which for standard correlation is $\kappa = n - 1$. Furthermore, the squared empirical correlation is distributed according to a Beta distribution

$$r^2 \sim \text{Beta}\left(\frac{1}{2}, \frac{\kappa - 1}{2}\right)$$

For partial correlation the null distribution of r^2 has the same form but with a different degree of freedom. Specifically, κ is reduced by the number of variables being conditioned on. If for d dimensions we condition on $d - 2$ variables the resulting degree of freedom is $\kappa = n - 1 - (d - 2) = n - d + 1$. For $d = 2$ we get back the degree of freedom for standard empirical correlation.

6.6.7 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 .

The test for conditional independence is done by statistical testing for vanishing partial correlation. Specifically, we compute the p -value assuming that the true underlying partial correlation is zero and then decide whether to reject the null assumption of zero partial correlation.

If there are many edges tested simultaneously we may need to adjust (i.e reduce) the test threshold, for example applying Bonferroni or FDR methods.

6.6.8 Example: exam score data

This is a data set from Mardia et al. (1979) and features $d = 5$ variables measured on $n = 88$ subjects.

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

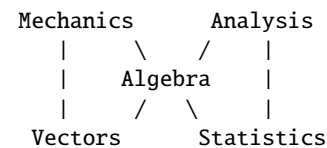
Note that there are no zero correlations but there are **four partial correlations close to 0**, indicating **conditional independence** between:

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

This can be verified by computing the normal p -values for the partial correlations (with $\kappa = 84$ as degree of freedom):

##	mechanics	vectors	algebra	analysis	statistics
## mechanics	NA	0.002	0.034	0.988	0.823
## vectors	NA	NA	0.009	0.477	0.854
## algebra	NA	NA	NA	0.000	0.001
## analysis	NA	NA	NA	NA	0.020
## statistics	NA	NA	NA	NA	NA

There are six edges with small p -value (smaller than say 0.05) and these correspond to the edges for which the null assumption of zero partial correlation can be rejected so that out of ten possible edges four are not statistically significant. Therefore the conditional independence graph looks as follows:



Chapter 7

Nonlinear and nonparametric models

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

Relevant textbooks:

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

- James et al. (2021) *An introduction to statistical learning with applications in R (2nd edition)*. 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!

The first book is also available in a version with examples in Python:

- James et al. (2023) *An introduction to statistical learning with applications in Python*. Springer.

7.1 Random forests

Another widely used approach for prediction in nonlinear settings is the method of random forests.

Relevant reading:

Please read: James et al. (2021) or James et al. (2023) **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

7.1.1 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* **16**:199–231. <https://doi.org/10.1214/ss/1009213726>

This paper has recently be revisited in the following discussion paper by Efron (2020) and discussants: Prediction, estimation, and attribution. *JASA* **115**:636–677. <https://doi.org/10.1080/01621459.2020.1762613>

7.1.2 Random forests

Proposed by Leo Breimann in 2001 as application of “bagging” (Breiman 1996) to decision trees.

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).

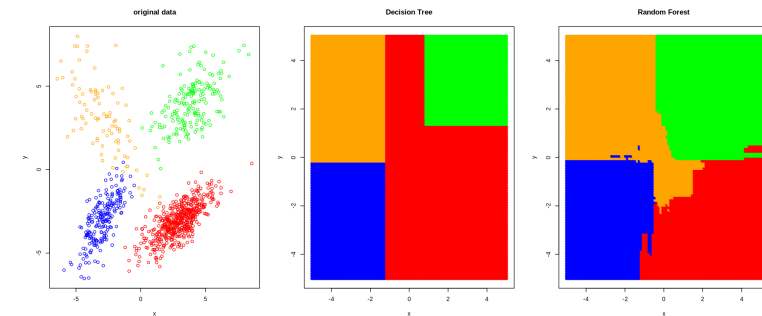
Variations: boosting, XGBoost (<https://xgboost.ai/>)

Random forests will be applied in Worksheet 11.

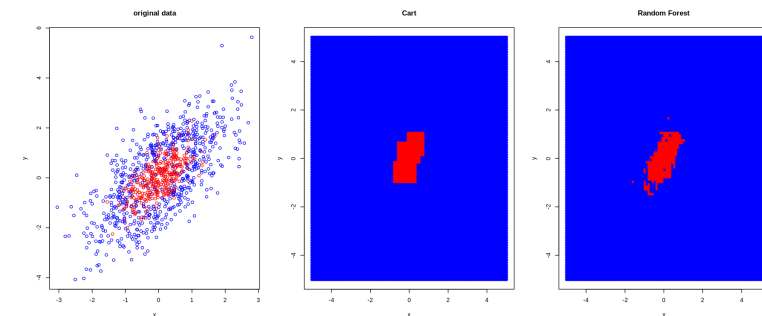
They are computationally expensive but typically perform very well!

7.1.3 Comparison of decision boundaries: decision tree vs. random forest

Non-nested case:



Nested case:



Compare also with the decision boundaries for LDA and QDA (Chapter 4).

7.2 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**.

7.2.1 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 a function)
- GPs use an infinitely dimensional multivariate normal distribution as prior

7.2.2 Conditional multivariate normal distribution

GPs make use of the fact that marginal and conditional distributions of a multivariate normal distribution are also multivariate normal.

Multivariate normal distribution:

$$\mathbf{z} \sim N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Assume:

$$\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

with

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

and

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_1 & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_2 \end{pmatrix}$$

with corresponding dimensions d_1 and d_2 and $d_1 + d_2 = d$.

Marginal distributions:

Any subset of \mathbf{z} is also multivariate normally distributed. Specifically,

$$\mathbf{z}_1 \sim N_{d_1}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$$

and

$$\mathbf{z}_2 \sim N_{d_2}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

Conditional multivariate normal:

The conditional distribution is also multivariate normal:

$$\mathbf{z}_1 | \mathbf{z}_2 = \mathbf{z}_{1|2} \sim N_{d_1}(\boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2})$$

with

$$\boldsymbol{\mu}_{1|2} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_2^{-1} (\mathbf{z}_2 - \boldsymbol{\mu}_2)$$

and

$$\boldsymbol{\Sigma}_{1|2} = \boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_{12}^T$$

$\mathbf{z}_{1|2}$ and $\boldsymbol{\mu}_{1|2}$ have dimension $d_1 \times 1$ and $\boldsymbol{\Sigma}_{1|2}$ has dimension $d_1 \times d_1$, i.e. the same dimension as the unconditioned variables.

You may recall the above formula in the context of linear regression, with $y = z_1$ and $\mathbf{x} = \mathbf{z}_2$ so that the conditional mean becomes

$$\begin{aligned} E(y|\mathbf{x}) &= \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_x^{-1} (\mathbf{x} - \boldsymbol{\mu}_x) \\ &= \beta_0 + \boldsymbol{\beta}^T \mathbf{x} \end{aligned}$$

with $\boldsymbol{\beta} = \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\Sigma}_{xy}$ and $\beta_0 = \boldsymbol{\mu}_y - \boldsymbol{\beta}^T \boldsymbol{\mu}_x$, and the corresponding conditional variance is

$$\text{Var}(y|\mathbf{x}) = \sigma_y^2 - \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\Sigma}_{xy}.$$

7.2.3 Covariance functions and kernels

The GP prior is an infinitely dimensional multivariate normal distribution with mean zero and the **covariance specified by a function** $k(\mathbf{x}, \mathbf{x}')$:

A widely used covariance function is

$$k(\mathbf{x}, \mathbf{x}') = \text{Cov}(\mathbf{x}, \mathbf{x}') = \sigma^2 e^{-\frac{(\mathbf{x}-\mathbf{x}')^2}{2l^2}}$$

This is known as the **squared-exponential kernel** or **Radial-basis function (RBF) kernel**.

Note that this kernel implies

- $k(\mathbf{x}, \mathbf{x}) = \text{Var}(\mathbf{x}) = \sigma^2$ and
- $\text{Cor}(\mathbf{x}, \mathbf{x}') = e^{-\frac{(\mathbf{x}-\mathbf{x}')^2}{2l^2}}.$

The parameter l in the RBF kernel is the length scale parameter and describes the “wigglyness” or “stiffness” of the resulting function. Small values of l correspond to more complex, more wiggly functions, and to low spatial correlation, as the

correlation decreases quicker with distance, and large values correspond to more rigid, stiffer functions, with longer range spatial correlation (note that in a time series context this would be called autocorrelation).

There are many other kernel functions, including linear, polynomial or periodic kernels.

7.2.4 GP model

Nonlinear regression in the GP approach is conceptually very simple:

- start with multivariate prior
- then condition on the observed data
- the resulting conditional multivariate normal can be used to predict the function values at any unobserved values
- the conditional variance can be used to compute credible intervals for predictions.

GP regression also provides a direct link with classical Bayesian linear regression (when using a linear kernel). Furthermore, GPs are also linked with neural networks as their limit in the case of an infinitely wide network (see section on neural networks).

Drawbacks of GPs: computationally expensive, typically $O(n^3)$ because of the matrix inversion. However, there are now variations of GPs that help to overcome this issue (e.g. sparse GPs).

7.2.5 Gaussian process example

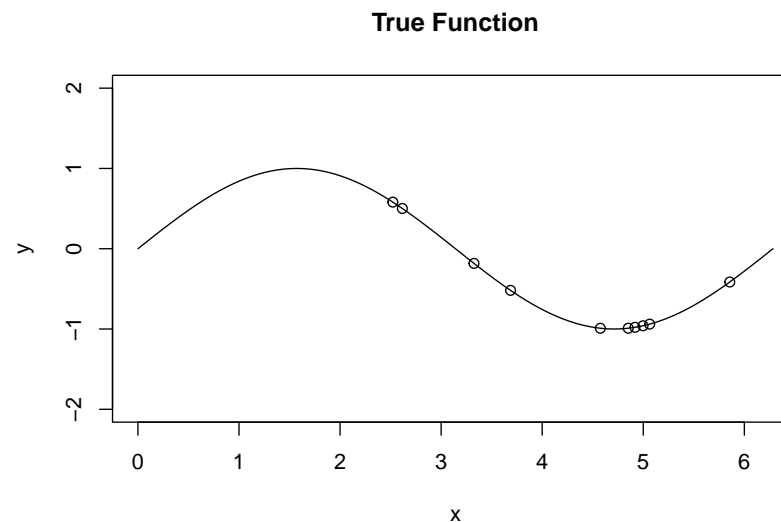
We now show how to apply Gaussian processes in R just using standard matrix calculations.

Our aim is to estimate the following nonlinear function from a number of observations. Note that initially we assume that there is no additional noise (so the observations lie directly on the curve):

```
truefunc = function(x) sin(x)
XLIM = c(0, 2*pi)
YLIM = c(-2, 2)

n2 = 10
x2 = runif(n2, min=XLIM[1], max=XLIM[2])
y2 = truefunc(x2) # no noise

curve( truefunc(x), xlim=XLIM, ylim=YLIM, xlab="x", ylab="y",
       main="True Function")
points(x2, y2)
```



Next we use the RBF kernel as the prior covariance and also assume that the prior has mean zero:

```
# RBF kernel
rbfkernel = function(xa, xb, s2=1, l=1/2) s2*exp(-1/2*(xa-xb)^2/l^2)
kfun.mat = function(xavec, xbvec, FUN=rbfkernel)
  outer(X=as.vector(xavec), Y=as.vector(xbvec), FUN=FUN)

# prior mean
mu.vec = function(x) rep(0, length(x))
```

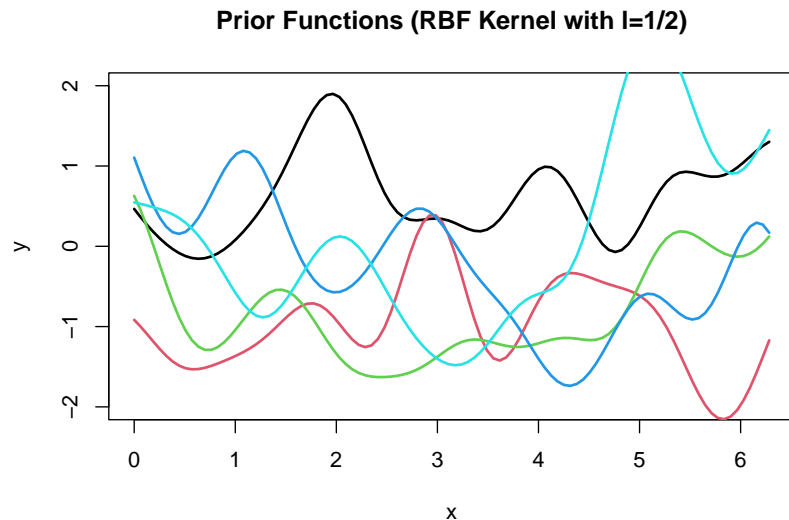
We can now visualise the functions sampled from the multivariate normal prior:

```
# grid of x-values
n1 = 100
x1 = seq(XLIM[1], XLIM[2], length.out=n1)

# unconditioned covariance and mean (unobserved samples x1)
K1 = kfun.mat(x1, x1)
m1 = mu.vec(x1)

## sample functions from GP prior
B = 5
library("MASS") # for mvrnorm
y1r = t(mvrnorm(B, mu = m1, Sigma=K1))
```

```
plot(x1, y1r[,1], type="l", lwd=2, ylab="y", xlab="x", ylim=YLIM,
     main="Prior Functions (RBF Kernel with l=1/2)")
for(i in 2:B)
  lines(x1, y1r[,i], col=i, lwd=2)
```



Now we compute the posterior mean and variance by conditioning on the observations:

```
# unconditioned covariance and mean (observed samples x2)
K2 = kfun.mat(x2, x2)
m2 = mu.vec(x2)
iK2 = solve(K2) # inverse

# cross-covariance
K12 = kfun.mat(x1, x2)

# Conditioning: x1 conditioned on x2

# conditional mean
m1.2 = m1 + K12 %*% iK2 %*% (y2 - m2)

# conditional variance
K1.2 = K1 - K12 %*% iK2 %*% t(K12)
```

Now we can plot the posterior mean and upper and lower bounds of a 95%

credible interval:

```
# upper and lower CI
upper.bound = m1.2 + 1.96*sqrt(diag(K1.2))

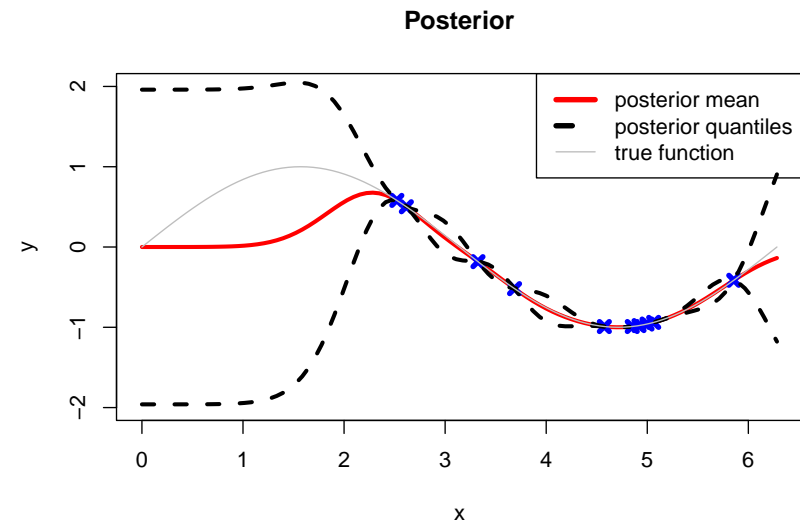
## Warning in sqrt(diag(K1.2)): NaNs wurden erzeugt
lower.bound = m1.2 - 1.96*sqrt(diag(K1.2))

## Warning in sqrt(diag(K1.2)): NaNs wurden erzeugt

plot(x1, m1.2, type="l", xlim=XLIM, ylim=YLIM, col="red", lwd=3,
     ylab="y", xlab="x", main="Posterior")

points(x2, y2, pch=4, lwd=4, col="blue")
lines(x1, upper.bound, lty=2, lwd=3)
lines(x1, lower.bound, lty=2, lwd=3)
curve(truefunc(x), xlim=XLIM, add=TRUE, col="gray")

legend(x="topright",
      legend=c("posterior mean", "posterior quantiles", "true function"),
      lty=c(1, 2, 1), lwd=c(4, 4, 1), col=c("red", "black", "gray"), cex=1.0)
```



Finally, we can take into account noise at the measured data points by adding an error term:

```
# add some noise
sdeps = 0.1
```

```

K2 = K2 + sdeps^2*diag(1,length(x2))

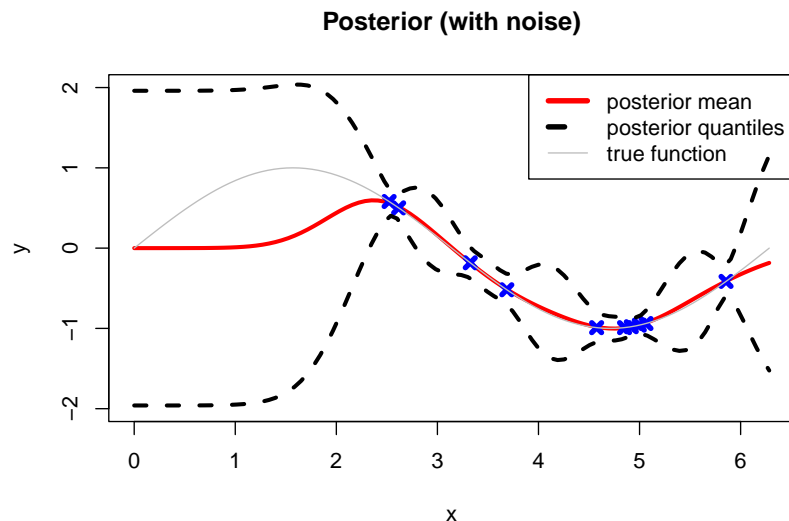
# update
iK2 = solve(K2) # inverse
m1.2 = m1 + K12 %*% iK2 %*% (y2 - m2)
K1.2 = K1 - K12 %*% iK2 %*% t(K12)
upper.bound = m1.2 + 1.96*sqrt(diag(K1.2))
lower.bound = m1.2 - 1.96*sqrt(diag(K1.2))

plot(x1, m1.2, type="l", xlim=XLIM, ylim=YLIM, col="red", lwd=3,
     ylab="y", xlab = "x", main = "Posterior (with noise)")

points(x2,y2,pch=4,lwd=4,col="blue")
lines(x1,upper.bound,lty=2,lwd=3)
lines(x1,lower.bound,lty=2,lwd=3)
curve(truefunc(x), xlim=XLIM, add=TRUE, col="gray")

legend(x="topright",
      legend=c("posterior mean", "posterior quantiles", "true function"),
      lty=c(1, 2, 1),lwd=c(4, 4, 1), col=c("red", "black", "gray"), cex=1.0)

```



Note that in the vicinity of data points the CIs are small and the further away from data the more uncertain the estimate of the underlying function becomes.

7.3 Neural networks

Another highly important class of models for nonlinear prediction (and nonlinear function approximation) are neural networks.

Relevant reading:

Please read: James et al. (2021) or James et al. (2023) **Chapter 10 “Deep Learning”**

7.3.1 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 (e.g. 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
- efficient automatic symbolic differentiation (e.g. autograd)
- stochastic gradient descent algorithms (e.g. Adam)
- use of GPUs and TPUs (e.g. for linear algebra)
- availability of packages for symbolic tensor computations and deep learning.

Currently the most popular frameworks are:

- [PyTorch](#) (PyTorch Foundation, formerly Meta/Facebook)
- [TensorFlow](#) (Google Research)
- [Flax / JAX](#) (Google Research)

and high-level wrappers:

- [skorch](#) (scikit-learn wrapper for PyTorch)
- [Keras 3](#) (for TensorFlow, JAX, and PyTorch)

7.3.2 Neural networks

Neural networks are essentially stacked systems of linear regressions, with nonlinear mappings between each layer, mapping the input to output via one or more layers of internal hidden nodes corresponding to internal latent variables:

- Each internal node is a nonlinear function of all or some of nodes in the previous layer

- Typically, the output of a node is computed using a **non-linear activation function**, such as the sigmoid function or a piecewise linear function (ReLU), from a linear combination of the input variables of that node.

A simple architecture is a feedforward network with a single hidden layer. More complex models are multilayer perceptrons and convolutional neural networks.

It can be shown that even simple network architectures can (with sufficient number of nodes) approximate any arbitrary non-linear function. This is called the **universal function approximation** property.

“Deep” neural networks have many layers, and the optimisation of their parameters requires advanced techniques (see above), with the objective function typically an empirical risk based on, e.g., squared error loss or cross-entropy loss. Neural networks are very highly parameterised models and therefore require lots of data for training, and typically also some form of regularisation (e.g. dropout).

As an extreme example, the neural network behind the [ChatGPT 4](#) language model that is trained on essentially the whole freely accessible text available on the internet has an estimated 1.76 trillion (!) parameters (1.76×10^{12}).

In the limit of an infinite width a single layer fully connected neural network becomes equivalent to a Gaussian process. This was first shown by R. M. Neal (1996)¹. More recently, this equivalence has also been demonstrated for other types of neural networks (with the kernel function of the GP being determined by the neural network architecture). This is formalised in the “neural tangent kernel” (NTK) framework.

Some of the statistical aspects of neural networks are not well understood. For example, there is the **paradox that neural networks typically overfit the training data but still generalise well** - this clearly violates the traditional understanding of bias-variance tradeoff for classical modelling in statistics and machine learning — see for example Belkin et al. (2019)². Some researchers argue that this contradiction can be resolved by better understanding the effective dimension of complex models. There is a lot of current research to explain this phenomenon of “multiple descent”, i.e. the decrease of prediction error for models with very many parameters. A further topic is robustness of the predictions, which is also caused by overfitting. It is well known that neural networks can sometimes be “fooled” by so-called adversarial examples, e.g., the classification of a sample may change if a small amount of noise is added to the test data.

¹Neal, R. M. 1996. Bayesian Learning for Neural Networks. Springer. <https://doi.org/10.1007/978-1-4612-0745-0>

²Belkin, M. et al. 2019. Reconciling modern machine-learning practice and the classical bias-variance trade-off. PNAS 116: 15849–15854. <https://doi.org/10.1073/pnas.1903070116>

7.3.3 Learning more about deep learning

A good place to to learn more about deep learning and actual application in computer code using various deep learning software frameworks is the online book “Dive into deep learning” by Zhang et al. (2023) available online at <https://d2l.ai/>

Appendix A

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.

A.1 Recommended reading

For multivariate statistics and machine learning:

- Izenman (2008) *Modern Multivariate Statistical Techniques*. Springer.
- Rogers and Girolami (2017) *A first course in machine learning (2nd Edition)*. Chapman and Hall / CRC.
- James et al. (2021) *An introduction to statistical learning with applications in R (2nd edition)*. Springer.
- James et al. (2023) *An introduction to statistical learning with applications in Python*. Springer.

A.2 Advanced reading

Additional (advanced) reference books for probabilistic machine learning are:

- Bishop (2006) *Pattern recognition and machine learning*. Springer.
- Murphy (2022) *Probabilistic Machine Learning: An Introduction*. MIT Press.
- Murphy (2023) *Probabilistic Machine Learning: Advanced Topics*. MIT Press.

You can find further suggestions on my list of [online textbooks in statistics and machine learning](#).

Bibliography

- Bishop, C. M. 2006. *Pattern Recognition and Machine Learning*. Springer. <https://www.microsoft.com/en-us/research/people/cmbishop/prml-book/>.
- Izenman, A. J. 2008. *Modern Multivariate Statistical Techniques*. New York: Springer. <https://doi.org/10.1007/978-0-387-78189-1>.
- James, G., D. Witten, T. Hastie, and R. Tibshirani. 2021. *An Introduction to Statistical Learning with Applications in R*. 2nd ed. Springer. <https://doi.org/10.1007/978-1-0716-1418-1>.
- James, G., D. Witten, T. Hastie, R. Tibshirani, and J. Taylor. 2023. *An Introduction to Statistical Learning with Applications in Python*. Springer. <https://doi.org/10.1007/978-3-031-38747-0>.
- Murphy, K. P. 2022. *Probabilistic Machine Learning: An Introduction*. MIT Press. <https://probml.github.io/pml-book/book1.html>.
- . 2023. *Probabilistic Machine Learning: Advanced Topic*. MIT Press. <https://probml.github.io/pml-book/book2.html>.
- Rogers, S., and M. Girolami. 2017. *A First Course in Machine Learning*. 2nd ed. Chapman; Hall / CRC. <https://doi.org/10.1201/9781315382159>.
- Zhang, A., Z. C. Lipton, M. Li, and A. J. Smola. 2023. *Dive into Deep Learning*. <https://d2l.ai>.