

# Statistical Methods: Likelihood, Bayes and Regression

Korbinian Strimmer

14 May 2021



# Contents

<b>Welcome</b>	<b>7</b>
License . . . . .	7
<b>Preface</b>	<b>9</b>
About the author . . . . .	9
About the module . . . . .	9
Acknowledgements . . . . .	10
<b>I Likelihood estimation and inference</b>	<b>11</b>
<b>1 Overview of statistical learning</b>	<b>13</b>
1.1 How to learn from data? . . . . .	13
1.2 Probability theory versus statistical learning . . . . .	14
1.3 Cartoon of statistical learning . . . . .	15
1.4 Likelihood . . . . .	16
<b>2 From entropy to maximum likelihood</b>	<b>19</b>
2.1 Entropy . . . . .	19
2.2 Kullback-Leibler divergence . . . . .	24
2.3 Local quadratic approximation and expected Fisher information .	26
2.4 Entropy learning and maximum likelihood . . . . .	28
<b>3 Maximum likelihood estimation</b>	<b>33</b>
3.1 Principle of maximum likelihood estimation . . . . .	33
3.2 Maximum likelihood estimation in practise . . . . .	36
3.3 Observed Fisher information . . . . .	41
<b>4 Quadratic approximation and normal asymptotics</b>	<b>47</b>
4.1 Multivariate statistics for random vectors . . . . .	47
4.2 Approximate distribution of maximum likelihood estimates . . .	50
4.3 Quantifying the uncertainty of maximum likelihood estimates . .	54
4.4 Example of a non-regular model . . . . .	60

<b>5</b>	<b>Likelihood-based confidence interval and likelihood ratio</b>	<b>63</b>
5.1	Likelihood-based confidence intervals and Wilks statistic . . . . .	63
5.2	Generalised likelihood ratio test (GLRT) . . . . .	69
<b>6</b>	<b>Optimality properties and conclusion</b>	<b>73</b>
6.1	Properties of maximum likelihood encountered so far . . . . .	73
6.2	Summarising data and the concept of minimal sufficiency . . . . .	74
6.3	Concluding remarks on maximum likelihood . . . . .	77
<b>II</b>	<b>Bayesian Statistics</b>	<b>79</b>
<b>7</b>	<b>Essentials of Bayesian statistics</b>	<b>81</b>
7.1	Conditional probability . . . . .	81
7.2	Bayes' theorem . . . . .	82
7.3	Principle of Bayesian learning . . . . .	82
7.4	What is exactly is the "Bayesian estimate"? . . . . .	83
7.5	Computer implementation of Bayesian learning . . . . .	84
7.6	Bayesian interpretation of probability . . . . .	85
7.7	Historical developments . . . . .	86
7.8	Connection with entropy learning . . . . .	87
<b>8</b>	<b>Beta-Binomial model for estimating a proportion</b>	<b>89</b>
8.1	Binomial likelihood . . . . .	89
8.2	Excursion: Properties of the Beta distribution . . . . .	90
8.3	Beta prior distribution . . . . .	91
8.4	Computing the posterior distribution . . . . .	91
<b>9</b>	<b>Properties of Bayesian learning</b>	<b>93</b>
9.1	Prior acting as pseudo-data . . . . .	93
9.2	Linear shrinkage of mean . . . . .	94
9.3	Conjugacy of prior and posterior distribution . . . . .	95
9.4	Large sample asymptotics . . . . .	95
9.5	Posterior variance for finite $n$ . . . . .	96
<b>10</b>	<b>Normal-Normal and Inverse-Gamma-Normal models for estimating the mean and the variance</b>	<b>99</b>
10.1	Normal-Normal model to estimate mean . . . . .	99
10.2	Inverse-Gamma-Normal model to estimate variance . . . . .	101
<b>11</b>	<b>Shrinkage estimation using empirical risk minimisation</b>	<b>103</b>
11.1	Linear shrinkage . . . . .	103
11.2	James-Stein estimator . . . . .	104
<b>12</b>	<b>Bayesian model comparison using Bayes factors and the BIC</b>	<b>105</b>
12.1	The Bayes factor . . . . .	105

12.2	Approximate computation of the marginal likelihood and of the log-Bayes factor . . . . .	108
<b>13</b>	<b>False discovery rates</b>	<b>111</b>
13.1	General setup . . . . .	111
13.2	Specificity and Sensitivity . . . . .	112
13.3	FDR and FNDR . . . . .	112
13.4	Bayesian perspective . . . . .	113
13.5	Software . . . . .	114
<b>14</b>	<b>Optimality properties and summary</b>	<b>115</b>
14.1	Bayesian statistics in a nutshell . . . . .	115
14.2	Frequentist properties of Bayesian estimators . . . . .	116
14.3	Specifying the prior — problem or advantage? . . . . .	117
14.4	Choosing a prior . . . . .	117
14.5	Optimality of Bayesian inference . . . . .	119
14.6	Conclusion . . . . .	119
<b>III</b>	<b>Regression</b>	<b>121</b>
<b>15</b>	<b>Overview over regression modelling</b>	<b>123</b>
15.1	General setup . . . . .	123
15.2	Objectives . . . . .	124
15.3	Regression as a form of supervised learning . . . . .	124
15.4	Various regression models used in statistics . . . . .	125
<b>16</b>	<b>Linear Regression</b>	<b>127</b>
16.1	The linear regression model . . . . .	127
16.2	Interpretation of regression coefficients and intercept . . . . .	128
16.3	Different types of linear regression: . . . . .	128
16.4	Distributional assumptions and properties . . . . .	128
16.5	Regression in data matrix notation . . . . .	130
16.6	Centering and vanishing of the intercept $\beta_0$ . . . . .	130
16.7	Regression objectives for linear model . . . . .	131
<b>17</b>	<b>Estimating regression coefficients</b>	<b>133</b>
17.1	Ordinary Least Squares (OLS) estimator of regression coefficients	133
17.2	Maximum likelihood estimation of regression coefficients . . . . .	135
17.3	Covariance plug-in estimator of regression coefficients . . . . .	137
17.4	Best linear predictor . . . . .	138
17.5	Regression by conditioning . . . . .	140
17.6	Standardised regression coefficients and relationship to correlation	141
<b>18</b>	<b>Squared multiple correlation and variance decomposition in linear regression</b>	<b>143</b>
18.1	Squared multiple correlation $\Omega^2$ and the $R^2$ coefficient . . . . .	143

18.2	Variance decomposition in regression . . . . .	145
18.3	Sample version of variance decomposition . . . . .	147
<b>19</b>	<b>Prediction and variable selection</b>	<b>149</b>
19.1	Prediction and prediction intervals . . . . .	149
19.2	Variable importance and prediction . . . . .	150
19.3	Regression <i>t</i> -scores. . . . .	152
19.4	Further approaches for variable selection . . . . .	154
<b>Appendix</b>		<b>157</b>
<b>A</b>	<b>Refresher</b>	<b>159</b>
A.1	Basic mathematical notation . . . . .	159
A.2	Vectors and matrices . . . . .	159
A.3	Functions . . . . .	160
A.4	Combinatorics . . . . .	162
A.5	Probability . . . . .	164
A.6	Distributions . . . . .	167
A.7	Statistics . . . . .	172
<b>B</b>	<b>Further study</b>	<b>179</b>
B.1	Recommended reading . . . . .	179
B.2	Additional references . . . . .	179
<b>Bibliography</b>		<b>181</b>

# Welcome

These are the lecture notes for MATH20802, a course in **Statistical Methods** for second year mathematics students at the [Department of Mathematics of the University of Manchester](#).

The course text was written by [Korbinian Strimmer](#) from 2019–2021. This version is from 13 May 2021.

The notes will be updated from time to time. To view the current version visit the [online MATH20802 lecture notes](#). You may also [download the MATH20802 lecture notes as PDF](#).

## 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 part 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 spring term 2019 at the University of Manchester, and subsequently also in 2020 and in 2021.

I hope you enjoy the course! If you have any questions, comments, or corrections then please email me at [korbinian.strimmer@manchester.ac.uk](mailto:korbinian.strimmer@manchester.ac.uk).

## About the module

### Topics covered

The MATH20802 module is designed to run over the course of 11 weeks. It has three parts:

1. Likelihood estimation and likelihood ratio tests (W1–W5)
2. Bayesian learning and inference (W6–W8)
3. Linear regression (W9–W11)

This module focuses on conceptual understanding and methods, not on theory. As such, the presentation in this course is non-technical. The aim is to offer insights how diverse statistical approaches are linked and to demonstrate that statistics offers a concise and coherent theory of information rather than being an adhoc collection of “recipes” for data analysis (a common but wrong perception of statistics).

### Prerequisites

For this module it is important that you refresh your knowledge in:

- Introduction to statistics

- Probability
- R data analysis and programming

In addition you will need to know matrix algebra and how to compute the gradient and the curvature of a function of several variables.

Check the Appendix for a brief refresher of the essential material.

## Additional support material

Accompanying these notes are

- [lecture videos](#) (visualiser style).

Furthermore, there is also an [MATH20802 online reading list](#) hosted by the University of Manchester library.

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 examples and solutions and R code, and
- exam papers of previous years.

## Acknowledgements

Many thanks to [Beatriz Costa Gomes](#) for her help in creating the 2019 version of the lecture notes and to [Kristijonas Raudys](#) for his extensive feedback on the 2020 version.

## **Part I**

# **Likelihood estimation and inference**



# Chapter 1

## Overview of statistical learning

### 1.1 How to learn from data?

A fundamental question is how to extract information from data in an optimal way, and to make predictions based on this information.

For this purpose, a number of competing **theories of information** have been developed. **Statistics** is the oldest science of information and is concerned with offering principled ways to learn from data and to extract and process information using probabilistic models. However, there are other theories of information (e.g. Vapnik-Chernov theory of learning, computational learning) that are more algorithmic than analytic and sometimes not even based on probability theory.

Furthermore, there are other disciplines, such computer science and machine learning that are closely linked with and also have substantial overlap with statistics. The field of “data science” today comprises both statistics and a machine learning and brings together mathematics, statistics and computer science. Also the growing field of so-called “artificial intelligence” makes substantial use of statistical and machine learning techniques.

The recent popular science book “The Master Algorithm” by Domingos (2015) provides an accessible informal overview over the various schools of science of information. It discusses the main algorithms used in machine learning and statistics:

- Starting as early as 1763, the **Bayesian school** of learning was started which later turned out to be closely linked with *likelihood inference* established in 1922 by [R.A. Fisher \(1890–1962\)](#) and generalised in 1951 to **entropy learning** by Kullback and Leibler.

- It was also in the 1950s that the concept of artificial **neural network** arises, essentially a nonlinear input-output map that works in a non-probabilistic way. This field saw another leap in the 1980 and further progress from 2010 onwards with the development of *deep learning*. It is now one of the most popular (and most effective) methods for analysis of imaging data. Even your mobile phone most likely has a dedicated computer chip with special neural network hardware, for example.
- Further advanced theories of information were developed in the 1960 under the term of **computational learning**, most notably the Vapnik-Chernov theory, with the most prominent example of the “support vector machine” (another non-probabilistic model).
- With the advent of large-scale genomic and other high-dimensional data there has been a surge of new and exciting developments in the field of high-dimensional (large dimension) and also big data (large dimension and large sample size), both in statistics and in machine learning.

**The connections between various fields of information is still not perfectly understood, but it is clear that an overarching theory will need to be based on probabilistic learning.**

## 1.2 Probability theory versus statistical learning

When you study statistics (or any other information theory) you need to be aware that there is a fundamental difference between probability theory and statistics, and that relates to the **distinction between “randomness” and “uncertainty”**.

Probability theory studies **randomness**, by developing mathematical models for randomness (such as probability distributions), and studying corresponding mathematical properties (including asymptotics etc). Probability theory may in fact be viewed as a branch of measure theory, and thus it belongs to the domain pure mathematics.

Probability theory provides probabilistic generative models for data, for simulation of data or for use in learning from data, i.e. inference about the model from observations. Methods and theory how to best learn from data is the domain of applied mathematics, specifically statistics and the related areas of machine learning and data science.

Note that statistics, in contrast to probability, is in fact not at all concerned with randomness. Instead, the focus is about measuring and elucidating the **uncertainty** of events, predictions, outcomes, parameters and this uncertainty measures the **state of knowledge**. Note that if new data or information becomes available, the state of knowledge and thus the uncertainty changes! Thus, **uncertainty is an epistemological property**.

The uncertainty most often is due to our ignorance of the true underlying pro-

cesses (on purpose or not), but not because the underlying process is actually random. The success of statistics is based on the fact that we can mathematically model the uncertainty without knowing any specifics of the underlying processes, and we still have procedures for optimal inference under uncertainty.

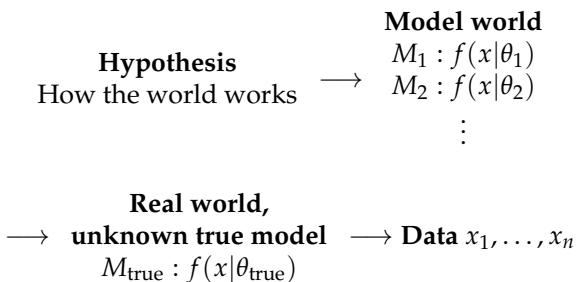
In short, statistics is about describing the state of knowledge of the world, which may be uncertain and incomplete, and to make decisions and prediction in the face of uncertainty, and this uncertainty sometimes derives from randomness but most often from our ignorance (and sometimes this ignorance even helps to create a simple yet effective model)!

## 1.3 Cartoon of statistical learning

We observe data  $x_1, \dots, x_n$  assumed to be generated by the underlying true model  $M_{\text{true}}$ .

To explain the data, and make prediction, we make hypotheses in the form of candidate models  $M_1, M_2, \dots$ . The true model  $M_{\text{true}}$  itself is unknown and cannot be observed. However, what we can observe is a finite amount of data from the model by measuring properties of objects interest (our observations from experiments). Sometimes we can also perturb the model and see what the effect is (interventional study).

The various candidate models  $M_1, M_2, \dots$  in the **model world** will never be perfect or correct as the true model  $M_{\text{true}}$  will only be among the candidate models in an idealised situation. However, even an imperfect candidate model will often provide a useful mathematical approximation and capture some important characteristics of the true model and thus will help to interpret observed data..



**The aim of statistical learning is to identify the model(s) that explain the current data and also predict future data (i.e. predict outcome of experiments that have not been conducted yet).**

Thus a good model provides a good fit to the current data (i.e. it explains current observations well) and also to future data (i.e. it generalises well).

A large proportion of statistical theory is devoted to finding these “good” models that avoid both *overfitting* (models being too complex and don’t generalise well) or *underfitting* (models being too simplistic and hence also don’t predict well).

Typically the aim is to find a model whose the **model complexity** matches the complexity of the unknown true model and also the complexity of the data observed from the unknown true model.

## 1.4 Likelihood

A core problem in statistics is how to find probabilistic models for explaining existing data and predicting new data. For this we need a measure of how good a hypothesis/candidate model  $M_k$  is as approximation for the (typically unknown) true data generating model  $M_{\text{true}}$ .

As you already know from the year 1 module MATH10282 “Introduction to Statistics”, one such measure is provided by the likelihood function which helps to choose among the various candidate models and estimate corresponding parameters by finding the model  $M$  that maximises the (log)-likelihood.

Given a probability distribution  $F_\theta$  with density or mass function  $f(x|\theta)$  where  $\theta$  is a parameter vector, and  $x_1, \dots, x_n$  is the observed iid data (i.e. independent and identically distributed), the **likelihood function** is defined as

$$L_n(\theta) = \prod_{i=1}^n f(x_i|\theta)$$

Typically, instead of the likelihood one uses the log-likelihood function:

$$\log L_n(\theta) = l_n(\theta) = \sum_{i=1}^n \log f(x_i|\theta)$$

Reasons for using log-likelihood (rather than likelihood) include that

- the log density is in fact the more “natural” and relevant quantity (this will become clear in the upcoming chapters) and that
- addition is numerically more stable than multiplication on a computer.

For discrete random variables for which  $f(x|\theta)$  is a probability mass function the likelihood is often interpreted as the probability to observe the data given the model with specified parameters  $\theta$ . In fact, this was indeed the way how the likelihood was historically introduced. However, this view is not strictly correct. First, given that the samples are iid and thus the ordering of the  $x_i$  is not important, an additional factor accounting for the possible permutations is needed in the above to obtain the actual probability of the data. Moreover, for continuous random variables this interpretation breaks down due to the use of densities rather than probability mass functions in the likelihood. Thus, the view of the likelihood being the probability of the data is in fact too simplistic.



In the next chapter we will see that the justification for using likelihood rather stems from its close link to the Kullback-Leibler information and cross-entropy. This also helps to understand why using likelihood for estimation is only optimal in the limit of large sample size.

In the first part of the MATH28082 “Statistical Methods” module we will study likelihood estimation and inference in much detail. We will provide links to related methods of inference and discuss its information-theoretic foundations. We will also discuss the optimality properties as well as the limitation of likelihood inference. Extensions of likelihood analysis, in particular Bayesian learning, which will be discussed in the second part module. In the third part of the module we will apply statistical learning to linear regression.



# Chapter 2

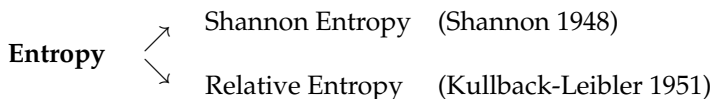
## From entropy to maximum likelihood

### 2.1 Entropy

#### 2.1.1 Overview

In this chapter we discuss various information criteria and their connection to maximum likelihood.

The modern definition of (relative) entropy, or “disorder”, was first discovered in 1875 by physicist [L. Boltzmann \(1844–1906\)](#) in the context of thermodynamics. In the 1940–1950’s the notion of entropy turned out to be central in information theory, a field pioneered by mathematicians such as [R. Hartley \(1988–1970\)](#), [S. Kullback \(1907–1994\)](#), [R. Leibler \(1914–2003\)](#), [A. Turing \(1912–1954\)](#), [I. J. Good \(1916–2009\)](#), [C. Shannon \(1916–2001\)](#), and [E. T. Jaynes \(1922–1998\)](#), and later further explored by [S. Amari \(1936–\)](#), [I. Ciszár \(1938–\)](#), [B. Efron \(1938–\)](#), [A. P. Dawid \(1946–\)](#) and many others.



Fisher information → Likelihood theory (Fisher 1922)

Mutual Information → Information theory (Shannon 1948, Lindley 1953)

### 2.1.2 Surprise, surprisal or Shannon information

The **surprise** to observe an event of probability  $p$  is **defined** as  $-\log(p)$ . This is also called **surprisal** or **Shannon information**.

Thus, the surprise to observe a certain event (with  $p = 1$ ) is zero, and conversely the surprise to observe an event that is certain not to happen (with  $p = 0$ ) is infinite.

The **log-odds ratio** can be viewed as the difference of the surprise of an event and the surprise of the complementary event:

$$\log\left(\frac{p}{1-p}\right) = -\log(1-p) - (-\log(p))$$

In this module we always use the *natural logarithm* by default, and will explicitly write  $\log_2$  and  $\log_{10}$  for logarithms with respect to base 2 and 10, respectively.

Surprise and entropy computed with the natural logarithm ( $\log$ ) is given in “nats” (=natural information units). Using  $\log_2$  leads to “bits” and using  $\log_{10}$  to “ban” or “Hartley”.

### 2.1.3 Shannon entropy

Assume we have a discrete distribution  $F$  with  $K$  classes and class probabilities  $p_1, \dots, p_K$  with  $\Pr(\text{"class k"}) = p_k$  and  $\sum_{k=1}^K p_k = 1$ . Let  $x$  indicate the selected class then define the PMF as  $f(x = \text{"class k"}) = p_k$ .

The **Shannon entropy** of the discrete distribution  $F$  is defined as the **expected surprise**, i.e. the negative expected log-probability

$$\begin{aligned} H(F) &= -E_F(\log f(x)) \\ &= -\sum_{i=1}^K p_i \log(p_i) \end{aligned}$$

As all  $p_k \in [0, 1]$  by construction Shannon entropy must be larger or equal to 0. Furthermore, it is bounded above by  $\log K$ . Hence for any discrete distribution  $F$  with  $K$  categories we have

$$\log K \geq H(F) \geq 0$$

**Example 2.1. Discrete uniform distribution  $U_K$ :** let  $p_1 = p_2 = \dots = p_K = \frac{1}{K}$ . Then

$$H(U_K) = -\sum_{i=1}^K \frac{1}{K} \log\left(\frac{1}{K}\right) = \log K$$

Note this is the largest value the Shannon entropy can assume with  $K$  classes.

**Example 2.2. Concentrated probability mass:** let  $p_1 = 1$  and  $p_2 = p_3 = \dots = p_K = 0$ . Using  $0 \times \log(0) = 0$  we obtain for the Shannon probability

$$H(F) = 1 \times \log(1) + 0 \times \log(0) + \dots = 0$$

Note that 0 is the smallest value that Shannon entropy can assume, and corresponds to maximum concentration.

Thus, **large entropy** implies that the **distribution is spread out** whereas **small entropy** means the **distribution is concentrated**.

Correspondingly, maximum entropy distributions can be considered minimally informative about a random variable.

This interpretation is also supported by the close link of Shannon entropy with multinomial coefficients counting the permutations of  $n$  items (samples) of  $K$  distinct types (classes).

**Example 2.3.** Large sample asymptotics of the multinomial coefficient and Shannon entropy:

The number of possible permutation of  $n$  items of  $K$  distinct types, with  $n_1$  of type 1,  $n_2$  of type 2 and so on, is given by the multinomial coefficient

$$\binom{n}{n_1, \dots, n_K} = \frac{n!}{n_1! \times n_2! \times \dots \times n_K!}$$

with  $\sum_{k=1}^K n_k = n$  and  $K \leq n$ .

Using the the Moivre-Sterling formula for large  $n$  the factorial can be approximated as

$$\log n! \approx n \log n - n$$

As a result

$$\begin{aligned} \log \binom{n}{n_1, \dots, n_K} &= \log n! - \sum_{k=1}^K \log n_k! \\ &\approx n \log n - n - \sum_{k=1}^K (n_k \log n_k - n_k) \\ &= n \log n - \sum_{k=1}^K n_k \log n_k \\ &= \sum_{k=1}^K n_k \log n - \sum_{k=1}^K n_k \log n_k \\ &= - \sum_{k=1}^K n_k \log \left( \frac{n_k}{n} \right) \end{aligned}$$

and thus

$$\begin{aligned} \frac{1}{n} \log \binom{n}{n_1, \dots, n_K} &\approx - \sum_{k=1}^K \frac{n_k}{n} \log \left( \frac{n_k}{n} \right) \\ &= - \sum_{k=1}^K \hat{p}_k \log (\hat{p}_k) \\ &= H(\hat{F}) \end{aligned}$$

where  $\hat{F}$  is the empirical discrete distribution with  $\hat{p}_k = \frac{n_k}{n}$ .

The combinatorial derivation of entropy is due to Boltzmann (1877) who discovered it in his research in statistical mechanics and thermodynamics.

### 2.1.4 Differential entropy

Shannon entropy is only defined for discrete random variables.

*Differential Entropy* results from applying the definition of Shannon entropy to a *continuous* random variable  $x$  with density  $f(x)$ :

$$H(F) = -E_F(\log f(x)) = - \int f(x) \log f(x) dx$$

Despite having essentially the same formula the different name is justified because differential entropy exhibits different properties compared to Shannon entropy, because the logarithm is taken of a density which in contrast to a probability can assume values larger than one. As a consequence, differential entropy is *not* bounded below by zero and can be negative.

**Example 2.4.** Consider the uniform distribution  $U(0, a)$  with  $a > 0$ , support from 0 to  $a$  and density  $f(x) = 1/a$ . As  $-\int_0^a f(x) \log f(x) dx = -\int_0^a \frac{1}{a} \log(\frac{1}{a}) dx = \log a$  the differential entropy is

$$H(U(0, a)) = \log a.$$

Note that for  $a < 1$  the differential entropy is negative.

**Example 2.5.** The density of the univariate normal  $N(\mu, \sigma^2)$  distribution is  $f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$  with  $\sigma^2 > 0$ . The corresponding differential entropy is

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

Note that it only depends on the variance and not on the mean, and that for  $\sigma^2 < 1/(2\pi e) \approx 0.0585$  the differential entropy is negative.

### 2.1.5 Maximum entropy principle to characterise distributions

Both maximum Shannon entropy and differential entropy are useful to characterise distributions:

- 1) The **discrete uniform distribution** is the **maximum entropy distribution** among all discrete distributions.
- 2) the maximum entropy distribution of a continuous random variable with support  $[-\infty, \infty]$  with a specific mean and variance is the normal distribution
- 3) the maximum entropy distribution among all continuous distributions supported in  $[0, \infty]$  with a specified mean is the exponential distribution.

The higher the entropy the more spread out (and more uninformative) is a distribution.

Using maximum entropy to characterise maximally uninformative distributions was advocated by E.T. Jaynes (who also proposed to use maximum entropy in the context of finding Bayesian priors).

A list of maximum entropy distribution is given here: [https://en.wikipedia.org/wiki/Maximum\\_entropy\\_probability\\_distribution](https://en.wikipedia.org/wiki/Maximum_entropy_probability_distribution).

### 2.1.6 Cross-entropy

If in the definition of Shannon entropy (and differential entropy) the expectation over the log-density (say  $g(x)$  of distribution  $G$ ) is with regard to a different distribution  $F$  over the same state space we arrive at the **cross-entropy**

$$H(F, G) = -\mathbb{E}_F(\log g(x))$$

Therefore, cross-entropy is a measure linking two distributions  $F$  and  $G$ .

Note that

- cross-entropy is not symmetric with regard to  $F$  and  $G$ , because the expectation is taken with reference to  $F$ .
- By construction  $H(F, F) = H(F)$ .

A crucial property of the cross-entropy  $H(F, G)$  is that it is bounded below by the entropy of  $F$ , therefore

$$H(F, G) \geq H(F)$$

with equality for  $F = G$ .

Equivalently we can write

$$H(F, G) - H(F) \geq 0$$

In fact, this recalibrated cross-entropy turns out to be more fundamental than

both cross-entropy and Shannon resp. differential entropy. It will be studied in detail in the next section.

**Example 2.6.** Cross-entropy between two normals:

Assume  $F_{\text{ref}} = N(\mu_{\text{ref}}, \sigma_{\text{ref}}^2)$  and  $F = N(\mu, \sigma^2)$ . Then the cross-entropy is

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

**Example 2.7.** If  $\mu_{\text{ref}} = \mu$  and  $\sigma_{\text{ref}}^2 = \sigma^2$  then the above cross-entropy  $H(F, G)$  degenerates to the differential entropy  $H(F_{\text{ref}}) = \frac{1}{2} (\log(2\pi\sigma_{\text{ref}}^2) + 1)$ .

## 2.2 Kullback-Leibler divergence

### 2.2.1 Definition

Also known as **relative entropy** and **discrimination information**.

The **relative entropy** measures the **divergence** of a distribution  $G$  from the distribution  $F$  and is defined as

$$\begin{aligned} D_{\text{KL}}(F, G) &= E_F \log \left( \frac{dF}{dG} \right) \\ &= E_F \log \left( \frac{f(x)}{g(x)} \right) \\ &= \underbrace{-E_F(\log g(x))}_{\text{cross-entropy}} - \underbrace{(-E_F(\log f(x)))}_{\text{(differential) entropy}} \\ &= H(F, G) - H(F) \end{aligned}$$

- $D_{\text{KL}}(F, G)$  measures the amount of information lost if  $G$  is used to approximate  $F$ .
- If  $F$  and  $G$  are identical (and no information is lost) then  $D_{\text{KL}}(F, G) = 0$ .

(Note: here “divergence” measures the dissimilarity between probability distributions. This type of divergence is not related and should not be confused with divergence (div) as used in vector analysis.)

The term divergence (rather than distance) implies also that the distributions  $F$  and  $G$  are not interchangeable in  $D_{\text{KL}}(F, G)$ .

In applications in statistics the typical roles of  $F$  and  $G$  are:

- $F$  as the (unknown) underlying true model for the data generating process
- $G$  as the approximating model (e.g. some parametric family)

In Bayesian statistics we use



- $F$  as posterior distribution
- $G$  as prior distribution

There exist various notations for KL divergence in the literature. Here we use  $D_{KL}(F, G)$  but often you can find  $KL(F||G)$  or  $I^{KL}(F; G)$  in other references.

Some authors (e.g. Efron) call twice the KL divergence  $2D_{KL}(F, G) = D(F, G)$  the **deviance** of  $G$  from  $F$ .

### 2.2.2 Properties of KL divergence

1.  $D_{KL}(F, G) \neq D_{KL}(G, F)$ , i.e., KL divergence is not symmetric,  $F$  and  $G$  cannot be interchanged.
2.  $D_{KL}(F, G) = 0$  if and only if  $F = G$ , i.e., the KL divergence is zero if and only if  $F$  and  $G$  are identical.
3.  $D_{KL}(F, G) \geq 0$ , proof via the **Jensen Inequality**.
4.  $D_{KL}(F, G)$  remains invariant under coordinate transformations, i.e. it is an invariant geometric quantity.

Note that in the KL divergence the expectation is taken over a ratio of densities (or ratio of probabilities for discrete random variables). This is what creates the transformation invariance.

For more details and proofs of properties 3 and 4 see Worksheet 1.

### 2.2.3 Examples

**Example 2.8.** KL divergence between two Bernoulli distributions  $\text{Ber}(p)$  and  $\text{Ber}(q)$ :

The “success” probabilities for the two distributions are  $p$  and  $q$ , respectively, and the complementary “failure” probabilities are  $1 - p$  and  $1 - q$ . With this we get for the KL divergence

$$D_{KL}(\text{Ber}(p), \text{Ber}(q)) = p \log \left( \frac{p}{q} \right) + (1 - p) \log \left( \frac{1 - p}{1 - q} \right)$$

**Example 2.9.** KL divergence between two univariate normals with different means and variances:

Assume  $F_{\text{ref}} = N(\mu_{\text{ref}}, \sigma_{\text{ref}}^2)$  and  $F = N(\mu, \sigma^2)$ . Then

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

**Example 2.10.** KL divergence between two univariate normals with different means and common variance:

An important special case of the previous Example 2.9 occurs if the variances are equal. Then we get

$$D_{\text{KL}}(N(\mu_{\text{ref}}, \sigma^2), N(\mu, \sigma^2)) = \frac{1}{2} \left( \frac{(\mu - \mu_{\text{ref}})^2}{\sigma^2} \right)$$

## 2.3 Local quadratic approximation and expected Fisher information

### 2.3.1 Definition of expected Fisher information

KL information measures the divergence of two distributions. We may thus use relative entropy to measure the divergence between two distributions in the same family, separated in parameter space only by some small  $\varepsilon$ .

Let  $h(\theta) = D_{\text{KL}}(F_{\theta_0}, F_{\theta}) = E_{\theta_0}(\log f(x|\theta_0) - \log f(x|\theta))$ . Note that the first distribution in the KL divergence is fixed at  $F_{\theta_0}$  and the second distribution is varied. Then  $h(\theta_0 + \varepsilon) = D_{\text{KL}}(F_{\theta_0}, F_{\theta_0 + \varepsilon})$ . Since the KL divergence vanishes only when the two arguments are identical  $h(\theta)$  reaches a minimum at  $\theta_0$  with  $h(\theta_0) = 0$  and flat gradient  $\nabla h(\theta_0) = 0$ .

We can therefore approximate  $h(\theta_0 + \varepsilon)$  by a quadratic function around  $\theta_0$

$$\begin{aligned} h(\theta_0 + \varepsilon) &\approx \frac{1}{2} \varepsilon^T \nabla^T \nabla h(\theta_0) \varepsilon \\ &= \frac{1}{2} \varepsilon^T \left( -E_{\theta_0} \nabla^T \nabla \log f(x|\theta_0) \right) \varepsilon \\ &= \frac{1}{2} \varepsilon^T \underbrace{\mathbf{I}^{\text{Fisher}}(\theta_0)}_{\text{expected Fisher information}} \varepsilon \end{aligned}$$

This yields the **expected Fisher information** at  $\theta_0$  as the negative expected Hessian matrix of the log-density at  $\theta_0$ . Since  $\theta_0$  is a minimum the expected Fisher information matrix must be positive definite!

Since there is no data involved the expected Fisher information is purely a property of the model, or more precisely of the space of the models indexed by  $\theta$ . In the next Chapter we will study a related quantity, the *observed Fisher information* that in contrast is a function of the observed data.

We can use the above approximation also to compute the divergence  $D_{\text{KL}}(F_{\theta_0 + \varepsilon}, F_{\theta_0})$  where the first argument varies and the second is kept fixed:

$$D_{\text{KL}}(F_{\theta_0 + \varepsilon}, F_{\theta_0}) \approx \frac{1}{2} \varepsilon^T \mathbf{I}^{\text{Fisher}}(\theta_0 + \varepsilon) \varepsilon$$

In a linear approximation  $\mathbf{I}^{\text{Fisher}}(\theta_0 + \varepsilon) \approx \mathbf{I}^{\text{Fisher}}(\theta_0) + \Delta_{\varepsilon}$  each element of the matrix  $\Delta_{\varepsilon}$  is the scalar product of  $\varepsilon$  and the gradient of the corresponding

element in  $\mathbf{I}^{\text{Fisher}}(\boldsymbol{\theta}_0)$  evaluated at  $\boldsymbol{\theta}_0$ . Therefore  $\boldsymbol{\varepsilon}^T \boldsymbol{\Delta}_{\boldsymbol{\varepsilon}} \boldsymbol{\varepsilon}$  is of *cubic order* in  $\boldsymbol{\varepsilon}$  and hence

$$\begin{aligned} D_{\text{KL}}(F_{\boldsymbol{\theta}_0 + \boldsymbol{\varepsilon}}, F_{\boldsymbol{\theta}_0}) &\approx \frac{1}{2} \boldsymbol{\varepsilon}^T \mathbf{I}^{\text{Fisher}}(\boldsymbol{\theta}_0 + \boldsymbol{\varepsilon}) \boldsymbol{\varepsilon} \\ &\approx \frac{1}{2} \boldsymbol{\varepsilon}^T \mathbf{I}^{\text{Fisher}}(\boldsymbol{\theta}_0) \boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}^T \boldsymbol{\Delta}_{\boldsymbol{\varepsilon}} \boldsymbol{\varepsilon} \\ &\approx \frac{1}{2} \boldsymbol{\varepsilon}^T \mathbf{I}^{\text{Fisher}}(\boldsymbol{\theta}_0) \boldsymbol{\varepsilon} \end{aligned}$$

keeping only terms quadratic in  $\boldsymbol{\varepsilon}$ .

### 2.3.2 Examples

**Example 2.11.** Expected Fisher information for the Bernoulli distribution:

The log-probability mass function of the Bernoulli  $\text{Ber}(p)$  distribution is

$$\log f(x|p) = x \log(p) + (1 - x) \log(1 - p)$$

where  $p$  is the proportion of “success”. The second derivative with regard to the parameter  $p$  is

$$\frac{d^2}{dp^2} \log f(x|p) = -\frac{x}{p^2} - \frac{1-x}{(1-p)^2}$$

Since  $E(x) = p$  we get as Fisher information

$$\begin{aligned} I^{\text{Fisher}}(p) &= -E \left( \frac{d^2}{dp^2} \log f(x|p) \right) \\ &= \frac{p}{p^2} + \frac{1-p}{(1-p)^2} \\ &= \frac{1}{p(1-p)} \end{aligned}$$

**Example 2.12.** Quadratic approximations of the KL divergence between two Bernoulli distributions:

From Example 2.8 we have as KL divergence

$$D_{\text{KL}}(\text{Ber}(p_1), \text{Ber}(p_2)) = p_1 \log \left( \frac{p_1}{p_2} \right) + (1 - p_1) \log \left( \frac{1 - p_1}{1 - p_2} \right)$$

and from Example 2.11 the corresponding expected Fisher information.

The quadratic approximation implies that

$$D_{\text{KL}}(\text{Ber}(p), \text{Ber}(p + \boldsymbol{\varepsilon})) \approx \frac{\boldsymbol{\varepsilon}^2}{2} I^{\text{Fisher}}(p) = \frac{\boldsymbol{\varepsilon}^2}{2p(1-p)}$$

and also that

$$D_{\text{KL}}(\text{Ber}(p + \varepsilon), \text{Ber}(p)) \approx \frac{\varepsilon^2}{2} I^{\text{Fisher}}(p) = \frac{\varepsilon^2}{2p(1-p)}$$

In Worksheet 1 this is verified by using a second order Taylor series applied to the KL divergence.

**Example 2.13.** Expected Fisher information for the normal distribution  $N(\mu, \sigma^2)$ .

The log-density is

$$\log f(x|\mu, \sigma^2) = -\frac{1}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (x - \mu)^2 - \frac{1}{2} \log(2\pi)$$

The gradient with respect to  $\mu$  and  $\sigma^2$  (!) is the row vector

$$\nabla \log f(x|\mu, \sigma^2) = \left( -\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4} (x - \mu)^2 \right)^T$$

Hint for calculating the gradient: replace  $\sigma^2$  by  $v$  and then take the partial derivative with regard to  $v$ , then substitute back.

The Hessian matrix is

$$\nabla^T \nabla \log f(x|\mu, \sigma^2) = \begin{pmatrix} -\frac{1}{\sigma^2} & -\frac{1}{\sigma^4} (x - \mu) \\ -\frac{1}{\sigma^4} (x - \mu) & \frac{1}{2\sigma^4} - \frac{1}{\sigma^6} (x - \mu)^2 \end{pmatrix}$$

As  $E(x) = \mu$  we have  $E(x - \mu) = 0$ . Furthermore, with  $E((x - \mu)^2) = \sigma^2$  we see that  $E\left(\frac{1}{\sigma^6} (x - \mu)^2\right) = \frac{1}{\sigma^4}$ . Therefore the expected Fisher information matrix is the negative expected Hessian matrix is

$$I^{\text{Fisher}}(\mu, \sigma^2) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix}$$

## 2.4 Entropy learning and maximum likelihood

### 2.4.1 The relative entropy between true model and approximating model

Assume we have observations  $x_1, \dots, x_n$ . The data is sampled from  $F$ , the true but unknown data generating distribution. We also specify models  $G_\theta$  indexed by  $\theta$  to approximate  $F$ .

The relative entropy  $D_{\text{KL}}(F, G_\theta)$  then measures the divergence of the approximation  $G_\theta$  from the unknown true model  $F$ . It can be written as:

$$\begin{aligned} D_{\text{KL}}(F, G_\theta) &= H(F, G_\theta) - H(F) \\ &= \underbrace{-\mathbb{E}_F \log g_\theta(x)}_{\text{cross-entropy}} - \underbrace{(\mathbb{E}_F \log f(x))}_{\text{entropy of } F, \text{ does not depend on } \theta} \end{aligned}$$

However, since we do not know  $F$  we cannot actually compute this divergence. Nonetheless, we may use the empirical distribution  $\hat{F}_n$  — a function of the observed data — as approximation for  $F$ , and in this way arrive at an approximation for  $D_{\text{KL}}(F, G_\theta)$  that becomes more and more accurate with growing sample size.

---

Recall the “Law of Large Numbers” :

- By the strong law of large numbers the empirical distribution  $\hat{F}_n$  converges to the true underlying distribution  $F$  as  $n \rightarrow \infty$  almost surely:

$$\hat{F}_n \xrightarrow{\text{a.s.}} F$$

- For  $n \rightarrow \infty$  the average  $\mathbb{E}_{\hat{F}_n}(h(X)) = \frac{1}{n} \sum_{i=1}^n h(x_i)$  converges to the expectation  $\mathbb{E}_F(h(X))$ .

---

Hence, for large sample size  $n$  we can approximate cross-entropy and as a result the KL divergence. The cross-entropy  $H(F, G_\theta)$  is approximated by the empirical cross-entropy where the expectation is taken with regard to  $\hat{F}_n$  rather than  $F$ :

$$\begin{aligned} H(F, G_\theta) &\approx H(\hat{F}_n, G_\theta) \\ &= -\mathbb{E}_{\hat{F}_n}(\log(x)) \\ &= -\frac{1}{n} \sum_{i=1}^n \log g(x_i | \theta) \\ &= -\frac{1}{n} l_n(\theta) \end{aligned}$$

This turns out to be equal to the negative log-likelihood standardised by the sample size  $n!$  Or in other words, the **likelihood** is the **negative empirical cross-entropy (times sample size  $n$ )**.

From the link of the multinomial coefficient with Shannon entropy (Example 2.3) we already know that for large sample size

$$H(\hat{F}) \approx \frac{1}{n} \log \binom{n}{n_1, \dots, n_K}$$

The KL divergence  $D_{\text{KL}}(F, G_{\theta})$  can therefore be approximated by

$$D_{\text{KL}}(F, G_{\theta}) \approx -\frac{1}{n} \left( \log \binom{n}{n_1, \dots, n_K} + l_n(\theta) \right)$$

Thus, with the KL divergence we obtain not just the log-likelihood (the cross-entropy part) but also the multiplicity factor taking account of the possible orderings of the data (the entropy part).

## 2.4.2 Minimum KL divergence and maximum likelihood

If we were to know  $F$  we would simply minimise  $D_{\text{KL}}(F, G_{\theta})$  to find the particular model  $G_{\theta}$  that is closest to the true model. Equivalently, we would minimise the cross-entropy  $H(F, G_{\theta})$ . However, since we actually don't know  $F$  this is not possible.

However, for large sample size  $n$  when the empirical distribution  $\hat{F}_n$  is a good approximation for  $F$ , we can use the results from the previous section. Thus, instead of minimising the KL divergence  $D_{\text{KL}}(F, G_{\theta})$  we simply minimise  $H(\hat{F}_n, G_{\theta})$  which is the same as maximising the likelihood  $l_n(\theta)$ . Note that the entropy of the true distribution  $F$  (and the corresponding empirical distribution  $\hat{F}$ ) that does not depend on the parameters  $\theta$  and hence it does not matter when minimising the divergence.

Conversely, this implies that maximising the likelihood with regard to the  $\theta$  is equivalent (asymptotically for large  $n$ !) to minimising the KL divergence of the approximating model and the unknown true model!

$$\begin{aligned} \hat{\theta}^{\text{ML}} &= \arg \max_{\theta} l_n(\theta) \\ &= \arg \min_{\theta} H(\hat{F}_n, G_{\theta}) \\ &\approx \arg \min_{\theta} D_{\text{KL}}(F, G_{\theta}) \end{aligned}$$

Therefore, the reasoning behind the method of **maximum likelihood** is that it minimises a **large sample approximation of the KL divergence** of the candidate model  $G_{\theta}$  from the unknown true model  $F$ .

As a consequence of the close link of maximum likelihood and relative entropy maximum likelihood inherits for large  $n$  (and only then!) all the optimality properties from KL divergence. These will be discussed in more detail later in the course.

## 2.4.3 Further connections

Since minimising KL divergence contains ML estimation as special case you may wonder whether there is a broader justification of relative entropy in the

context of statistical data analysis?

Indeed, KL divergence has strong geometrical interpretation that forms the basis of *information geometry*. In this field the manifold of distributions is studied using tools from differential geometry. The expected Fisher information plays an important role as [metric tensor in the space of distributions](#).

Furthermore, it is also linked to probabilistic forecasting. In the framework of so-called [scoring rules](#), the only local proper scoring rule is the negative log-probability (“surprise”). The expected “surprise” is the cross-entropy and relative entropy is the corresponding natural divergence connected with the log scoring rule.

Furthermore, another intriguing property of KL divergence is that the relative entropy  $D_{\text{KL}}(F, G)$  is the *only divergence measure* that is both a Bregman and an  $f$ -divergence. Note that [f-divergences](#) and [Bregman-divergences](#) (in turn related to proper scoring rules) are two large classes of measures of similarity and divergence between two probability distributions.

Finally, not only the likelihood estimation but also the Bayesian update rule (as discussed later in this module) is another special case of entropy learning.





## Chapter 3

# Maximum likelihood estimation

### 3.1 Principle of maximum likelihood estimation

#### 3.1.1 Outline

The starting points in an ML analysis are

- the observed  $n$  data samples  $x_1, \dots, x_n$ , iid (=independent and identically distributed), with the ordering irrelevant, and a
- model  $F_\theta$  with corresponding probability density or probability mass function  $f(x|\theta)$  with parameters  $\theta$

From this we construct the likelihood function:

- $L_n(\theta|x_1, \dots, x_n) = \prod_{i=1}^n f(x_i|\theta)$

Historically, the likelihood is also often interpreted as the probability of the data given the model. However, this is not strictly correct. First this interpretation only applies to discrete random variables. Second, since the samples are iid even in this case one would still need to add a factor accounting for the multiplicity of possible orderings of the samples to obtain the correct probability of the data. Third, the interpretation of likelihood as probability of the data completely breaks down for continuous random variables because then  $f(x)$  is a density, not a probability.

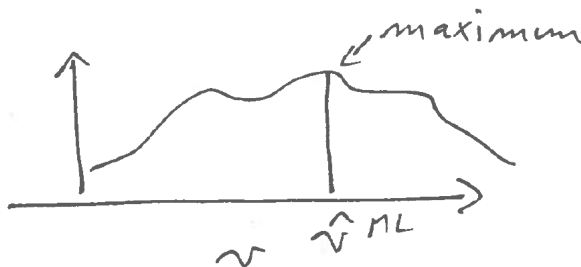
As we have seen in the previous chapter the origin of the likelihood function lies in its connection to relative entropy. Specifically, the log-likelihood function

- $l_n(\theta|x_1, \dots, x_n) = \sum_{i=1}^n \log f(x_i|\theta)$

divided by sample size  $n$  is a large sample approximation of the cross-entropy between the unknown true data generating model and the approximating model  $F_\theta$ . Note that the log-likelihood is additive over the samples  $x_i$ .

The maximum likelihood point estimate  $\hat{\theta}^{ML}$  is then given by maximising the (log)-likelihood

$$\hat{\theta}^{ML} = \arg \max_{\theta} l_n(\theta | x_1, \dots, x_n)$$



### 3.1.2 Obtaining MLEs for a regular model

In regular situations, i.e. when

- the log-likelihood function is smooth and twice differentiable,
- the second derivative is negative and not zero, and for more than one parameter the Hessian matrix is negative definite and not singular,
- the parameters of the model are all identifiable (in particular the model is not overparameterised), and
- the true parameter values lie inside the support and not on the border,

then in order to maximise  $l_n$  one may use the **score function**  $S(\theta)$  which is the first order derivative of the log-likelihood function:

$$S_n(\theta) = \frac{dl_n(\theta | x_1, \dots, x_n)}{d\theta} \quad \text{scalar parameter: first derivative of log-likelihood function}$$

$$S_n(\theta) = \nabla l_n(\theta | x_1, \dots, x_n) \quad \text{gradient if } \theta \text{ is a vector (i.e. if there's more than one parameter)}$$

A necessary (but not sufficient) condition for the MLE is that

$$S_n(\hat{\theta}_{ML}) = 0$$

To demonstrate that the log-likelihood function actually achieves a maximum at  $\hat{\theta}_{ML}$  the curvature at the MLE must be negative, i.e. that the log-likelihood must be locally concave at the MLE.

In the case of a single parameter (scalar  $\theta$ ) this requires to check that the second derivative of the log-likelihood function is negative:

$$\frac{d^2 l_n(\hat{\theta}_{ML})}{d\theta^2} < 0$$

In the case of a parameter vector (multivariate  $\theta$ ) you need to compute the Hessian matrix (matrix of second order derivatives) at the MLE:

$$\nabla^T \nabla l_n(\hat{\theta}_{ML})$$

and then verify that this matrix is negative definite (i.e. all its eigenvalues must be negative).

As we will see later the second order derivatives of the log-likelihood function also play an important role for assessing the uncertainty of the MLE.

### 3.1.3 Invariance property of the maximum likelihood

Maximisation is a procedure that is invariant against coordinate transformations of the argument. Suppose  $x_{\max} = \arg \max h(x)$  and  $y = g(x)$  where  $g$  is an invertible function. Then  $y_{\max} = \arg \max h(g^{-1}(y)) = g(x_{\max})$ . The achieved maximum itself remains invariant:  $h(x_{\max}) = h(g^{-1}(y_{\max}))$ .

With regard to maximum likelihood estimation this implies the following **invariance property** of the maximum likelihood:

- Suppose that  $\hat{\theta}_{ML}$  is the MLE of  $\theta$ .
- We transform the parameter to  $\theta^* = g(\theta)$  where  $g$  is an invertible function.
- Then  $g(\hat{\theta}_{ML}) = \hat{\theta}^*$  is the MLE of  $\theta^*$ .
- The value of the achieved maximum likelihood is the same in both cases, i.e. it is invariant against transformation of the parameters.

The invariance property can be very useful in practise because it may be easier to perform the maximisation required for finding the MLE in a particular coordinate system.

See Worksheet 2 for an example application of the invariance principle.

### 3.1.4 Consistency of maximum likelihood estimates

One important property of maximum likelihood is that it produces **consistent estimates**.

Specifically, if the true underlying model  $F_{\text{true}}$  with parameter  $\theta_{\text{true}}$  is contained in the set of specified candidate models  $F_{\theta}$

$$\underbrace{F_{\text{true}}}_{\text{true model}} \subset \underbrace{F_{\theta}}_{\text{specified models}}$$

then

$$\hat{\theta}_{ML} \xrightarrow{\text{large } n} \theta_{\text{true}}$$

This is a consequence of  $D_{\text{KL}}(F_{\text{true}}, F_{\theta}) \rightarrow 0$  for  $F_{\theta} \rightarrow F_{\text{true}}$ , and that maximisation of the likelihood function is for large  $n$  equivalent to minimising the relative entropy.

Thus given sufficient data the MLE will converge to the true value. As a consequence, **MLEs are asymptotically unbiased**. As we will see in the examples they can still be biased in finite samples.

Note that even if the candidate model  $F_{\theta}$  is misspecified (i.e. it does not contain the actual true model) the MLE is still optimal in the sense in that it will find the closest possible model.

It is possible to find inconsistent MLEs, but this occurs only in situations where the dimension of the model / number of parameters increases with sample size, or when the MLE is at a boundary or when there are singularities in the likelihood function.

## 3.2 Maximum likelihood estimation in practise

### 3.2.1 Worked examples

In this section we now provide a number of worked example how ML estimation works in practise.

**Example 3.1.** Estimation of a proportion:

We aim to estimate the true proportion  $p$  in a Bernoulli experiment with binary outcomes, say the proportion of “successes” vs. “failures” or of “heads” vs. “tails” in a coin tossing experiment.

- Bernoulli model  $\text{Ber}(p)$ :  $\Pr(\text{“success”}) = p$  and  $\Pr(\text{“failure”}) = 1 - p$ .
- The “success” is indicated by outcome  $x = 1$  and the “failure” by  $x = 0$ .
- We conduct  $n$  trials and record  $n_1$  successes and  $n - n_1$  failures.
- Parameter:  $p$ : probability of “success”.

What is the MLE of  $p$ ?

- the data  $x_1, \dots, x_n$  take on values 0 or 1.
- the average of the data points is  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{n_1}{n}$ .
- the probability mass function (PMF) of the Bernoulli distribution  $\text{Ber}(p)$  is:

$$f(x) = p^x(1-p)^{1-x} = \begin{cases} p & \text{if } x = 1 \\ 1-p & \text{if } x = 0 \end{cases}$$

- log-PMF:

$$\log f(x) = x \log(p) + (1 - x) \log(1 - p)$$

- log-likelihood function:

$$\begin{aligned} l_n(p) &= \sum_{i=1}^n \log f(x_i) \\ &= n_1 \log p + (n - n_1) \log(1 - p) \\ &= n (\bar{x} \log p + (1 - \bar{x}) \log(1 - p)) \end{aligned}$$

Note how the log-likelihood depends on the data only through  $\bar{x}$ ! This is an example of a *sufficient statistic* for the parameter  $p$  (in fact it is also a *minimally sufficient statistic*). This will be discussed in more detail later.

- Score function:

$$S_n(p) = \frac{dl_n(p)}{dp} = n \left( \frac{\bar{x}}{p} - \frac{1 - \bar{x}}{1 - p} \right)$$

- Maximum likelihood estimate: Setting  $S_n(\hat{p}_{ML}) = 0$  yields as solution

$$\hat{p}_{ML} = \bar{x} = \frac{n_1}{n}$$

With  $\frac{dS_n(p)}{dp} = -n \left( \frac{\bar{x}}{p^2} + \frac{1 - \bar{x}}{(1 - p)^2} \right) < 0$  the optimum corresponds indeed to the maximum of the (log-)likelihood function as this is negative for  $\hat{p}_{ML}$  (and indeed for any  $p$ ).

The maximum likelihood estimator of  $p$  is therefore identical to the frequency of the successes among all observations.

Note that to analyse the coin tossing experiment and to estimate  $p$  we may equally well use the Binomial distribution  $\text{Bin}(n, p)$  as model for the number of successes. In this case we then have only a single observation, namely the observed  $k$ . This results in the same MLE for  $p$  but the likelihood function based on the Binomial PMF includes the Binomial coefficient  $\binom{n}{k}$ . However, as this factor does not depend on  $p$  it disappears in the score function and has no influence in the derivation of the MLE.

**Example 3.2.** Normal distribution with unknown mean and known variance:

- $x \sim N(\mu, \sigma^2)$  with  $E(x) = \mu$  and  $\text{Var}(x) = \sigma^2$
- the parameter to be estimated is  $\mu$  whereas  $\sigma^2$  is known.

What's the MLE of parameter  $\mu$ ?

- the data  $x_1, \dots, x_n \in [-\infty, \infty]$  are real values.
- the average  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  is real as well.

- Density:

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

- Log-Density:

$$\log f(x) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(x-\mu)^2}{2\sigma^2}$$

- Log-likelihood function:

$$\begin{aligned} l_n(\mu) &= \sum_{i=1}^n \log f(x_i) \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \quad \underbrace{-\frac{n}{2} \log(2\pi\sigma^2)}_{\text{constant term, does not depend on } \mu, \text{ can be removed}} \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i^2 - 2x_i\mu + \mu^2) + C \\ &= \frac{n}{\sigma^2} (\bar{x}\mu - \frac{1}{2}\mu^2) \quad \underbrace{-\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2}_{\text{another constant term}} + C \end{aligned}$$

Note how the non-constant terms of the log-likelihood depend on the data only through  $\bar{x}$ !

- Score function:

$$S_n(\mu) = \frac{n}{\sigma^2} (\bar{x} - \mu)$$

- Maximum likelihood estimate:

$$S_n(\hat{\mu}_{ML}) = 0 \Rightarrow \hat{\mu}_{ML} = \bar{x}$$

- With  $\frac{dS_n(\mu)}{d\mu} = -\frac{n}{\sigma^2} < 0$  the optimum is indeed the maximum

The constant term  $C$  in the log-likelihood function collects all terms that do not depend on the parameter. After taking the first derivative with regard to the parameter this term disappears thus  **$C$  is not relevant for finding the MLE of the parameter. In the future we will often omit such constant terms from the log-likelihood function without further mention.**

**Example 3.3.** Normal distribution with mean and variance both unknown:

- $x \sim N(\mu, \sigma^2)$  with  $E(x) = \mu$  and  $\text{Var}(x) = \sigma^2$
- both  $\mu$  and  $\sigma^2$  need to be estimated.

What's the MLE of the parameter vector  $\theta = (\mu, \sigma^2)^T$ ?

- the data  $x_1, \dots, x_n \in [-\infty, \infty]$  are real values.
- the average  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  is real as well.
- the average of the squared data  $\overline{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2 \geq 0$  is non-negative.
- Density:

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

- Log-Density:

$$\log f(x) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(x-\mu)^2}{2\sigma^2}$$

- Log-likelihood function:

$$\begin{aligned} l_n(\theta) &= \sum_{i=1}^n \log f(x_i) \\ &= -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \quad \underbrace{-\frac{n}{2} \log(2\pi)}_{\text{constant not depending on } \mu \text{ or } \sigma^2} \\ &= -\frac{n}{2} \log(\sigma^2) - \frac{n}{2\sigma^2} (\overline{x^2} - 2\bar{x}\mu + \mu^2) + C \end{aligned}$$

Note how the log-likelihood function depends on the data only through  $\bar{x}$  and  $\overline{x^2}$ !

- Score function  $S$  (row vector!), gradient of  $l_n(\theta)$ :

$$\begin{aligned} S(\theta) &= \nabla l_n(\theta) \\ &= \left( -\frac{n}{2\sigma^2} + \frac{n}{2\sigma^4} \left( \overline{x^2} - 2\bar{x}\mu + \mu^2 \right) \right)^T \end{aligned}$$

Note that to obtain the second component of the score function the partial derivative needs to be taken with regard to the variance parameter  $\sigma^2$  — not with regard to  $\sigma$ ! Hint: replace  $\sigma^2 = v$  in the log-likelihood function, then take the partial derivative with regard to  $v$ , then backsubstitute  $v = \sigma^2$  in the result.

- Maximum likelihood estimate:

$$\begin{aligned} S(\hat{\theta}_{ML}) &= 0 \Rightarrow \\ \hat{\theta}_{ML} &= \begin{pmatrix} \hat{\mu}_{ML} \\ \hat{\sigma}_{ML}^2 \end{pmatrix} = \begin{pmatrix} \bar{x} \\ \overline{x^2} - \bar{x}^2 \end{pmatrix} \end{aligned}$$

The ML estimate of the variance we can also write  $\hat{\sigma}_{ML}^2 = \overline{x^2} - \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ .

- To confirm that we actually have maximum we need to verify that the eigenvalues of the Hessian matrix are all negative. This is indeed the case, for details see Example 3.6.

### 3.2.2 Relationship with least squares estimation

In Example 3.2 the form of the log-likelihood function is a function of the sum of squared differences. Maximising  $l_n(\mu) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$  is equivalent to *minimising*  $\sum_{i=1}^n (x_i - \mu)^2$ . Hence, finding the mean by **maximum likelihood assuming a normal model is equivalent to least-squares estimation!**

Note that least-squares estimation has been in use at least since the early 1800s and thus predates maximum likelihood (1924). Due to its simplicity it is still very popular in particular in regression and the link with maximum likelihood and normality allows to understand why it usually works well!

### 3.2.3 Bias and maximum likelihood

Example 3.3 is interesting because it shows that maximum likelihood can result in both biased and as well as unbiased estimators.

Recall that  $x \sim N(\mu, \sigma^2)$ . As a result

$$\hat{\mu}_{ML} = \bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

with  $E(\hat{\mu}_{ML}) = \mu$  and

$$\hat{\sigma}_{ML}^2 \sim \frac{\sigma^2}{n} \chi_{n-1}^2$$

with  $E(\hat{\sigma}_{ML}^2) = \frac{n-1}{n} \sigma^2$ .

Therefore, the MLE of  $\mu$  is unbiased as

$$\text{Bias}(\hat{\mu}_{ML}) = E(\hat{\mu}_{ML}) - \mu = 0$$

In contrast, however, the MLE of  $\sigma^2$  is negatively biased because

$$\text{Bias}(\hat{\sigma}_{ML}^2) = E(\hat{\sigma}_{ML}^2) - \sigma^2 = -\frac{1}{n} \sigma^2$$

Thus, in the case of the variance parameter of the normal distribution the MLE is *not* recovering the well-known unbiased estimator of the variance

$$\hat{\sigma}_{UB}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n}{n-1} \hat{\sigma}_{ML}^2$$

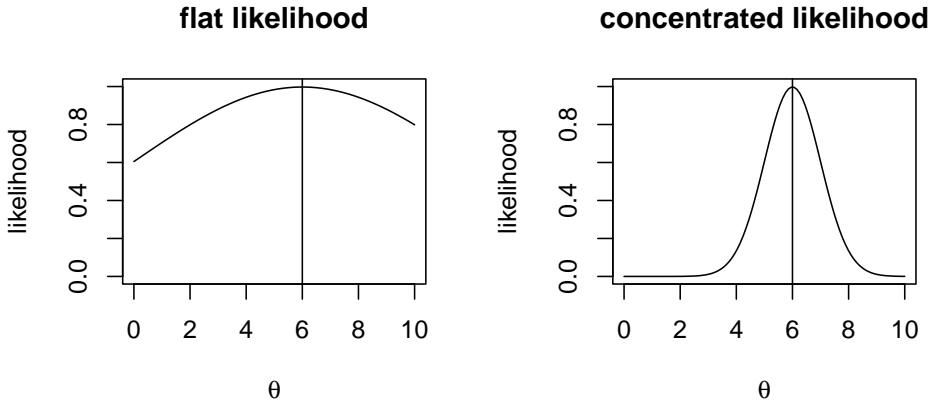
Conversely, the unbiased estimator is not a maximum likelihood estimate!



Therefore it is worth keeping in mind that maximum likelihood can result in biased estimates for finite  $n$ . For large  $n$ , however, the bias disappears as MLEs are consistent.

### 3.3 Observed Fisher information

#### 3.3.1 Motivation and definition



By inspection of some log-likelihood curves it is apparent that the log-likelihood function contains more information about the parameter  $\theta$  than just the maximum point  $\hat{\theta}_{ML}$ .

In particular the **curvature** of the log-likelihood function at the MLE must be somehow related the accuracy of  $\hat{\theta}_{ML}$ : if the likelihood surface is flat near the maximum (low curvature) then it is more difficult to find the optimal parameter (also numerically!). Conversely, if the likelihood surface is peaked (strong curvature) then the maximum point is clearly defined.

The curvature is described by the second-order derivatives (Hessian matrix) of the log-likelihood function.

For univariate  $\theta$  the Hessian is a scalar:

$$\frac{d^2 l_n(\theta)}{d\theta^2}$$

For multivariate parameter vector  $\theta$  of dimension  $d$  the Hessian is a matrix of size  $d \times d$ :

$$\nabla^T \nabla l_n(\theta)$$

By construction the Hessian is negative definite at the MLE (i.e. its eigenvalues are all negative) to ensure the the function is concave at the MLE (i.e. peak shaped).

The **observed Fisher information** (matrix) is defined as the negative curvature at the MLE  $\hat{\theta}_{ML}$ :

$$J_n(\hat{\theta}_{ML}) = -\nabla^T \nabla l_n(\hat{\theta}_{ML})$$

Sometimes this is simply called the “observed information”. To avoid confusion with the expected Fisher information introduced earlier

$$I^{\text{Fisher}}(\theta) = -E_{F_\theta} \left( \nabla^T \nabla \log f(x|\theta) \right)$$

it is necessary to always use the qualifier “observed” when referring to  $J_n(\hat{\theta}_{ML})$ .

### 3.3.2 Examples of observed Fisher information

**Example 3.4.** Bernoulli model  $\text{Ber}(p)$ :

We continue Example 3.1. Recall that  $\hat{p}_{ML} = \bar{x} = \frac{n_1}{n}$  and the score function  $S_n(p) = n \left( \frac{\bar{x}}{p} - \frac{1-\bar{x}}{1-p} \right)$ . The negative second derivative of the log-likelihood function is

$$-\frac{dS_n(p)}{dp} = n \left( \frac{\bar{x}}{p^2} + \frac{1-\bar{x}}{(1-p)^2} \right)$$

The observed Fisher information is therefore

$$\begin{aligned} J_n(\hat{p}_{ML}) &= n \left( \frac{\bar{x}}{\hat{p}_{ML}^2} + \frac{1-\bar{x}}{(1-\hat{p}_{ML})^2} \right) \\ &= n \left( \frac{1}{\hat{p}_{ML}} + \frac{1}{1-\hat{p}_{ML}} \right) \\ &= \frac{n}{\hat{p}_{ML}(1-\hat{p}_{ML})} \end{aligned}$$

The inverse of the observed Fisher information is:

$$J_n(\hat{p}_{ML})^{-1} = \frac{\hat{p}_{ML}(1-\hat{p}_{ML})}{n}$$

Compare this with  $\text{Var} \left( \frac{x}{n} \right) = \frac{p(1-p)}{n}$  for  $x \sim \text{Bin}(n, p)$ .

**Example 3.5.** Normal distribution with unknown mean and known variance:

This is the continuation of Example 3.2. Recall the MLE for the mean  $\hat{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$  and the score function  $S_n(\mu) = \frac{n}{\sigma^2} (\bar{x} - \mu)$ . The negative second derivative of the score function is

$$-\frac{dS_n(\mu)}{d\mu} = \frac{n}{\sigma^2}$$

The observed Fisher information at the MLE is therefore

$$J_n(\hat{\mu}_{ML}) = \frac{n}{\sigma^2}$$

and the inverse of the observed Fisher information is

$$J_n(\hat{\nu}_{ML})^{-1} = \frac{\sigma^2}{n}$$

For  $x_i \sim N(\mu, \sigma^2)$  we have  $\text{Var}(x_i) = \sigma^2$  and hence  $\text{Var}(\bar{x}) = \frac{\sigma^2}{n}$ , which is equal to the inverse observed Fisher information.

**Example 3.6.** Normal distribution with mean and variance parameter:

This is the continuation of Example 3.3. Recall the MLE for the mean and variance:

$$\begin{aligned}\hat{\mu}_{ML} &= \frac{1}{n} \sum_{i=1}^n x_i = \bar{x} \\ \hat{\sigma}_{ML}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \overline{x^2} - \bar{x}^2\end{aligned}$$

with score function

$$S_n(\mu, \sigma^2) = \nabla l_n(\mu, \sigma^2) = \begin{pmatrix} \frac{n}{\sigma^2}(\bar{x} - \mu) \\ -\frac{n}{2\sigma^2} + \frac{n}{2\sigma^4}(\overline{x^2} - 2\mu\bar{x} + \mu^2) \end{pmatrix}^T$$

The Hessian matrix of the log-likelihood function is

$$\nabla^T \nabla l_n(\mu, \sigma^2) = \begin{pmatrix} -\frac{n}{\sigma^2} & -\frac{n}{\sigma^4}(\bar{x} - \mu) \\ -\frac{n}{\sigma^4}(\bar{x} - \mu) & \frac{n}{2\sigma^4} - \frac{n}{\sigma^6}(\overline{x^2} - 2\mu\bar{x} + \mu^2) \end{pmatrix}$$

The negative Hessian at the MLE, i.e. at  $\hat{\mu}_{ML} = \bar{x}$  and  $\hat{\sigma}_{ML}^2 = \overline{x^2} - \bar{x}^2$  yields the **observed Fisher information matrix**:

$$J_n(\hat{\mu}_{ML}, \hat{\sigma}_{ML}^2) = \begin{pmatrix} \frac{n}{\hat{\sigma}_{ML}^2} & 0 \\ 0 & \frac{n}{2(\hat{\sigma}_{ML}^2)^2} \end{pmatrix}$$

Note that the observed Fisher information matrix is diagonal with positive entries. Therefore its eigenvalues are all positive as required for a maximum, because for a diagonal matrix the eigenvalues are simply the entries on the diagonal.

The inverse of the observed Fisher information matrix is

$$J_n(\hat{\mu}_{ML}, \hat{\sigma}_{ML}^2)^{-1} = \begin{pmatrix} \frac{\hat{\sigma}_{ML}^2}{n} & 0 \\ 0 & \frac{2(\hat{\sigma}_{ML}^2)^2}{n} \end{pmatrix}$$

Recall that  $x \sim N(\mu, \sigma^2)$  and therefore

$$\hat{\mu}_{ML} = \bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Hence  $\text{Var}(\hat{\mu}_{ML}) = \frac{\sigma^2}{n}$ . If you compare this with the first diagonal entry of the inverse observed Fisher information matrix you see that this is essentially the same expression (apart from the “hat”).

The empirical variance  $\hat{\sigma}_{ML}^2$  follows a scaled chi-squared distribution

$$\hat{\sigma}_{ML}^2 \sim \frac{\sigma^2}{n} \chi_{n-1}^2$$

with variance  $\text{Var}(\hat{\sigma}_{ML}^2) = \frac{n-1}{n} \frac{2\sigma^4}{n}$ . For large  $n$  this becomes  $\text{Var}(\hat{\sigma}_{ML}^2) \stackrel{a}{=} \frac{2\sigma^4}{n}$  which is essentially (apart from the “hat”) the second diagonal entry of the inverse observed Fisher information matrix.

### 3.3.3 Relationship between observed and expected Fisher information

The observed Fisher information  $J_n(\hat{\theta}_{ML})$  and the expected Fisher information  $I^{\text{Fisher}}(\theta)$  are related but also two clearly different entities:

- Both types of Fisher information are based on computing the second order derivative (Hessian matrix), thus are based on the curvature of a function.
- The observed Fisher information is computed from the log-likelihood function. Therefore it takes the observed data into account. It explicitly depends on the sample size  $n$ . It contains estimates of the parameters but not the parameters themselves. While the curvature of the log-likelihood function may be computed for any point the the observed Fisher information specifically refers to the MLE  $\hat{\theta}_{ML}$ . It is linked to the (asymptotic) variance of the MLE as we have seen in the examples and will discuss in more detail later.
- In contrast, the expected Fisher information is derived directly from the log-density. It does not depend on the observed data, and thus does not have dependency on sample size. It can be computed for any value of the parameters. It describes the geometry of the space of the models, and is the local approximation of relative entropy.
- Asymptotically, for large sample size  $n$  the MLE converges to  $\hat{\theta}_{ML} \rightarrow \theta_0$ . It follows from the construction of the observed Fisher information and the law of large numbers that correspondingly  $J_n(\hat{\theta}_{ML}) \rightarrow nI^{\text{Fisher}}(\theta_0)$ .
- In a very important class of models, namely **in the exponential family**, we find that  $J_n(\hat{\theta}_{ML}) = nI^{\text{Fisher}}(\hat{\theta}_{ML})$  also for finite sample size  $n$ . This is in fact the case in all the examples discussed above (e.g. see Examples 2.11 and 3.4 for the Bernoulli and Examples 2.13 and 3.6 for the normal distribution).
- However, this is an exception. In a general model  $J_n(\hat{\theta}_{ML}) \neq nI^{\text{Fisher}}(\hat{\theta}_{ML})$  for finite sample size  $n$ . An example is provided by the Cauchy distribution

with median parameter  $\theta$ . It is not part of the exponential family and has expected Fisher information  $I^{\text{Fisher}}(\theta) = \frac{1}{2}$  regardless of the choice of the median parameter whereas the observed Fisher information  $J_n(\hat{\theta}_{ML})$  depends on the MLE  $\hat{\theta}_{ML}$  of the median parameter and is not simply  $\frac{n}{2}$ .



## Chapter 4

# Quadratic approximation and normal asymptotics

### 4.1 Multivariate statistics for random vectors

#### 4.1.1 Covariance and correlation

Assume a scalar random variable  $x$  with mean  $E(x) = \mu$ . The corresponding variance is given by

$$\begin{aligned}\text{Var}(x) &= E\left((x - \mu)^2\right) \\ &= E\left((x - \mu)(x - \mu)\right) \\ &= E(x^2) - \mu^2\end{aligned}$$

For a random vector  $\mathbf{x} = (x_1, x_2, \dots, x_d)^T$  the mean  $E(\mathbf{x}) = \boldsymbol{\mu}$  is simply comprised of the means of its components, i.e.  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)^T$ . Thus, the mean of a random vector of dimension  $d$  is a vector of the same length.

The variance of a random vector of length  $d$ , however, is not a vector but a matrix of size  $d \times d$ . This matrix is called the **covariance matrix**:

$$\begin{aligned}\text{Var}(\mathbf{x}) &= \underbrace{\boldsymbol{\Sigma}}_{d \times d} = (\sigma_{ij}) = \begin{pmatrix} \sigma_{11} & \dots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \dots & \sigma_{dd} \end{pmatrix} \\ &= E\left(\underbrace{(\mathbf{x} - \boldsymbol{\mu})}_{d \times 1} \underbrace{(\mathbf{x} - \boldsymbol{\mu})^T}_{1 \times d}\right) \\ &= E(\mathbf{x}\mathbf{x}^T) - \boldsymbol{\mu}\boldsymbol{\mu}^T\end{aligned}$$

The entries of the covariance matrix  $\sigma_{ij} = \text{Cov}(x_i, x_j)$  describe the covariance between the random variables  $x_i$  and  $x_j$ . The covariance matrix is symmetric, hence  $\sigma_{ij} = \sigma_{ji}$ . The diagonal entries  $\sigma_{ii} = \text{Cov}(x_i, x_i) = \text{Var}(x_i) = \sigma_i^2$  correspond to the variances of the components of  $\mathbf{x}$ . The covariance matrix is **positive semi-definite**, i.e. the eigenvalues of  $\Sigma$  are all positive or equal to zero. However, in practise one aims to use non-singular covariance matrices, with all eigenvalues positive, so that they are invertible.

A covariance matrix can be factorised into the product

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

where  $V$  is a diagonal matrix containing the variances

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

and the matrix  $P$  ("capital rho") is the symmetric **correlation matrix**

$$P = (\rho_{ij}) = \begin{pmatrix} 1 & \dots & \rho_{1d} \\ \vdots & \ddots & \vdots \\ \rho_{d1} & \dots & 1 \end{pmatrix} = V^{-\frac{1}{2}} \Sigma V^{-\frac{1}{2}}$$

Thus, the correlation between  $x_i$  and  $x_j$  is defined as

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

For univariate  $x$  and scalar constant  $a$  the variance of  $ax$  equals  $\text{Var}(ax) = a^2 \text{Var}(x)$ . For a random vector  $\mathbf{x}$  of dimension  $d$  and matrix  $A$  of dimension  $m \times d$  this generalises to  $\text{Var}(A\mathbf{x}) = A \text{Var}(\mathbf{x}) A^T$ .

### 4.1.2 Multivariate normal distribution

The density of a normally distributed scalar variable  $x \sim N(\mu, \sigma^2)$  with mean  $E(x) = \mu$  and variance  $\text{Var}(x) = \sigma^2$  is

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

The univariate normal distribution for a scalar  $x$  generalises to the **multivariate normal distribution** for a vector  $\mathbf{x} = (x_1, x_2, \dots, x_d)^T \sim N_d(\boldsymbol{\mu}, \Sigma)$  with with mean



$E(\mathbf{x}) = \boldsymbol{\mu}$  and covariance matrix  $\text{Var}(\mathbf{x}) = \boldsymbol{\Sigma}$ . The corresponding density is

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

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

For  $d = 1$  we have  $\mathbf{x} = x$ ,  $\boldsymbol{\mu} = \mu$  and  $\boldsymbol{\Sigma} = \sigma^2$  so that the multivariate normal density reduces to the univariate normal density.

**Example 4.1.** Maximum likelihood estimates of the parameters of the multivariate normal distribution:

Maximising the log-likelihood based on the multivariate normal density yields the MLEs for  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . These are generalisations of the MLEs for the mean  $\mu$  and variance  $\sigma^2$  of the univariate normal as encountered in Example 3.3.

The estimates can be written in three different ways:

**a) data vector notation**

with  $\mathbf{x}_1, \dots, \mathbf{x}_n$  the  $n$  vector-valued observations from the multivariate normal:

MLE for the mean:

$$\hat{\boldsymbol{\mu}}_{ML} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k = \bar{\mathbf{x}}$$

MLE for the covariance:

$$\hat{\boldsymbol{\Sigma}}_{ML} = \frac{1}{n} \sum_{k=1}^n \underbrace{(\mathbf{x}_k - \bar{\mathbf{x}})}_{d \times 1} \underbrace{(\mathbf{x}_k - \bar{\mathbf{x}})^T}_{1 \times d}$$

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

With  $\overline{\mathbf{x}\mathbf{x}^T} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k \mathbf{x}_k^T$  we can also write

$$\hat{\boldsymbol{\Sigma}}_{ML} = \overline{\mathbf{x}\mathbf{x}^T} - \bar{\mathbf{x}}\bar{\mathbf{x}}^T$$

**b) data component notation**

with  $x_{ki}$  the  $i$ -th component of the  $k$ -th sample  $\mathbf{x}_k$ :

$$\hat{\mu}_i = \frac{1}{n} \sum_{k=1}^n x_{ki} \text{ with } \hat{\boldsymbol{\mu}} = \begin{pmatrix} \hat{\mu}_1 \\ \vdots \\ \hat{\mu}_d \end{pmatrix}$$

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

### c) data matrix notation

with  $X = \begin{pmatrix} x_1^T \\ \dots \\ x_n^T \end{pmatrix}$  as a data matrix containing the samples in its rows. Note that this is the *statistics convention* — in much of the engineering and computer science literature the data matrix is often transposed and samples are stored in the columns. Thus, the formulas below are only correct assuming the statistics convention.

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

Here  $\mathbf{1}_n$  is a vector of length  $n$  containing 1 at each component.

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

To simplify the expression for the estimate of the covariance matrix one often assumes that the data matrix is centered, i.e. that  $\hat{\mu} = 0$ .

Because of the ambiguity in convention (machine learning vs statistics convention) and the often implicit use of centered data matrices the matrix notation is often confusing. Hence, using the other two notations is generally preferable.

## 4.2 Approximate distribution of maximum likelihood estimates

### 4.2.1 Quadratic log-likelihood resulting from normal model

Assume we observe a single sample  $x \sim N(\mu, \Sigma^2)$  with known covariance. The corresponding log-likelihood for  $\mu$  is

$$l_1(\mu) = C - \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)$$

where  $C$  is a constant that does not depend on  $\mu$ . Note that the log-likelihood is exactly quadratic and the maximum lies at  $(x, C)$ .

### 4.2.2 Quadratic approximation of a log-likelihood function

Now consider the quadratic approximation of the log-likelihood function  $l_n(\theta)$  for a general model around the MLE  $\hat{\theta}_{ML}$ .



We assume the model is regular with  $\nabla l_n(\hat{\theta}_{ML}) = 0$ . The Taylor series approximation of scalar-valued function  $f(x)$  around  $x_0$  is

$$f(x) = f(x_0) + \nabla f(x_0)(x - x_0) + \frac{1}{2}(x - x_0)^T \nabla^2 f(x_0)(x - x_0) + \dots$$

Applied to the log-likelihood function this yields

$$l_n(\theta) \approx l_n(\hat{\theta}_{ML}) - \frac{1}{2}(\hat{\theta}_{ML} - \theta)^T J_n(\hat{\theta}_{ML})(\hat{\theta}_{ML} - \theta)$$

This is a quadratic function with maximum at  $(\hat{\theta}_{ML}, l_n(\hat{\theta}_{ML}))$ . Note the natural appearance of the observed Fisher information  $J_n(\hat{\theta}_{ML})$  in the quadratic term. There is no linear term because of the vanishing gradient at the MLE.

Crucially, we realise that the approximation has the same form as if  $\hat{\theta}_{ML}$  was a sample from a multivariate normal distribution with mean  $\theta$  and with covariance given by the *inverse* observed Fisher information! Note that this requires a positive definite observed Fisher information matrix so that  $J_n(\hat{\theta}_{ML})$  is actually invertible!

**Example 4.2.** Quadratic approximation of the log-likelihood for a proportion:

From Example 3.1 we have the log-likelihood

$$l_n(p) = n(\bar{x} \log p + (1 - \bar{x}) \log(1 - p))$$

and the MLE

$$\hat{p}_{ML} = \bar{x}$$

and from Example 3.4 the observed Fisher information

$$J_n(\hat{p}_{ML}) = \frac{n}{\bar{x}(1 - \bar{x})}$$

The log-likelihood at the MLE is

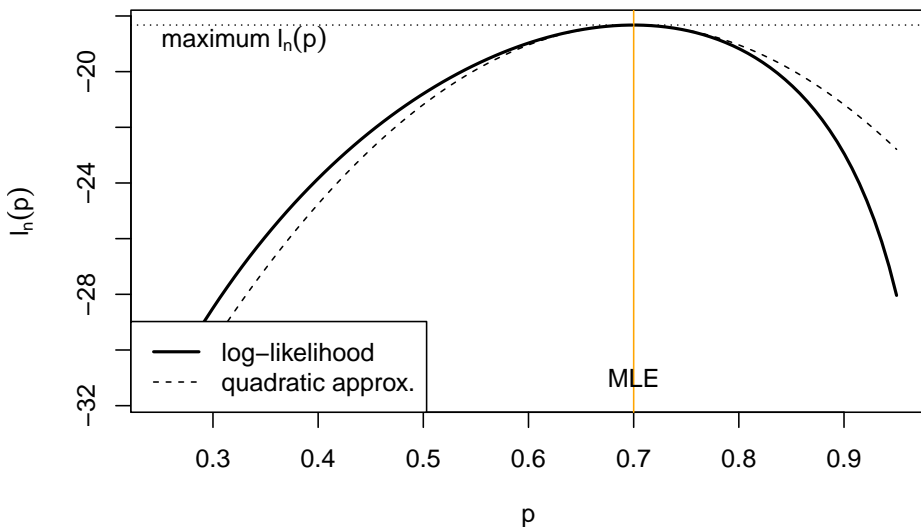
$$l_n(\hat{p}_{ML}) = n(\bar{x} \log \bar{x} + (1 - \bar{x}) \log(1 - \bar{x}))$$

This allows us to construct the quadratic approximation of the log-likelihood around the MLE as

$$\begin{aligned}
 l_n(p) &\approx l_n(\hat{p}_{ML}) - \frac{1}{2} J_n(\hat{p}_{ML})(p - \hat{p}_{ML})^2 \\
 &= n \left( \bar{x} \log \bar{x} + (1 - \bar{x}) \log(1 - \bar{x}) - \frac{(p - \bar{x})^2}{2\bar{x}(1 - \bar{x})} \right) \\
 &= C + \frac{\bar{x}p - \frac{1}{2}p^2}{\bar{x}(1 - \bar{x})/n}
 \end{aligned}$$

The constant  $C$  does not depend on  $p$ , its only purpose is to match the approximate log-likelihood at the MLE with that of the corresponding original log-likelihood. The approximate log-likelihood takes on the form of a normal log-likelihood (Example 3.2) for one observation of  $\hat{p}_{ML} = \bar{x}$  from  $N\left(p, \frac{\bar{x}(1-\bar{x})}{n}\right)$ .

The following figure shows the above log-likelihood function and its quadratic approximation for example data with  $n = 30$  and  $\bar{x} = 0.7$ :



### 4.2.3 Asymptotic normality of maximum likelihood estimates

Intuitively, it makes sense to associate large amount of curvature of the log-likelihood at the MLE with low variance of the MLE (and conversely, low amount of curvature with high variance).

From the above we see that

- normality implies a quadratic log-likelihood,
- conversely, taking a quadratic approximation of the log-likelihood implies approximate normality, and

- in the quadratic approximation **the inverse observed Fisher information plays the role of the covariance** of the MLE.

This suggests the following theorem: **Asymptotically, the MLE is normally distributed around the true parameter and with covariance equal to the inverse of the observed Fisher information:**

$$\hat{\theta}_{ML} \overset{a}{\sim} \underbrace{N_d}_{\text{multivariate normal}} \left( \underbrace{\theta}_{\text{mean vector}}, \underbrace{J_n(\hat{\theta}_{ML})^{-1}}_{\text{covariance matrix}} \right)$$

This theorem about the distributional properties of MLEs greatly enhances the usefulness of the method of maximum likelihood. It implies that in regular settings maximum likelihood is not just a method for obtaining point estimates but also also provides estimates of their uncertainty.

However, we need to clarify what “asymptotic” actually means in the context of the above theorem:

- 1) Primarily, it means to have sufficient sample size so that the log-likelihood  $l_n(\theta)$  is sufficiently well approximated by a quadratic function around  $\hat{\theta}_{ML}$ . The better the local quadratic approximation the better the normal approximation!
- 2) In a regular model with positive definite observed Fisher information matrix this is guaranteed for large sample size  $n \rightarrow \infty$  thanks to the central limit theorem).
- 3) However,  $n$  going to infinity is in fact not always required for the normal approximation to hold! Depending on the particular model a good local fit to a quadratic log-likelihood may be available also for finite  $n$ . As a trivial example, for the normal log-likelihood it is valid for any  $n$ .
- 4) In the other hand, in non-regular models (with nondifferentiable log-likelihood at the MLE and/or a singular Fisher information matrix) no amount of data, not even  $n \rightarrow \infty$ , will make the quadratic approximation work.

Remarks:

- The technical details of the above considerations are worked out in the theory of [locally asymptotically normal \(LAN\) models](#) pioneered in 1960 by [Lucien LeCam \(1924–2000\)](#).
- There are also methods to obtain higher-order (higher than quadratic and thus non-normal) asymptotic approximations. These relate to so-called [saddle point approximations](#).

### 4.2.4 Asymptotic optimal efficiency

Assume now that  $\hat{\theta}$  is an arbitrary and unbiased estimator for  $\theta$  and the underlying data generating model is regular with density  $f(x|\theta)$ .

H. Cramér (1893–1985), C. R. Rao (1920–) and others demonstrated in 1945 the so-called **information inequality**,

$$\text{Var}(\hat{\theta}) \geq \frac{1}{n} \mathbf{I}^{\text{Fisher}}(\theta)^{-1}$$

which puts a lower bound on the variance of an estimator for  $\theta$ . (Note for  $d > 1$  this is a matrix inequality, meaning that the difference matrix is positive semidefinite).

For large sample size with  $n \rightarrow \infty$  and  $\hat{\theta}_{ML} \rightarrow \theta$  the observed Fisher information becomes  $J_n(\hat{\theta}) \rightarrow n \mathbf{I}^{\text{Fisher}}(\theta)$  and therefore we can write the asymptotic distribution of  $\hat{\theta}_{ML}$  as

$$\hat{\theta}_{ML} \stackrel{a}{\sim} N_d \left( \theta, \frac{1}{n} \mathbf{I}^{\text{Fisher}}(\theta)^{-1} \right)$$

This means that for large  $n$  in regular models  $\hat{\theta}_{ML}$  achieves the lowest variance possible according to the Cramér–Rao information inequality. In other words, for large sample size maximum likelihood is optimally efficient and thus the best available estimator will in fact be the MLE!

However, as we will see later this does not hold for small sample size where it is indeed possible (and necessary) to improve over the MLE (e.g. via Bayesian estimation or regularisation).

## 4.3 Quantifying the uncertainty of maximum likelihood estimates

### 4.3.1 Estimating the variance of MLEs

In the previous section we saw that MLEs are asymptotically normally distributed, with the inverse Fisher information (both expected and observed) linked to the asymptotic variance.

This leads to the question whether to use the observed Fisher information  $J_n(\hat{\theta}_{ML})$  or the expected Fisher information at the MLE  $n \mathbf{I}^{\text{Fisher}}(\hat{\theta}_{ML})$  to estimate the variance of the MLE?

- Clearly, for  $n \rightarrow \infty$  both can be used interchangeably.
- However, they can be very different for finite  $n$  in particular for models outside the exponential family.
- Also normality may occur well before  $n$  goes to  $\infty$ .

### 4.3. QUANTIFYING THE UNCERTAINTY OF MAXIMUM LIKELIHOOD ESTIMATES

Therefore one needs to choose between the two, considering also that

- the expected Fisher information at the MLE is the average curvature at the MLE, whereas the observed Fisher information is the actual observed curvature, and
- the observed Fisher information naturally occurs in the quadratic approximation of the log-likelihood.

All in all, the observed Fisher information as estimator of the variance is more appropriate as it is based on the actual observed data and also works for large  $n$  (in which case it yields the same result as using expected Fisher information):

$$\widehat{\text{Var}}(\hat{\theta}_{ML}) = J_n(\hat{\theta}_{ML})^{-1}$$

and its square-root as the estimate of the standard deviation

$$\widehat{\text{SD}}(\hat{\theta}_{ML}) = J_n(\hat{\theta}_{ML})^{-1/2}$$

Note that in the above we use *matrix inversion* and the (inverse) *matrix square root*.

The reasons for preferring observed Fisher information are made mathematically precise in a classic paper by Efron and Hinkley (1978).

**Example 4.3.** Estimated variance and distribution of the MLE of a proportion:

From Examples 3.1 and 3.4 we know the MLE

$$\hat{p}_{ML} = \bar{x} = \frac{k}{n}$$

and the corresponding observed Fisher information

$$J_n(\hat{p}_{ML}) = \frac{n}{\hat{p}_{ML}(1 - \hat{p}_{ML})}$$

The estimated variance of the MLE is therefore

$$\widehat{\text{Var}}(\hat{p}_{ML}) = \frac{\hat{p}_{ML}(1 - \hat{p}_{ML})}{n}$$

and the corresponding asymptotic normal distribution is

$$\hat{p}_{ML} \overset{a}{\sim} N\left(p, \frac{\hat{p}_{ML}(1 - \hat{p}_{ML})}{n}\right)$$

**Example 4.4.** Estimated variance and distribution of the MLE of the mean parameter for the normal distribution with known variance:

From Examples 3.2 and 3.5 we know that

$$\hat{\mu}_{ML} = \bar{x}$$

and that the corresponding observed Fisher information at  $\hat{\mu}_{ML}$  is

$$J_n(\hat{\mu}_{ML}) = \frac{n}{\sigma^2}$$

The estimated variance of the MLE is therefore

$$\widehat{\text{Var}}(\hat{\mu}_{ML}) = \frac{\sigma^2}{n}$$

and the corresponding asymptotic normal distribution is

$$\hat{\mu}_{ML} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Note that in this case the distribution is not asymptotic but is **exact**, i.e. valid also for small  $n$  (as long as the data  $x_i$  are actually from  $N(\mu, \sigma^2)$ !).

### 4.3.2 Wald statistic

Centering the MLE  $\hat{\theta}_{ML}$  with  $\theta_0$  followed by standardising with  $\widehat{\text{SD}}(\hat{\theta}_{ML})$  yields the **Wald statistic** (named after [Abraham Wald, 1902–1950](#)):

$$\begin{aligned} \mathbf{t}(\theta_0) &= \widehat{\text{SD}}(\hat{\theta}_{ML})^{-1}(\hat{\theta}_{ML} - \theta_0) \\ &= J_n(\hat{\theta}_{ML})^{1/2}(\hat{\theta}_{ML} - \theta_0) \end{aligned}$$

The **squared Wald statistic** is a scalar defined as

$$\begin{aligned} t(\theta_0)^2 &= \mathbf{t}(\theta_0)^T \mathbf{t}(\theta_0) \\ &= (\hat{\theta}_{ML} - \theta_0)^T J_n(\hat{\theta}_{ML})(\hat{\theta}_{ML} - \theta_0) \end{aligned}$$

Note that in the literature both  $\mathbf{t}(\theta_0)$  and  $t(\theta_0)^2$  are commonly referred to as Wald statistics. In this text we use the qualifier “squared” if we refer to the latter.

We now assume that the true underlying parameter is  $\theta_0$ . Since the MLE is asymptotically normal the Wald statistic is asymptotically **standard normal** distributed:

$$\begin{aligned} \mathbf{t}(\theta_0) &\stackrel{a}{\sim} N_d(0, \mathbf{I}_d) && \text{for vector } \theta \\ t(\theta_0) &\stackrel{a}{\sim} N(0, 1) && \text{for scalar } \theta \end{aligned}$$

Correspondingly, the **squared** Wald statistic is chi-squared distributed:

$$\begin{aligned} t(\theta_0)^2 &\stackrel{a}{\sim} \chi_d^2 && \text{for vector } \theta \\ t(\theta_0)^2 &\stackrel{a}{\sim} \chi_1^2 && \text{for scalar } \theta \end{aligned}$$

The degree of freedom of the chi-squared distribution is the dimension  $d$  of the parameter vector  $\theta$ .



### 4.3. QUANTIFYING THE UNCERTAINTY OF MAXIMUM LIKELIHOOD ESTIMATES

**Example 4.5.** Wald statistic for a proportion:

We continue from Example 4.3. With  $\hat{p}_{ML} = \bar{x}$  and  $\widehat{\text{Var}}(\hat{p}_{ML}) = \frac{\hat{p}_{ML}(1-\hat{p}_{ML})}{n}$  and thus  $\widehat{\text{SD}}(\hat{p}_{ML}) = \sqrt{\frac{\hat{p}_{ML}(1-\hat{p}_{ML})}{n}}$  we get as **Wald statistic**:

$$t(p_0) = \frac{\bar{x} - p_0}{\sqrt{\bar{x}(1-\bar{x})/n}} \stackrel{a}{\sim} N(0, 1)$$

The **squared Wald statistic** is:

$$t(p_0)^2 = n \frac{(\bar{x} - p_0)^2}{\bar{x}(1-\bar{x})} \stackrel{a}{\sim} \chi_1^2$$

**Example 4.6.** Wald statistic for the mean parameter of a normal distribution with known variance:

We continue from Example 4.4. With  $\hat{\mu}_{ML} = \bar{x}$  and  $\widehat{\text{Var}}(\hat{\mu}_{ML}) = \frac{\sigma^2}{n}$  and thus  $\widehat{\text{SD}}(\hat{\mu}_{ML}) = \frac{\sigma}{\sqrt{n}}$  we get as **Wald statistic**:

$$t(\mu_0) = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \sim N(0, 1)$$

Note this is the one sample  $t$ -statistic with given  $\sigma$ . The **squared Wald statistic** is:

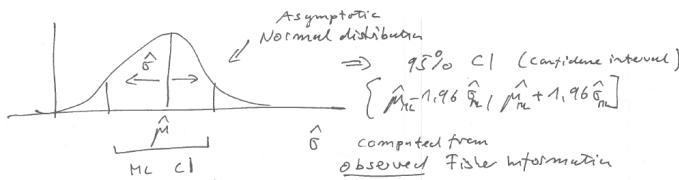
$$t(\mu_0)^2 = \frac{(\bar{x} - \mu_0)^2}{\sigma^2/n} \sim \chi_1^2$$

Again, in this instance this is the exact distribution, not just the asymptotic one.

Using the Wald statistic or the squared Wald statistic we can test whether a particular  $\mu_0$  can be rejected as underlying true parameter, and we can also construct corresponding confidence intervals.

#### 4.3.3 Normal confidence intervals using the Wald statistic

The asymptotic normality of MLEs derived from regular models enables us to construct a corresponding normal confidence interval (CI):



For example, to construct the asymptotic normal CI for the MLE of a scalar parameter  $\theta$  we use the MLE  $\hat{\theta}_{ML}$  as estimate of the mean and its standard deviation  $\widehat{SD}(\hat{\theta}_{ML})$  computed from the observed Fisher information:

$$CI = [\hat{\theta}_{ML} \pm c_{normal} \widehat{SD}(\hat{\theta}_{ML})]$$

$c_{normal}$  is a critical value for the standard-normal symmetric confidence interval chosen to achieve the desired nominal coverage- The critical values are computed using the inverse standard normal distribution function via  $c_{normal} = \Phi^{-1}\left(\frac{1+\kappa}{2}\right)$  (cf. refresher section in the Appendix).

coverage $\kappa$	Critical value $c_{normal}$
0.9	1.64
0.95	1.96
0.99	2.58

For example, for a CI with 95% coverage one uses the factor 1.96 so that

$$CI = [\hat{\theta}_{ML} \pm 1.96 \widehat{SD}(\hat{\theta}_{ML})]$$

The normal CI can be expressed using Wald statistic as follows:

$$CI = \{\theta_0 : |t(\theta_0)| < c_{normal}\}$$

Similary, it can also be expressed using the squared Wald statistic:

$$CI = \{\theta_0 : t(\theta_0)^2 < c_{chisq}\}$$

Note that this form facilitates the construction of normal confidence intervals for a parameter vector  $\theta_0$ .

The following lists containst the critical values resulting from the chi-squared distribution with degree of freedom  $m = 1$  for the three most common choices of coverage  $\kappa$  for a normal CI for a univariate parameter:

coverage $\kappa$	Critical value $c_{chisq} (m = 1)$
0.9	2.71
0.95	3.84
0.99	6.63

**Example 4.7.** Asymptotic normal confidence interval for a proportion:

We continue from Examples 4.3 and 4.5. Assume we observe  $n = 30$  measurements with average  $\bar{x} = 0.7$ . Then  $\hat{p}_{ML} = \bar{x} = 0.7$  and  $\widehat{SD}(\hat{p}_{ML}) = \sqrt{\frac{\bar{x}(1-\bar{x})}{n}} \approx 0.084$ .

The symmetric asymptotic normal CI for  $p$  with 95% coverage is given by  $\hat{p}_{ML} \pm 1.96 \widehat{SD}(\hat{p}_{ML})$  which for the present data results in the interval  $[0.536, 0.864]$ .

**Example 4.8.** Normal confidence interval for the mean:

We continue from Examples 4.4 and 4.6. Assume that we observe  $n = 25$  measurements with average  $\bar{x} = 10$ , from a normal with unknown mean and variance  $\sigma^2 = 4$ .

Then  $\hat{\mu}_{ML} = \bar{x} = 10$  and  $\widehat{SD}(\hat{\mu}_{ML}) = \sqrt{\frac{\sigma^2}{n}} = \frac{2}{5}$ .

The symmetric asymptotic normal CI for  $p$  with 95% coverage is given by  $\hat{\mu}_{ML} \pm 1.96 \widehat{SD}(\hat{\mu}_{ML})$  which for the present data results in the interval  $[9.216, 10.784]$ .

#### 4.3.4 Normal tests using the Wald statistic

Finally, recall the **duality between confidence intervals and statistical tests**. Specifically, a confidence interval with coverage  $\kappa$  can be also used for testing as follows.

- for every  $\theta_0$  inside the CI the data do not allow to reject the hypothesis that  $\theta_0$  is the true parameter with significance level  $1 - \kappa$ .
- Conversely, all values  $\theta_0$  outside the CI can be rejected to be the true parameter with significance level  $1 - \kappa$ .

Hence, in order to test whether  $\theta_0$  is the true underlying parameter value we can compute the corresponding (squared) Wald statistic, find the desired critical value and then decide on rejection.

**Example 4.9.** Asymptotic normal test for a proportion:

We continue from Example 4.7.

We now consider two possible values ( $p_0 = 0.5$  and  $p_0 = 0.8$ ) as potentially true underlying proportion.

The value  $p_0 = 0.8$  lies inside the 95% confidence interval  $[0.536, 0.864]$ . This implies we cannot reject the hypothesis that this is the true underlying parameter on 5% significance level. In contrast,  $p_0 = 0.5$  is outside the confidence interval so we can indeed reject this value. In other words, data plus model exclude this value as statistically implausible.

This can be verified more directly by computing the corresponding (squared) Wald statistics (see Example 4.5) and comparing them with the relevant critical value (3.84 from chi-squared distribution for 5% significance level):

- $t(0.5)^2 = 5.71 > 3.84$  hence  $p_0 = 0.5$  can be rejected.

- $t(0.8)^2 = 1.43 < 3.84$  hence  $p_0 = 0.8$  cannot be rejected.

Note that the squared Wald statistic at the boundaries of the normal confidence interval is equal to the critical value.

**Example 4.10.** Normal confidence interval and test for the mean:

We continue from Example 4.8.

We now consider two possible values ( $\mu_0 = 9.5$  and  $\mu_0 = 11$ ) as potentially true underlying mean parameter.

The value  $\mu_0 = 9.5$  lies inside the 95% confidence interval  $[9.216, 10.784]$ . This implies we cannot reject the hypothesis that this is the true underlying parameter on 5% significance level. In contrast,  $\mu_0 = 11$  is outside the confidence interval so we can indeed reject this value. In other words, data plus model exclude this value as a statistically implausible.

This can be verified more directly by computing the corresponding (squared) Wald statistics (see Example 4.6) and comparing them with the relevant critical values:

- $t(9.5)^2 = 1.56 < 3.84$  hence  $\mu_0 = 9.5$  cannot be rejected.
- $t(11)^2 = 6.25 > 3.84$  hence  $\mu_0 = 11$  can be rejected.

The squared Wald statistic at the boundaries of the confidence interval equals the critical value.

Note that this is the standard one-sample test of the mean, and that it is exact, not an approximation.

## 4.4 Example of a non-regular model

Not all models allow a quadratic approximation of the log-likelihood function around the MLE. This is the case when the log-likelihood function is not differentiable at the MLE. These models are called non-regular and for those models the normal approximation is not available.

**Example 4.11.** Uniform distribution with upper bound  $\theta$ :

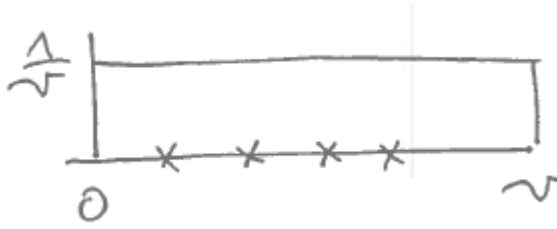
$$x_1, \dots, x_n \sim U(0, \theta)$$

With  $x_{[i]}$  we denote the *ordered* observations with  $0 \leq x_{[1]} < x_{[2]} < \dots < x_{[n]} \leq \theta$  and  $x_{[n]} = \max(x_1, \dots, x_n)$ .

We would like to obtain both the maximum likelihood estimator  $\hat{\theta}_{ML}$  and its distribution.

The probability density function of  $U(0, \theta)$  is

$$f(x|\theta) = \begin{cases} \frac{1}{\theta} & \text{if } x \in [0, \theta] \\ 0 & \text{otherwise.} \end{cases}$$



and on the log-scale

$$\log f(x|\theta) = \begin{cases} -\log \theta & \text{if } x \in [0, \theta] \\ -\infty & \text{otherwise.} \end{cases}$$

Since all observed data  $x_1, \dots, x_n$  lie in the interval  $[0, \theta]$  we get as log-likelihood function

$$l_n(\theta) = \begin{cases} -n \log \theta & \text{for } x_{[n]} \leq \theta \\ -\infty & \text{otherwise} \end{cases}$$

Obtaining the MLE of  $\theta$  is straightforward:  $-n \log \theta$  is monotonically decreasing therefore the log-likelihood function has a maximum at  $\hat{\theta}_{ML} = x_{[n]}$ .

However, there is a discontinuity in  $l_n(\theta)$  at  $x_{[n]}$  and therefore  $l_n(\theta)$  **is not differentiable** at  $\hat{\theta}_{ML}$ . Thus, **there is no quadratic approximation around  $\hat{\theta}_{ML}$**  and the **observed Fisher information cannot be computed**. Hence, the normal approximation for the distribution of  $\hat{\theta}_{ML}$  is not valid regardless of sample size, i.e. not even asymptotically for  $n \rightarrow \infty$ .

Nonetheless, we can in fact still obtain the sampling distribution of  $\hat{\theta}_{ML} = x_{[n]}$ . However, *not* via asymptotic arguments but instead by understanding that  $x_{[n]}$  is an order statistic (see [https://en.wikipedia.org/wiki/Order\\_statistic](https://en.wikipedia.org/wiki/Order_statistic)) with the following properties:

$$x_{[n]} \sim \theta \text{Beta}(n, 1) \quad \text{"n-th order statistic"}$$

$$E(x_{[n]}) = \frac{n}{n+1}\theta$$

$$\text{Var}(x_{[n]}) = \frac{n}{(n+1)^2(n+2)}\theta^2 \approx \frac{\theta^2}{n^2}$$

Note that the variance decreases with  $\frac{1}{n^2}$  which is much faster than the usual  $\frac{1}{n}$  of an "efficient" estimator. Correspondingly,  $\hat{\theta}_{ML}$  is a so-called "super efficient" estimator.



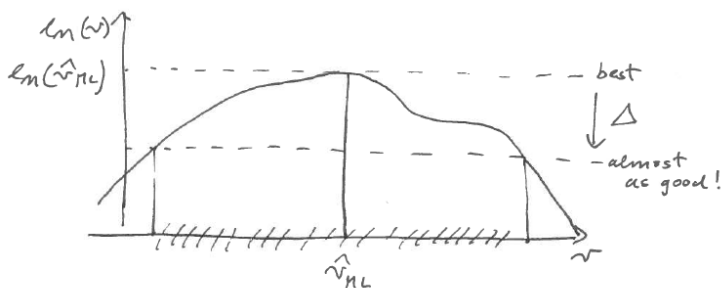
# Chapter 5

## Likelihood-based confidence interval and likelihood ratio

### 5.1 Likelihood-based confidence intervals and Wilks statistic

#### 5.1.1 General idea and definition of Wilks statistic

Instead of relying on normal / quadratic approximation, we can also use the log-likelihood directly to find the so called **likelihood confidence intervals**:



Idea: find all  $\theta_0$  that have a log-likelihood that is almost as good as  $l_n(\hat{\theta}_{ML})$ .

$$CI = \{\theta_0 : l_n(\hat{\theta}_{ML}) - l_n(\theta_0) \leq \Delta\}$$

Here  $\Delta$  is the tolerated deviation from the maximum log-likelihood. We will see below how to determine a suitable  $\Delta$  further below.

The above leads naturally to the **Wilks log likelihood ratio statistic**  $W(\theta_0)$

defined as:

$$\begin{aligned} W(\theta_0) &= 2 \log \left( \frac{L(\hat{\theta}_{ML})}{L(\theta_0)} \right) \\ &= 2(l_n(\hat{\theta}_{ML}) - l_n(\theta_0)) \end{aligned}$$

With its help we can write the likelihood CI follows:

$$CI = \{\theta_0 : W(\theta_0) \leq 2\Delta\}$$

The Wilks statistic is named after [Samuel S. Wilks \(1906–1964\)](#).

Advantages of using a likelihood-based CI:

- not restricted to be symmetric
- enables to construct multivariate CIs for parameter vector easily even in non-normal cases
- contains normal CI as special case

**Question:** how to choose  $\Delta$ , i.e how to calibrate the likelihood interval? Essentially, by comparing with normal CI!

**Example 5.1.** Wilks statistic for the proportion:

The log-likelihood for the parameter  $p$  is (cf. Example 3.1)

$$l_n(p) = n(\bar{x} \log p + (1 - \bar{x}) \log(1 - p))$$

Hence the Wilks statistic is

$$\begin{aligned} W(p_0) &= 2(l_n(\hat{p}_{ML}) - l_n(p_0)) \\ &= 2n \left( \bar{x} \log \left( \frac{\bar{x}}{p_0} \right) + (1 - \bar{x}) \log \left( \frac{1 - \bar{x}}{1 - p_0} \right) \right) \end{aligned}$$

Comparing with Example 2.8 we see that in this case the Wilks statistic is essentially (apart from a scale factor  $2n$ ) the KL divergence between two Bernoulli distributions:

$$W(p_0) = 2n D_{KL}(\text{Ber}(\hat{p}_{ML}), \text{Ber}(p_0))$$

**Example 5.2.** Wilks statistic for the mean parameter of a normal model:

The Wilks statistic is

$$W(\mu_0)^2 = \frac{(\bar{x} - \mu_0)^2}{\sigma^2/n}$$

See Worksheet 4 for a derivation of the Wilks statistic directly from the log-likelihood function.

Note this is the same as the squared Wald statistic discussed in Example 4.6.



Comparing with Example 2.10 we see that in this case the Wilks statistic is essentially (apart from a scale factor  $2n$ ) the KL divergence between two normal distributions with different means and variance equal to  $\sigma^2$ :

$$W(p_0) = 2nD_{\text{KL}}(N(\hat{\mu}_{ML}, \sigma^2), N(\mu_0, \sigma^2))$$

### 5.1.2 Quadratic approximation of Wilks statistic and squared Wald statistic

Recall the *quadratic approximation* (= second order Taylor series around the MLE  $\hat{\theta}_{ML}$ ) applied the log-likelihood function  $l_n(\theta_0)$ :

$$l_n(\theta_0) \approx l_n(\hat{\theta}_{ML}) - \frac{1}{2}(\theta_0 - \hat{\theta}_{ML})^T J_n(\hat{\theta}_{ML})(\theta_0 - \hat{\theta}_{ML})$$

With this we can then approximate the Wilks statistic:

$$\begin{aligned} W(\theta_0) &= 2(l_n(\hat{\theta}_{ML}) - l_n(\theta_0)) \\ &\approx (\theta_0 - \hat{\theta}_{ML})^T J_n(\hat{\theta}_{ML})(\theta_0 - \hat{\theta}_{ML}) \\ &= t(\theta_0)^2 \end{aligned}$$

Thus the quadratic approximation of the Wilks statistic yields the squared Wald statistic!

Conversely, the Wilks statistic can be understood a generalisation of the squared Wald statistic.

**Example 5.3.** Quadratic approximation of the Wilks statistic for a proportion (continued from Example 5.1):

A Taylor series of second order (for  $p_0$  around  $\bar{x}$ ) yields

$$\log\left(\frac{\bar{x}}{p_0}\right) \approx -\frac{p_0 - \bar{x}}{\bar{x}} + \frac{(p_0 - \bar{x})^2}{2\bar{x}^2}$$

and

$$\log\left(\frac{1 - \bar{x}}{1 - p_0}\right) \approx \frac{p_0 - \bar{x}}{1 - \bar{x}} + \frac{(p_0 - \bar{x})^2}{2(1 - \bar{x})^2}$$

With this we can approximate the Wilks statistic of the proportion as

$$\begin{aligned} W(p_0) &\approx 2n \left( -(p_0 - \bar{x}) + \frac{(p_0 - \bar{x})^2}{2\bar{x}} + (p_0 - \bar{x}) + \frac{(p_0 - \bar{x})^2}{2(1 - \bar{x})} \right) \\ &= n \left( \frac{(p_0 - \bar{x})^2}{\bar{x}} + \frac{(p_0 - \bar{x})^2}{(1 - \bar{x})} \right) \\ &= n \left( \frac{(p_0 - \bar{x})^2}{\bar{x}(1 - \bar{x})} \right) \\ &= t(p_0)^2. \end{aligned}$$

This verifies that the quadratic approximation of the Wilks statistic leads back to the squared Wald statistic of Example 4.5.

**Example 5.4.** Quadratic approximation of the Wilks statistic for the mean parameter of a normal model (continued from Example 5.2):

The normal log-likelihood is already quadratic in the mean parameter (cf. Example 3.2). Correspondingly, the Wilks statistic is quadratic in the mean parameter as well. Hence in this particular case the quadratic “approximation” is in fact exact and the Wilks statistic and the squared Wald statistic are identical!

Correspondingly, confidence intervals and tests based on the Wilks statistic are identical to those obtained using the Wald statistic.

### 5.1.3 Distribution of the Wilks statistic

The connection with the squared Wald statistic implies that both have asymptotically the same distribution.

Hence, under  $\theta_0$  the Wilks statistic is distributed asymptotically as

$$W(\theta_0) \overset{a}{\sim} \chi_d^2$$

where  $d$  is the number of parameters in  $\theta$ , i.e. the dimension of the model.

For scalar  $\theta$  (i.e. single parameter and  $d = 1$ ) this becomes

$$W(\theta_0) \overset{a}{\sim} \chi_1^2$$

This fact is known as **Wilks’ theorem**.

### 5.1.4 Cutoff values for the likelihood CI

coverage $\kappa$	$\Delta = \frac{c_{chisq}}{2} \ (m = 1)$
0.9	1.35
0.95	1.92
0.99	3.32

The asymptotic distribution for  $W$  is useful to choose a suitable  $\Delta$  for the likelihood CI — note that  $2\Delta = c_{chisq}$  where  $c_{chisq}$  is the critical value for a specified coverage  $\kappa$ . This yields the above table for scalar parameter

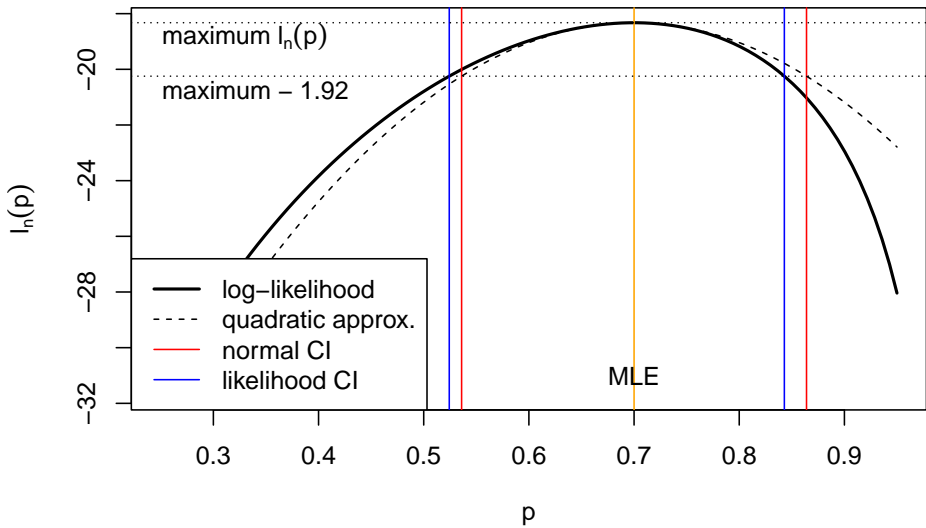
**Example 5.5.** Likelihood confidence interval for a proportion:

We continue from Example 5.1, and as in Example 4.7 we assume we have data with  $n = 30$  and  $\bar{x} = 0.7$ .

This yields (via numerical root finding) as the 95% likelihood confidence interval the interval  $[0.524, 0.843]$ . It is similar but not identical to the corresponding asymptotic normal interval  $[0.536, 0.864]$  obtained in Example 4.7.

The following figure illustrate the relationship between the normal CI, the likelihood CI and also shows the role of the quadratic approximation (see also Example 4.2). Note that:

- the normal CI is symmetric around the MLE whereas the likelihood CI is not symmetric
- the normal CI is identical to the likelihood CI when using the quadratic approximation!



### 5.1.5 Likelihood ratio test (LRT) using Wilks statistic

As in the normal case (with Wald statistic and normal CIs) one can also construct a test using the Wilks statistic:

$$\begin{array}{lll}
 H_0 : \theta = \theta_0 & \text{True model is } \theta_0 & \text{Null hypothesis} \\
 H_1 : \theta \neq \theta_0 & \text{True model is not } \theta_0 & \text{Alternative hypothesis}
 \end{array}$$

As test statistic we use the Wilks log likelihood ratio  $W(\theta_0)$ . Extreme values of this test statistic imply evidence against  $H_0$ .

Note that the null model is “simple” (= a single parameter value) whereas the alternative model is “composite” (= a set of parameter values).

#### Remarks:

- The composite alternative  $H_1$  is represented by a single point (the MLE).

- **Reject  $H_0$  for large values of  $W(\theta_0)$**
- under  $H_0$  and for large  $n$  the statistic  $W(\theta_0)$  is chi-squared distributed, i.e.  $W(\theta_0) \stackrel{a}{\sim} \chi_d^2$ . This allows to compute critical values (i.e. thresholds to declared rejection under a given significance level) and also  $p$ -values corresponding to the observed test statistics.
- Models **outside** the CI are **rejected**
- Models **inside** the CI **cannot be rejected**, i.e. they can't be statistically distinguished from the best alternative model.

A statistic equivalent to  $W(\theta_0)$  is the **likelihood ratio**

$$\Lambda(\theta_0) = \frac{L(\theta_0)}{L(\hat{\theta}_{ML})}$$

The two statistics can be transformed into each other by  $W(\theta_0) = -2 \log \Lambda(\theta_0)$  and  $\Lambda(\theta_0) = e^{-W(\theta_0)/2}$ . We **reject  $H_0$  for small values of  $\Lambda$** .

It can be shown that the likelihood ratio test to compare two simple model is optimal in the sense that for any given specified type I error (=probability of wrongly rejecting  $H_0$ , i.e. the significance level) it will maximise the power (=1- type II error, probability of correctly accepting  $H_1$ ). This is known as the **Neyman-Pearson theorem**.

**Example 5.6.** Likelihood test for a proportion:

We continue from Example 5.5 with 95% likelihood confidence interval  $[0.524, 0.843]$ .

The value  $p_0 = 0.5$  is outside the CI and hence can be rejected whereas  $p_0 = 0.8$  is inside the CI and hence cannot be rejected on 5% significance level.

The Wilks statistic for  $p_0 = 0.5$  and  $p_0 = 0.8$  take on the following values:

- $W(0.5)^2 = 4.94 > 3.84$  hence  $p_0 = 0.5$  can be rejected.
- $W(0.8)^2 = 1.69 < 3.84$  hence  $p_0 = 0.8$  cannot be rejected.

Note that the Wilks statistic at the boundaries of the likelihood confidence interval is equal to the critical value (3.84 corresponding to 5% significance level for a chi-squared distribution with 1 degree of freedom).

## 5.1.6 Origin of likelihood ratio statistic

The likelihood ratio statistic is asymptotically linked to differences in the KL divergences of the two compared models with the underlying true model.

Assume that  $F$  is the true (and unknown) data generating model  $G_\theta$  is a family of models and we would like to compare two candidate models  $G_A$  and  $G_B$  corresponding to parameters  $\theta_A$  and  $\theta_B$  on the basis of observed data  $x_1, \dots, x_n$ . The KL divergences  $D_A = D_{KL}(F, G_A)$  and  $D_B = D_{KL}(F, G_B)$  indicate how close each of the models  $G_A$  and  $G_B$  fit the true  $F$ . The difference  $D_B - D_A$  is

thus a way to measure the relative fit of the two models, and can be computed as

$$D_B - D_A = D_{\text{KL}}(F, G_B) - D_{\text{KL}}(F, G_A) = E_F \log \frac{g_A(x)}{g_B(x)}$$

Replacing  $F$  by the empirical distribution  $\hat{F}_n$  leads to the large sample approximation

$$2n(D_B - D_A) \approx 2(l_n(\theta_A) - l_n(\theta_B))$$

Hence, the difference in the log-likelihoods provides an estimate of the difference in the KL divergence of the two models involved.

The Wilks log likelihood ratio statistic

$$W(\theta_0) = 2(l_n(\hat{\theta}_{ML}) - l_n(\theta_0)) \approx 2n(D_{F_{\theta_0}} - D_{F_{\hat{\theta}_{ML}}})$$

thus compares the best-fit distribution with  $\hat{\theta}_{ML}$  as the parameter to the distribution with parameter  $\theta_0$ .

For some specific models the Wilks statistic can also be written in the form of the KL divergence:

$$W(\theta_0) = 2nD_{\text{KL}}(F_{\hat{\theta}_{ML}}, F_{\theta_0})$$

This is the case for the examples 5.1 and 5.2 and also more generally for exponential family models, but it is not true in general.

## 5.2 Generalised likelihood ratio test (GLRT)

Also known as **maximum likelihood ratio test (MLRT)**. The Generalised Likelihood Ratio Test (GLRT) works like the standard likelihood ratio test with the difference that now the null model  $H_0$  is a composite model. This means that in the denominator in the test statistic needs to be optimised as well.

$$\begin{array}{ll} H_0 : \theta \in \omega_0 \subset \Omega & \text{True model lies in restricted model space} \\ H_1 : \theta \in \omega_1 = \Omega \setminus \omega_0 & \text{True model is not the restricted model space} \end{array}$$

Both  $H_0$  and  $H_1$  are now composite hypotheses.  $\Omega$  represents the unrestricted model space with dimension (=number of free parameters)  $d = |\Omega|$ . The constrained space  $\omega_0$  has degree of freedom  $d_0 = |\omega_0|$  with  $d_0 < d$ . Note that in the standard LRT the set  $\omega_0$  is a simple point with  $d_0 = 0$  as the null model is a simple distribution. Thus, LRT is contained in GLRT as special case!

The corresponding generalised (log) likelihood ratio statistic is given by

$$W = 2 \log \left( \frac{L(\hat{\theta}_{ML})}{L(\hat{\theta}_{ML}^0)} \right) \text{ and } \Lambda = \frac{\max_{\theta \in \omega_0} L(\theta)}{\max_{\theta \in \Omega} L(\theta)}$$

where  $L(\hat{\theta}_{ML})$  is the maximised likelihood assuming the full model (with parameter space  $\Omega$ ) and  $L(\hat{\theta}_{ML}^0)$  is the maximised likelihood for the restricted model (with parameter space  $\omega_0$ ).

**Remarks:**

- MLE in the restricted model space  $\omega_0$  is taken as a representative of  $H_0$ .
- The likelihood is **maximised in both numerator and denominator**.
- The restricted model is a special case of the full model (i.e. the two models are nested).
- The asymptotic distribution of  $W$  is chi-squared with degree of freedom depending on both  $d$  and  $d_0$ :

$$W \stackrel{a}{\sim} \chi_{d-d_0}^2$$

- This result is due to Wilks (1938). Note that it assumes that the true model is contained among the investigated models.
- If  $H_0$  is a simple hypothesis (i.e.  $d_0 = 0$ ) then the standard LRT (and corresponding CI) is recovered as special case of the GLRT.

**Example 5.7.** GLRT example:

*Case-control study:* (e.g. “healthy” vs. “disease”)

we observe normal data from two groups with sample size  $n_1$  and  $n_2$  (and  $n = n_1 + n_2$ ):

$$x_1, \dots, x_{n_1} \sim N(\mu_1, \sigma^2)$$

and

$$x_{n_1+1}, \dots, x_n \sim N(\mu_2, \sigma^2)$$

Question: are the two means  $\mu_1$  and  $\mu_2$  the same in the two groups?

$$H_0 : \mu_1 = \mu_2 \text{ (with variance unknown nuisance parameter)}$$

$$H_1 : \mu_1 \neq \mu_2$$

*Restricted and full models:*

$\omega_0$ : restricted model with two parameters  $\mu_0$  and  $\sigma_0^2$  (so that  $x_1, \dots, x_n \sim N(\mu_0, \sigma_0^2)$ ).

$\Omega$ : full model with three parameters  $\mu_1, \mu_2, \sigma^2$ .

*Corresponding log-likelihood functions:*

Restricted model  $\omega_0$ :

$$\log L(\mu_0, \sigma_0^2) = -\frac{n}{2} \log(\sigma_0^2) - \frac{1}{2\sigma_0^2} \sum_{i=1}^n (x_i - \mu_0)^2$$

Full model  $\Omega$ :

$$\begin{aligned}\log L(\mu_1, \mu_2, \sigma^2) &= \left( -\frac{n_1}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^{n_1} (x_i - \mu_1)^2 \right) + \\ &\quad \left( -\frac{n_2}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=n_1+1}^n (x_i - \mu_1)^2 \right) \\ &= -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \left( \sum_{i=1}^{n_1} (x_i - \mu_1)^2 + \sum_{i=n_1+1}^n (x_i - \mu_2)^2 \right)\end{aligned}$$

Corresponding MLEs:

$$\begin{aligned}\omega_0 : \quad \hat{\mu}_0 &= \frac{1}{n} \sum_{i=1}^n x_i & \hat{\sigma}_0^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_0)^2 \\ \Omega : \quad \hat{\mu}_1 &= \frac{1}{n_1} \sum_{i=1}^{n_1} x_i & \hat{\sigma}^2 &= \frac{1}{n} \left\{ \sum_{i=1}^{n_1} (x_i - \hat{\mu}_1)^2 + \sum_{i=n_1+1}^n (x_i - \hat{\mu}_2)^2 \right\} \\ \hat{\mu}_2 &= \frac{1}{n_2} \sum_{i=n_1+1}^n x_i\end{aligned}$$

We note that the two estimated variances are related by

$$\begin{aligned}\hat{\sigma}_0^2 &= \hat{\sigma}^2 + \frac{n_1 n_2}{n^2} (\hat{\mu}_1 - \hat{\mu}_2)^2 \\ &= \hat{\sigma}^2 \left( 1 + \frac{1}{n} \frac{(\hat{\mu}_1 - \hat{\mu}_2)^2}{\frac{n}{n_1 n_2} \hat{\sigma}^2} \right) \\ &= \hat{\sigma}^2 \left( 1 + \frac{t_{ML}^2}{n} \right)\end{aligned}$$

with

$$t_{ML} = \frac{\hat{\mu}_1 - \hat{\mu}_2}{\sqrt{\left( \frac{1}{n_1} + \frac{1}{n_2} \right) \hat{\sigma}^2}}$$

This is an example of a variance decomposition, with  $\hat{\sigma}_0^2$  being the estimated total variance and  $\hat{\sigma}^2$  the estimated within-group variance.

Corresponding maximised log-likelihood:

Restricted model:

$$\log L(\hat{\mu}_0, \hat{\sigma}_0^2) = -\frac{n}{2} \log(\hat{\sigma}_0^2) - \frac{n}{2}$$

Full model:

$$\log L(\hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}^2) = -\frac{n}{2} \log(\hat{\sigma}^2) - \frac{n}{2}$$

*Likelihood ratio statistic:*

$$\begin{aligned} W &= 2 \log \left( \frac{L(\hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}^2)}{L(\hat{\mu}_0, \hat{\sigma}_0^2)} \right) \\ &= 2 \log L(\hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}^2) - 2 \log L(\hat{\mu}_0, \hat{\sigma}_0^2) \\ &= n \log \left( \frac{\hat{\sigma}_0^2}{\hat{\sigma}^2} \right) \\ &= n \log \left( 1 + \frac{t_{ML}^2}{n} \right) \end{aligned}$$

The last step uses the decomposition for the total variance  $\hat{\sigma}_0^2$ . If an unbiased total variance estimate is used the

$$W = n \log \left( 1 + \frac{1}{n-2} t^2 \right)$$

with

$$t = \frac{\hat{\mu}_1 - \hat{\mu}_2}{\sqrt{\left( \frac{1}{n_1} + \frac{1}{n_2} \right) \frac{n}{n-2} \hat{\sigma}^2}}$$

→ the GRLT is a monotone function of the (squared) two-sample  $t$ -statistic!

**It can be shown that all standard tests with normal distributions can be interpreted as GLRTs!**

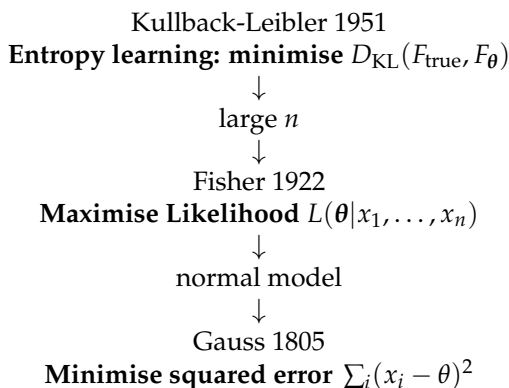


## Chapter 6

# Optimality properties and conclusion

### 6.1 Properties of maximum likelihood encountered so far

1. MLE is a special case of relative entropy minimisation *valid for large samples*.
2. MLE can be seen as generalisation of least squares (and conversely, least squares is a special case of ML).



3. Given a model, derivation of the MLE is basically automatic (only optimisation required)!
4. MLEs are **consistent**, i.e. if the true underlying model  $F_{\text{true}}$  with parameter  $\theta_{\text{true}}$  is contained in the set of specified candidate models  $F_{\theta}$  then the MLE will converge to the true model.

5. Correspondingly, **MLEs are asymptotically unbiased**.
6. However, MLEs are *not* necessarily unbiased in finite samples (e.g. the MLE of the variance parameter in the normal distribution).
7. The maximum likelihood is invariant against parameter transformations.
8. In regular situations (when local quadratic approximation is possible) MLEs are **asymptotically normally distributed**, with the asymptotic variance determined by the observed Fisher information.
9. In regular situations and for large sample size MLEs are **asymptotically optimally efficient** (Cramer-Rao theorem): For large samples the MLE achieves the lowest possible variance possible in an estimator — this is the so-called Cramer-Rao lower bound. The variance decreases to zero with  $n \rightarrow \infty$  typically with rate  $1/n$ .
10. The likelihood ratio can be used to construct optimal tests (in the sense of the Neyman-Pearson theorem).

## 6.2 Summarising data and the concept of minimal sufficiency

Another important concept in statistics and likelihood theory (especially when applied to the exponential family) is that of a **minimally sufficient statistic** to optimally summarise the information available in the data about a parameter in a model.

Generally, a **statistic**  $T(x_1, \dots, x_n) = T(x_i)$  is function of the data  $x_1, \dots, x_n$ . In the following we write  $x_i$  as a shorthand for the complete data set with  $n$  observations. The statistic  $T(x_i)$  can be of any type and value (scalar, vector, matrix etc. — even a function).  $T(x_i)$  is called a *summary statistic* if it describes important aspects of the data such as location (e.g. the average  $\text{avg}(x_i) = \bar{x}$ , the median) or scale (e.g. standard deviation, interquartile range).

A statistic  $T(x_i)$  is said to be **sufficient** for a parameter  $\theta$  in a model if the corresponding likelihood function can be written in terms of  $T(x_i)$  so that

$$L(\theta|x_i) = h(T(x_i), \theta) k(x_i),$$

where  $h(x)$  and  $k(x)$  are positive-valued functions, and or equivalently on log-scale

$$l_n(\theta) = \log h(T(x_i), \theta) + \log k(x_i).$$

This is known as the **Fisher-Pearson factorisation**. By construction, estimation and inference about  $\theta$  based on the factorised likelihood  $L(\theta)$  is mediated through the sufficient statistic  $T(x_i)$  and does not require the original data  $x_i$ . Instead, the sufficient statistic  $T(x_i)$  contains all the information in  $x_i$  required to learn about the parameter  $\theta$ . Therefore, if the MLE  $\hat{\theta}_{ML}$  of  $\theta$  exists and is unique

then **the MLE is a unique function of the sufficient statistic**  $T(x_i)$ . If the MLE is not unique then it can be chosen to be function of  $T(x_i)$ . Note that **a sufficient statistic always exists** since the data  $x_i$  are themselves sufficient statistics, with  $T(x_i) = x_i$ . Furthermore, sufficient statistics are **not unique** since applying a one-to-one transformation to  $T(x_i)$  yields another sufficient statistic.

Every sufficient statistic  $T(x_i)$  induces a partitioning of the space of data sets by clustering all hypothetical outcomes for which the statistic  $T(x_i)$  assumes the same value  $t$ :

$$\mathcal{X}_t = \{x_i : T(x_i) = t\}$$

The **data sets in  $\mathcal{X}_t$  are equivalent in terms of the sufficient statistic  $T(x_i)$** . Note that the dimensions of  $T(x_i)$  may be much smaller than those of  $x_i$ . Instead of  $n$  data points as few as one or two summaries may be sufficient to fully convey all the information in the data about the model parameters. Thus, transforming data  $x_i$  using a sufficient statistic  $T(x_i)$  may result in substantial **data reduction**.

Data sets  $x_i$  and  $y_i$  for which the ratio of the likelihoods  $L(\theta|x_i)/L(\theta|y_i)$  does not depend on  $\theta$  (so the two likelihoods are proportional to each other by a constant) are called **likelihood equivalent** because a likelihood-based procedure to learn about  $\theta$  will draw identical conclusions from  $x_i$  and  $y_i$ . For data sets  $x_i, y_i \in \mathcal{X}_t$  equivalent with respect to a sufficient statistic  $T(x_i)$  it follows directly from the Fisher-Pearson factorisation that the ratio

$$L(\theta|x_i)/L(\theta|y_i) = k(x_i)/k(y_i)$$

and thus is constant with regard to  $\theta$ . Consequently, all **data sets in  $\mathcal{X}_t$  are also likelihood equivalent**. However, the converse is not true: depending on the sufficient statistics there usually will be many likelihood equivalent data sets that are not part of the same set  $\mathcal{X}_t$ .

Of particular interest is therefore to find those sufficient statistics that achieve the coarsest partitioning of the sample space and thus may allow the highest data reduction. Specifically, a **minimal sufficient statistic** is a sufficient statistic  $T(x_i)$  for which all likelihood equivalent data sets also are equivalent under  $T(x_i)$ . Therefore, to check whether a sufficient statistic  $T(x_i)$  is minimally sufficient we verify whether for any two likelihood equivalent data sets  $x_i$  and  $y_i$  it also follows that  $T(x_i) = T(y_i)$ . If this holds true then  $T(x_i)$  is a minimally sufficient statistic.

An equivalent non-operational definition is that a minimal sufficient statistic  $T(x_i)$  is a sufficient statistic that can be computed from any other sufficient statistic  $S(x_i)$ . This follows from the above directly: assume any sufficient statistic  $S(x_i)$ , this defines a corresponding set  $\mathcal{X}_s$  of likelihood equivalent data sets. By implication any  $x_i, y_i \in \mathcal{X}_s$  will necessarily also be in  $\mathcal{X}_t$ , thus whenever  $S(x_i) = S(y_i)$  we also have  $T(x_i) = T(y_i)$ , and therefore  $T(x_i)$  is a function of  $S(x_i)$ .

A trivial but **important example of a minimal sufficient statistic is the likelihood function itself** since by definition it can be computed from any set of

sufficient statistics. Thus the likelihood function  $L(\theta)$  captures all information about  $\theta$  that is available in the data. In other words, it provides an *optimal summary* of the observed data with regard to a model. Note that in Bayesian statistics (to be discussed in Part 2 of the module) the likelihood function is used as proxy/summary of the data.

**Example 6.1.** Sufficient statistics for the parameters of the normal distribution:

The normal model  $N(\mu, \sigma^2)$  with parameter vector  $\theta = (\mu, \sigma^2)^T$  and log-likelihood

$$l_n(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

One possible set of minimal sufficient statistics for  $\theta$  are  $\bar{x}$  and  $\bar{x}^2$ , and with these we can rewrite the log-likelihood function without any reference to the original data  $x_i$  as follows

$$l_n(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{n}{2\sigma^2} (\bar{x}^2 - 2\bar{x}\mu + \mu^2)$$

An alternative set of minimal sufficient statistics for  $\theta$  consists of  $s^2 = \bar{x}^2 - \bar{x}^2 = \hat{\sigma}_{ML}^2$  as and  $\bar{x} = \hat{\mu}_{ML}$ . The log-likelihood written in terms of  $s^2$  and  $\bar{x}$  is

$$l_n(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{n}{2\sigma^2} (s^2 + (\bar{x} - \mu)^2)$$

Note that in this example the dimension of the parameter vector  $\theta$  equals the dimension of the minimal sufficient statistic, and furthermore, that the MLEs of the parameters are in fact minimal sufficient!

The conclusion from Examples 6.1 holds true more generally: in the exponential family (which contains the normal distribution as special case) the MLEs of the natural parameters are minimal sufficient statistics. Thus, there will typically be substantial dimension reduction from the raw data to the sufficient statistics.

However, outside the exponential family the MLE is not necessarily a minimal sufficient statistic, and may not even be a sufficient statistic. This is because **a (minimal) sufficient statistic of the same dimension as the parameters does not always exist**. A classic example is the Cauchy distribution for which the minimal sufficient statistics are the ordered observations, thus the MLE of the parameters do not constitute sufficient statistics, let alone minimal sufficient statistics. However, the MLE is of course still a function of the minimal sufficient statistic.

In summary, the likelihood function acts as perfect data summariser (i.e. as minimally sufficient statistic), and in exponential families (e.g. Normal distribution) the MLEs of the parameters  $\hat{\theta}_{ML}$  are minimally sufficient.

Finally, while sufficiency is clearly a useful concept for data reduction one needs to keep in mind that this is always in reference to a specific model. Therefore,

unless one strongly believes in a certain model it is generally a good idea to keep (and not discard!) the original data.

## 6.3 Concluding remarks on maximum likelihood

### 6.3.1 Remark on KL divergence

Finding the model  $F_\theta$  that best approximates the underlying true model  $F_0$  is done by minimising the relative entropy  $D_{\text{KL}}(F_0, F_\theta)$ . For large sample size  $n$  we may approximate  $F_0$  by the empirical distribution  $\hat{F}_0$ , and minimising  $D_{\text{KL}}(\hat{F}_0, F_\theta)$  then yields the method of maximum likelihood.

However, since the KL divergence is not symmetric there are in fact two ways to minimise the divergence between a fixed  $F_0$  and the optimised  $F_\theta$ , each with different properties:

- a) **forward KL, approximation KL**:  $\min_\theta D_{\text{KL}}(F_0, F_\theta)$

This is also called an “M (Moment) projection”. It has a **zero avoiding** property:  $f_\theta(x) > 0$  whenever  $f_0(x) > 0$

- b) **reverse KL, inference KL**:  $\min_\theta D_{\text{KL}}(F_\theta, F_0)$

This is also called an “I (Information) projection”. It has a **zero forcing** property:  $f_\theta(x) = 0$  whenever  $f_0(x) = 0$

Maximum likelihood is based on “forward KL”, whereas Bayesian updating and Variational Bayes approximations use “reverse KL”.

### 6.3.2 What happens if $n$ is small?

From the long list of optimality properties of ML it is clear that for large sample size  $n$  the best estimator will typically be the MLE.

However, for **small sample size it is indeed possible (and necessary) to improve over the MLE** (e.g. via Bayesian estimation or regularisation). Some of these ideas will be discussed in Part II.

- Likelihood will *overfit*!

Alternative methods need to be used:

- regularised/penalised likelihood
- Bayesian methods

which are essentially two sides of the same coin.

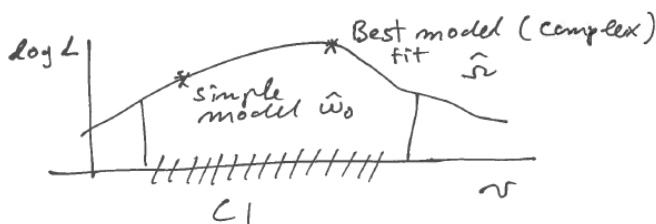
Classic example of a simple non-ML estimator that is better than the MLE: **Stein’s example / Stein paradox** (C. Stein, 1955):

- Problem setting: estimation of the mean in multivariate case

- Maximum likelihood estimation breaks down!  $\rightarrow$  average (=MLE) is worse in terms of MSE than Stein estimator.
- For small  $n$  the asymptotic distributions for the MLE and for the LRT are not accurate, so for inference in these situations the distributions may need to be obtained by simulation (e.g. parametric or nonparametric bootstrap).

### 6.3.3 Model selection

- CI are sets of models that are not statistically distinguishable from the best ML model
- in doubt, choose the simplest model compatible with data
- better prediction, avoids overfitting
- Useful for model exploration and model building.



- Note that, by construction, the model with more parameters always has a higher likelihood, implying likelihood favours complex models
- Complex model may overfit!
- For comparison of models penalised likelihood or Bayesian approaches may be necessary
- Model selection in small samples and high dimension is challenging
- Recall that the aim in statistics is **not** about rejecting models (this is easy as for large sample size any model will be rejected!)
- Instead, the aim is model building, i.e. to find a model that **explains the data well** and that **predicts well**!
- Typically, this will not be the best-fit ML model, but rather a simpler model that is close enough to the best / most complex model.

**Part II**

**Bayesian Statistics**





# Chapter 7

## Essentials of Bayesian statistics

### 7.1 Conditional probability

Assume we have two random variables  $x$  and  $y$  with a **joint density** (or joint PMF)  $p(x, y)$ . By definition  $\int \int_{x,y} p(x, y) dx dy = 1$ .

The **marginal densities** for the individual  $x$  and  $y$  are given by  $p(x) = \int_y p(x, y) dy$  and  $p(y) = \int_x p(x, y) dx$ . Thus, when computing the marginal densities a variable is removed from the joint density by integrating over all possible states of that variable. It follows also that  $\int_x p(x) dx = 1$  and  $\int_y p(y) dy = 1$ , i.e. the marginal densities also integrate to 1.

As alternative to integrating out a random variable in the joint density  $p(x, y)$  we may wish to keep it fixed at some value, say keep  $y$  fixed at  $y_0$ . In this case  $p(x, y = y_0)$  is proportional to the **conditional density** (or PMF) given by the ratio

$$p(x|y = y_0) = \frac{p(x, y = y_0)}{p(y = y_0)}$$

The denominator  $p(y = y_0) = \int_x p(x, y = y_0) dx$  is needed to ensure that  $\int_x p(x|y = y_0) dx = 1$ , thus it renormalises  $p(x, y = y_0)$  so that it is a proper density.

To simplify notation, the specific value on which a variable is conditioned is often left out.

The mean and variance of the conditional distribution are called conditional mean and conditional variance.

Rearranging the above we see that the joint density can be written as the product of marginal and conditional density in two different ways:

$$p(x, y) = p(x|y)f(y) = p(y|x)p(x)$$

## 7.2 Bayes' theorem

Bayesian statistical learning is linked with the name of [Thomas Bayes \(1701-1761\)](#) who was the first to state [Bayes' theorem \(1763\)](#) on conditional probability. Interestingly, this work published only after Bayes' death by [Richard Price \(1723-1791\)](#):

$$p(A|B) = p(B|A) \frac{p(A)}{p(B)}$$

This theorem relates the two possible conditional densities (or conditional probability mass functions) for two events  $A$  and  $B$ .

It follows directly from the product rule linking the joint density with the marginal and conditional densities.

## 7.3 Principle of Bayesian learning

Ingredients:

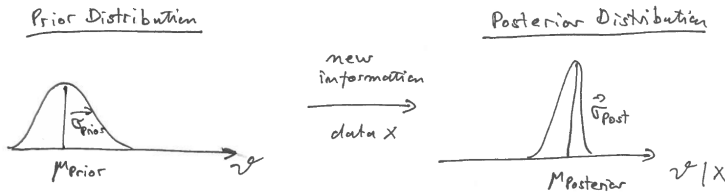
- $\theta$  parameter of interest, unknown and fixed.
- prior distribution with density  $p(\theta)$  describing the *uncertainty* (not randomness!) about  $\theta$
- data generating process  $p(x|\theta)$  (likelihood!)

Question: new information in the form of new observation  $x$  arrives - how does the uncertainty about  $\theta$  change?

Answer: use Bayes' theorem to **update prior distribution to posterior distribution**.

$$\underbrace{p(\theta|x)}_{\text{posterior}} = \frac{p(x|\theta)}{p(x)} \underbrace{p(\theta)}_{\text{prior}}$$

Note that this update procedure can be repeated again and again: we can use the posterior as our new prior and then update it with further data.



For the denominator in Bayes formula we need to compute  $p(x)$ . This is obtained by

$$\begin{aligned}
 p(x) &= E_{F_\theta} p(x|\theta) \\
 &= \int_{\theta} p(x|\theta) p(\theta) d\theta \\
 &= \int_{\theta} p(x, \theta) d\theta
 \end{aligned}$$

i.e. by marginalisation of the parameter  $\theta$  from the joint distribution of  $\theta$  and  $x$ . (For discrete  $\theta$  replace the integral by a sum).

Depending on the context this quantity is either called *marginal likelihood* (of the underlying model) or *prior predictive distribution* (for the data).

Intriguingly, to conduct a Bayesian statistical analysis typically require integration and/or averaging (e.g. to compute the marginal likelihood), in contrast to maximum likelihood that requires optimisation (to find the maximum likelihood).

## 7.4 What is exactly is the “Bayesian estimate”?

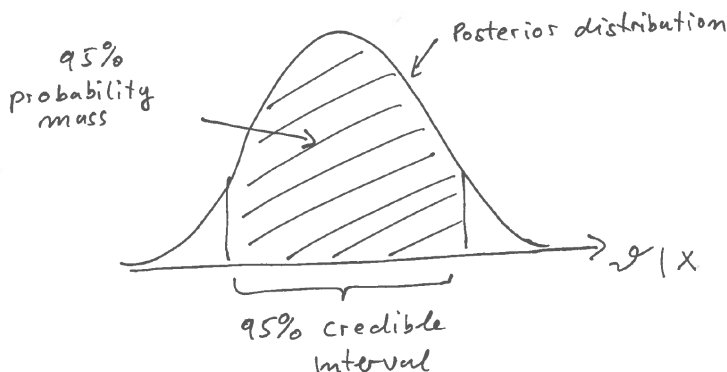
**The Bayesian estimate is the full complete posterior distribution!**

However, it is useful to summarise aspects of the posterior distribution:

- Posterior mean  $E(\theta|x)$
- Posterior variance  $\text{Var}(\theta|x)$
- Posterior mode etc.

In particular the mean of the posterior distribution is often taken as a *Bayesian point estimate*.

The posterior distribution also allows to define **credible regions** or **credible intervals**. These are the **Bayesian equivalent to confidence intervals** and are constructed by finding the areas of highest probability mass (say 95%) in the posterior distribution.



Bayesian credible intervals (unlike their frequentist confidence counterparts) are thus very easy to interpret - they simply correspond to the area in the parameter space in which we can find the parameter with a given specified probability. In contrast, in frequentist statistics it does not make sense to assign a probability to a parameter value!

Note that there are typically many credible intervals with the given specified coverage  $\alpha$  (say 95%). Therefore, we may need further criteria to construct these intervals.

In the univariate case a **two-sided equal-tail credible interval** is obtained by finding the corresponding lower  $1 - \alpha/2$  and upper  $\alpha/2$  quantiles.

A **highest posterior density (HPD)** interval of coverage  $\alpha$  is found by identifying the shortest interval (i.e. with smallest support) for the given  $\alpha$  probability mass. Any point within such an interval has higher density resp. probability than outside the credible interval, and the density / probability at the boundaries are all equal. Thus a Bayesian HPD credible interval is constructed similar like a likelihood based confidence interval. When the posterior has multiple modes this means the the HPD interval may be disjoint.

## 7.5 Computer implementation of Bayesian learning

As we have seen Bayesian learning is *conceptually straightforward*:

- 1) Specify prior uncertainty  $p(\theta)$  about the parameters of interest  $\theta$ .
- 2) Specify the data generating process for a specified parameter:  $p(x|\theta)$ .
- 3) Apply Bayes' theorem to update prior uncertainty in the light of the new data.

In practise, however, computing the posterior distribution can be *computationally very demanding*, especially for complex models.

For this reason specialised software packages have been developed for computational Bayesian modelling, for example:

- Bayesian statistics in R: <https://cran.r-project.org/web/views/Bayesian.html>
- Stan probabilistic programming language (can be used with R and Python) — <https://mc-stan.org/>
- Bayesian statistics in Python: [PyMC3](#) using Theano, [Pyro](#) using PyTorch, [NumPyro](#) using JAX, [TensorFlow Probability](#) using Tensorflow
- Bayesian statistics in Julia: [Turing.jl](#)
- Bayesian hierarchical modelling with [BUGS](#), [JAGS](#) and [NIMBLE](#).

In addition to numerical procedures to sample from the posterior distribution there are also many procedures aiming to approximate the Bayesian posterior, employing the Laplace approximation, integrated nested Laplace approximation (INLA), variational Bayes etc.

## 7.6 Bayesian interpretation of probability

### 7.6.1 What makes you “Bayesian”?

If you use Bayes’ theorem are you therefore automatically a Bayesian? No!!

Bayes’ theorem is a mathematical fact from probability theory. Hence, Bayes’ theorem is valid for everyone, whichever form for statistical learning you are subscribing (such as frequentist ideas, likelihood methods, entropy learning, Bayesian learning).

As we discuss now the key difference between Bayesian and frequentist statistical learning lies in the differences in *interpretation of probability*, not in the mathematical formalism for probability (which includes Bayes’ theorem).

### 7.6.2 Mathematics of probability

The mathematics of probability in its modern foundation was developed by [Andrey Kolmogorov \(1903–1987\)](#). In this book [Foundations of the Theory of Probability \(1933\)](#) he establishes probability in terms of set theory/ measure theory. This theory provides a coherent mathematical framework to work with probabilities.

However, Kolmogorov’s theory does *not* provide an interpretation of probability!

→ The Kolmogorov framework is the basis for both the frequentist and the Bayesian interpretation of probability.

### 7.6.3 Interpretations of probability

Essentially, there are two major commonly used interpretation of probability in statistics - the **frequentist interpretation** and the **Bayesian interpretation**.

**A: Frequentist interpretation**

probability = frequency (of an event in a long-running series of identically repeated experiments)

This is the *ontological view* of probability (i.e. probability “exists” and is identical to something that can be observed.).

It is also a very restrictive view of probability. For example, frequentist probability cannot be used to describe events that occur only a single time. Frequentist probability thus can only be applied asymptotically, for large samples!

### **B: Bayesian probability**

“Probability does not exist” — famous quote by [Bruno de Finetti \(1906–1985\)](#), a Bayesian statistician.

What does this mean?

Probability is a **description of the state of knowledge** and of **uncertainty**.

Probability is thus an *epistemological quantity* that is assigned and that changes rather than something that is an inherent property of an object.

Note that this does not require any repeated experiments. The Bayesian interpretation of probability is valid regardless of sample size or the number or repetitions of an experiment.

**Hence, the key difference between frequentist and Bayesian approaches is not the use of Bayes’ theorem. Rather it is whether you consider probability as ontological (frequentist) or epistemological entity (Bayesian).**

## **7.7 Historical developments**

- Thomas Bayes (1701-1761) the father of Bayesian statistics  
Only after his death his paper on Bayes’ theorem was published (1763).
- Laplace (from 1800) was actually the first to use Bayes’ theorem for statistical calculations. This activity was then called “inverse probability”.
- Between 1900 and 1940 classical mathematical statistics was developed and the field was heavily influenced and dominated by R.A. Fisher (who invented likelihood theory and ANOVA, among other things - he also was working in population genetics). Fisher himself was very much opposed to Bayesian theory.
- 1931 de Finetti publishes his “representation theorem”. This shows that the joint distribution of a sequence of exchangeable events (i.e. where the ordering can be permuted) can be represented by a mixture distribution that can be constructed via Bayes’ theorem. (Note that exchangeability is a weaker condition than i.i.d.) This theorem is often used as a justification Bayesian statistics (along with the so-called Dutch book argument, also by de Finetti).

- 1933 publication of Kolmogorov's book on probability theory.
- 1946 Cox theorem ([Richard T. Cox \(1898–1991\)](#)): the aim to generalise classical logic (from TRUE/FALSE to continuous measures of uncertainty) inevitably leads to probability theory and Bayesian learning! This justification of Bayesian statistics was later popularised by [Edwin T. Jaynes \(1922–1998\)](#) in various books (1959, 2003).
- 1955 Stein Paradox - [Charles M. Stein \(1920–2016\)](#) publishes paper on the Stein estimator - an estimator of the mean that dominates ML estimator. His estimator is always better in terms of MSE than the ML estimator, and this was very puzzling at that time!

From 1970 onwards Bayesian learning has become more pervasive!

- Computers allow to do the complex computations needed in Bayesian statistics
- Metropolis-Hastings algorithm published
- A lot of work on interpreting Stein estimators as empirical Bayes estimators (Efron and Morris 1975) and on development of regularised estimation techniques such as penalised likelihood in regression (e.g. ridge regression)
- regularisation originally was only meant to make singular systems/matrices invertible - but then it turned out regularisation has a simple Bayesian interpretation!
- work on reference priors (Bernado 1979)
- penalised likelihood via KL divergence for model selection (Akaike 1973)

Another boost was in the 1990/2000s when in science (e.g. genomics) many complex and high-dimensional data set were becoming widely available. Classical statistical methods cannot be used in this setting (overfitting!) so many new methods were developed for high-dimensional data analysis, many with direct link to Bayesian statistics:

- 1996 lasso regression (Tibshirani)
- Machine learning etc (many Bayesians in this field!)

## 7.8 Connection with entropy learning

### 7.8.1 Zero forcing property

It is easy to see that if in Bayes rule the prior probability for an event is set to 0, then the posterior probability for that event will remain at 0, regardless of the data! This **zero-forcing property** of the Bayes update rule has been called **Cromwell's rule** by D. Lindley. Therefore, assigning prior probability 0 to an event should be avoided.

Note that this implies that assigning prior probability 1 to an event should be avoided, too, since this means assigning 0 to all other alternative events.

## 7.8.2 Connection with entropy learning

The *Bayesian update rule* is a very general form of learning when the *new information arrives in the form of data*.

But actually there is an even a more general principle: the **principle of minimal information update** (e.g. Jaynes 1959, 2003) or **principle of minimum information discrimination (MDI)** (Kullback 1959):

- **Change your beliefs only as much as necessary to be coherent with new evidence!**

This is also called **entropy learning** since the KL divergence ( $F_{\theta|\text{new information}}; F_{\theta}$ ) is employed to measure the divergence from the updated distribution to the distribution prior to the arrival of the information.

Note that this update is based on an *I*-projection (see Part I, Likelihood), which also does have the zero forcing property (hinting that Bayes rule is a special case).

Thus, when new information arrives then the uncertainty about the parameter is only minimally adjusted, just as much as needed to account for the new information (“inertia of beliefs”).

There are three main special cases that follow from the entropy learning rule:

- 1) if information arrives in form of data  $\rightarrow$  update by T. Bayes’ theorem (1763)
- 2) if information is in the form of another distribution  $\rightarrow$  update using R. Jeffrey’s rule (1965)
- 3) if the information is in form of constraints  $\rightarrow$  Kullback’s principle of minimum MDI (1959), E. T. Jaynes MaxEnt principle (1957)

Since 1) is by far the most common situation it is clear why it is important to study Bayesian learning!

This shows (again) how fundamentally important KL divergence is in statistics - it not only leads to likelihood inference but also to Bayesian learning, as well as to other forms of information updating! Furthermore, relative entropy is useful to choose priors (e.g. reference priors) and in experimental design.



## Chapter 8

# Beta-Binomial model for estimating a proportion

In this chapter we discuss how to estimate a portion in the Bayesian framework.

### 8.1 Binomial likelihood

In order to apply Bayes' theorem we first need to find a suitable likelihood based on modeling the data generating process. Here we follow the Binomial model as used previously in Part I:

Repeated Bernoulli experiment (Binomial model):

$x \in \{0, 1\}$  (e.g. "tails" vs. "heads")

probability mass function (pmf):  $\Pr(x = 1) = p, \Pr(x = 0) = 1 - p$

Mean:  $E(x) = p$

Variance  $\text{Var}(x) = p(1 - p)$

$\text{Bin}(n, p)$  (sum of  $n$  Bernoulli experiments)

$x \in \{0, 1, \dots, n\}$

Mean:  $E(x) = np$

Variance:  $\text{Var}(x) = np(1 - p)$

Standardised Binomial (average of  $n$  Bernoulli experiments):

$\frac{x}{n} \in \{0, \frac{1}{n}, \dots, 1\}$

Mean:  $E(\frac{x}{n}) = p$

Variance:  $\text{Var}(\frac{x}{n}) = \frac{p(1-p)}{n}$

From part I (likelihood theory) we know that the *maximum likelihood estimate* of the proportion is the frequency  $\hat{p}_{ML} = \frac{x}{n}$  given  $x$  (number of "heads") is

observed in  $n$  repeats.

## 8.2 Excursion: Properties of the Beta distribution

The density of the Beta distribution  $\text{Beta}(\alpha, \beta)$  for  $x \in [0, 1]$  and  $\alpha > 0$  and  $\beta > 0$  is

$$f(x|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

The mean is  $E(x) = \mu = \frac{\alpha}{\alpha+\beta}$  and the variance  $\text{Var}(x) = \frac{\mu(1-\mu)}{\alpha+\beta+1}$ .

The density depends on the Beta function  $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$  which in turn is defined via Euler's Gamma function

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$$

Note that  $\Gamma(x) = (x-1)!$  for any positive integer  $x$

A useful reparameterisation of the Beta distribution is in terms of the parameters  $\mu \in [0, 1]$  and  $m > 0$ , yielding the original parameters via  $\alpha = \mu m$  and  $\beta = (1-\mu)m$ . Conversely,  $m = \alpha + \beta$  and  $\mu = \frac{\alpha}{\alpha+\beta}$ .

The Beta distribution is very flexible and can assume a number of different shapes, depending on the value of  $\alpha$  and  $\beta$ :



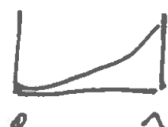
$\alpha = \beta = 1$  "uniform"  
"flat"



$\alpha > 1$   
 $\beta > 1$  "hill"



$\alpha < 1$   
 $\beta < 1$  "U shape"



$\alpha > 1$   
 $\beta < 1$  "slope"

## 8.3 Beta prior distribution

In Bayesian learning we need to make explicit our uncertainty about  $p$ .

$p$  has support  $[0, 1] \rightarrow$  we use the **Beta distribution**  $\text{Beta}(\alpha, \beta)$  as prior for  $p$  with parameters  $\alpha \geq 0$  and  $\beta \geq 0$ :

$$p \sim \text{Beta}(\alpha, \beta)$$

Note this does not actually mean that  $p$  is random! It only means that we model the uncertainty about  $p$  using a Beta random variable!

The flexibility of the Beta distribution allows to accomodate a large variety of possible scenarios for our prior knowledge.

The prior mean is

$$E(p) = \frac{\alpha}{m} = \mu_{\text{prior}}$$

and the prior variance

$$\text{Var}(p) = \frac{\mu_{\text{prior}}(1 - \mu_{\text{prior}})}{m + 1}$$

where  $m = \alpha + \beta$ .

Note the similarity to the moments of the standardised Binomial above!

## 8.4 Computing the posterior distribution

Bayes' theorem for continuous random variables to compute posterior density:

$$f(p|x) = \frac{f(x|p)f(p)}{\int_{p'} f(x|p')f(p')dp'}$$

We use in our analysis the Beta-Binomial model:

a) **Beta prior:**

$$p \sim \text{Beta}(\alpha, \beta)$$

$$f(p) = \frac{1}{B(\alpha, \beta)} p^{\alpha-1} (1-p)^{\beta-1}$$

b) **Binomial likelihood:**

$$x|p \sim \text{Bin}(n, p)$$

$$f(x|p) = \binom{n}{x} p^x (1-p)^{(n-x)}$$

Applying Bayes' theorem results in

c) **Beta posterior distribution**

$$p|x \sim \text{Beta}(\alpha + x, \beta + n - x)$$

$$f(p|x) = \frac{1}{B(\alpha + x, \beta + n - x)} p^{\alpha+x-1} (1-p)^{\beta+n-x-1}$$

(for a proof see Worksheet 5!)

The posterior can be summarised by its first two moments (mean and variance):

Posterior mean:

$$\mu_{\text{posterior}} = E(p|x) = \frac{x + \alpha}{n + m}$$

Posterior variance:

$$\sigma_{\text{posterior}}^2 = \text{Var}(p|x) = \frac{\mu_{\text{posterior}}(1 - \mu_{\text{posterior}})}{n + m + 1}$$

## Chapter 9

# Properties of Bayesian learning

The Beta-Binomial models allows to observe a number of intriguing features and properties of Bayesian learning. Many of these extend also to other models as we will see later.

### 9.1 Prior acting as pseudo-data

In the expression for the posterior mean and variance you can see that  $m = \alpha + \beta$  behaves like an implicit sample size connected with prior information!

Specifically,  $\alpha$  and  $\beta$  act as **pseudo-counts** that influence both the posterior mean and the posterior variance, exactly in the same way as coventional data.

For example, larger  $m$  (and thus  $\alpha$  and  $\beta$ ) the smaller is the posterior variance, with variance decreasing proportional to the inverse of  $m$ . If the prior is highly concentrated, i.e. if it has low variance and large precision (=inverse variance) then the implicit data size  $m$  is large. Conversely, if the prior has a large variance, then the prior is vague and the implicit data size  $m$  is small.

Hence, a prior has the same effect as if one would add data – but without actually adding data! This is precisely this why a prior acts as a regulariser and prevents overfitting, because it increases effective sample size.

Another interpretation is that any prior summarises data that may have been available previously as observations.

## 9.2 Linear shrinkage of mean

The posterior mean  $\mu_{\text{posterior}}$  is a linearly adjusted  $\hat{\mu}_{ML}$ . This becomes evident by writing  $\mu_{\text{posterior}}$  as

$$\mu_{\text{posterior}} = \lambda \mu_{\text{prior}} + (1 - \lambda) \hat{\mu}_{ML}$$

with weight  $\lambda \in [0, 1]$

$$\lambda = \frac{m}{m + n}.$$

**The posterior mean is a convex combination (i.e. the weighted average) of the ML estimate and the prior mean.** The factor  $\lambda$  is called the **shrinkage intensity** — note that it is the ratio of the “prior sample size” ( $m$ ) and the “effective overall sample size” ( $m + n$ ).

1. This is called *shrinkage*, because the ML estimator is “shrunk” towards the prior mean (which is often called the “target”, and sometimes the target is zero, and then the terminology “shrinking” makes most sense).
2. If the shrinkage intensity is zero ( $\lambda = 0$ ) then the ML point estimator is recovered. This implies  $\alpha = 0$  and  $\beta = 0$ , or  $n \rightarrow \infty$ .

Note that using maximum likelihood to estimate of the proportion  $p$  (for moderate or small  $n$ ) is the same as Bayesian estimation using the Beta-Binomial model with prior  $\alpha = 0$  and  $\beta = 0$ . This prior is extremely “u-shaped” and the implicit prior for the ML estimation. (Would you would use such a prior intentionally?)

3. If the shrinkage intensity is large ( $\lambda \rightarrow 1$ ) then the posterior mean corresponds to the prior. This happens if  $n = 0$  or if  $m$  is very large (implying that the prior is sharply concentrated around the prior mean).
4. Since the ML estimate (=frequency) is unbiased the Bayesian point estimate is biased (for finite  $n$ )! And the bias is in fact the prior mean! So Bayesian statistics produces by default biased estimators (but asymptotically they will be unbiased like in ML).
5. That the posterior mean is a linear combination of the ML estimate and the prior mean is not a coincidence. In fact, this is true for all distributions in the exponential family (see e.g. Diaconis and Ylvisaker, 1979). Furthermore, it is possible (and indeed quite useful for computational reasons!) to formulate Bayes theory completely in terms of linear shrinkage (e.g. Hartigan 1969). The resulting theory is called “Bayes linear statistics” (Goldstein and Wooff, 2007).

## 9.3 Conjugacy of prior and posterior distribution

In the Beta-Binomial model for estimating the proportion  $p$  the choice of the **Beta distribution as prior distribution** along with the Binomial likelihood resulted in having the **Beta distribution as posterior distribution** as well.

If the prior and posterior belong to the same distributional family the prior is called a **conjugate prior**. This will be the case if the prior has the same functional form as the likelihood.

In the Beta-Binomial the likelihood is based on the Binomial distribution and has the following form (only terms depending on the parameter  $p$  are shown):

$$p^x(1-p)^{n-x}$$

The form of the Beta prior is (again, only showing terms depending on  $p$ ):

$$p^{\alpha-1}(1-p)^{\beta-1}$$

Since the posterior is proportional to the product of prior and likelihood the posterior will have exactly the same form as the prior:

$$p^{\alpha+x-1}(1-p)^{\beta+n-x-1}$$

Choosing the prior distribution from a family conjugate to the likelihood greatly simplifies Bayesian analysis since the Bayes formula can then be written in form of an update formula for the parameters of the Beta distribution:

$$\alpha \rightarrow \alpha + x$$

$$\beta \rightarrow \beta + n - x$$

Thus, conjugate prior distributions are very convenient choices. However, in their application it must be ensured that the prior distribution is flexible enough to encapsulate all prior information that may be available. In cases where this is not the case alternative priors should be used (and most likely this will then require to compute the posterior distribution numerically rather than analytically).

## 9.4 Large sample asymptotics

### 9.4.1 Large sample limits of mean and variance

If  $n$  is large and  $n \gg \alpha, \beta$  the posterior mean and variance become asymptotically

$$\mu_{\text{posterior}} \stackrel{a}{=} \frac{x}{n} = \hat{\mu}_{ML}$$

and

$$\sigma_{\text{posterior}}^2 \stackrel{a}{=} \frac{\hat{\mu}_{ML}(1 - \hat{\mu}_{ML})}{n}$$

Thus, if sample size is large the Bayes' estimator turns into the ML estimator! Specifically, the posterior mean becomes the ML point estimate, and the posterior variance is equal to the asymptotic variance computed via the observed Fisher information!

Thus, for large  $n$  the data dominate and any details about the prior (such as values of  $\alpha$  and  $\beta$  become irrelevant!

## 9.4.2 Asymptotic Normality of the Posterior distribution

Also known as **Bayesian Central Limit Theorem (CLT)**.

Under some regularity conditions (such as regular likelihood and positive prior probability for all parameter values, finite number of parameters, etc.) for large sample size the Bayesian posterior distribution converges to a Normal distribution centered around the MLE and with the variance of the MLE:

$$\text{for large } n: p(\theta|x_1, x_2, \dots, x_n) \rightarrow N(\hat{\theta}_{ML}, \text{Var}(\hat{\theta}_{ML}))$$

So not only are the posterior mean and variance converging to the MLE and the variance of the MLE for large sample size, but also the posterior distribution itself converges to the sampling distribution!

This holds generally in many regular cases, not just in our example of the Beta-Bernoulli model.

The Bayesian CLT is generally known as the **Bernstein-van Mises theorem** (who discovered it at around 1920-30), but special cases were already known as by Laplace.

In the Worksheet 5 the asymptotic convergence of the posterior distribution to a normal distribution is demonstrated graphically.

## 9.5 Posterior variance for finite $n$

In the previous chapter we have derived a Bayesian point estimate for the proportion  $p$  as the posterior mean

$$E(p|x) = \frac{x + \alpha}{n + m} = \hat{p}_{\text{Bayes}}$$

with posterior variance

$$\text{Var}(p|x) = \frac{\hat{p}_{\text{Bayes}}(1 - \hat{p}_{\text{Bayes}})}{n + m + 1}$$



Asymptotically, we have seen that for large  $n$  the posterior becomes the ML estimator, and the posterior variance becomes the asymptotic variance of the MLE. Thus, the Bayesian estimate will be indistinguishable from the MLE for large  $n$  and shares its favourable properties.

In addition, for finite sample size the posterior variance will typically be *smaller* than both the asymptotic posterior variance (for large  $n$ ) and the prior variance, showing that combining the information in the prior and in the data leads to a more efficient estimate.



## Chapter 10

# Normal-Normal and Inverse-Gamma-Normal models for estimating the mean and the variance

### 10.1 Normal-Normal model to estimate mean

#### 10.1.1 Normal likelihood

For the **likelihood** we assume as data-generating model the normal distribution with known fixed variance  $\sigma^2$

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

This yields as the MLE  $\hat{\mu}_{ML} = \bar{x}$ .

#### 10.1.2 Normal prior distribution

To model the uncertainty about  $\mu$  we use the normal distribution  $N(\mu, \sigma^2/k)$  parameterised by the two parameters  $\mu$  and  $k$  (remember  $\sigma^2$  is fixed).

With  $\mu = \mu_0$  and  $k = m$  we get the **normal prior**

$$\mu \sim N(\mu_0, \sigma^2/m)$$

with prior mean  $E(\mu) = \mu_0$  and prior variance  $\text{Var}(\mu) = \frac{\sigma^2}{m}$  where  $m$  is the implied sample size from the prior. Note that  $m$  does not need to be an integer value!

### 10.1.3 Normal posterior distribution

The **posterior distribution** after observing  $n$  samples  $x_1, \dots, x_n$  is normal with  $\mu = \mu_1$  and  $k = m + n$

$$\mu | x_1, \dots, x_n \sim N(\mu_1, \sigma^2 / (m + n))$$

with posterior mean

$$E(\mu | x_1, \dots, x_n) = \mu_1 = \frac{m\mu_0 + n\bar{x}}{n + m} = \lambda\mu_0 + (1 - \lambda)\hat{\mu}_{ML}$$

with  $\lambda = \frac{m}{n+m}$ . Note the linear shrinkage of  $\hat{\mu}_{ML}$  towards  $\mu_0$ !

The corresponding posterior variance is

$$\text{Var}(\mu | x_1, \dots, x_n) = \frac{\sigma^2}{n + m}$$

Thus, the **normal distribution is the conjugate distribution to the mean parameter in the normal likelihood**.

### 10.1.4 Large sample asymptotics and Stein paradox

For  $n$  large and  $n \gg m$  we get

$$E(\mu | x_1, \dots, x_n) \stackrel{a}{=} \hat{\mu}_{ML}$$

$$\text{Var}(\mu | x_1, \dots, x_n) \stackrel{a}{=} \frac{\sigma^2}{n}$$

i.e. the MLE and its asymptotic variance!

Note that the posterior variance  $\frac{\sigma^2}{n+m}$  is smaller than the asymptotic variance  $\frac{\sigma^2}{n}$  and the prior variance  $\frac{\sigma^2}{m}$ .

When studying the frequentist properties of the posterior mean  $\mu_1$  it turns out that by an appropriate choice of  $m$  (or  $\lambda$ ) it is possible to construct an estimator that will outperform the MLE for finite  $n$  in terms of MSE (with the reduced variance compensating for the increase in bias)! Charles Stein was one of the first to present such an estimator (see next chapter), and by many of his contemporaries it was considered very puzzling to have any estimator outperform the MLE, hence this effect is called **Stein paradox**.

## 10.2 Inverse-Gamma-Normal model to estimate variance

### 10.2.1 Inverse Gamma distribution

Next, we study a common Bayesian model for estimating the variance parameter of the normal distribution. For this we use the inverse Gamma distribution:

$$x \sim \text{Inv-Gam}(\alpha, \beta)$$

This distribution is closely linked with the Gamma distribution — the inverse of  $x$  is Gamma-distributed with inverted scale parameter:

$$\frac{1}{x} \sim \text{Gam}(\alpha, \beta^{-1})$$

For use as prior and posterior we employ a different parameterisation with  $k = 2(\alpha - 1)$  and  $v = \beta / (\alpha - 1)$ :

$$x \sim \text{Inv-Gam}(1 + \frac{k}{2}, \frac{k}{2}v)$$

The first two moments of the inverse Gamma distribution are

$$E(x) = \frac{\beta}{\alpha - 1} = v$$

and

$$\text{Var}(x) = \frac{\beta^2}{(\alpha - 1)^2(\alpha - 2)} = \frac{2v^2}{k - 2}$$

### 10.2.2 Normal likelihood

As data likelihood / generating model we use normal distribution  $N(\mu, \sigma^2)$  with given fixed mean  $\mu$ .

This yields as MLE  $\hat{\sigma}_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$

### 10.2.3 Inverse Gamma prior distribution

For the prior distribution we use the inverse Gamma distribution with  $k = m$  and  $v = \sigma_0^2$

$$\sigma^2 \sim \text{Inv-Gam}(k = m, v = \sigma_0^2)$$

The corresponding prior mean is

$$E(\sigma^2) = \sigma_0^2$$

and the prior variance is

$$\text{Var}(\sigma^2) = \frac{2\sigma_0^4}{m - 2}$$

(note that  $m > 2$ )

### 10.2.4 Inverse Gamma posterior distribution

As the inverse Gamma distribution is conjugate to the normal likelihood the posterior distribution is inverse Gamma as well:

$$\sigma^2 | x_1, \dots, x_n \sim \text{Inv-Gam}(k = m + n, v = \sigma_1^2)$$

$$\text{with } \sigma_1^2 = \frac{\sigma_0^2 m + n \hat{\sigma}_{ML}^2}{m + n}.$$

The posterior mean is

$$E(\sigma^2 | x_1, \dots, x_n) = \sigma_1^2$$

and the posterior variance

$$\text{Var}(\sigma^2 | x_1, \dots, x_n) = \frac{2\sigma_1^4}{m + n - 2}$$

The update formula for the posterior mean of the variance follows the usual linear shrinkage rule:

$$\sigma_1^2 = \lambda \sigma_0^2 + (1 - \lambda) \hat{\sigma}_{ML}^2$$

$$\text{with } \lambda = \frac{m}{m + n}.$$

### 10.2.5 Large sample asymptotics

For  $n$  large and  $n \gg m$  we get

$$E(\sigma^2 | x_1, \dots, x_n) \stackrel{a}{=} \hat{\sigma}_{ML}^2$$

$$\text{Var}(\sigma^2 | x_1, \dots, x_n) \stackrel{a}{=} \frac{2\sigma^4}{n}$$

which is indeed the MLE of  $\sigma^2$  and its asymptotic variance!

### 10.2.6 Estimating precision

Instead of estimating the variance it is actually a bit simpler to estimate the precision (i.e. the inverse variance). For this one would then use a Gamma prior and a normal likelihood, resulting in a Gamma posterior.

### 10.2.7 Joint estimation of mean and variance

It is possible to combine the Normal-Normal for the mean and the Inverse-Gamma-Normal model into a joint model for the mean and variance.

This implies having a joint prior and a joint posterior for  $\mu$  and  $\sigma^2$ .

Details are not shown here but the resulting joint point estimators are identical to the above individual estimators.

# Chapter 11

## Shrinkage estimation using empirical risk minimisation

### 11.1 Linear shrinkage

In the examples for Bayesian estimation we have seen so far the posterior mean of the parameter of interest was obtained by linear shrinkage

$$\hat{\theta}_{\text{shrink}} = E(\theta|x_1, \dots, x_n) = \lambda\theta_0 + (1 - \lambda)\hat{\theta}_{\text{ML}}$$

of the MLE  $\hat{\theta}_{\text{ML}}$  towards the prior mean  $\theta_0$ , with shrinkage intensity  $\lambda = \frac{m}{m+n}$  determined by the pseudo-sample size  $m$  (which in turn is linked the precision of the prior) and the sample size  $n$ .

The resulting point estimate  $\hat{\theta}_{\text{shrink}}$  is called *shrinkage estimate* and is a convex combination of  $\theta_0$  and  $\hat{\theta}_{\text{ML}}$ . The prior mean  $\theta_0$  is also called the “target”.

In a Bayesian estimation the parameter  $m$  and hence  $\lambda$  is given a priori, but it turns out that it is possible and useful to find an optimal value for  $\lambda$  by minimising the mean squared error of the estimator  $\hat{\theta}_{\text{shrink}}$ .

In particular, by construction, the target  $\theta_0$  has zero variance but substantial bias, whereas the MLE  $\hat{\theta}_{\text{ML}}$  will have low or zero bias but a non-vanishing variance. By combining these two estimators with opposite properties the aim is to achieve a *bias-variance tradeoff* so that the resulting estimator  $\hat{\theta}_{\text{shrink}}$  has lower MSE than either  $\theta_0$  and  $\hat{\theta}_{\text{ML}}$ .

Specifically, the aim is to find

$$\lambda^* = \arg \min_{\lambda} E \left( (\theta - \hat{\theta}_{\text{shrink}})^2 \right)$$

It turns out that this can be minimised without knowing the actual true value of  $\theta$  and the result for an unbiased  $\hat{\theta}_{\text{ML}}$  is

$$\lambda^* = \frac{\text{Var}(\hat{\theta}_{\text{ML}})}{\text{E}((\hat{\theta}_{\text{ML}} - \theta_0)^2)}$$

Hence, the shrinkage intensity will be small if the variance of the MLE is small and/or if the target and the MLE differ substantially. On the other hand, if the variance of the MLE is large and/or the target is close to the MLE the shrinkage intensity will be large.

## 11.2 James-Stein estimator

We can now use empirical risk minimisation to estimate the shrinkage parameter the Normal-Normal model.

In 1955 James and Stein propose the following estimate for the multivariate mean  $\mu$  of using a single sample  $x$  drawn from the multivariate normal  $N_d(\mu, I)$ :

$$\hat{\mu}_{JS} = (1 - \frac{d-2}{\|x\|^2})x$$

Here, we recognise  $\hat{\mu}_{\text{ML}} = x$ ,  $\mu_0 = 0$  and shrinkage intensity  $\lambda^* = \frac{d-2}{\|x\|^2}$ .

Efron and Morris (1972) and Lindley and Smith (1972) generalised this shrinkage estimator to the case of multiple observations  $x_1, \dots, x_n$  and target  $\mu_0$ , yielding an empirical Bayes estimate of  $\mu$  based on the Normal-Normal model.



## Chapter 12

# Bayesian model comparison using Bayes factors and the BIC

### 12.1 The Bayes factor

We would like to compare two models  $M_1$  and  $M_2$ . Before seeing data  $D$  we can check their **Prior odds** (= ratio of prior probabilities of the models  $M_1$  and  $M_2$ ):

$$\frac{\Pr(M_1)}{\Pr(M_2)}$$

After seeing data  $D$  we arrive at the **Posterior odds** (= ratio of posterior probabilities):

$$\frac{\Pr(M_1|D)}{\Pr(M_2|D)}$$

Using Bayes Theorem  $\Pr(M_i|D) = \Pr(M_i) \frac{p(D|M_i)}{p(D)}$  we can rewrite the posterior odds as

$$\underbrace{\frac{\Pr(M_1|D)}{\Pr(M_2|D)}}_{\text{posterior odds}} = \underbrace{\frac{p(D|M_1)}{p(D|M_2)}}_{\text{Bayes factor } B_{12}} \underbrace{\frac{\Pr(M_1)}{\Pr(M_2)}}_{\text{prior odds}}$$

The **Bayes factor** is the multiplicative factor that updates the prior odds to the posterior odds, and is the ratio of the (marginal) likelihoods of the two models:

$$B_{12} = \frac{p(D|M_1)}{p(D|M_2)}$$

The **log-Bayes factor**  $\log B_{12}$  is also called the **weight of evidence** for  $M_1$  over  $M_2$ . Therefore, we see that

$$\text{log-posterior odds} = \text{weight of evidence} + \text{log-prior odds}$$

### 12.1.1 Connection with relative entropy

The *expected* weight of evidence, with expectation taken with regard to one of the two models, is in fact the **KL divergence** between the two models (plus a minus sign depending on direction):

$$E_{M_1}(\log B_{12}) = KL(M_1||M_2)$$

$$E_{M_2}(\log B_{12}) = -E_{M_2}(\log B_{21}) = -KL(M_2||M_1)$$

### 12.1.2 Interpretation of and scale for Bayes factor

Following Harold Jeffreys (1961) one may interpret the strength of the Bayes factor as follows:

$B_{12}$	$\log B_{12}$	evidence in favour of $M_1$ versus $M_2$
$> 100$	$> 4.6$	decisive
10 to 100	2.3 to 4.6	strong
3.2 to 10	1.16 to 2.3	substantial
1 to 3.2	0 to 1.16	not worth more than a bare mention

More recently, Kass and Raftery (1995) proposed to use the following slightly modified scale:

$B_{12}$	$\log B_{12}$	evidence in favour of $M_1$ versus $M_2$
$> 150$	$> 5$	very strong
20 to 150	3 to 5	strong
3 to 20	1 to 3	positive
1 to 3	0 to 1	not worth more than a bare mention

### 12.1.3 Computing $p(D|M)$ for simple and composite models

In the Bayes factor we need to compute  $p(D|M)$ , and it turns out that this is different for simple and composite models.

A model is called “simple” if it directly corresponds to a specific distribution, say, a Normal with fixed mean and variance, or a Binomial distribution with a

set probability for the two classes. Thus, a simple model is a point in the model space described by the parameters of a distribution family (e.g.  $\mu$  and  $\sigma^2$  for the normal family  $N(\mu, \sigma^2)$ ). For a simple model  $M$  the density  $p(D|M)$  corresponds to standard likelihood of  $M$ .

On the other hand, a model is “composite” if it is composed of simple models. This can be a finite set, or it can be comprised of infinite number of models. For example, a Normal with a given mean but unspecified variance, or a Binomial model with unspecified parameter  $p$ , is a composite model.

If  $M$  is a composite model, with the underlying simple models indexed by a parameter  $\theta$ , the probability of the data given the model is obtained by marginalisation over  $\theta$ :

$$\begin{aligned} p(D|M) &= \int_{\theta} p(D|\theta, M)p(\theta|M)d\theta \\ &= \int_{\theta} p(D, \theta|M)d\theta \end{aligned}$$

i.e. we *integrate* over all parameter values  $\theta$ . The resulting probability is called the *marginal likelihood* of the model  $M$ . Note the marginal likelihood appears also in the denominator of Bayes formula! The marginal distribution for  $D$  is also called the prior predictive distribution given  $M$ .

If the distribution over  $\theta$  is strongly concentrated around a specific value then the composite model degenerates to a simple point model.

A worked example (in the form of the Beta-Binomial distribution) is discussed in more detail in the Worksheet 6, Question 3.

### 12.1.4 Bayes factor versus likelihood ratio

**If both  $M_1$  and  $M_2$  are simple models then the Bayes factor is identical to the likelihood ratio of the two models.**

However, if one of the two models is composite then the Bayes factor and the generalised likelihood ratio differ: In the Bayes factor the representative of a composite model is the **model average** of the simple models indexed by  $\theta$ , with weights taken from the prior distribution over the simple models contained in  $M$ . In contrast, in contrast in the generalised likelihood ratio statistic the representative of a composite model is chosen by *maximisation*!

Thus, **for composite models, the Bayes factor does *not* equal the corresponding generalised likelihood ratio statistic.** As we will see next when studying the BIC approximation, the key difference is that the Bayes factor takes into account the dimension of the composite models.

## 12.2 Approximate computation of the marginal likelihood and of the log-Bayes factor

The marginal likelihood and the Bayes factor can be difficult to compute in practise. Therefore, a number of approximations for Bayesian modeling and model selection have been developed. The most important is the so-called BIC approximation.

### 12.2.1 Schwarz (1978) approximation of log-marginal likelihood

The logarithm of the marginal likelihood of a model can be approximated using the so-called BIC approximation (Schwarz 1978) as follow:

$$\log p(D|M) \approx l_n^M(\hat{\theta}_{ML}^M) - \frac{1}{2}d_M \log n$$

where  $d_M$  is the dimension of the model  $M$  (number of parameters in  $\theta$  belonging to  $M$ ) and  $n$  is the sample size and  $\hat{\theta}_{ML}^M$  is the MLE. For a simple model  $d_M = 0$  so then there is no approximation as in this case the marginal likelihood equals the likelihood.

The above formula can be obtained by quadratic approximation of the likelihood **assuming large  $n$**  and that the prior is uniform around the MLE.

Note that the approximation is the maximum log-likelihood minus a penalty that depends on the model complexity (as measured by dimension  $d$ ), thus this is an example of penalised ML! Also note that the distribution over the parameter  $\theta$  is not required in the approximation.

### 12.2.2 Bayesian information criterion (BIC)

The BIC (Bayesian information criterion) of the model  $M$  is the approximated log-marginal likelihood times the factor -2:

$$BIC(M) = -2l_n^M(\hat{\theta}_{ML}^M) + d_M \log n$$

Thus, when comparing models one aims to maximise the marginal likelihood or, as approximation, minimise the BIC.

The reason for the factor “-2” is simply to have a quantity that is on the same scale as the Wilks log likelihood ratio. Some people / software packages also use the factor “2”.

### 12.2.3 Approximating the weight of evidence (log-Bayes factor) with BIC

Using BIC (twice) the log-Bayes factor can be approximated as

$$\begin{aligned} 2 \log B_{12} &\approx -BIC(M_1) + BIC(M_2) \\ &= 2 \left( l_n^{M_1}(\hat{\theta}_{ML}^{M_1}) - l_n^{M_2}(\hat{\theta}_{ML}^{M_2}) \right) - \log(n)(d_{M_1} - d_{M_2}) \end{aligned}$$

i.e. it is the penalised log-likelihood ratio of model  $M_1$  vs.  $M_2$ .

### 12.2.4 Model complexity and Occams razor

As demonstrated above the averaging over  $\theta$  in the marginal likelihood has the effect of automatically penalising complex models.

Therefore, when comparing models using marginal likelihood, as in the Bayes factor, a complex model may be ranked below simpler models. In contrast, when selecting a model by maximum likelihood directly, without averaging, the model with the highest number of parameters always wins over simpler models.

Thus, the penalisation implicit in the marginal likelihood is very much desired as it prevents the overfitting of maximum likelihood. **The principle of preferring a less complex model is called “Occam’s razor”,** and it is a natural property of the Bayes factor.

Note that when comparing models a simpler model is often preferably over a more complex model, because the simpler model is typically better suited to both explaining the currently observed data as well as future data, whereas a complex model will only excel in fitting the current data but will then perform poorly in prediction.



# Chapter 13

## False discovery rates

### 13.1 General setup

#### 13.1.1 Overview

In this chapter we introduce False Discovery Rates (FDR) as a Bayesian method to distinguish a null model from an alternative model. This is closely linked with classical frequentist multiple testing procedures.

#### 13.1.2 Choosing between $H_0$ and $H_A$

We consider two models:

$H_0$  : null model, with density  $f_0(x)$  and distribution  $F_0(x)$

$H_A$  : alternative model, with density  $f_A(x)$  and distribution  $F_A(x)$

Aim: given observations  $x_1, \dots, x_n$  we would like to decide for each  $x_i$  whether it belongs to  $H_0$  or  $H_A$ .

This is done by a critical decision threshold  $x_c$ : if  $x_i > x_c$  then  $x_i$  is called “significant” and otherwise called “not significant”.

In classical statistics one of the the most widely used approach to find the decision threshold is by computing  $p$ -values from the  $x_i$  (this uses only the null model but not the alternative model), and then thresholding the  $p$ -values a a certain level (say 5%). If  $n$  is large then often the test is modified by adjusting the  $p$ -values or the threshold (e.g. if Bonferroni correction).

Note that this procedure ignores any information we may have about the alternative model!

### 13.1.3 True and false positives and negatives

For any decision threshold  $x_c$  we can distinguish the following errors:

- False positives (FP), “false alarm”, type I error:  $x_i$  belongs to null but is called “significant”
- False negative (FN), “miss”, type II error:  $x_i$  belongs to alternative, but is called “not significant”

In addition we have:

- True positives (TP), “hits”: belongs to alternative and is called “significant”
- True negatives (TN), “correct rejections”: belongs to null and is called “not significant”

## 13.2 Specificity and Sensitivity

From counts of TP, TN, FN, FP we can derive further quantities:

- True Negative Rate TNR, **specificity**:  $TNR = \frac{TN}{TN+FP} = 1 - FPR$  with  $FPR = \text{False Positive Rate} = 1 - \alpha_I$
- True Positive Rate TPR, **sensitivity, power**, recall:  $TPR = \frac{TP}{TP+FN} = 1 - FNR$  with  $FNR = \text{False negative rate} = 1 - \alpha_{II}$
- Accuracy:  $ACC = \frac{TP+TN}{TP+TN+FP+FN}$

Another common way to choose the decision threshold  $x_d$  in classical statistics is to balance sensitivity/power vs. specificity (maximising both power and specificity, or equivalently, minimising both false positive and false negative rates). ROC curves plot TPR/sensitivity vs.  $FPR = 1 - \text{specificity}$ .

### 13.3 FDR and FNDR

It is possible to link the above with the observed counts of TP, FP, TN, FN:

- False Discovery Rate (FDR):  $FDR = \frac{FP}{FP+TP}$
- False Nondiscovery Rate (FNDR):  $FNDR = \frac{FN}{TN+FN}$
- Positive predictive value (PPV), True Discovery Rate (TDR), precision:  $PPV = \frac{TP}{FP+TP} = 1 - FDR$
- Negative predictive value (NPV):  $NPV = \frac{TN}{TN+FN} = 1 - FNDR$

In order to choose the decision threshold it is natural to balance FDR and FNDR (or PPV and NPV), by minimising both FDR and FNDR or maximising both PPV and NPV.

In machine learning it is common to use “precision-recall plots” that plot precision (=PPV, TDR) vs. recall (=power, sensitivity).



## 13.4 Bayesian perspective

### 13.4.1 Two component mixture model

In the Bayesian perspective the problem of choosing the decision threshold is related to computing the posterior probability

$$\Pr(H_0|x_i),$$

i.e. probability of the null model given the observation  $x_i$ , or equivalently computing

$$\Pr(H_A|x_i) = 1 - \Pr(H_0|x_i)$$

the probability of the alternative model given the observation  $x_i$ .

This is done by assuming a mixture model

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

where  $\pi_0 = \Pr(H_0)$  is the prior probability of  $H_0$  and.  $\pi_A = 1 - \pi_0 = \Pr(H_A)$  the prior probability of  $H_A$ .

Note that the weights  $\pi_0$  can in fact be estimated from the observations by fitting the mixture distribution to the observations  $x_1, \dots, x_n$  (which implies that this yields a form of empirical Bayes method).

### 13.4.2 Local FDR

The posterior probability of the null model given a data point is then given by

$$\Pr(H_0|x_i) = \frac{\pi_0 f_0(x_i)}{f(x_i)} = LFDR(x_i)$$

This quantity is also known as the **local FDR** or **local False Discovery Rate**.

In the given one-sided setup the local FDR is large (close to 1) for small  $x$ , and will become close to 0 for large  $x$ . A common decision rule is given by thresholding local false discovery rates: if  $LFDR(x_i) < 0.1$  the  $x_i$  is called significant.

### 13.4.3 q-values

In correspondence to  $p$ -values one can also define tail-area based false discovery rates:

$$Fdr(x_i) = \Pr(H_0|X > x_i) = \frac{\pi_0 F_0(x_i)}{F(x_i)}$$

These are called **q-values**, or simply **False Discovery Rates (FDR)**. Intriguingly, these also have a frequentist interpretation as adjusted  $p$ -values (using a Benjamini-Hochberg adjustment procedure).

## 13.5 Software

There are a number of R packages to compute (local) FDR values:

For example:

- locfdr
- qvalue
- fdrtool

and many more.

Using FDR values for screening is especially useful in high-dimensional settings (e.g. when analysing genomic and other high-throughput data).

FDR values have both a Bayesian as well as frequentist interpretation, providing further evidence that good classical statistical methods do have a Bayesian interpretation.

## Chapter 14

# Optimality properties and summary

### 14.1 Bayesian statistics in a nutshell

- Bayesian statistics explicitly models the uncertainty about the parameters of interests by probability
- In the light of new evidence (observed data) the uncertainty is updated, i.e. the prior distribution is combined with the likelihood to form the posterior distribution

Example: Beta-Binomial model

- Binomial likelihood
- $n$  observations:  $x$  “heads”,  $n - x$  “tails”
- Frequency  $\hat{\theta}_{ML} = \frac{x}{n}$
- Beta prior  $\theta \sim \text{Beta}(\alpha_0, \beta_0)$  with mean  $\theta_0 = \frac{\alpha_0}{m}$  and  $m = \alpha_0 + \beta_0$
- Beta posterior  $\theta|x, n \sim \text{Beta}(\alpha_1, \beta_1)$  with mean  $\theta_1 = \frac{\alpha_1}{\alpha_1 + \beta_1}$  and  $\alpha_1 = \alpha_0 + x$  and  $\beta_1 = \beta_0 + n - x$
- Update of prior mean to posterior mean by shrinkage of MLE:

$$\theta_1 = \lambda \theta_0 + (1 - \lambda) \hat{\theta}_{ML}$$

with shrinkage intensity  $\lambda = \frac{m}{n+m}$

- $m$  can be interpreted as prior sample size

#### 14.1.1 Remarks

- If posterior in same family as prior  $\rightarrow$  conjugate prior
- In the exponential family the Bayesian update of the mean is always expressible as linear shrinkage of the MLE

- For sample size  $n \rightarrow \infty$  then  $\lambda \rightarrow 0$  and  $\theta_1 \rightarrow \hat{\theta}_{ML}$  (for large samples posterior mean = maximum likelihood estimator)
- For  $n \rightarrow 0$  then  $\lambda \rightarrow 1$  and  $\theta_1 \rightarrow \hat{\theta}_0$  (if no data is available fall back to prior)
- Note that the Bayesian estimator is biased for finite  $n$  by construction (but asymptotically unbiased like the MLE).

### 14.1.2 Advantages

- adding prior information has regularisation properties. This is very important in more complex models with many parameters, e.g., in estimation of a covariance matrix (to avoid singularity).
- improves small-sample accuracy (e.g. MSE)
- that Bayesian estimators tend to be better than MLE is not surprising - they use the data plus extra information!
- Bayesian credible intervals are conceptually much more simple than frequentist confidence intervals

## 14.2 Frequentist properties of Bayesian estimators

A Bayesian point estimator (e.g. the posterior mean) can also be assessed by its frequentist properties.

- First, we know that, by construction, the Bayesian estimator  $\hat{p}_{\text{Bayes}}$  will be biased for finite  $n$  even if the MLE is unbiased (with the bias being the posterior mean in this case).
- Second, intriguingly it turns out that the sampling variance of the Bayes point estimator (not to be confused with the posterior variance!) can be smaller than the variance of the MLE. This depends on the choice of the shrinkage parameter  $\lambda$  that also determines the posterior variance.

As a result, Bayesian estimators may have smaller MSE (=squared bias + variance) than the ML estimator for finite  $n$ .

In statistical decision theory this is called the theorem of **admissibility of Bayes rules**. It states that under mild conditions every admissible estimation rule (i.e. one that dominates all other estimators with regard to some expected loss, such as the MSE) is in fact a Bayes estimator with some prior.

Unfortunately, this theorem does not tell which prior is needed to achieve optimality, however an optimal estimator with minimum MSE can often be found by tuning  $\lambda$ .

## 14.3 Specifying the prior — problem or advantage?

In Bayesian statistics the analyst needs to be very explicit about the modeling assumptions:

Model = data generating process (likelihood) + prior uncertainty (prior distribution)

Note that alternative statistical methods can often be interpreted as Bayesian methods assuming a specific *implicit* prior!

For example, likelihood estimation for the Binomial model is equivalent to Bayes estimation using the Beta-Binomial model with a Beta(0,0) prior (=Haldane prior).

However, when choosing a prior explicitly for this model, interestingly most analysts would rather use a flat prior Beta(1, 1) (=Laplace prior) with implicit sample size  $m = 2$  or a transformation-invariant prior Beta(1/2, 1/2) (=Jeffreys prior) with implicit sample size  $m = 1$  than the Haldane prior!

→ be aware about the implicit priors!!

Better to acknowledge that a prior is being used (even if implicit!)

Writing down all your assumptions is enforced by the Bayesian approach.

Specifying a prior is thus best understood as an intrinsic part of model specification. It helps to improve inference and it may only be ignored if there is lots of data.

## 14.4 Choosing a prior

It is **essential in a Bayesian analysis to specify your prior uncertainty about the model parameters**. Note that this is simply **part of the modeling process**!

Typically, the location of the prior determines the amount of bias, and the precision (inverse variance) of the prior is proportional to the implied sample size of the prior.

As we have seen before for large  $n$  the Bayesian estimate converges to the ML estimate, so for large  $n$  you may ignore specifying a prior.

However, for small  $n$  it is essential that a prior is specified. In non-Bayesian approaches (if interpreted from Bayesian perspective) this prior is still there but it is implicit (e.g. uniform prior for likelihood estimation).

### 14.4.1 Some guidelines

So the question remains what are good ways to choose a prior? Two useful ways (among many others) are:

1. Use a weakly informative prior (cf. Gelman). This means that you have a vague idea about the suitable values of the parameter of interest, and you use a corresponding prior (with moderate variance) to model the uncertainty. This acknowledges that there are no uninformative priors and aims to ensure that the prior will not dominate the likelihood.
2. Empirical Bayes methods can often be used to determine one or all of the hyperparameters (i.e. the parameters in the prior). There are several ways to do this, one of them is to tune the shrinkage parameter  $\lambda$  to achieve minimum MSE. We discuss this further below.

In contrast, there also exists many proposals advocating to select so-called “uninformative priors”. However, as it is easily shown, there are no true uninformative priors, since a prior that looks uninformative (i.e. “flat”) in one coordinate system can be informative in another — this is a simple consequence of the rule for transformation of probability densities. Furthermore, often these priors are improper, i.e. are not actually probability distributions. For this (and many other reasons) the search for “uninformative” priors is not just futile but in fact also undesirable (e.g. the prior typically also needs to act as regulariser)!

**Instead, specifying the prior needs to be viewed as part of the modelling process, with specification of the prior as integral as the specification of the likelihood.**

### 14.4.2 Jeffreys prior

In order to complement the discussion on non-informative priors we now look (briefly) at the proposal by Jeffreys (1946).

Specifically, this prior is constructed from the expected Fisher observation using the log-likelihood function and thus promises automatic construction of objective uninformative priors:

$$p(\theta) \propto \sqrt{\det I^{\text{Fisher}}(\theta)}$$

The reasoning underlying this prior is **invariance against transformation of the coordinate system of the parameters**.

For the Beta-Binomial model the Jeffreys prior corresponds to  $\text{Beta}(\frac{1}{2}, \frac{1}{2})$ .

For the Normal-Normal model it corresponds to the flat improper prior  $p(\mu) = 1$ .

For the Inverse-Gamma-Normal model the Jeffreys prior is the improper prior  $p(\sigma^2) = \frac{1}{\sigma^2}$ .

This already illustrates the main problem with this type of prior – namely that it often is an improper prior.

Another issue is that Jeffreys priors are usually not conjugate which complicates the update from the prior to the posterior. An alternative to Jeffreys prior is the **reference prior** developed by Bernardo (1979).

## 14.5 Optimality of Bayesian inference

The optimality of Bayesian model making use of full model specification (likelihood plus prior) can be shown from a number of different perspectives. Correspondingly, there are many theorems that prove (or at least indicate) this optimality:

- 1) Richard Cox's theorem: the aim to generalise classic logic inevitably leads to Bayesian inference.
- 2) Entropy perspective: Bayesian inference is a consequence of minimal information update where new information arrives in form of observations
- 3) de Finetti's representation theorem: joint distribution of exchangeable sequences can be viewed as posterior distributions computed by Bayes theorem)
- 4) Frequentist decision theory: all admissible decision rules are Bayes rules! (admissible = always better than all other methods!)

Remark: the above also excludes a few other (somewhat esoteric) suggestions for propagating uncertainty (e.g. Fuzzy Logic, imprecise probabilities, etc).

## 14.6 Conclusion

Bayesian statistics offers a coherent framework for statistical learning from data, with methods for

- estimation
- testing
- model building

There are a number of theorems that show that "optimal" estimators (defined in various ways) are all Bayesian.

It is conceptually very simple — but can be computationally very involved!

It provides a coherent generalisation of classical TRUE/FALSE logic (and therefore does not suffer from some of the inconsistencies prevalent in frequentist statistics).

Bayesian statistics is a non-asymptotic theory, it works for any sample size. Asymptotically (large  $n$ ) it is consistent and converges to the true model (like ML!). But Bayesian reasoning can also be applied to events that take place

only once — no assumption of hypothetical infinitely many repetitions as in frequentist statistics is needed.

Moreover, many classical (frequentist) procedures may be viewed as *approximations* to Bayesian methods and estimators, so using classical approaches in the correct application domain is perfectly in line with the Bayesian framework.

Bayesian estimation and inference also automatically regularises (via the prior) which is important for complex models and when there is the problem of overfitting.

### 14.6.1 Current directions of research

For example: connection between Bayesian models and algorithmic models widely used in machine learning (such as neural networks, deep learning, convolutional networks, ensemble methods, XGBoost etc).

Are these models optimal (as in the Bayesian sense)? Can we learn something about highly complex, non-parametric statistical models?

How do we do effective Bayesian learning for these parameter-rich models? Both in terms of computational and statistical efficiency.



# **Part III**

## **Regression**



## Chapter 15

# Overview over regression modelling

### 15.1 General setup



- $y$ : **response variable**, also known as **outcome** or **label**
- $x_1, x_2, x_3, \dots, x_d$ : **predictor variables**, also known as **covariates** or **covariates**
- The relationship between the outcomes and the predictor variables is assumed to follow

$$y = f(x_1, x_2, \dots, x_d) + \varepsilon$$

where  $f$  is the **regression function** (not a density) and  $\varepsilon$  represents **noise**.

## 15.2 Objectives

1. **Understand the relationship** between the response  $y$  and the predictor variables  $x_i$  by **learning the regression function**  $f$  from observed data (training data). The estimated regression function is  $\hat{f}$ .
2. **Prediction of outcomes**

$$\underbrace{\hat{y}}_{\substack{\text{predicted response} \\ \text{using fitted } \hat{f}}} = \hat{f}(x_1, x_2, \dots, x_d)$$

If instead of the fitted function  $\hat{f}$  the known regression function  $f$  is used we denote this by

$$\underbrace{y^*}_{\substack{\text{predicted response} \\ \text{using known } f}} = f(x_1, x_2, \dots, x_d)$$

### 3. Variable importance

- which covariates are most relevant in predicting the outcome?
- allows to better understand the data and model  
→ variable selection (to build simpler model with same predictive capability)

## 15.3 Regression as a form of supervised learning

Regression modeling is a special case of **supervised learning**.

In supervised learning we make use of labeled data, i.e.  $x_i$  has an associated label  $y_i$ . Thus, the data consists of pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ .

The *supervision* part of in supervised learning refers to the fact that the labels are given.

In **regression** typically the label  $y_i$  is continuous and called the *response*.

On the other hand, if the label  $y_i$  is discrete/categorical then supervised learning is called **classification**.

Supervised Learning	→ Discrete $y$	→ Classification Methods
	→ Continuous $y$	→ Regression Methods

Another important type of statistical learning is **unsupervised learning** where labels  $y$  are inferred from the data  $x$  (this is also known as **clustering**). Furthermore, there is also *semi-supervised learning* with labels only partly known.

Note that there are regression models (e.g. logistic regression) with discrete response that are performing classification, so one may argue that “supervised learning”=“generalised regression”.

## 15.4 Various regression models used in statistics

In this course we only study linear multiple regression. However, you should be aware that the linear model is in fact just a special cases of some much more general regression approaches.

General regression model:

$$y = f(x_1, \dots, x_d) + \text{"noise"}$$

- The function  $f$  is estimated nonparametrically - splines - Gaussian processes
- Generalised Additive Models (GAM): - the function  $f$  is assumed to be the sum of individual functions  $f_i(x_i)$
- Generalised Linear Models (GLM): -  $f$  is a transformed linear predictor  $h(\sum b_i x_i)$ , noise is assumed from exponential family
- Linear Model (LM): - linear predictor  $\sum b_i x_i$ , normal noise

In R the linear model is implemented in the function `lm()`, and generalised linear models in the function `glm()`. Generalised additive models are available in the package “mgcv”.

In the following we focus on the linear regression model with continuous response.



# Chapter 16

## Linear Regression

### 16.1 The linear regression model

In this module we assume that  $f$  is a linear function:

$$f(x_1, \dots, x_d) = \beta_0 + \sum_{j=1}^d \beta_j x_j = y^*$$

In vector notation:

$$f(\mathbf{x}) = \beta_0 + \boldsymbol{\beta}^T \mathbf{x} = y^*$$

with  $\boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_d \end{pmatrix}$  and  $\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$

Therefore, the linear regression model is

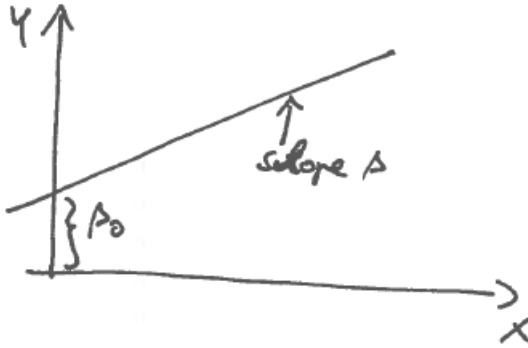
$$\begin{aligned} y &= \beta_0 + \sum_{j=1}^d \beta_j x_j + \varepsilon \\ &= \beta_0 + \boldsymbol{\beta}^T \mathbf{x} + \varepsilon \\ &= y^* + \varepsilon \end{aligned}$$

where:

- $\beta_0$  is the **intercept**
- $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d)^T$  are the **regression coefficients**
- $\mathbf{x} = (x_1, \dots, x_d)^T$  is the predictor vector containing the **predictor variables**

## 16.2 Interpretation of regression coefficients and intercept

- The regression coefficient  $\beta_i$  corresponds to the slope (first partial derivative) of the regression function in the direction of  $x_i$ . In other words, the gradient of  $f(x)$  are the regression coefficients:  $\nabla f(x) = \beta$
- The intercept  $\beta_0$  is the offset at the origin ( $x_1 = x_2 = \dots = x_d = 0$ ):



## 16.3 Different types of linear regression:

- **Simple linear regression:**  $y = \beta_0 + \beta x + \varepsilon$  (=single predictor)
- **Multiple linear regression:**  $y = \beta_0 + \sum_{j=1}^d \beta_j x_j + \varepsilon$  (= multiple predictor variables)
- **Multivariate regression:** multivariate response  $y$

## 16.4 Distributional assumptions and properties

*General assumptions:*

- We treat  $y$  and  $x_1, \dots, x_d$  as the primary observables that can be described by random variables.
- $\beta_0, \beta$  are parameters to be inferred from the observations on  $y$  and  $x_1, \dots, x_d$ .
- Specifically, will we assume that response and predictors have a mean and a (cov)variance:
  - Response:
 
$$E(y) = \mu_y$$

$$\text{Var}(y) = \sigma_y^2$$



The **variance of the response**  $\text{Var}(y)$  is also called the **total variation**.

ii. Predictors:

$$E(x_i) = \mu_{x_i} \text{ (or } E(\mathbf{x}) = \boldsymbol{\mu}_x)$$

$$\text{Var}(x_i) = \sigma_{x_i}^2 \text{ and } \text{Cor}(x_i, x_j) = \rho_{ij} \text{ (or } \text{Var}(\mathbf{x}) = \boldsymbol{\Sigma}_x)$$

The **signal variance**  $\text{Var}(y^*) = \text{Var}(\beta_0 + \boldsymbol{\beta}^T \mathbf{x}) = \boldsymbol{\beta}^T \boldsymbol{\Sigma}_x \boldsymbol{\beta}$  is also called the **explained variation**.

- We assume that  $y$  and  $x$  are jointly distributed with correlation  $\text{Cor}(y, x_j) = \rho_{y, x_j}$  between each predictor variable  $x_j$  and the response  $y$ .
- In contrast to  $y$  and  $x$  the noise variable  $\varepsilon$  is only indirectly observed via the difference  $\varepsilon = y - y^*$ . We denote the mean and variance of the noise by  $E(\varepsilon)$  and  $\text{Var}(\varepsilon)$ .  
The **noise variance**  $\text{Var}(\varepsilon)$  is also called the **unexplained variation**.

*Identifiability assumptions:*

In a statistical analysis we would like to be able to separate signal ( $y^*$ ) from noise ( $\varepsilon$ ). To achieve this we require some **distributional assumptions to ensure identifiability** and avoid confounding:

- 1) **Assumption 1:**  $\varepsilon$  and  $y^*$  are independent. This implies  $\text{Var}(y) = \text{Var}(y^*) + \text{Var}(\varepsilon)$ , or equivalently  $\text{Var}(\varepsilon) = \text{Var}(y) - \text{Var}(y^*)$ .

Thus, this assumption implies the **decomposition of variance**, i.e. that the **total variation**  $\text{Var}(y)$  equals the sum of the **explained variation**  $\text{Var}(y^*)$  and the **unexplained variation**  $\text{Var}(\varepsilon)$ .

- 2) **Assumption 2:**  $E(\varepsilon) = 0$ . This allows to identify the intercept  $\beta_0$  and implies  $E(y) = E(y^*)$ .

*Optional assumptions (often but not always):*

- The noise  $\varepsilon$  is normally distributed
- The response  $y$  and the predictor variables  $x_i$  are continuous variables
- The response and predictor variables are jointly normally distributed

*Further properties:*

- As a result of the independence assumption 1) we can only choose two out of the three variances freely:
  - i. in a generative perspective we will choose signal variance  $\text{Var}(y^*)$  (or equivalently the variances  $\text{Var}(x_j)$ ) and the noise variance  $\text{Var}(\varepsilon)$ , then the variance of the response  $\text{Var}(y)$  follows.
  - ii. in an observational perspective we will observe the variance of the response  $\text{Var}(y)$  and the variances  $\text{Var}(x_j)$ , and then the error variance  $\text{Var}(\varepsilon)$  follows.
- As we will see later the regression coefficients  $\beta_j$  depend on the correlations between the response  $y$  and the predictor variables  $x_j$ . Thus, the choice

of regression coefficients implies a specific correlation pattern, and vice versa (in fact, we will use this correlation pattern to infer the regression coefficient from data!).

## 16.5 Regression in data matrix notation

We can also write the regression in terms of actual observed data (rather than random variables):

Data matrix for the predictors:

$$\mathbf{X} = \begin{pmatrix} x_{11} & \dots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nd} \end{pmatrix}$$

Note the statistics convention: the  $n$  rows of  $\mathbf{X}$  contain the samples, and the  $d$  columns contain variables.

Response data vector:  $(y_1, \dots, y_n)^T = \mathbf{y}$

Then the regression equation is written in data matrix notation:

$$\underbrace{\mathbf{y}}_{n \times 1} = \underbrace{\mathbf{1}_n}_{n \times 1} \beta_0 + \underbrace{\mathbf{X}}_{n \times d} \underbrace{\boldsymbol{\beta}}_{d \times 1} + \underbrace{\boldsymbol{\varepsilon}}_{\substack{n \times 1 \\ \text{residuals}}}$$

where  $\mathbf{1}_n = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$  is a column vector of length  $n$  (size  $n \times 1$ ).

Note that here the regression coefficients are now multiplied *after* the data matrix (compare with the original vector notation where the *transpose* of regression coefficients come *before* the vector of the predictors).

The **observed noise** values (i.e. realisations of  $\varepsilon$ ) are called the **residuals**.

## 16.6 Centering and vanishing of the intercept $\beta_0$

If  $x$  and  $y$  are centered, i.e.  $E(x) = \mu_x = 0$  and  $E(y) = \mu_y = 0$  then the intercept  $\beta_0$  disappears:

The regression equation is

$$y = \beta_0 + \boldsymbol{\beta}^T \mathbf{x} + \varepsilon$$

with  $E(\varepsilon)$ . Taking the expectation on both sides we get  $\mu_y = \beta_0 + \boldsymbol{\beta}^T \boldsymbol{\mu}_x$  and therefore

$$\beta_0 = \mu_y - \boldsymbol{\beta}^T \boldsymbol{\mu}_x$$

This is equal to zero if the means of the response and predictors vanish. Conversely, if we assume that the intercept vanishes ( $\beta_0 = 0$ ) this is only possible for general  $\boldsymbol{\beta}$  if both  $\boldsymbol{\mu}_x = 0$  and  $\mu_y = 0$ .

Thus, in the linear model is always possible to transform  $y$  and  $x$  (or data  $\mathbf{y}$  and  $\mathbf{X}$ ) so that the intercept vanishes!

$\Rightarrow$  we will therefore often set  $\beta_0 = 0$ .

## 16.7 Regression objectives for linear model

1. Understand functional relationship: find estimates of intercept ( $\hat{\beta}_0$ ) and regression coefficients ( $\hat{\beta}_j$ )
2. Prediction:
  - Known coefficients  $\beta_0$  and  $\boldsymbol{\beta}$ :  $y^* = \beta_0 + \boldsymbol{\beta}^T \mathbf{x}$
  - Estimated coefficients  $\hat{\beta}_0$  and  $\hat{\boldsymbol{\beta}}$  (note the “hat”!):  $\hat{y} = \hat{\beta}_0 + \sum_{j=1}^d \hat{\beta}_j x_j = \hat{\beta}_0 + \hat{\boldsymbol{\beta}}^T \mathbf{x}$

Also find the **corresponding prediction errors!**

3. Variable importance: Which predictors  $x_j$  are most relevant?
  - $\rightarrow$  test whether  $\beta_j = 0$
  - $\rightarrow$  find measures of variable importance

Remark: as we will see  $\beta_j$  or  $\hat{\beta}_j$  itself is **not** a measure of importance!



# Chapter 17

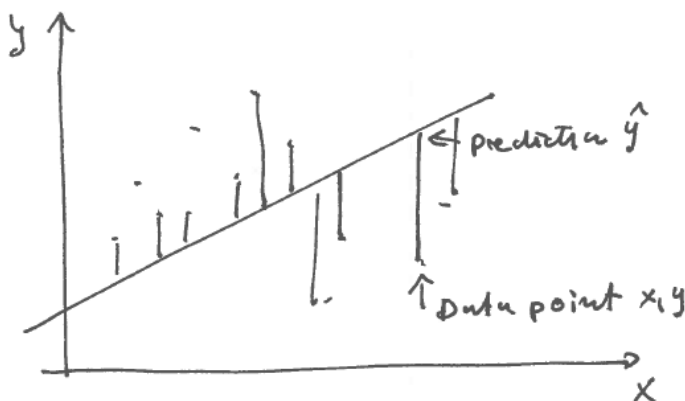
## Estimating regression coefficients

In this chapter we discuss various ways to estimate the regression coefficients. First, we discuss estimation by Ordinary Least Squares (OLS) by minimising the residual sum of squares. This yields the famous Gauss estimator. Second, we derive estimates of the regression coefficients using the methods of maximum likelihood assuming normal errors. This also leads to the Gauss estimator. Third, we show that the coefficients in linear regression can be written and interpreted in terms of two covariance matrices and that the Gauss estimator of the regression coefficients is a plug-in estimator using the MLEs of these covariance matrices. Furthermore, we show that the (population version) of the Gauss estimator can also be derived by finding the best linear predictor and by conditioning. Finally, we discuss special cases of regression coefficients and their relationship to marginal correlation.

### 17.1 Ordinary Least Squares (OLS) estimator of regression coefficients

Now we show the classic way (Gauss 1809; Legendre 1805) to **estimate regression coefficients** by the method of **ordinary least squares (OLS)**.

*Idea:* choose regression coefficients such as to *minimise* the *squared error* between observations and the prediction.



In data matrix notation (note we assume  $\beta_0 = 0$  and thus *centered data*  $X$  and  $y$ ):

$$\text{RSS}(\beta) = (y - X\beta)^T (y - X\beta)$$

RSS is an abbreviation for “Residual Sum of Squares” which is a function of  $\beta$ . Minimising RSS yields the OLS estimate:

$$\hat{\beta}_{\text{OLS}} = \arg \min_{\beta} \text{RSS}(\beta)$$

$$\text{RSS}(\beta) = y^T y - 2\beta^T X^T y + \beta^T X^T X \beta$$

Gradient:

$$\nabla \text{RSS}(\beta) = -2X^T y + 2X^T X \beta$$

$$\nabla \text{RSS}(\hat{\beta}) = 0 \longrightarrow X^T y = X^T X \hat{\beta}$$

$$\implies \hat{\beta}_{\text{OLS}} = (X^T X)^{-1} X^T y$$

Note the similarities in the procedure to maximum likelihood (ML) estimation (with minimisation instead of maximisation)! In fact, as we see next this is not by chance as OLS *is* indeed a special case of ML! This also implies that OLS is generally a good method — but only if sample size  $n$  is large!

The above Gauss’ estimator is fundamental in statistics so it is worthwhile to memorise it!

## 17.2 Maximum likelihood estimation of regression coefficients

We now show how to estimate regression coefficients using the method of maximum likelihood. This is a second method to derive  $\hat{\beta}$ .

We recall the basic regression equation

$$y = \beta_0 + \beta^T x + \varepsilon$$

with  $E(\varepsilon) = 0$  and observed data  $y_1, \dots, y_n$  and  $x_1, \dots, x_n$ . The intercept is identified as

$$\beta_0 = \mu_y - \beta^T \mu_x$$

so that we can solve for the noise variable

$$\varepsilon = (y - \mu_y) - \beta^T (x - \mu_x)$$

Assuming joint (multivariate) normality for the response  $y$  and  $x$  we get as the MLEs for the respective means and (co)variances:

- $\hat{\mu}_y = \hat{E}(y) = \frac{1}{n} \sum_{i=1}^n y_i$
- $\hat{\sigma}_y^2 = \widehat{\text{Var}}(y) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu}_y)^2$
- $\hat{\mu}_x = \hat{E}(x) = \frac{1}{n} \sum_{i=1}^n x_i$
- $\hat{\Sigma}_{xx} = \widehat{\text{Var}}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_x)(x_i - \hat{\mu}_x)^T$
- $\hat{\Sigma}_{xy} = \widehat{\text{Cov}}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_x)(y_i - \hat{\mu}_y)$

The noise  $\varepsilon \sim N(0, \sigma_\varepsilon^2)$  is normally distributed with mean 0 and variance  $\text{Var}(\varepsilon) = \sigma_\varepsilon^2$ . Having obtained MLEs  $\hat{\mu}_y$  and  $\hat{\mu}_x$  corresponding (indirect) observations are given by

$$(y_i - \hat{\mu}_y) - \beta^T (x_i - \hat{\mu}_x)$$

which leads to the normal log-likelihood function

$$\log L(\beta, \sigma_\varepsilon^2) = -\frac{n}{2} \log \sigma_\varepsilon^2 - \frac{1}{2\sigma_\varepsilon^2} \sum_{i=1}^n \left( (y_i - \hat{\mu}_y) - \beta^T (x_i - \hat{\mu}_x) \right)^2$$

We now only need to maximise the log-likelihood to obtain MLEs of  $\sigma_\varepsilon^2$  and  $\beta$ !

Note that the residual sum of squares appears in the log-likelihood function (with a minus sign), which implies that ML assuming normal distribution will recover the OLS estimator for the regression coefficients! So OLS is a special case of ML !

### 17.2.1 Detailed derivation of the MLEs

The gradient with regard to  $\beta$  is

$$\begin{aligned}\nabla_{\beta} \log L(\beta, \sigma_{\varepsilon}^2) &= \frac{1}{\sigma_{\varepsilon}^2} \sum_{i=1}^n \left( (x_i - \hat{\mu}_x)(y_i - \hat{\mu}_y) - (x_i - \hat{\mu}_x)(x_i - \hat{\mu}_x)^T \beta \right) \\ &= \frac{n}{\sigma_{\varepsilon}^2} (\hat{\Sigma}_{xy} - \hat{\Sigma}_{xx} \beta)\end{aligned}$$

Setting this equal to zero yields the Gauss estimator

$$\hat{\beta} = \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy}$$

By plugin we get the MLE of  $\beta_0$  as

$$\hat{\beta}_0 = \hat{\mu}_y - \hat{\beta}^T \hat{\mu}_x$$

Taking the derivative of  $\log L(\hat{\beta}, \sigma_{\varepsilon}^2)$  with regard to  $\sigma_{\varepsilon}^2$  yields

$$\frac{\partial}{\partial \sigma_{\varepsilon}^2} \log L(\hat{\beta}, \sigma_{\varepsilon}^2) = -\frac{n}{2\sigma_{\varepsilon}^2} + \frac{1}{2\sigma_{\varepsilon}^4} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

with  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}^T x_i$  and the residuals  $y_i - \hat{y}_i$  resulting from the fitted linear model. This leads to the MLE of the noise variance

$$\hat{\sigma}_{\varepsilon}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Note that the MLE  $\hat{\sigma}_{\varepsilon}^2$  is a biased estimate of  $\sigma_{\varepsilon}^2$ . The unbiased estimate is  $\frac{1}{n-d-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ , where  $d$  is the dimension of  $\beta$  (i.e. the number of predictors).

### 17.2.2 Asymptotics

The advantage of using maximum likelihood is that we also get the (asymptotic) variance associated with each estimator and typically can also assume asymptotic normality.

Specifically, for  $\hat{\beta}$  we get via the observed Fisher information at the MLE an asymptotic estimator of its variance

$$\widehat{\text{Var}}(\hat{\beta}) = \frac{1}{n} \hat{\sigma}_{\varepsilon}^2 \hat{\Sigma}_{xx}^{-1}$$

Similarly, for  $\hat{\beta}_0$  we have

$$\widehat{\text{Var}}(\hat{\beta}_0) = \frac{1}{n} \hat{\sigma}_{\varepsilon}^2 (1 + \hat{\mu}^T \hat{\Sigma}_{xx}^{-1} \hat{\mu})$$



For finite sample size  $n$  with known  $\text{Var}(\varepsilon)$  one can show that the variances are

$$\text{Var}(\hat{\boldsymbol{\beta}}) = \frac{1}{n} \sigma_\varepsilon^2 \hat{\boldsymbol{\Sigma}}_{xx}^{-1}$$

and

$$\text{Var}(\hat{\beta}_0) = \frac{1}{n} \sigma_\varepsilon^2 (1 + \hat{\boldsymbol{\mu}}_x^T \hat{\boldsymbol{\Sigma}}_{xx}^{-1} \hat{\boldsymbol{\mu}}_x)$$

and that the regression coefficients and the intercept are normally distributed according to

$$\hat{\boldsymbol{\beta}} \sim N_d(\boldsymbol{\beta}, \text{Var}(\hat{\boldsymbol{\beta}}))$$

and

$$\hat{\beta}_0 \sim N(\beta_0, \text{Var}(\hat{\beta}_0))$$

We may use this to test whether whether  $\beta_j = 0$  and  $\beta_0 = 0$ .

## 17.3 Covariance plug-in estimator of regression coefficients

We now try to understand regression coefficients in terms of covariances (thus obtaining a third way to compute and estimate them).

We recall that the Gauss regression coefficients are given by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

where  $\mathbf{X}$  is the  $n \times d$  data matrix (in statistics convention)

$$\mathbf{X} = \begin{pmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nd} \end{pmatrix}$$

Note that we assume that the data matrix  $\mathbf{X}$  is centered (i.e. column sums  $\mathbf{X}^T \mathbf{1}_n = \mathbf{0}$  are zero).

Likewise  $\mathbf{y} = (y_1, \dots, y_n)^T$  is the response data vector (also centered with  $\mathbf{y}^T \mathbf{1}_n = 0$ ).

Noting that

$$\hat{\boldsymbol{\Sigma}}_{xx} = \frac{1}{n} (\mathbf{X}^T \mathbf{X})$$

is the MLE of covariance matrix among  $\mathbf{x}$  and

$$\hat{\boldsymbol{\Sigma}}_{xy} = \frac{1}{n} (\mathbf{X}^T \mathbf{y})$$

is the MLE of the covariance between  $x$  and  $y$  we see that the OLS estimate of the regression coefficients can be expressed as

$$\hat{\beta} = (\hat{\Sigma}_{xx})^{-1} \hat{\Sigma}_{xy}$$

We can write down a population version (with no hats!):

$$\beta = \Sigma_{xx}^{-1} \Sigma_{xy}$$

Thus, OLS regression coefficients can be interpreted as plugin estimator using MLEs of covariances! In fact, we may also use the unbiased estimates since the scale factor ( $1/n$  or  $1/(n-1)$ ) cancels out so it does not matter which one you use!

### 17.3.1 Importance of positive definiteness of estimated covariance matrix

Note that  $\hat{\Sigma}_{xx}$  is inverted in  $\hat{\beta} = (\hat{\Sigma}_{xx})^{-1} \hat{\Sigma}_{xy}$ .

- Hence, the estimate  $\hat{\Sigma}_{xx}$  needs to be positive definite!
- But  $\hat{\Sigma}_{xx}^{\text{MLE}}$  is only positive definite if  $n > d$ !

Therefore we can use the ML estimate (empirical estimator) only for large  $n > d$ , otherwise we need to employ a different (regularised) estimation approach (e.g. Bayes or a penalised ML)!

Remark: writing  $\hat{\beta}$  explicitly based on covariance estimates has the advantage that we can construct plug-in estimators of regression coefficient based on regularised covariance estimators that improve over ML for small sample size. This leads to the so-called SCOUT method (=covariance-regularized regression by Witten and Tibshirani, 2008).

## 17.4 Best linear predictor

The **best linear predictor** is a fourth way to arrive at the linear model. This is closely related to OLS and minimising squared residual error.

Without assuming normality the above multiple regression model can be shown to be optimal linear predictor under the minimum mean squared prediction error:

Assumptions:

- $y$  and  $x$  are random variables
- we construct a new variable (the linear predictor)  $y^{**} = b_0 + \mathbf{b}^T \mathbf{x}$  to optimally approximate  $y$

Aim:

- choose  $b_0$  and  $\mathbf{b}$  such to minimize the mean squared prediction error  $E((y - y^{**})^2)$

### 17.4.1 Result:

The mean squared prediction error  $MSPE$  in dependence of  $(b_0, \mathbf{b})$  is

$$\begin{aligned}
 E((y - y^{**})^2) &= \text{Var}(y - y^{**}) + E(y - y^{**})^2 \\
 &= \text{Var}(y - b_0 - \mathbf{b}^T \mathbf{x}) + (E(y) - b_0 - \mathbf{b}^T E(\mathbf{x}))^2 \\
 &= \sigma_y^2 + \text{Var}(\mathbf{b}^T \mathbf{x}) + 2 \text{Cov}(y, -\mathbf{b}^T \mathbf{x}) + (\mu_y - b_0 - \mathbf{b}^T \boldsymbol{\mu}_x)^2 \\
 &= \sigma_y^2 + \mathbf{b}^T \boldsymbol{\Sigma}_{xx} \mathbf{b} - 2 \mathbf{b}^T \boldsymbol{\Sigma}_{xy} + (\mu_y - b_0 - \mathbf{b}^T \boldsymbol{\mu}_x)^2 \\
 &= MSPE(b_0, \mathbf{b})
 \end{aligned}$$

We look for

$$(\beta_0, \boldsymbol{\beta}) = \arg \min_{b_0, \mathbf{b}} MSPE(b_0, \mathbf{b})$$

In order to find the minimum we compute the gradient with regard to  $(b_0, \mathbf{b})$

$$\nabla MSPE = \begin{pmatrix} -2(\mu_y - b_0 - \mathbf{b}^T \boldsymbol{\mu}_x) \\ 2 \boldsymbol{\Sigma}_{xx} \mathbf{b} - 2 \boldsymbol{\Sigma}_{xy} - 2 \boldsymbol{\mu}_x (\mu_y - b_0 - \mathbf{b}^T \boldsymbol{\mu}_x) \end{pmatrix}$$

and setting this equal to zero yields

$$\begin{pmatrix} \beta_0 \\ \boldsymbol{\beta} \end{pmatrix} = \begin{pmatrix} \mu_y - \boldsymbol{\beta}^T \boldsymbol{\mu}_x \\ \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy} \end{pmatrix}$$

Thus, the optimal values for  $b_0$  and  $\mathbf{b}$  in the best linear predictor correspond to the previously derived coefficients  $\beta_0$  and  $\boldsymbol{\beta}$ !

### 17.4.2 Irreducible Error

The minimum achieved  $MSPE$  (=irreducible error) is

$$MSPE(\beta_0, \boldsymbol{\beta}) = \sigma_y^2 - \boldsymbol{\beta}^T \boldsymbol{\Sigma}_{xx} \boldsymbol{\beta} = \sigma_y^2 - \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}$$

With the abbreviation  $\Omega^2 = \mathbf{P}_{yx} \mathbf{P}_{xx}^{-1} \mathbf{P}_{xy} = \sigma_y^{-2} \boldsymbol{\Sigma}_{yx} \boldsymbol{\Sigma}_{xx}^{-1} \boldsymbol{\Sigma}_{xy}$  we can simplify this to

$$MSPE(\beta_0, \boldsymbol{\beta}) = \sigma_y^2 (1 - \Omega^2) = \text{Var}(\varepsilon)$$

Writing  $b_0 = \beta_0 + \Delta_0$  and  $\mathbf{b} = \boldsymbol{\beta} + \boldsymbol{\Delta}$  it is easy to see that the mean squared predictive error is a quadratic function around the minimum:

$$MSPE(\beta_0 + \Delta_0, \boldsymbol{\beta} + \boldsymbol{\Delta}) = \text{Var}(\varepsilon) + \Delta_0^2 + \boldsymbol{\Delta}^T \boldsymbol{\Sigma}_{xx} \boldsymbol{\Delta}$$

Note that usually  $y^* = \beta_0 + \beta^T x$  does not perfectly approximate  $y$  so there *will* be an irreducible error (= noise variance)

$$\text{Var}(\varepsilon) = \sigma_y^2(1 - \Omega^2) > 0$$

which implies  $\Omega^2 < 1$ .

The quantity  $\Omega^2$  has a further interpretation of the population version of as the squared multiple correlation coefficient between the response and the predictors and plays a vital role in decomposition of variance, as discussed later.

## 17.5 Regression by conditioning

**Conditioning** is a fifth way to arrive at the linear model. This is also the most general way and can be used to derive many other regression models (not just the simple linear model).

### 17.5.1 General idea:

- two random variables  $y$  (response, scalar) and  $x$  (predictor variables, vector)
- we assume that  $y$  and  $x$  have a joint distribution  $F_{y,x}$
- compute *conditional* random variable  $y|x$  and the corresponding distribution  $F_{y|x}$

### 17.5.2 Multivariate normal assumption

Now we assume that  $y$  and  $x$  are (jointly) multivariate normal. Then the conditional distribution  $F_{y|x}$  is a univariate normal with the following moments (you can verify this by looking up the general conditional multivariate normal distribution):

#### a) Conditional expectation:

$$E(y|x) = y^* = \beta_0 + \beta^T x$$

with coefficients  $\beta = \Sigma_{xx}^{-1} \Sigma_{xy}$  and intercept  $\beta_0 = \mu_y - \beta^T \mu_x$ .

Note that as  $y^*$  depends on  $x$  it is a random variable itself with mean

$$E(y^*) = \beta_0 + \beta^T \mu_x = \mu_y$$

and variance

$$\begin{aligned}
 \text{Var}(y^*) &= \text{Var}(E(y|x)) \\
 &= \beta^T \Sigma_{xx} \beta = \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \\
 &= \sigma_y^2 P_{yx} P_{xx}^{-1} P_{xy} \\
 &= \sigma_y^2 \Omega^2
 \end{aligned}$$

**b) Conditional variance:**

$$\begin{aligned}
 \text{Var}(y|x) &= \sigma_y^2 - \beta^T \Sigma_{xx} \beta \\
 &= \sigma_y^2 - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \\
 &= \sigma_y^2 (1 - \Omega^2)
 \end{aligned}$$

Note this is a constant so  $E(\text{Var}(y|x)) = \sigma_y^2 (1 - \Omega^2)$  as well.

## 17.6 Standardised regression coefficients and relationship to correlation

First we note that we can decompose regression coefficients into the product of marginal correlations and correlations among predictors.

Using the variance-correlation decompositions  $\Sigma_{xx} = V_x^{1/2} P_{xx} V_x^{1/2}$  and  $\Sigma_{xy} = V_x^{1/2} P_{xy} \sigma_y$  we rewrite the regression coefficients as

$$\beta = \underbrace{V_x^{-1/2}}_{\text{(inverse) scale of } x_i} \underbrace{P_{xx}^{-1}}_{\text{(inverse) correlation among predictors}} \underbrace{P_{xy}}_{\text{marginal correlations}} \underbrace{\sigma_y}_{\text{scale of } y}$$

Thus the regression coefficients  $\beta$  contain the scale of the variables, and take into account the correlations among the predictors ( $P_{xx}$ ) in addition to the marginal correlations between the response  $y$  and the predictors  $x_i$  ( $P_{xy}$ ).

This decomposition allows to understand a number special cases when the regression coefficient simplify further:

- a) If the response and the predictors are standardised to have variance one, i.e.  $\text{Var}(y) = 1$  and  $\text{Var}(x_i)$ , then  $\beta$  becomes equal to the **standardised regression coefficients**

$$\beta_{\text{std}} = P_{xx}^{-1} P_{xy}$$

Note that standardised regression coefficients do not make use of variances and thus are scale-independent.

- b) If there is no correlation among the predictors , i.e.  $P_{xx} = I$  the the regression coefficients reduce to

$$\beta = V_x^{-1} \Sigma_{xy}$$

where  $V_x$  is a diagonal matrix containing the variances of the predictors. This is also called **marginal regression**. Note that the inversion of  $V_x$  is trivial since you only need to invert each diagonal element individually.

- c) If both a) and b) apply simultaneously (i.e. there is no correlation among predictors and response and predictors and predictors are standardised) then the regression coefficients simplify even further to

$$\beta = P_{xy}$$

Thus, in this very special case the regression coefficients are identical to the correlations between the response and the predictors!

## Chapter 18

# Squared multiple correlation and variance decomposition in linear regression

In this chapter we first introduce the (squared) multiple correlation and the multiple and adjusted  $R^2$  coefficients as estimators. Subsequently we discuss variance decomposition.

### 18.1 Squared multiple correlation $\Omega^2$ and the $R^2$ coefficient

In the previous chapter we encountered the following quantity:

$$\Omega^2 = P_{yx}P_{xx}^{-1}P_{xy} = \sigma_y^{-2}\Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}$$

With  $\beta = \Sigma_{xx}^{-1}\Sigma_{xy}$  and  $\beta_0 = \mu_y - \beta^T\mu_x$  it is straightforward to verify the following:

- the cross-covariance between  $y$  and  $y^*$  is

$$\begin{aligned}\text{Cov}(y, y^*) &= \Sigma_{yx}\beta = \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy} \\ &= \sigma_y^2 P_{yx}P_{xx}^{-1}P_{xy} = \sigma_y^2 \Omega^2\end{aligned}$$

- the (signal) variance of  $y^*$  is

$$\begin{aligned}\text{Var}(y^*) &= \beta^T \Sigma_{xx} \beta = \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \\ &= \sigma_y^2 P_{yx} P_{xx}^{-1} P_{xy} = \sigma_y^2 \Omega^2\end{aligned}$$

hence the correlation  $\text{Cor}(y, y^*) = \frac{\text{Cov}(y, y^*)}{\text{SD}(y)\text{SD}(y^*)} = \Omega$  with  $\Omega \geq 0$ .

This helps to understand the  $\Omega$  and  $\Omega^2$  coefficients:

- $\Omega$  is the linear correlation between the response ( $y$ ) and prediction  $y^*$ .
- $\Omega^2$  is called the **squared multiple correlation** between the scalar  $y$  and the vector  $x$ .
- Note that if we only have one predictor (if  $x$  is a scalar) then  $P_{xx} = 1$  and  $P_{yx} = \rho_{yx}$  so the multiple squared correlation coefficient reduces to squared correlation  $\Omega^2 = \rho_{yx}^2$  between two scalar random variables  $y$  and  $x$ .

### 18.1.1 Estimation of $\Omega^2$ and the multiple $R^2$ coefficient

The multiple squared correlation coefficient  $\Omega^2$  can be estimated by plug-in of empirical estimates for the corresponding correlation matrices:

$$R^2 = \hat{P}_{yx} \hat{P}_{xx}^{-1} \hat{P}_{xy} = \hat{\sigma}_y^{-2} \hat{\Sigma}_{yx} \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy}$$

This estimator of  $\Omega^2$  is called the **multiple  $R^2$  coefficient**.

If the same scale factor  $1/n$  or  $1/(n-1)$  is used in estimating the variance  $\sigma_y^2$  and the covariances  $\Sigma_{xx}$  and  $\Sigma_{yx}$  then this factor will cancel out.

Above we have seen that  $\Omega^2$  is directly linked with the noise variance via

$$\text{Var}(\varepsilon) = \sigma_y^2(1 - \Omega^2).$$

so we can express the squared multiple correlation as

$$\Omega^2 = 1 - \text{Var}(\varepsilon)/\sigma_y^2$$

The **maximum likelihood estimate** of the noise variance  $\text{Var}(\varepsilon)$  (also called **residual variance**) can be computed from the residual sum of squares  $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$  as follows:

$$\widehat{\text{Var}}(\varepsilon)_{ML} = \frac{RSS}{n}$$

whereas the **unbiased estimate** is obtained by

$$\widehat{\text{Var}}(\varepsilon)_{UB} = \frac{RSS}{n - d - 1} = \frac{RSS}{df}$$

where the **degree of freedom** is  $df = n - d - 1$  and  $d$  is the number of predictors.

Similarly, we can find the maximum likelihood estimate  $v_y^{ML}$  for  $\sigma_y^2$  (with factor  $1/n$ ) as well as an unbiased estimate  $v_y^{UB}$  (with scale factor  $1/(n-1)$ )



The **multiple  $R^2$  coefficient** can then be written as

$$R^2 = 1 - \widehat{\text{Var}}(\varepsilon)_{ML} / v_y^{ML}$$

Note we use MLEs.

In contrast, the so-called **adjusted multiple  $R^2$  coefficient** is given by

$$R_{\text{adj}}^2 = 1 - \widehat{\text{Var}}(\varepsilon)_{UB} / v_y^{UB}$$

where the unbiased variances are used.

Both  $R^2$  and  $R_{\text{adj}}^2$  are estimates of  $\Omega^2$  and are related by

$$1 - R^2 = (1 - R_{\text{adj}}^2) \frac{df}{n - 1}$$

### 18.1.2 R commands

In R the command `lm()` fits the linear regression model.

In addition to the regression coefficients (and derived quantities) the R function `lm()` also lists

- the multiple R-squared  $R^2$ ,
- the adjusted R-squared  $R_{\text{adj}}^2$ ,
- the degrees of freedom  $df$  and
- the residual standard error  $\sqrt{\widehat{\text{Var}}(\varepsilon)_{UB}}$  (computed from the unbiased variance estimate).

See also Worksheet 9 which provides R code to reproduce the exact output of the native `lm()` R function.

## 18.2 Variance decomposition in regression

The squared multiple correlation coefficient is useful also because it plays an important role in the decomposition of the total variance:

- total variance:  $\text{Var}(y) = \sigma_y^2$
- unexplained variance (irreducible error):  $\sigma_y^2(1 - \Omega^2) = \text{Var}(\varepsilon)$
- the explained variance is the complement:  $\sigma_y^2\Omega^2 = \text{Var}(y^*)$

In summary:

$$\text{Var}(y) = \text{Var}(y^*) + \text{Var}(\varepsilon)$$

becomes

$$\underbrace{\sigma_y^2}_{\text{total variance}} = \underbrace{\sigma_y^2\Omega^2}_{\text{explained variance}} + \underbrace{\sigma_y^2(1 - \Omega^2)}_{\text{unexplained variance}}$$

The unexplained variance measures the fit after introducing predictors into the model (smaller means better fit). The total variance measures the fit of the model without any predictors. The explained variance is the difference between total and unexplained variance, it indicates the increase in model fit due to the predictors.

## 18.2.1 Law of total variance and variance decomposition

The law of total variance is

$$\underbrace{\text{Var}(y)}_{\text{total variance}} = \underbrace{\text{Var}(\mathbb{E}(y|x))}_{\text{explained variance}} + \underbrace{\mathbb{E}(\text{Var}(y|x))}_{\text{unexplained variance}}$$

provides a very general decomposition in explained and unexplained parts of the variance that is valid regardless of the form of the distributions  $F_{y,x}$  and  $F_{y|x}$ .

In regression it connects variance decomposition and conditioning. If you plug-in the conditional expectations for the multivariate normal model (cf. previous chapter) we recover

$$\underbrace{\sigma_y^2}_{\text{total variance}} = \underbrace{\sigma_y^2 \Omega^2}_{\text{explained variance}} + \underbrace{\sigma_y^2 (1 - \Omega^2)}_{\text{unexplained variance}}$$

## 18.2.2 Related quantities

Using the above three quantities (total variance, explained variance, and unexplained variance) we can construct a number of scores:

1) **coefficient of determination, squared multiple correlation:**

$$\frac{\text{explained var}}{\text{total var}} = \frac{\sigma_y^2 \Omega^2}{\sigma_y^2} = \Omega^2$$

(range 0 to 1, with 1 indicating perfect fit)

2) **coefficient of non-determination, coefficient of alienation:**

$$\frac{\text{unexplained var}}{\text{total var}} = \frac{\sigma_y^2 (1 - \Omega^2)}{\sigma_y^2} = 1 - \Omega^2$$

(range 0 to 1, with 0 indicating perfect fit)

3) **F score,  $t^2$  score:**

$$\frac{\text{explained var}}{\text{unexplained var}} = \frac{\sigma_y^2 \Omega^2}{\sigma_y^2 (1 - \Omega^2)} = \frac{\Omega^2}{1 - \Omega^2} = \mathcal{F} = \frac{\tau^2}{n}$$

(range 0 to  $\infty$ , with  $\infty$  indicating perfect fit)

Note that the  $\mathcal{F}$  and  $\tau^2$  scores are population versions of the  $F$  and  $t^2$  statistics.

Also note that  $\Omega^2 = \frac{\tau^2}{\tau^2 + n} = \frac{\mathcal{F}}{\mathcal{F} + 1}$  links squared correlation with squared  $t$ -scores and  $F$ -scores.

## 18.3 Sample version of variance decomposition

If  $\Omega^2$  and  $\sigma_y^2$  are replaced by their MLEs this can be written in a sample version as follows using data points  $y_i$ , predictions  $\hat{y}_i$  and  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{total sum of squares (TSS)}} = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\text{explained sum of squares (ESS)}} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\text{residual sum of squares (RSS)}}$$

Note that TSS, ESS and RSS all scale with  $n$ . Using data vector notation the sample-based variance decomposition can be written in form of the Pythagorean theorem:

$$\underbrace{\|\mathbf{y} - \bar{y}\mathbf{1}\|^2}_{\text{total sum of squares (TSS)}} = \underbrace{\|\hat{\mathbf{y}} - \bar{y}\mathbf{1}\|^2}_{\text{explained sum of squares (ESS)}} + \underbrace{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}_{\text{residual sum of squares (RSS)}}$$

### 18.3.1 Geometric interpretation of regression as orthogonal projection:

The above equation can be further simplified to

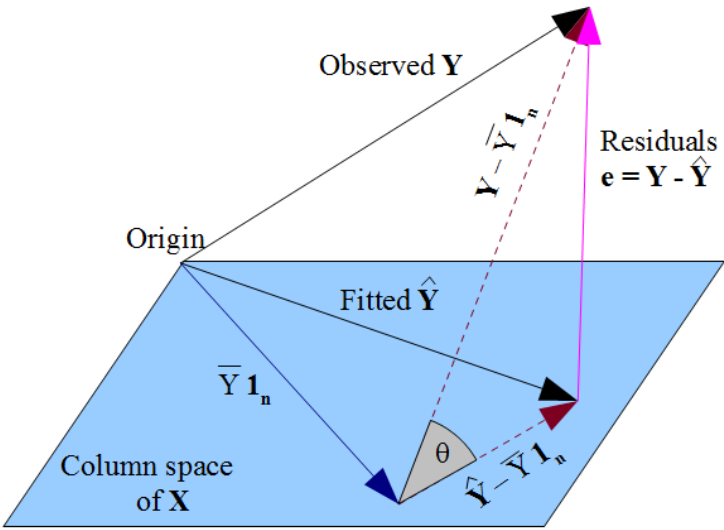
$$\|\mathbf{y}\|^2 = \|\hat{\mathbf{y}}\|^2 + \underbrace{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}_{\text{RSS}}$$

Geometrically speaking, this implies  $\hat{\mathbf{y}}$  is an orthogonal projection of  $\mathbf{y}$ , since the residuals  $\mathbf{y} - \hat{\mathbf{y}}$  and the predictions  $\hat{\mathbf{y}}$  are orthogonal (by construction!).

This also valid for the centered versions of the vectors, i.e.  $\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n$  is an orthogonal projection of  $\mathbf{y} - \bar{y}\mathbf{1}_n$  (see Figure).

Also note that the angle  $\theta$  between the two centered vectors is directly related to the (estimated) multiple correlation, with  $R = \cos(\theta) = \frac{\|\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n\|}{\|\mathbf{y} - \bar{y}\mathbf{1}_n\|}$ , or  $R^2 =$

$$\cos(\theta)^2 = \frac{\|\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n\|^2}{\|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2} = \frac{\text{ESS}}{\text{TSS}}.$$



Source of Figure: [Stack Exchange](#)

## Chapter 19

# Prediction and variable selection

In this chapter we discuss how to compute (lower bounds) of the prediction error and how to select variables relevant for prediction

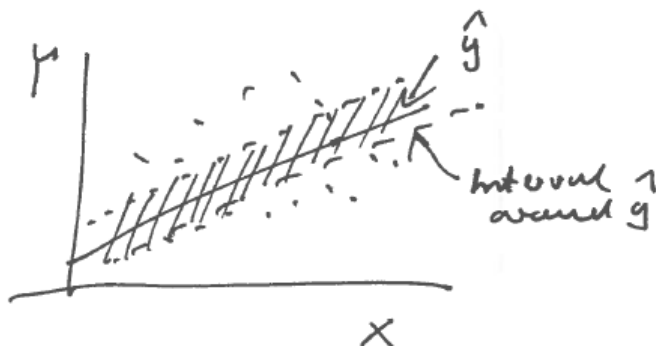
### 19.1 Prediction and prediction intervals

Learning the regression function from (training) data is only the first step in application of regression models.

The next step is to actually make **prediction** of future outcomes  $y^{\text{test}}$  given test data  $x^{\text{test}}$ :

$$y^{\text{test}} = \hat{y}(x^{\text{test}}) = \hat{f}_{\hat{\beta}_0, \hat{\beta}}(x^{\text{test}})$$

Note that  $\hat{y}^{\text{test}}$  is a point estimator. Is it possible also to construct a corresponding interval estimate?



The answer is yes, and leads back to the conditioning approach:

$$y^* = E(y|x) = \beta_0 + \beta^T x$$

$$\text{Var}(\varepsilon) = \text{Var}(y|x) = \sigma_y^2(1 - \Omega^2)$$

We know that the mean squared prediction error for  $y^*$  is  $E((y - y^*)^2) = \text{Var}(\varepsilon)$  and that this is the minimal irreducible error. Hence, we may use  $\text{Var}(\varepsilon)$  as the *minimum* variability for the prediction.

The corresponding prediction interval is

$$[y^*(x^{\text{test}}) \pm c \times \text{SD}(\varepsilon)]$$

where  $c$  is some suitable constant (e.g. 1.96 for symmetric 95% normal intervals).

However, please note that the prediction interval constructed in this fashion will be an *underestimate*. The reason is that this assumes that we employ  $y^* = \beta_0 + \beta^T x$  but in reality we actually use  $\hat{y} = \hat{\beta}_0 + \hat{\beta}^T x$  for prediction — note the estimated coefficients! We recall from an earlier chapter (best linear predictor) that this leads to increase of MSPE compared with using the optimal  $\beta_0$  and  $\beta$ .

Thus, for better prediction intervals we would need to consider the mean squared prediction error of  $\hat{y}$  that can be written as  $E((y - \hat{y})^2) = \text{Var}(\varepsilon) + \delta$  where  $\delta$  is an **additional error term due to using an estimated rather than the true regression function**.  $\delta$  typically declines with  $1/n$  but can be substantial for small  $n$  (in particular as it usually depends on the number of predictors  $d$ ).

For more details on this we refer to later modules on regression.

## 19.2 Variable importance and prediction

Another key question in regression modeling is to find out predictor variables  $x_1, x_2, \dots, x_d$  are actually important for predicting the outcome  $y$ .

→ We need to study variable importance measures (VIM).

### 19.2.1 How to quantify variable importance?

A variable  $x_i$  is **important** if it **improves prediction** of the response  $y$ .

Recall variance decomposition:

$$\text{Var}(y) = \sigma_y^2 = \underbrace{\sigma_y^2 \Omega^2}_{\text{explained variance}} + \underbrace{\sigma_y^2(1 - \Omega^2)}_{\text{unexplained/residual variance} = \text{Var}(\varepsilon)}$$

- $\Omega^2$  squared multiple correlation  $\in [0, 1]$
- $\Omega^2$  large  $\rightarrow$  1 predictor variables explain most of  $\sigma_y^2$
- $\Omega^2$  small  $\rightarrow$  0 linear model fails and predictors do not explain variability
- $\Rightarrow$  If a predictor helps to increase explained variance  
decrease unexplained variance then it is important!
- $\Omega^2 = \mathbf{P}_{yx} \mathbf{P}_{xx}^{-1} \mathbf{P}_{xy} \triangleq$  a function of the  $X$ !

VIM: which predictors contribute most to  $\Omega^2$

### 19.2.2 Some candidates for VIMs

#### 1. The regression coefficients $\beta$

- $\beta = \Sigma_{xx}^{-1} \Sigma_{xy} = V_x^{-1/2} \mathbf{P}_{xx}^{-1} \mathbf{P}_{xy} \sigma_y$
- Not a good VIM since  $\beta$  contains the scale!
- Large  $\hat{\beta}_i$  does not indicate that  $x_i$  is important.
- Small  $\hat{\beta}_i$  does not indicate that  $x_i$  is not important.

#### 2. Standardised regression coefficients $\beta_{\text{std}}$

- $\beta_{\text{std}} = \mathbf{P}_{xx}^{-1} \mathbf{P}_{xy}$
- implies  $\text{Var}(y) = 1, \text{Var}(x_i) = 1$
- These do not contain the scale (so better than  $\hat{\beta}$ )
- But still unclear how this relates to decomposition of variance

#### 3. Squared marginal correlations $\rho_{y,x_i}^2$

Consider case of uncorrelated predictors:  $\mathbf{P}_{xx} = \mathbf{I}$  (no correlation among  $x_i$ )

$$\Rightarrow \Omega^2 = \mathbf{P}_{yx} \mathbf{P}_{xy} = \sum_{i=1}^d \rho_{y,x_i}^2$$

$\rho_{y,x_i}^2 = \text{Cor}(y, x_i)$  is the marginal correlation between  $y$  and  $x_i$ , and  $\Omega^2$  is (for uncorrelated predictors) the sum of squared marginal correlations.

- If  $\mathbf{P}_{xx} = \mathbf{I}$ , then *ranking* predictors by  $\rho_{y,x_i}^2$  is optimal!
- The predictor with largest marginal correlation reduces the unexplained variance most!
- good news: even if there is weak correlation among predictors the marginal correlations are still good as VIM (but then they will not perfectly add up to  $\Omega^2$ )
- Advantage: very simple but often also very effective.
- Caution! If there is strong correlation in  $\mathbf{P}_{xx}$ , then there is **colinearity** (in this case it is often best to remove one of the strongly correlated variables, or to merge the correlated variables).

Often, ranking predictors by their squared marginal correlations is done as a prefiltering step (independence screening).

## 19.3 Regression $t$ -scores.

### 19.3.1 Wald statistic for regression coefficients

So far, we discussed three obvious candidates for variable importance measures (regression coefficients, standardised regression coefficients, marginal correlations).

In this section we consider a further quantity, the **regression  $t$ -score**:

Recall that ML estimation of the regression coefficients yields

- a point estimate  $\hat{\beta}$
- the (asymptotic) variance  $\widehat{\text{Var}}(\hat{\beta})$
- the asymptotic normal distribution  $\hat{\beta} \stackrel{a}{\sim} N_d(\beta, \widehat{\text{Var}}(\hat{\beta}))$

Corresponding to each predictor  $x_i$  we can construct from the above a  $t$ -score

$$t_i = \frac{\hat{\beta}_i}{\widehat{\text{SD}}(\hat{\beta}_i)}$$

where the standard deviations are computed by  $\widehat{\text{SD}}(\hat{\beta}_i) = \text{Diag}(\widehat{\text{Var}}(\hat{\beta}))_i$ . This corresponds to the **Wald statistic** to test that the underlying true regression coefficient is zero ( $\beta_i = 0$ ).

Correspondingly, under the null hypothesis that  $\beta_i = 0$  asymptotically for large  $n$  the regression  $t$ -score is standard normal distributed:

$$t_i \stackrel{a}{\sim} N(0, 1)$$

This allows to compute (symmetric)  $p$ -values  $p = 2\Phi(-|t_i|)$  where  $\Phi$  is the standard normal distribution function.

For finite  $n$ , assuming normality of the observation and using the unbiased estimate for variance when computing  $t_i$ , the exact distribution of  $t_i$  is given by the Student- $t$  distribution:

$$t_i \sim t_{n-d-1}$$

Regression  $t$ -scores can thus be used to test whether a regression coefficient is zero. A large magnitude of the  $t_i$  score indicates that the hypothesis  $\beta_i = 0$  can be rejected. Thus, a small  $p$ -value (say smaller than 0.05) signals that the regression coefficient is non-zero and hence that the corresponding predictor variable should be included in the model.

This allows rank predictor variables by  $|t_i|$  or the corresponding  $p$ -values with regard to their relevance in the linear model. Typically, in order to simplify a



model, predictors with the largest  $p$ -values (and thus smallest absolute  $t$ -scores) may be removed from a model. However, note that having a  $p$ -value say larger than 0.05 by itself is not sufficient to declare a regression coefficient to be zero (because in classical statistical testing you can only reject the null hypothesis, but not accept it!).

Note that by construction the regression  $t$ -scores do not depend on the scale, so when the original data are rescaled this will not affect the corresponding regression  $t$ -scores. Furthermore, if  $\widehat{SD}(\hat{\beta}_i)$  is small, then the regression  $t$ -score  $t_i$  can still be large even if  $\hat{\beta}_i$  is small!

### 19.3.2 Computing

When you perform regression analysis in R (or another statistical software package) the computer will return the following:

$\hat{\beta}_i$	$\widehat{SD}(\hat{\beta}_i)$	$t_i = \frac{\hat{\beta}_i}{\widehat{SD}(\hat{\beta}_i)}$	p-values	Indicator of
Estimated	Error of	t-score	for $t_i$	Significance
repression	$\hat{\beta}_i$	computed from	based on t-distribution	* 0.9
coefficient		first two columns		** 0.95
				*** 0.99

In the `lm()` function in R the standard deviation is the square root of the unbiased estimate of the variance (but note that it itself is not unbiased!).

### 19.3.3 Connection with partial correlation

The deeper reason why ranking predictors by regression  $t$ -scores and associated  $p$ -values is useful is their link with **partial correlation**.

In particular, the (squared) regression  $t$ -score can be 1:1 transformed into the (estimated) (squared) partial correlation

$$\hat{\rho}_{y,x_i|x_{j \neq i}}^2 = \frac{t_i^2}{t_i^2 + df}$$

with  $df = n - d - 1$ , and it can be shown that the  $p$ -values for testing that  $\beta_i = 0$  are exactly the same as the  $p$ -values for testing that the partial correlation  $\rho_{y,x_i|x_{j \neq i}}$  vanishes!

Therefore, ranking the predictors  $x_i$  by regression  $t$ -scores leads to exactly the same ranking and  $p$ -values as partial correlation!

### 19.3.4 Squared Wald statistic and the $F$ statistic

In the above we looked at individual regression coefficients. However, we can also construct a Wald test using the complete vector  $\hat{\beta}$ . The squared Wald statistic to test that  $\beta = 0$  is given by

$$\begin{aligned} t^2 &= \hat{\beta}^T \widehat{\text{Var}}(\hat{\beta}^{-1}) \hat{\beta} \\ &= \left( \hat{\Sigma}_{yx} \hat{\Sigma}_{xx}^{-1} \right) \left( \frac{n}{\widehat{\sigma}_\varepsilon^2} \hat{\Sigma}_{xx} \right) \left( \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy} \right) \\ &= \frac{n}{\widehat{\sigma}_\varepsilon^2} \hat{\Sigma}_{yx} \hat{\Sigma}_{xx}^{-1} \hat{\Sigma}_{xy} \\ &= \frac{n}{\widehat{\sigma}_\varepsilon^2} \hat{\sigma}_y^2 R^2 \end{aligned}$$

With  $\widehat{\sigma}_\varepsilon^2 / \hat{\sigma}_y^2 = 1 - R^2$  we finally get the related  $F$  statistic

$$\frac{t^2}{n} = \frac{R^2}{1 - R^2} = F$$

which is a function of  $R^2$ . If  $R^2 = 0$  then  $F = 0$ . If  $R^2$  is large ( $< 1$ ) then  $F$  is large as well ( $< \infty$ ) and the null hypothesis  $\beta = 0$  can be rejected, which implies that at least one regression coefficient is non-zero. Note that the squared Wald statistic  $t^2$  is asymptotically  $\chi_d^2$  distributed which is useful to find critical values and to compute  $p$ -values.

## 19.4 Further approaches for variable selection

In addition to ranking by marginal and partial correlation, there are many other approaches for variable selection in regression!

a) Search-based methods:

- search through subsets of linear models for  $d$  variables, ranging from full model (including all predictors) to the empty model (includes no predictor) and everything inbetween.
- Problem: exhaustive search not possible even for relatively small  $d$  as space of models is very large!
- Therefore heuristic approaches such as forward selection (adding predictors), backward selection (removing predictors), or monte-carlo random search are employed.
- Problem: maximum likelihood cannot be used for choosing among the models - since ML will always pick the best model. Therefore, penalised ML criteria such as AIC or Bayesian criteria are often employed instead.

b) Integrative estimation and variable selection:

- there are methods that fit the regression model and perform variable selection *simultaneously*.
- the most well-known approach of this type is “lasso” regression (Tibshirani 1996)
- This applies a (generalised) linear model with ML plus L1 penalty.
- Alternative: Bayesian variable selection and estimation procedures

c) Entropy-based variable selection:

As seen above, two of the most popular approaches in linear models are based on correlation, either marginal correlation or partial correlation (via regression  $t$ -scores).

Correlation measures can be generalised to non-linear settings. One very popular measure is the **mutual information** which is computed using the KL divergence. In case of two variables  $x$  and  $y$  with joint normal distribution and correlation  $\rho$  the mutual information is a function of the correlation:

$$\text{MI}(x, y) = \frac{1}{2} \log(1 - \rho^2)$$

In regression the mutual information between the response  $y$  and predictor  $x_i$  is  $\text{MI}(y, x_i)$ , and this is widely used for feature selection, in particular in machine learning.



# Appendix



# Appendix A

## Refresher

Statistics is a mathematical science that requires practical use of tools from probability, vector and matrices, analysis etc.

Here we briefly list some essentials that are needed for “Statistical Methods”. Please familiarise yourself (again) with these topics.

### A.1 Basic mathematical notation

Summation:

$$\sum_{i=1}^n x_i = x_1 + x_2 + \dots + x_n$$

Multiplication:

$$\prod_{i=1}^n x_i = x_1 \times x_2 \times \dots \times x_n$$

### A.2 Vectors and matrices

Vector and matrix notation.

Vector algebra.

Eigenvectors and eigenvalues for a real symmetric matrix.

Eigenvalue (spectral) decomposition of a real symmetric matrix.

Positive and negative definiteness of a real symmetric matrix (containing only positive or only negative eigenvalues).

Singularity of a real symmetric matrix (containing one or more eigenvalues identical to zero).

Singular value decomposition of a real matrix.

## A.3 Functions

### A.3.1 Gradient

The **nabla operator** (also known as **del operator**) is the *row* vector

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

containing the first order partial derivative operators.

The **gradient** of a scalar-valued function  $h(\mathbf{x})$  with vector argument  $\mathbf{x} = (x_1, \dots, x_d)^T$  is also a *row* vector (with  $d$  columns) and can be expressed using the nabla operator

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

Note the various notations for the gradient.

**Example A.1.** Examples for the gradient:

- $h(\mathbf{x}) = \mathbf{a}^T \mathbf{x} + b$ . Then  $\nabla h(\mathbf{x}) = \frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{a}^T$ .
- $h(\mathbf{x}) = \mathbf{x}^T \mathbf{x}$ . Then  $\nabla h(\mathbf{x}) = \frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} = 2\mathbf{x}^T$ .
- $h(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$ . Then  $\nabla h(\mathbf{x}) = \frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T)$ .

### A.3.2 Hessian matrix

The matrix of all second order partial derivatives of scalar-valued function with vector-valued argument is called the **Hessian matrix** and is computed by double application of the nabla operator:

$$\nabla^T \nabla h(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 h(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 h(\mathbf{x})}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 h(\mathbf{x})}{\partial x_1 \partial x_d} \\ \frac{\partial^2 h(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 h(\mathbf{x})}{\partial x_2^2} & \dots & \frac{\partial^2 h(\mathbf{x})}{\partial x_2 \partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 h(\mathbf{x})}{\partial x_d \partial x_1} & \frac{\partial^2 h(\mathbf{x})}{\partial x_d \partial x_2} & \dots & \frac{\partial^2 h(\mathbf{x})}{\partial x_d^2} \end{pmatrix} = \left( \frac{\partial h(\mathbf{x})}{\partial x_i \partial x_j} \right) = \left( \frac{\partial}{\partial \mathbf{x}} \right)^T \frac{\partial h(\mathbf{x})}{\partial \mathbf{x}}.$$

By construction it is square and symmetric.

### A.3.3 Convex and concave functions

A function  $h(\mathbf{x})$  is convex if the second derivative  $h''(\mathbf{x}) \geq 0$  for all  $\mathbf{x}$ . More generally, a function  $h(\mathbf{x})$  is convex if the Hessian matrix  $\nabla^T \nabla h(\mathbf{x})$  is positive definite, i.e. if it contains only positive eigenvalues.



If  $h(\mathbf{x})$  is convex, then  $-h(\mathbf{x})$  is *concave*. A function is concave if the Hessian matrix is negative definite.

**Example A.2.** The logarithm  $\log(x)$  is an example of a concave function whereas  $x^2$  is a convex function.

To memorise, a *valley* is convex.

### A.3.4 Linear and quadratic approximation

Taylor series of first / second order.

Applied to scalar-valued function of a scalar:

$$h(x) \approx h(x_0) + h'(x_0)(x - x_0) + \frac{1}{2}h''(x_0)(x - x_0)^2$$

With  $x = x_0 + \varepsilon$  this can be written as

$$h(x_0 + \varepsilon) \approx h(x_0) + h'(x_0)\varepsilon + \frac{1}{2}h''(x_0)\varepsilon^2$$

Applied to scalar-valued function of a vector:

$$h(\mathbf{x}) \approx h(\mathbf{x}_0) + \nabla h(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \nabla^T \nabla h(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)$$

With  $\mathbf{x} = \mathbf{x}_0 + \varepsilon$  this can be written as

$$h(\mathbf{x}_0 + \varepsilon) \approx h(\mathbf{x}_0) + \nabla h(\mathbf{x}_0)\varepsilon + \frac{1}{2}\varepsilon^T \nabla^T \nabla h(\mathbf{x}_0)\varepsilon$$

**Example A.3.** Commonly occurring Taylor series approximations of second order are for example

$$\log(x_0 + \varepsilon) \approx \log(x_0) + \frac{\varepsilon}{x_0} - \frac{\varepsilon^2}{2x_0^2}$$

and

$$\frac{x_0}{x_0 + \varepsilon} \approx 1 - \frac{\varepsilon}{x_0} + \frac{\varepsilon^2}{x_0^2}$$

### A.3.5 Conditions for local optimum of a function

To check if  $x_0$  or  $\mathbf{x}_0$  is a local maximum or minimum we can use the following conditions:

For a function of a single variable:

- i) First derivative is zero at optimum  $h'(x_0) = 0$ .

- ii) If the second derivative  $h''(x_0) < 0$  at the optimum is negative the function is locally concave and the optimum is a maximum.
- iii) If the second derivative  $h''(x_0) > 0$  is positive at the optimum the function is locally convex and the optimum is a minimum.

For a function of several variables:

- i) Gradient vanishes at maximum,  $\nabla h(x_0) = 0$ .
- ii) If the Hessian  $\nabla^T \nabla h(x_0)$  is negative definite (= all eigenvalues of Hessian matrix are negative) then the function is locally concave and the optimum is a maximum.
- iii) If the Hessian is positive definite (= all eigenvalues of Hessian matrix are positive) then the function is locally convex and the optimum is a minimum.

Around the local optimum  $x_0$  we can approximate the function quadratically using

$$h(x_0 + \epsilon) \approx h(x_0) + \frac{1}{2} \epsilon^T \nabla^T \nabla h(x_0) \epsilon$$

Note the linear term is missing due to the gradient being zero at  $x_0$ .

### A.3.6 Functions of matrices

Matrix inverse, matrix square root etc. of symmetric real matrices.

Computation via eigenvalue decomposition i.e. apply function such as inverse, sqrt etc. on the eigenvalues.

In this course we do not actually compute matrix functions, but we will use matrix notation for matrix square roots, so you do need to know that it exists and that it is not the same as taking the square root of the matrix entries.

Trace and determinant of a square matrix.

Connection with eigenvalues (trace = sum of eigenvalues, determinant = product of eigenvalues).

## A.4 Combinatorics

### A.4.1 Number of permutations

The number of possible orderings, or permutations, of  $n$  distinct items is the number of ways to put  $n$  items in  $n$  bins with exactly one item in each bin. It is given by the factorial

$$n! = \prod_{i=1}^n i = 1 \times 2 \times \dots \times n$$

where  $n$  is a positive integer. For  $n = 0$  the factorial is defined as

$$0! = 1$$

as there is exactly one permutation of zero objects.

The factorial can also be obtained using the [Gamma function](#)

$$n! = \Gamma(n + 1)$$

which can be viewed as continuous version of the factorial.

### A.4.2 Multinomial and binomial coefficient

The number of possible permutation of  $n$  items of  $K$  distinct types, with  $n_1$  of type 1,  $n_2$  of type 2 and so on, equals the number of ways to put  $n$  items into  $K$  bins with  $n_1$  items in the first bin,  $n_2$  in the second and so on. It is given by the **multinomial coefficient**

$$\binom{n}{n_1, \dots, n_K} = \frac{n!}{n_1! \times n_2! \times \dots \times n_K!}$$

with  $\sum_{k=1}^K n_k = n$  and  $K \leq n$ . Note that it equals the number of permutation of all items divided by the number of permutations of the items in each bin (or of each type).

If all  $n_k = 1$  and hence  $K = n$  the multinomial coefficient reduces to the factorial.

If there are only two bins / types ( $K = 2$ ) the multinomial coefficients becomes the **binomial coefficient**

$$\binom{n}{n_1} = \binom{n}{n_1, n - n_1} = \frac{n!}{n_1!(n - n_1)!}$$

which counts the number of ways to choose  $n_1$  elements from a set of  $n$  elements.

### A.4.3 De Moivre-Sterling approximation of the factorial

The factorial is frequently approximated by the following formula derived by [Abraham de Moivre \(1667–1754\)](#) and [James Stirling \(1692–1770\)](#)

$$n! \approx \sqrt{2\pi n} n^{n+\frac{1}{2}} e^{-n}$$

or equivalently on logarithmic scale

$$\log n! \approx \left(n + \frac{1}{2}\right) \log n - n + \frac{1}{2} \log(2\pi)$$

The approximation is good for small  $n$  (but fails for  $n = 0$ ) and becomes more and more accurate with increasing  $n$ . For large  $n$  the approximation can be simplified to

$$\log n! \approx n \log n - n$$

## A.5 Probability

### A.5.1 Random variables

A **random variable** describes a random experiment. The set of possible outcomes is the **sample space** or **state space** and is denoted by  $\Omega = \{\omega_1, \omega_2, \dots\}$ . The outcomes  $\omega_i$  are the **elementary events**. The sample space  $\Omega$  can be finite or infinite. Depending on type of outcomes the random variable is **discrete** or **continuous**.

An event  $A \subseteq \Omega$  is subset of  $\Omega$  and thus itself a set of elementary events  $A = \{a_1, a_2, \dots\}$ . This includes as special cases the full set  $A = \Omega$ , the empty set  $A = \emptyset$ , and the elementary events  $A = \omega_i$ . The complementary event  $A^C$  is the complement of the set  $A$  in the set  $\Omega$  so that  $A^C = \Omega \setminus A = \{\omega_i \in \Omega : \omega_i \notin A\}$ .

The probability of an event is denoted by  $\Pr(A)$ . We assume that

- $\Pr(A) \geq 0$ , probabilities are positive,
- $\Pr(\Omega) = 1$ , the certain event has probability 1, and
- $\Pr(A) = \sum_{a_i \in A} \Pr(a_i)$ , the probability of an event equals the sum of its constituting elementary events  $a_i$ .

This implies

- $\Pr(A) \leq 1$ , i.e. probabilities all lie in the interval  $[0, 1]$
- $\Pr(A^C) = 1 - \Pr(A)$ , and
- $\Pr(\emptyset) = 0$

Assume now we have two events  $A$  and  $B$ . The probability of the event “ $A$  and  $B$ ” is then given by the probability of the set intersection  $\Pr(A \cap B)$ . Likewise the probability of the event “ $A$  or  $B$ ” is given by the probability of the set union  $\Pr(A \cup B)$ .

From the above it is clear that probability theory is closely linked to set theory, and in particular to measure theory. This allows for an unified treatment of discrete and continuous random variables (an elegant framework but not needed for this module).

### A.5.2 Probability mass and density function and distribution and quantile function

To describe a random variable  $x$  we need to assign probabilities to the corresponding elementary outcomes  $x \in \Omega$ . For convenience we use the same name to denote the random variable and the elementary outcomes.

For a discrete random variable we employ a probability mass function (PMF). We denote the it by a lower case  $f$  but occasionally we also use  $p$  or  $q$ . In the discrete case we can define the event  $A = \{x : x = a\} = \{a\}$  and obtain the

probability directly from the PMF:

$$\Pr(A) = \Pr(x = a) = f(a).$$

The PMF has the property that  $\sum_{x \in \Omega} f(x) = 1$  and that  $f(x) \in [0, 1]$ .

For continuous random variables we need to use a probability density function (PDF) instead. We define the event  $A = \{x : a < x \leq a + da\}$  as an infinitesimal interval and then assign the probability

$$\Pr(A) = \Pr(a < x \leq a + da) = f(a)da.$$

The PDF has the property that  $\int_{x \in \Omega} f(x)dx = 1$  but in contrast to a PMF the density  $f(x) \geq 0$  may take on values larger than 1.

Assuming an ordering we can define the event  $A = \{x : x \leq a\}$  and compute its probability

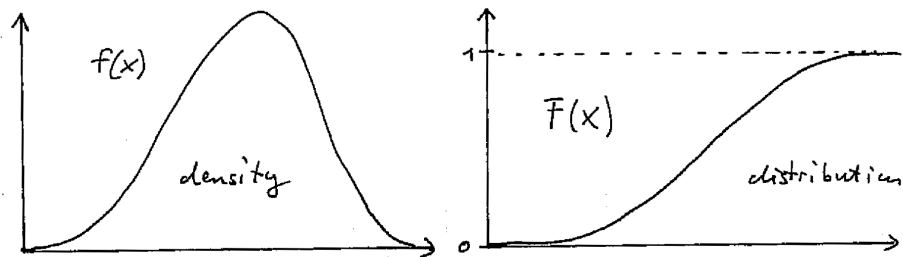
$$F(a) = \Pr(A) = \Pr(x \leq a) = \begin{cases} \sum_{x \in A} f(x) & \text{discrete case} \\ \int_{x \in A} f(x)dx & \text{continuous case} \end{cases}$$

This is known as the **distribution function**, or **cumulative distribution function** (CDF) and is denoted by upper case  $F$  if the corresponding PDF/PMF is  $f$  (or  $P$  and  $Q$  if the corresponding PDF/PMF are  $p$  and  $q$ ). By construction the distribution function is monotonically increasing and its value ranges from 0 to 1. With its help we can compute the probability of general interval sets such as

$$\Pr(a < x \leq b) = F(b) - F(a).$$

The inverse of the distribution function  $y = F(x)$  is the **quantile function**  $x = F^{-1}(y)$ . The 50% quantile  $F^{-1}\left(\frac{1}{2}\right)$  is the **median**.

If the random variable  $x$  has distribution function  $F$  we write  $x \sim F$ .



### A.5.3 Expection and variance of a random variable

The expected value  $E(x)$  of a random variable is defined as weighted average over all possible outcomes:

$$E(x) = \begin{cases} \sum_{x \in \Omega} x f(x) & \text{discrete case} \\ \int_{x \in \Omega} x f(x)dx & \text{continuous case} \end{cases}$$

The expectation is not necessarily always defined for a continuous random variable as the integral can diverge.

The expected value of a function of a random variable  $h(x)$  is obtained similarly:

$$E(h(x)) = \begin{cases} \sum_{x \in \Omega} h(x)f(x) & \text{discrete case} \\ \int_{x \in \Omega} h(x)f(x)dx & \text{continuous case} \end{cases}$$

This is called the “**law of the unconscious statistician**”, or short LOTUS.

For an event  $A$  we can define a corresponding **indicator function**

$$1_A(x) = \begin{cases} 1 & x \in A \\ 0 & x \notin A \end{cases}$$

Intriguingly,

$$E(1_A(x)) = \Pr(A)$$

i.e. the expectation of the indicator variable for  $A$  is the probability of  $A$ .

The moments of random variables are also defined by expectation:

- Zeroth moment:  $E(x^0) = 1$  by definition of PDF and PMF,
- First moment:  $E(x^1) = E(x) = \mu$ , the mean,
- Second moment:  $E(x^2)$
- The variance is the second moment centered about the mean  $\mu$ :

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

- The variance can also be computed by  $\text{Var}(x) = E(x^2) - E(x)^2$ .

A distribution does not necessarily need to have any finite first or higher moments. An example is the **Cauchy distribution** that does not have a mean or variance (or any other higher moment).

### A.5.4 Transformation of random variables

Linear transformation of random variables: if  $a$  and  $b$  are constants and  $x$  is a random variable, then the random variable  $y = a + bx$  has mean  $E(y) = a + bE(x)$  and variance  $\text{Var}(y) = b^2\text{Var}(x)$ .

For a general invertible coordinate transformation  $y = h(x) = y(x)$  the back-transformation is  $x = h^{-1}(y) = x(y)$ .

The transformation of the infinitesimal volume element is  $dy = \left| \frac{dy}{dx} \right| dx$ .

The transformation of the density is  $f_y(y) = \left| \frac{dx}{dy} \right| f_x(x(y))$ .

Note that  $\left| \frac{dx}{dy} \right| = \left| \frac{dy}{dx} \right|^{-1}$ .

### A.5.5 Law of large numbers:

- By the strong law of large numbers the empirical distribution  $\hat{F}_n$  converges to the true underlying distribution  $F$  as  $n \rightarrow \infty$  almost surely:

$$\hat{F}_n \xrightarrow{a.s.} F$$

The Glivenko–Cantelli theorem asserts that the convergence is uniform. Since the strong law implies the weak law we also have convergence in probability:

$$\hat{F}_n \xrightarrow{P} F$$

- Correspondingly, for  $n \rightarrow \infty$  the average  $E_{\hat{F}_n}(h(X)) = \frac{1}{n} \sum_{i=1}^n h(x_i)$  converges to the expectation  $E_F(h(X))$ .

### A.5.6 Jensen’s inequality

$$E(h(x)) \geq h(E(x))$$

for a *convex* function  $h(x)$ .

Recall: a convex function (such as  $x^2$ ) has the shape of a “valley”.

## A.6 Distributions

### A.6.1 Bernoulli and Binomial distribution

The Bernoulli distribution  $\text{Ber}(p)$  is simplest distribution possible. It is named after [Jacob Bernoulli \(1655-1705\)](#) who also invented the law of large numbers.

It describes a discrete binary random variable with two states  $x = 0$  (“failure”) and  $x = 1$  (“success”), where the parameter  $p \in [0, 1]$  is the probability of “success”. Often the Bernoulli distribution is also referred to as “coin tossing” model with the two outcomes “heads” and “tails”.

Correspondingly, the probability mass function of  $\text{Ber}(p)$  is

$$f(x = 0) = \Pr(\text{"failure"}) = 1 - p$$

and

$$f(x = 1) = \Pr(\text{"success"}) = p$$

A compact way to write the PMF of the Bernoulli distribution is

$$f(x|p) = p^x(1 - p)^{1-x}$$

If a random variable  $x$  follows the Bernoulli distribution we write

$$x \sim \text{Ber}(p).$$

The expected value is  $E(x) = p$  and the variance is  $\text{Var}(x) = p(1 - p)$ .

Closely related to the Bernoulli distribution is the Binomial distribution  $\text{Bin}(m, p)$  which results from repeating a Bernoulli experiment  $m$  times and counting the number of successes among the  $m$  trials (without keeping track of the ordering of the experiments).

Its probability mass function is:

$$f(x|p) = \binom{m}{x} p^x (1 - p)^{m-x}$$

for  $x = 0, 1, 2, \dots, m$ . The Binomial coefficient  $\binom{m}{x}$  is needed to account for the multiplicity of ways (orderings of samples) in which we can observe  $x$  successes.

The expected value is  $E(x) = mp$  and the variance is  $\text{Var}(x) = mp(1 - p)$ .

If a random variable  $x$  follows the Binomial distribution we write

$$x \sim \text{Bin}(m, p)$$

For  $m = 1$  it reduces to the Bernoulli distribution  $\text{Ber}(p)$ .

In R the PMF of the Binomial distribution is called `dbinom()`. The Binomial coefficient itself is computed by `choose()`.

## A.6.2 Normal distribution

Univariate normal distribution:

$$x \sim N(\mu, \sigma^2) \text{ with } E(x) = \mu \text{ and } \text{Var}(x) = \sigma^2.$$

Probability density function (PDF):

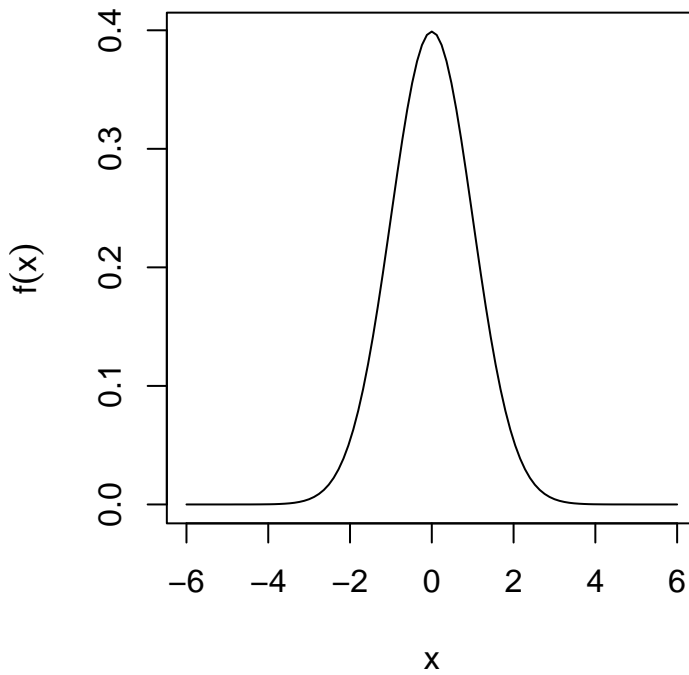
$$f(x|\mu, \sigma^2) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

In R the density function is called `dnorm()`.

The standard normal distribution is  $N(0, 1)$  with mean 1 and variance 1.

Plot of the PDF of the standard normal:

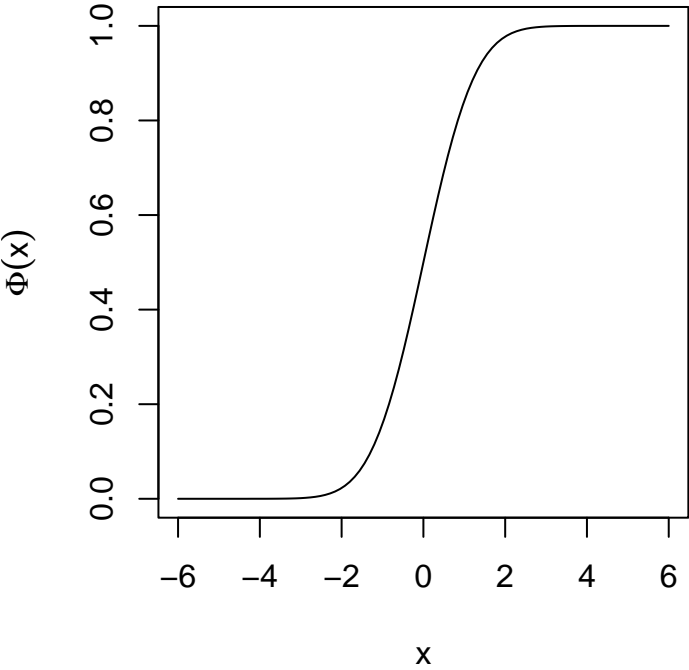




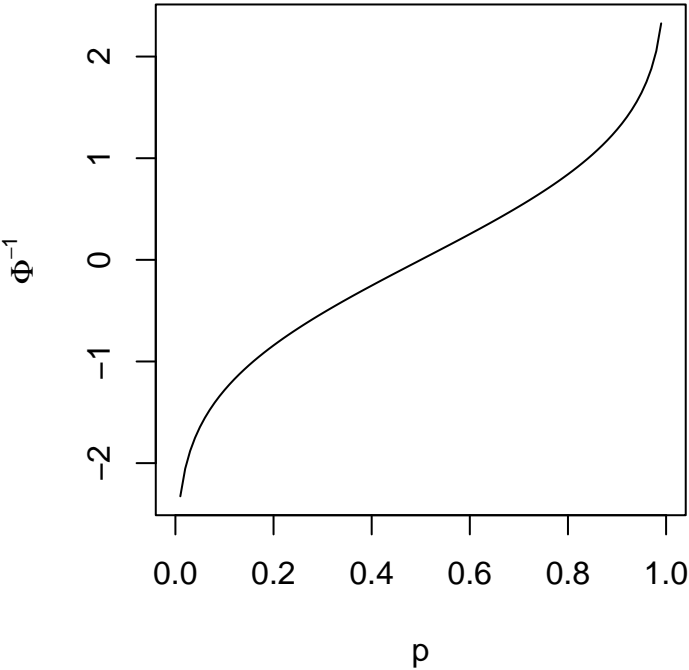
The cumulative distribution function (CDF) of the standard normal  $N(0,1)$  is

$$\Phi(x) = \int_{-\infty}^x f(x' | \mu = 0, \sigma^2 = 1) dx'$$

There is no analytic expression for  $\Phi(x)$ . In R the function is called `pnorm()`.



The inverse  $\Phi^{-1}(p)$  is called the quantile function of the standard normal. In R the function is called `qnorm()`.



The sum of two normal random variables is also normal (with the appropriate mean and variance).

### A.6.3 Scaled chi-squared / Wishart / gamma distribution and exponential distribution

Assume  $m$  independent normal random variables

$$z_1, z_2, \dots, z_m \sim N(0, \sigma^2)$$

Then the sum of the squares

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

follows a **scaled chi-squared distribution**

$$x \sim \sigma^2 \chi_m^2$$

with degree of freedom  $m$  and  $x \geq 0$ . The mean and variance of a scaled chi-squared distributed variable is  $E(x) = m\sigma^2$  and  $\text{Var}(x) = 2m\sigma^4$ .

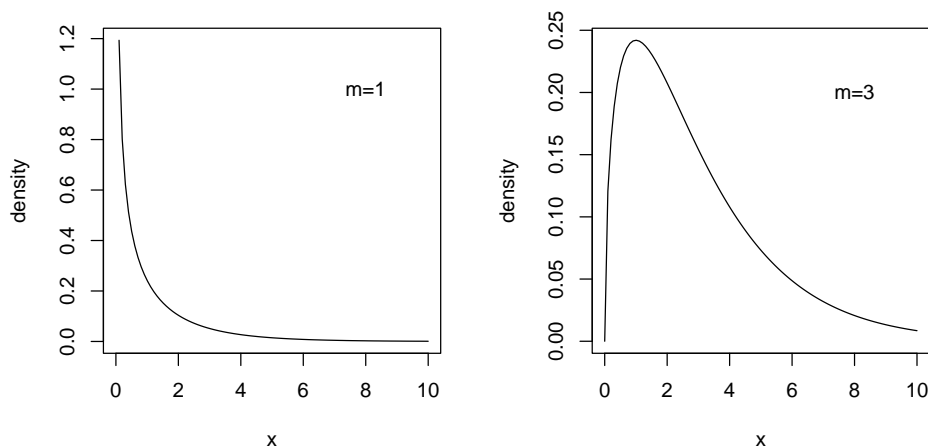
Another name for the scaled chi-squared distribution is **univariate Wishart distribution**  $W_1(\sigma^2, m)$  which uses the same parameters.

The **gamma distribution**  $\text{Gam}(\alpha, \beta)$  is a further variant of the scaled chi-squared distribution which uses a different parameterisation in terms of a shape parameter  $\alpha$  and a scale parameter  $\beta$ . The scaled chi-squared distribution  $\sigma^2 \chi_m^2$  equals  $\text{Gam}(\frac{m}{2}, 2\sigma^2)$ . The mean of  $\text{Gam}(\alpha, \beta)$  is  $\alpha\beta$  and its variance is  $\alpha\beta^2$ .

The **chi-squared distribution** is a special case with  $\sigma^2 = 1$  with mean  $E(x) = m$  and variance  $\text{Var}(x) = 2m$ . The chi-squared distribution  $\chi_m^2$  equals  $\text{Gam}(\frac{m}{2}, 2)$ .

The **exponential distribution**  $\text{Exp}(\beta)$  with scale parameter  $\beta$  (and mean  $\beta$  and variance  $\beta^2$ ) is another special case of the gamma distribution with shape parameter  $\alpha = 1$ . Instead of the scale parameter the exponential distribution is also often specified using a rate parameter  $\lambda = \frac{1}{\beta}$ .

Here is a plot of the density of the chi-squared distribution for degrees of freedom  $m = 1$  and  $m = 3$ :



In R the density of the chi-squared distribution is given by `dchisq()`. The cumulative density function is `pchisq()` and the quantile function is `qchisq()`.

The density of the gamma distribution (aka scaled chi-squared distribution) is available in the R function `dgamma()`. The cumulative density function is `pgamma()` and the quantile function is `qgamma()`.

## A.7 Statistics

### A.7.1 Statistical learning

The aim in statistics - data science - machine learning is to learn from data (from experiments, observations, measurements) to learn about and understand the world.

Specifically, to identify the best model(s) for the data in order to

- to explain the current data, and
- to enable good prediction of future data

Note that it is easy to get models that only explain the data but do not predict well!

This is called *overfitting* the data and happens in particular if the model is overparameterized for the amount of data available.

Specifically, we have data  $x_1, \dots, x_n$  and models  $f(x|\theta)$  that are indexed the parameter  $\theta$ .

Often (but not always)  $\theta$  can be interpreted and/or is associated with some property of the model.

If there is only a single parameter we write  $\theta$  (scalar parameter). For a parameter vector we write  $\theta$  (in bold type).

### A.7.2 Point and interval estimation

- There is a parameter  $\theta$  of interest in a model
- we are uncertain about this parameter (i.e. we don't know the exact value)
- we would like to learn about this parameter by observing data  $x_1, \dots, x_n$  from the model

Estimation:

- An **estimator for  $\theta$**  is a function  $\hat{\theta}(x_1, \dots, x_n)$  that maps the data (input) to a “guess” (output) about  $\theta$ .
- A **point estimator** provides a single number for each parameter
- An **interval estimator** provides a set of possible values for each parameter.

### A.7.3 Sampling properties of a point estimator $\hat{\theta}$

A point estimator  $\hat{\theta}$  depends on the data, hence it has sampling variation (i.e. estimate will be different for a new set of observations)

Thus  $\hat{\theta}$  can be seen as a random variable, and its distribution is called sampling distribution (across different experiments).

Properties of this distribution can be used to evaluate how far the estimator deviates (on average across different experiments) from the true value:

$$\begin{aligned}
 \text{Bias:} \quad \text{Bias}(\hat{\theta}) &= E(\hat{\theta}) - \theta \\
 \text{Variance:} \quad \text{Var}(\hat{\theta}) &= E((\hat{\theta} - E(\hat{\theta}))^2) \\
 \text{Mean squared error:} \quad \text{MSE}(\hat{\theta}) &= E((\hat{\theta} - \theta)^2) \\
 &= \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta})^2
 \end{aligned}$$

The last identity about MSE follows from  $E(X^2) = \text{Var}(X) + E(X)^2$ .

At first sight it seems desirable to focus on unbiased (for finite  $n$ ) estimators. However, requiring strict unbiasedness is not always a good idea!

In many situations it is better to allow for some small bias and in order to achieve a smaller variance and an overall total smaller MSE. This is called *bias-variance tradeoff* — as more bias is traded for smaller variance (or, conversely, less bias is traded for higher variance)

### A.7.4 Asymptotics

Typically, Bias, Var and MSE all decrease with increasing sample size so that with more data  $n \rightarrow \infty$  the errors become smaller and smaller.

The typical rate of decrease of variance of a good estimator is  $\frac{1}{n}$ . Thus, when sample size is doubled the variance is divided by 2 (and the standard deviation is divided by  $\sqrt{2}$ ).

Consistency:  $\hat{\theta}$  is called consistent if

$$\text{MSE}(\hat{\theta}) \longrightarrow 0 \text{ with } n \rightarrow \infty$$

Consistency implies we recover the true model in the limit of infinite data and if the model class contains the true model.

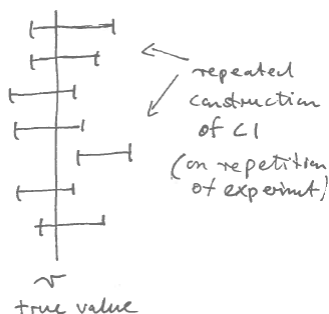
Consistency is a *minimum* essential requirement for any reasonable estimator! Of all consistent estimators we typically prefer the estimator that is most efficient (i.e. with fastest decrease in MSE) and that thus has smallest variance and/or MSE for given finite  $n$ .

Note that if the model class does not contain the true model then strict consistency cannot be achieved but we still wish to get at least as close as possible to the true model.

### A.7.5 Confidence intervals

- A **confidence interval (CI)** is an **interval estimate** with a **frequentist** interpretation.
- Definition of **coverage**  $\kappa$  of a CI: how often (in repeated identical experiment) does the estimated CI overlap the true parameter value  $\theta$ 
  - Eg.: Coverage  $\kappa = 0.95$  (95%) means that in 95 out of 100 case the estimated CI will contain the (unknown) true value (i.e. it will “cover”  $\theta$ ).

Illustration of the repeated construction of a CI for  $\theta$ :



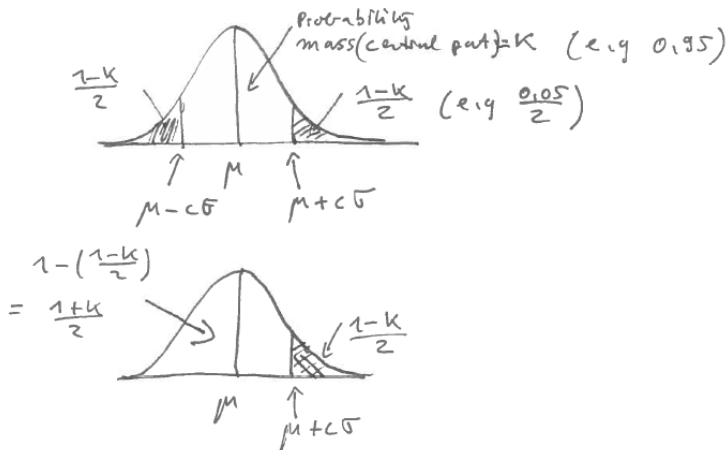
- Note that a CI is actually an **estimate**:  $\widehat{\text{CI}}(x_1, \dots, x_n)$ , i.e. it depends on data and has a random (sampling) variation.

- A good CI has high coverage and is compact.

**Note:** the coverage probability is **not** the probability that the true value is contained in a given estimated interval (that would be the Bayesian *Credible* Interval).

### A.7.6 Symmetric normal confidence interval

For a normally distributed univariate random variable it is straightforward to construct a symmetric two-sided CI with a given desired coverage  $\kappa$ .



For a normal random variable  $X \sim N(\mu, \sigma^2)$  with mean  $\mu$  and variance  $\sigma^2$  and density function  $f(x)$  we can compute the probability

$$\Pr(X \leq \mu + c\sigma) = \int_{-\infty}^{\mu + c\sigma} f(x) dx = \Phi(c) = \frac{1 + \kappa}{2}$$

Note  $\Phi(c)$  is the cumulative distribution function (CDF) of the standard normal  $N(0, 1)$ :

From the above we obtain the critical point  $c$  from the quantile function, i.e. by inversion of  $\Phi$ :

$$c = \Phi^{-1}\left(\frac{1 + \kappa}{2}\right)$$

The following table lists  $c$  for the three most commonly used values of  $\kappa$  - it is useful to memorise these values!

Coverage $\kappa$	Critical value $c$
0.9	1.64

Coverage $\kappa$	Critical value $c$
0.95	1.96
0.99	2.58

A **symmetric standard normal CI** with nominal coverage  $\kappa$  for

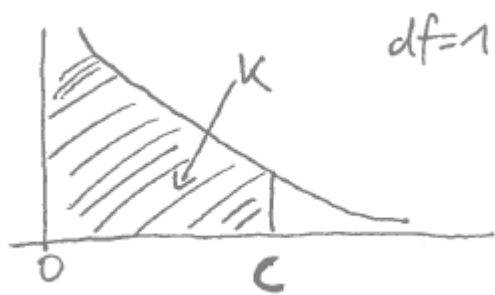
- a scalar parameter  $\theta$
- with normally distributed estimate  $\hat{\theta}$  and
- with estimated standard deviation  $\widehat{SD}(\hat{\theta}) = \hat{\sigma}$

is then given by

$$\widehat{CI} = [\hat{\theta} \pm c\hat{\sigma}]$$

where  $c$  is chosen for desired coverage level  $\kappa$ .

**A.7.7 Confidence interval for chi-squared distribution**



As for the normal CI we can compute critical values but for the chi-squared distribution we use a one-sided interval:

$$\Pr(X \leq c) = \kappa$$

As before we get  $c$  by the quantile function, i.e. by inverting the CDF of the chi-squared distribution.

The following list the critical values for the three most common choice of  $\kappa$  for  $m = 1$  (one degree of freedom):

Coverage $\kappa$	Critical value $c$ ( $m = 1$ )
0.9	2.71
0.95	3.84
0.99	6.63



A one-sided CI with nominal coverage  $\kappa$  is then given by  $[0, c]$ .



# Appendix B

## Further study

In this module we can only touch the surface of likelihood and Bayes inference. As a starting point for further reading the following text books are recommended.

### B.1 Recommended reading

- Held and Bové (2014) *Applied Statistical Inference: Likelihood and Bayes*. Springer.
- Faraway (2015) *Linear Models with R (second edition)*. Chapman and Hall/CRC.

### B.2 Additional references

- Wood (2015) *Core Statistics*. Cambridge University Press.
- Gelman et al. (2014) *Bayesian data analysis (3rd edition)*. CRC Press.



# Bibliography

Domingos, P. 2015. *The Master Algorithm: How the Quest for the Ultimate Learning Machine Will Remake Our World*. Basic Books.

Efron, B., and D. V. Hinkley. 1978. "Assessing the Accuracy of the Maximum Likelihood Estimator: Observed Versus Expected Fisher Information." *Biometrika* 65: 457–82. <https://doi.org/10.1093/biomet/65.3.457>.

Faraway, J. J. 2015. *Linear Models with R*. 2nd ed. Chapman; Hall/CRC.

Gelman, A., J. B. Carlin, H. A. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. 2014. *Bayesian Data Analysis*. 3rd ed. CRC Press.

Held, L., and D. S. Bové. 2014. *Applied Statistical Inference: Likelihood and Bayes*. Springer.

Wilks, S. S. 1938. "The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses." *Ann. Math. Statist.* 9: 60–62. <https://doi.org/10.1214/aoms/1177732360>.

Wood, S. 2015. *Core Statistics*. Cambridge University Press.