



THE UNIVERSITY
of EDINBURGH

Lecture Notes

Data Mining and Exploration

Michael Gutmann and Arno Onken

School of Informatics

Spring Semester 2021

January 2021

Contents

Preface	vii
1 First Steps in Exploratory Data Analysis	1
1.1 Numerical Data Description	1
1.1.1 Location	1
1.1.2 Scale	3
1.1.3 Shape	4
1.1.4 Multivariate Measures	5
1.2 Data Visualisation	10
1.2.1 Bar Plot	10
1.2.2 Box Plot	11
1.2.3 Scatter Plot	12
1.2.4 Histogram	12
1.2.5 Kernel Density Plot	12
1.2.6 Violin Plot	15
1.3 Data Pre-Processing	15
1.3.1 Standardisation	15
1.3.2 Outlier Detection and Removal	17
2 Principal Component Analysis	19
2.1 PCA by Variance Maximisation	19
2.1.1 First Principal Component Direction	19
2.1.2 Subsequent Principal Component Directions	21
2.1.3 Simultaneous Variance Maximisation	23
2.2 PCA by Minimisation of Approximation Error	23
2.3 PCA by Low Rank Matrix Approximation	25
2.3.1 Approximating the Data Matrix	25
2.3.2 Approximating the Sample Covariance Matrix	27
2.3.3 Approximating the Gram Matrix	28
2.4 Probabilistic PCA	29
2.4.1 Probabilistic Model	29
2.4.2 Joint, Conditional and Observation Distributions	29
2.4.3 Maximum Likelihood	31
2.4.4 Relation to PCA	33

3	Dimensionality Reduction	35
3.1	Linear Dimensionality Reduction	35
3.1.1	From Data Points	35
3.1.2	From Inner Products	36
3.1.3	From Distances	37
3.1.4	Example	38
3.2	Dimensionality Reduction by Kernel PCA	39
3.2.1	Idea	39
3.2.2	Kernel Trick	40
3.2.3	Example	41
3.3	Multidimensional Scaling	42
3.3.1	Metric MDS	42
3.3.2	Nonmetric MDS	45
3.3.3	Classical MDS	45
3.3.4	Example	46
3.4	Isomap	46
3.5	UMAP	47
4	Predictive Modelling and Generalisation	51
4.1	Prediction and Training Loss	51
4.1.1	Prediction Loss	51
4.1.2	Training Loss	52
4.1.3	Example	53
4.2	Generalisation Performance	55
4.2.1	Generalisation for Prediction Functions and Algorithms	56
4.2.2	Overfitting and Underfitting	57
4.2.3	Example	59
4.3	Estimating the Generalisation Performance	60
4.3.1	Methods for Estimating the Generalisation Performance	61
4.3.2	Hyperparameter Selection and Performance Evaluation	64
4.4	Loss Functions in Predictive Modelling	66
4.4.1	Loss Functions in Regression	67
4.4.2	Loss Functions in Classification	67
A	Linear Algebra	73
A.1	Matrices	73
A.2	Vectors	74
A.3	Matrix Operations as Operations on Column Vectors	75
A.4	Orthogonal Basis	77
A.5	Subspaces	78
A.6	Orthogonal Projections	79
A.7	Singular Value Decomposition	79
A.8	Eigenvalue Decomposition	80
A.9	Positive Semi-definite and Definite Matrices	81
A.10	Matrix Approximations	81

B Proofs Related to PCA	87
B.1 Sequential Maximisation Yields Simultaneous Maximisation	87
B.2 Equivalence to PCA by Variance Maximisation	89

Preface

With the massive progress in automation and computerisation, more data than ever are constantly generated. Sources of data range from scientific measurements to internet searches and connected devices. The aim of the course “Data Mining and Exploration” is to introduce and discuss modern techniques for analysing, interpreting, visualising and exploiting the data that are captured in scientific and commercial environments.

In the Large Hadron Collider, for example, 600 terabyte of raw data are generated per second when high energy particles colliding.¹ Not all data are of interest, and it is major challenge to extract useful information from that huge amount of raw data.

Figure 1 illustrates the amount of data generated in the digital world, which includes more than 500 hours of video uploaded to YouTube or more than 400'000 hours of Netflix video streamed *per minute* in 2020. As another example, sensors in modern cars produce a steady stream of data, with connected cars estimated to produce 25 gigabytes of data per hour.² As in the physics example, the raw data itself will not be useful but it contains valuable information that some companies may wish to extract for diagnostics and product improvement, or because of its monetary value for third-parties such as insurance companies.

The pervasiveness of data generation and collection raises a number of important issues related to privacy and fairness.³ These societal issues need to be taken into account when analysing data, interpreting the results and ultimately making decisions, and are explored in the second half of the course as part of paper presentations and potentially in the mini-projects. These lecture notes, and the associated computer labs, deal with the technical aspects of data analysis as taught in the first half of the course.

Figure 2 provides an overview of the overall data analysis process, showing the key steps (green ovals) and associated questions addressed. The diagram highlights that the data analysis process is iterative. On the one hand, we may want to go back a step when, for example, sanity checks flag potential issues. On the other hand, unless data collection has been carefully designed, the data collection and analysis are often coupled and this feedback loop needs to be kept in mind. In recommendation systems, for instance, users may not choose certain items solely because they have not been recommended to them and are unable to choose them. Absence of selected items in the data should thus not be considered

¹<https://home.cern/science/computing/processing-what-record>

²<https://www.statista.com/chart/8018/connected-car-data-generation/>

³See, for example, *Understanding artificial intelligence ethics and safety* by D. Leslie, <https://arxiv.org/abs/1906.05684>

as evidence of any disliking. Ignoring the feedback loop between data analysis and collection can lead to “filter bubbles” and associated biases in, for example, media consumption or predictive policing.⁴

The lecture notes focus on the three steps on the right side of the data analysis cycle in Figure 2—exploratory data analysis, pre-processing, and performance evaluation in (predictive) model building & fitting. The other topics are covered or explored in the paper presentations and the mini-projects in the second half of the course.

The lecture notes were first written by Michael Gutmann when he revived the course “Data Mining and Exploration” in academic year 2016-17 after it has not been taught for some years. They were expanded and partly re-organised by Arno Onken when he was teaching the course in the following years, adding sections on probabilistic PCA and UMAP, re-organising and adding new material to the first chapter and delegating proofs to the appendix. The text was revised by Michael Gutmann, adding and clarifying some minor explanations, when he was teaching the course again in academic year 2020-21 as part of the COVID-19 teaching adjustments of the School of Informatics.

Michael Gutmann and Arno Onken
Edinburgh, January 2021

⁴See, for example, https://en.wikipedia.org/wiki/Filter_bubble, or *Ethical principles in machine learning and artificial intelligence: cases from the field and possible ways forward*, Humanit Soc Sci Commun 7, 9, 2020, <https://www.nature.com/articles/s41599-020-0501-9>



Figure 1: Data generated in the digital world in 2020. Figure adapted from <https://www.domo.com/learn/data-never-sleeps-8>

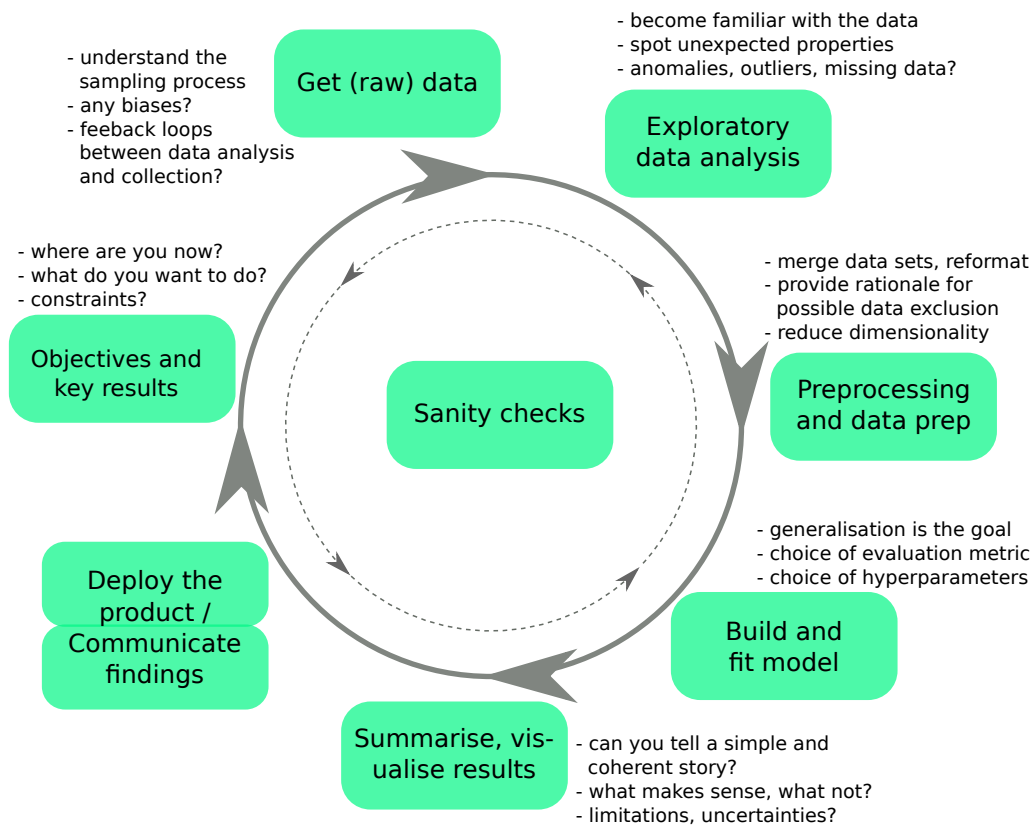


Figure 2: Overview of the data analysis process.

Chapter 1

First Steps in Exploratory Data Analysis

This chapter is about the first steps of exploratory data analysis. It is assumed that we have available n data points $\mathbf{x}_1, \dots, \mathbf{x}_n$ each containing d attributes and that the data have been transformed to numbers. Each data point \mathbf{x}_i thus is a d dimensional vector. We first explain numerical summaries of the data, and then methods to visualise their properties. This is followed by some elements of preprocessing for further analysis of the data. The presented methods are simple and inexpensive to compute. They provide invaluable descriptive statistics of the data that should generally always be first computed before any further data analysis is performed.

1.1 Numerical Data Description

This section covers a range of numerical measures to characterise a dataset. The measures provide answers to questions such as the following ones: Where are the data located? Where is the centre of the distribution? How scattered are the data? The measures vary in how robust they are to measurement errors and outliers.

1.1.1 Location

Location measures provide answers to questions about the overall location of the data. The *sample mean*, also simply known as *arithmetic mean* or *average* is the most well-known and most commonly used location measure. It is defined as

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i. \quad (1.1)$$

Assuming that the data were drawn from a random variable x with probability density function p , the sample mean \bar{x} of the data is an estimate of the mean or expected value of x ,

$$\mathbb{E}[x] = \int \xi p(\xi) d\xi. \quad (1.2)$$

2 First Steps in Exploratory Data Analysis

For multivariate data, the vector-valued mean is the element-wise mean:

$$\bar{\mathbf{x}} = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_d \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad (1.3)$$

Another common location measure is the *median*. It splits the samples into two chunks of equal size and can be thought of as the centre of the distribution. It can be found by ordering the samples and then taking the value of the sample in the middle if the number of samples is odd, or the mean of the middle two samples if the number of samples is even. More formally, let $x_{(1)}, x_{(2)}, \dots, x_{(n)}$ denote the ordered samples, so that $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$. The median of these samples is defined as:

$$\text{median}(x) = \begin{cases} x_{((n+1)/2)} & \text{if } n \text{ is odd} \\ \frac{1}{2}(x_{(n/2)} + x_{(n/2+1)}) & \text{if } n \text{ is even} \end{cases} \quad (1.4)$$

In contrast to the mean, the median is robust to corrupted observations (i.e. outliers, see Section 1.3.2). The precise location of most data points does not affect the median. Let us assume that an observation is recorded at $x_i + \delta$ rather than x_i because of a malfunction of the measurement device. The mean then changes from \bar{x} to $\bar{x} + \delta/n$ which can be arbitrarily large, while the median changes at most to a neighbouring observation.

For instance, the first 10 training digits of the MNIST dataset (a large dataset of handwritten digits that is commonly used for training various image processing systems), sorted from lowest to highest are:

$$(0, 1, 1, 1, 2, 3, 4, 4, 5, 9) \quad (1.5)$$

The median is $(2 + 3)/2 = 2.5$ whereas the arithmetic mean is 3.0.

If we change the last element from 9 to 9000, thus obtaining the list

$$(0, 1, 1, 1, 2, 3, 4, 4, 5, 9000), \quad (1.6)$$

then the median is still 2.5 whereas the mean changes dramatically from 3.0 to 902.1.

Another measure of the location of the data that is more robust than the average is the trimmed average. It is the average of the data when leaving out the smallest and largest $k < n/2$ values,

$$\frac{1}{n - 2k} \sum_{i=k+1}^{n-k} x_{(i)}. \quad (1.7)$$

If $k = 0$, the trimmed average is the usual sample average. As k approaches $n/2$, the trimmed average approaches the median.

A further important measure is the *mode*. It is the value that occurs most frequently. The mode is not necessarily unique: there can be more than one mode if several different values occur equally often. The mode is always applicable and

meaningful even if the data attributes are just categorical and unordered (e.g. nationality and gender). In the example above, the mode is 1.

Besides the median, we can define additional robust quantities that are invariant to the precise location of most data points. The α -th sample *quantile* q_α is roughly the data point with a proportion α of the ordered data $x_{(i)}$ to its left, i.e. $q_\alpha \approx x_{(\lceil n\alpha \rceil)}$. For example, the minimum and the maximum are the 0 and 1 quantiles q_0 and q_1 and the median is the 0.5 quantile $q_{0.5}$. Like the median, quantiles are often computed by interpolation if αn is not an integer.

The *quartiles* Q_1 , Q_2 and Q_3 are given by $q_{0.25}$, $q_{0.5}$ and $q_{0.75}$ respectively. To obtain the first and the third quartiles, we can split the ordered dataset into a lower half L_1 and an upper half L_2 , where we include the median in both L_1 and L_2 if the number of samples n is odd. The first quartile Q_1 is then the median of L_1 and the third quartile Q_3 is the median of L_2 .

1.1.2 Scale

Scale measures answer questions about the spread of the data. The *sample variance* is the mean squared difference from the sample mean:

$$\text{var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \quad (1.8)$$

It is an estimator of the variance of the random variable x ,

$$\text{Var}[x] = \int (\xi - \mathbb{E}[x])^2 p(\xi) d\xi = \mathbb{E}[x^2] - \mathbb{E}[x]^2. \quad (1.9)$$

The variance is the expected value of the squared deviation of the random variable from the mean. This is also known as the second central moment. “Central” because it refers to the deviation from the mean, and “second”, because it refers to the squared deviation, i.e. the power two. Other powers also provide useful measures which we will encounter in the following sections.

In our definition of the sample variance we divide by n . Other variance estimators divide by $n - 1$ (known as Bessel’s correction¹) rather than n .

Note that the variance does *not* have the same unit as the samples: the variance unit is the squared sample unit. The *sample standard deviation* is given by $\text{std}(x) = \sqrt{\text{var}(x)}$ which has same unit as the samples.

Because of the squaring, the sample variance is more affected by outliers than the sample mean. The median can be used to obtain a more robust measure of the scale of the data: instead of measuring the average squared deviation from the average, we measure the *median absolute deviation* from the median,

$$\text{MAD}(x) = \text{median}(|x_i - \text{median}(x)|) \quad (1.10)$$

This measure has the same units as the x_i themselves.

The range $x_{(n)} - x_{(1)}$, i.e. the difference between the largest and the smallest value, is another measure of the scale of the data, but it is not robust. A more

¹The reason for this correction is a bias that arises from the sample mean. To see this, calculate $\mathbb{E}[\text{var}(x)]$ by expanding \bar{x} and compare to $\text{Var}[x]$. For large n the bias is small, so for conceptual simplicity in later chapters we use the more intuitive biased sample variance.

robust quantity is the difference between the upper and lower end of what contains the central 50% of the data. The *interquartile range* (IQR) is defined as the difference between the first and the third quartile:

$$\text{IQR} = Q_3 - Q_1 \tag{1.11}$$

We will later use the IQR to detect unusual data points.

Returning to our MNIST example,

$$(0, 1, \underbrace{1, 2, 3, 4, 4, 5, 9}_{\text{IQR}}) \tag{1.12}$$

the sample standard deviation is $\sqrt{6.4} \approx 2.53$, the MAD is 1.5 and the IQR is $4 - 1 = 3$.

1.1.3 Shape

The *sample skewness* measures the asymmetry of the data. It is defined as

$$\text{skew}(x) = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\text{std}(x)} \right)^3 \tag{1.13}$$

Subtraction of the mean and division by the standard deviation normalises the terms in the brackets: the variable $z = \frac{x - \bar{x}}{\text{std}(x)}$ has zero mean and unit standard deviation. Location and scale are therefore not taken into account by skewness.

For unimodal distributions, positive skewness means that the distribution has a longer right tail, that is, the mass is concentrated on the left of the distribution and elongated on the right. For negative skewness, it is the other way around: the distribution has a longer left tail. Data that are symmetric around their mean have zero skewness. The converse, however, is not true: zero skewness does not necessarily mean that the data are symmetric around their mean.

Due to the third power, sample skewness is sensitive to outliers. A more robust measure can be obtained by means of the quartiles,

$$\text{Galton's measure of skewness} = \frac{(Q_3 - Q_2) - (Q_2 - Q_1)}{Q_3 - Q_1}. \tag{1.14}$$

The denominator is the interquartile range and normalises the skewness measure like the standard deviation in (1.13). By definition of the quartiles, both $Q_3 - Q_2$ and $Q_2 - Q_1$ are positive. The first term measures the range of the third quarter while the second term measures the range of the second quarter. Galton's skewness thus computes the difference between the ranges of the two quarters in a normalised way. It is positive if the range of the third quarter is larger than the range of the first quarter, and conversely. Figure 1.1 shows an example.

If we take the sample skewness equation with the fourth power instead of the third, then we get a measure called the *sample kurtosis*:

$$\text{kurt}(x) = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\text{std}(x)} \right)^4 \tag{1.15}$$

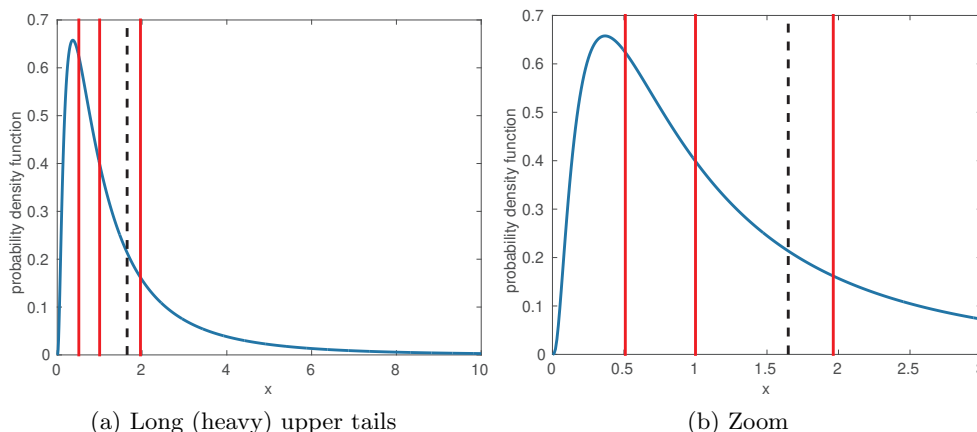


Figure 1.1: Example of positive skewness. The distribution has skewness equal to 6.18 according to (1.13), its interquartile range is 1.45, and $Q_3 - Q_2 = 0.96$ while $Q_2 - Q_1 = 0.49$, so that Galton’s measure of skewness is positive. The black dashed line indicates the mean, while the three red lines show from left to right Q_1 , Q_2 , and Q_3 .

Due to the fourth power, kurtosis is insensitive to the symmetry of the distribution of x . It measures how often x takes on values that are considerably larger or smaller than its standard deviation; it is said to measure how heavy the tails of the distribution of x are. Figure 1.2 shows the function $u \mapsto u^4$. The function is relatively flat for $-1 < u < 1$ so that kurtosis basically ignores the behaviour of x within one standard deviation around its mean. The function then grows rapidly and values larger than two standard deviations away from the mean contribute strongly to the value of kurtosis.

A robust version of the kurtosis can be defined using quantiles:

$$\text{robust kurtosis}(x) = \frac{(q_{7/8} - q_{5/8}) + (q_{3/8} - q_{1/8})}{Q_3 - Q_1} \quad (1.16)$$

This measure estimates the lengths of the upper and lower tails and normalises with the IQR. Further robust measures of kurtosis and skewness are discussed by Kim and White (2004).

1.1.4 Multivariate Measures

We now consider measures of dependency between two or more variables. The *sample covariance* between observations of two variables $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ measures the strength of linear association between them. It is given by

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \quad (1.17)$$

and is an estimator of the *covariance* of two random variables x and y

$$\text{Cov}[x, y] = \mathbb{E}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])] = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]. \quad (1.18)$$

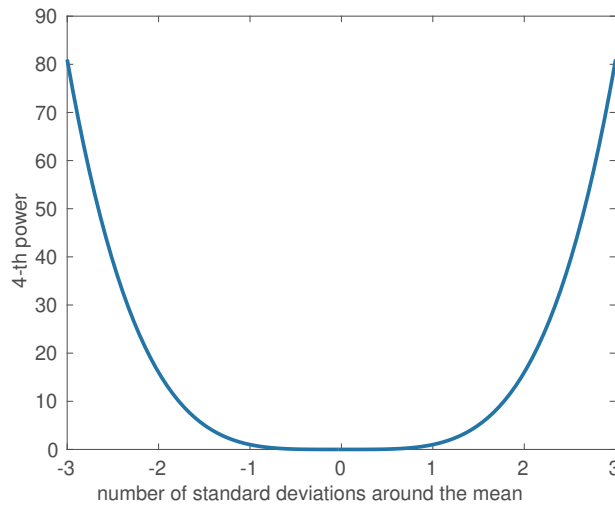


Figure 1.2: The figure shows the function $u \mapsto u^4$ that occurs in the definition of the kurtosis. It is relatively flat on $(-1, 1)$ and grows quickly for values outside the interval.

Like for the variance, other estimators of the covariance divide by $n - 1$ rather than n . It holds $\text{Cov}[x, x] = \text{Var}[x]$, $\text{Cov}[x, y] = \text{Cov}[y, x]$, and $\text{Cov}[ax + b, y] = a \text{Cov}[x, y]$. The value of the covariance thus depends on the scale of x and y which is often undesirable.

Pearson's correlation coefficient (also simply known as correlation coefficient) normalises the covariance using the product of standard deviations:

$$\rho(x, y) = \frac{\text{cov}(x, y)}{\text{std}(x)\text{std}(y)} \quad (1.19)$$

The measure takes values between -1 and 1 and is also known as linear correlation coefficient, because it measures linear relation. To see this, suppose that $y = ax + b$ where $a \neq 0$ and b are constants. From the linearity of the mean, it follows that $\bar{y} = a\bar{x} + b$. For the standard deviation, we obtain

$$\text{std}(y) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2} \quad (1.20)$$

$$= \sqrt{\frac{1}{n} \sum_{i=1}^n (ax_i + b - a\bar{x} - b)^2} \quad (1.21)$$

$$= \sqrt{a^2} \text{std}(x) \quad (1.22)$$

and for the covariance

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \quad (1.23)$$

$$= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(ax_i + b - a\bar{x} - b) \quad (1.24)$$

$$= a \text{var}(x). \quad (1.25)$$

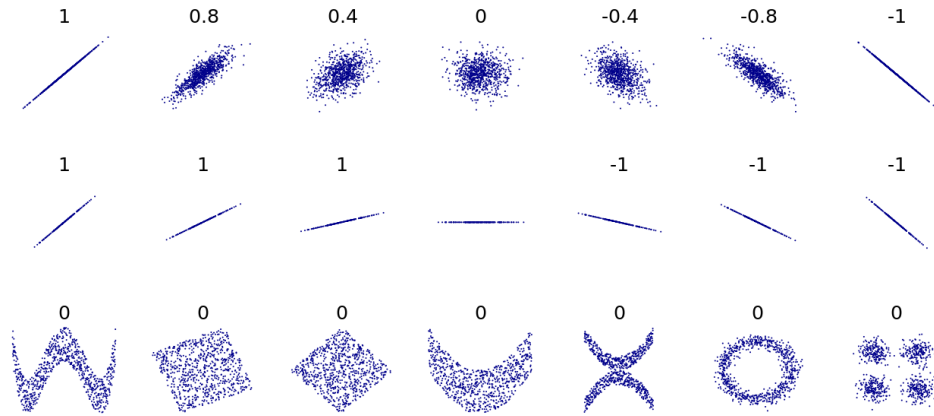


Figure 1.3: Correlation coefficients for different data sets. Figure from https://en.wikipedia.org/wiki/Pearson_product-moment_correlation_coefficient.

The correlation coefficient therefore is given by

$$\rho(x, y) = \frac{\text{avar}(x)}{\text{std}(x)\sqrt{a^2}\text{std}(x)} = \frac{a}{\sqrt{a^2}} = \begin{cases} 1 & \text{if } a > 0 \\ -1 & \text{if } a < 0 \end{cases}. \quad (1.26)$$

Thus, for linear relationships, ρ is either 1 or -1 with 1 indicating positive linear relationships and -1 indicating negative linear relationships.

If $\rho = 0$, then we refer to the variables as uncorrelated. It means $\frac{1}{n} \sum_{i=1}^n (x_i y_i) = \bar{x}\bar{y}$, that is, there is no linear relationship between the variables, but $\rho = 0$ does *not* mean that the variables are statistically independent which is a much stronger statement (see Figure 1.3 for uncorrelated but dependent examples).

For d attributes x_1, \dots, x_d with n observations for each attribute, the *sample covariance matrix* is given by

$$\text{cov}(x_1, \dots, x_d) = \begin{pmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) & \dots & \text{cov}(x_1, x_d) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) & \dots & \text{cov}(x_2, x_d) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(x_d, x_1) & \text{cov}(x_d, x_2) & \dots & \text{cov}(x_d, x_d) \end{pmatrix} \quad (1.27)$$

The matrix is symmetric because $\text{cov}(x_i, x_j) = \text{cov}(x_j, x_i)$. Moreover, the diagonal elements of the sample covariance matrix are the sample variances because $\text{cov}(x_i, x_i) = \text{var}(x_i)$.

The sample covariance matrix is an estimator of the *covariance matrix*

$$\text{Cov}[\mathbf{x}] = \mathbb{E} \left[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\top \right] \quad (1.28)$$

where \mathbf{x} denotes a d -dimensional random variable. This follows immediately from the properties of the outer product \mathbf{ab}^\top between two vectors \mathbf{a} and \mathbf{b} , see e.g. Section A.2, because $(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\top$ is a $d \times d$ matrix where the (i, j) -th element is $(x_i - \mathbb{E}[x_i])(x_j - \mathbb{E}[x_j])$. The covariance matrix is symmetric and,

moreover, positive semi-definite, because by linearity of expectation,

$$\mathbf{w}^\top \text{Cov}[\mathbf{x}]\mathbf{w} = \mathbb{E} \left[\underbrace{\mathbf{w}^\top (\mathbf{x} - \mathbb{E}[\mathbf{x}])}_{\text{scalar}} \underbrace{(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\top \mathbf{w}}_{\text{scalar}} \right] = \mathbb{E} \left[(\mathbf{w}^\top (\mathbf{x} - \mathbb{E}[\mathbf{x}]))^2 \right] \geq 0. \quad (1.29)$$

It thus has an eigenvalue decomposition $\text{Cov}[\mathbf{x}] = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$, where $\mathbf{\Lambda}$ is a diagonal matrix containing the eigenvalues, and \mathbf{U} is an orthogonal matrix with the eigenvectors as columns (see e.g. Appendix A for a linear algebra refresher).

The total variance of the d random variables x_i is the sum of all eigenvalues: with the definition of the trace of a matrix, see e.g. (A.6), we have

$$\sum_{i=1}^d \text{Var}[x_i] = \text{trace}(\text{Cov}[\mathbf{x}]) \quad (1.30)$$

$$= \text{trace}(\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top) \quad (1.31)$$

$$\stackrel{\text{(A.8)}}{=} \text{trace}(\mathbf{\Lambda}\mathbf{U}^\top\mathbf{U}) \quad (1.32)$$

$$= \text{trace}(\mathbf{\Lambda}) \quad (1.33)$$

$$= \sum_{i=1}^d \lambda_i, \quad (1.34)$$

where we have used that, for an orthogonal matrix \mathbf{U} , $\mathbf{U}^\top\mathbf{U}$ is the identity matrix.

The linearly transformed random variable $\mathbf{A}\mathbf{x} + \mathbf{b}$ has covariance matrix $\mathbf{A}\text{Cov}[\mathbf{x}]\mathbf{A}^\top$. This is due to the linearity of expectation,

$$\text{Cov}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbb{E} \left[(\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbb{E}[\mathbf{A}\mathbf{x} + \mathbf{b}])(\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbb{E}[\mathbf{A}\mathbf{x} + \mathbf{b}])^\top \right] \quad (1.35)$$

$$= \mathbb{E} \left[(\mathbf{A}\mathbf{x} - \mathbb{E}[\mathbf{A}\mathbf{x}])(\mathbf{A}\mathbf{x} - \mathbb{E}[\mathbf{A}\mathbf{x}])^\top \right] \quad (1.36)$$

$$= \mathbb{E} \left[(\mathbf{A}(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{A}(\mathbf{x} - \mathbb{E}[\mathbf{x}]))^\top \right] \quad (1.37)$$

$$= \mathbb{E} \left[\mathbf{A}(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\top \mathbf{A}^\top \right] \quad (1.38)$$

$$= \mathbf{A} \mathbb{E} \left[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\top \right] \mathbf{A}^\top \quad (1.39)$$

$$= \mathbf{A} \text{Cov}[\mathbf{x}]\mathbf{A}^\top. \quad (1.40)$$

Following this scheme, we can transform the sample covariance matrix to a *sample correlation matrix* in the following way:

$$\rho(\mathbf{x}) = \text{diag} \left(\frac{1}{\text{std}(\mathbf{x})} \right) \text{cov}(\mathbf{x}) \text{diag} \left(\frac{1}{\text{std}(\mathbf{x})} \right) \quad (1.41)$$

where $\text{diag}(\mathbf{x})$ denotes the diagonal matrix with \mathbf{x} on the diagonal and 0 everywhere else. The correlation matrix has all ones on the diagonal.

A simple way to measure nonlinear relationships between observations of two random variables x and y is to compute their covariance or correlation after transforming them nonlinearly, i.e. to compute

$$\rho(g(x), g(y)) = \frac{\text{cov}(g(x), g(y))}{\text{std}(g(x))\text{std}(g(y))} \quad (1.42)$$

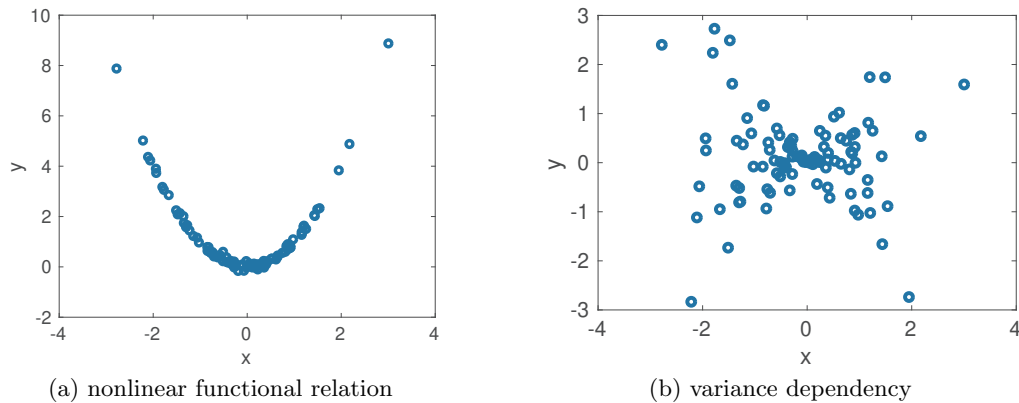


Figure 1.4: Measuring nonlinear relationships between two variables. The linear correlation coefficient is small in both (a) and (b). The correlation of the absolute values, however, captures the relation between x and y . (a) $\rho(|x|, |y|) = 0.93$, (b) $\rho(|x|, |y|) = 0.68$.

where g denotes a nonlinear function. Different nonlinearities g can be used to measure different properties of the dependencies. The absolute value, for example, can be used to measure variance dependencies.

Figure 1.4 shows two examples. In Figure 1.4(a) there is a clear functional relation between x and y but the (linear) correlation coefficient is -0.15 , wrongly indicating a negative correlation between x and y . Computing the correlation between the absolute values $|x|$ and $|y|$, however, yields a correlation coefficient of 0.93 . In Figure 1.4(b), the variance of y depends on the magnitude of x . The linear correlation is practically zero while the absolute values have a correlation coefficient of 0.68 .

Similar to the median, we can also define robust correlation measures that do not depend on the relative distances between samples. Such correlation measures are known as *rank correlation*. One example is Kendall's τ . To define this measure, we need the concept of concordant and discordant observations. A pair of observations (x_i, y_i) and (x_j, y_j) with $i \neq j$ is said to be *concordant* if the order of elements is consistent, that is if both $x_i > x_j$ and $y_i > y_j$ or if both $x_i < x_j$ and $y_i < y_j$. The pair is said to be *discordant* if both $x_i > x_j$ and $y_i < y_j$ or if both $x_i < x_j$ and $y_i > y_j$. Visually, a line connecting two concordant observations has a positive slope, while the slope is negative for discordant observations. If $x_i = x_j$ or $y_i = y_j$ then the pair is said to be neither concordant nor discordant. Kendall's τ is given by

$$\tau(x, y) = \frac{n_c(x, y) - n_d(x, y)}{n(n-1)/2} \quad (1.43)$$

where $n_c(x, y)$ denotes the total number of concordant pairs, $n_d(x, y)$ denotes the total number of discordant pairs and n denotes the total number of (bivariate) observations. Note that the denominator $n(n-1)/2$ is the number of (unordered) pairs that you can form with n data points, so that Kendall's τ is the difference between the fraction of concordant and discordant pairs of observations. Figure 1.5 illustrates the calculation of Kendall's τ .

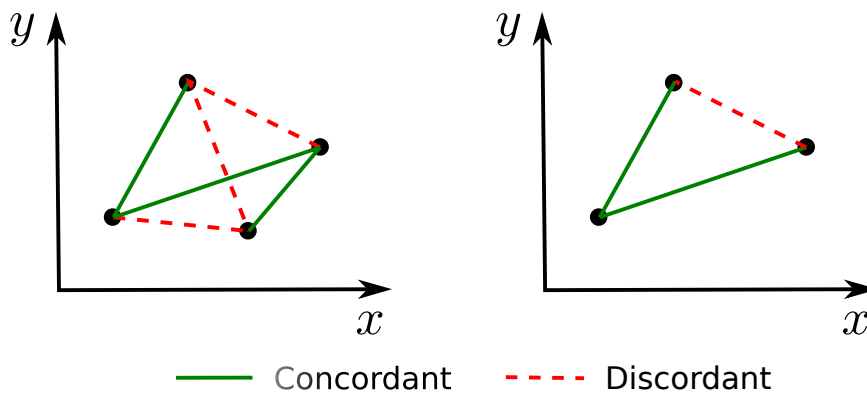


Figure 1.5: Kendall's τ is given by the normalised difference between the number of concordant (green solid) and discordant (red dashed) data pairs. Concordant data pairs are connected by lines with positive slope while discordant pairs are connected by line with negative slope. For the example on the left, $\tau = 0$, for the example on the right, $\tau = (2 - 1)/3 = 1/3$.

Like Pearson's correlation coefficient, Kendall's τ takes values between -1 and 1. Unlike Pearson's correlation coefficient, Kendall's τ does not take relative distances between neighbouring samples into account with -1 indicating that all pairs are discordant and 1 indicating that all pairs are concordant. Like the median, Kendall's τ is robust to outliers.

1.2 Data Visualisation

The data visualisations that we are discussing in this section are means to describe frequencies of occurrences. The various plots take different characteristics of the data into account and, like numerical data descriptions, vary in their robustness.

Before visualising the data, it is useful to determine what values the data attributes can take. If an attribute can take a finite or countably infinite number of values (e.g. nationality, number of students in this course), then we call it *discrete*. If, on the other hand, the attribute takes values on a continuum (e.g. time point), then we call the attribute *continuous*. Which visualisation is appropriate depends on the type of the data under consideration.

1.2.1 Bar Plot

A *bar plot* shows values characterising particular attributes of a population. This is done by means of a collection of bars with different heights and labels denoting the attributes. The heights of the bars are proportional to the values that they represent, for instance the number of occurrences n_j of the attribute v_j in data x_1, \dots, x_n .

When formally counting observations that have certain attributes, we will often make use of the indicator function $\mathbb{1}$. For a given set A , this function is defined as

$$\mathbb{1}_A(\xi) = \begin{cases} 1 & \text{if } \xi \in A \\ 0 & \text{otherwise} \end{cases} \quad (1.44)$$

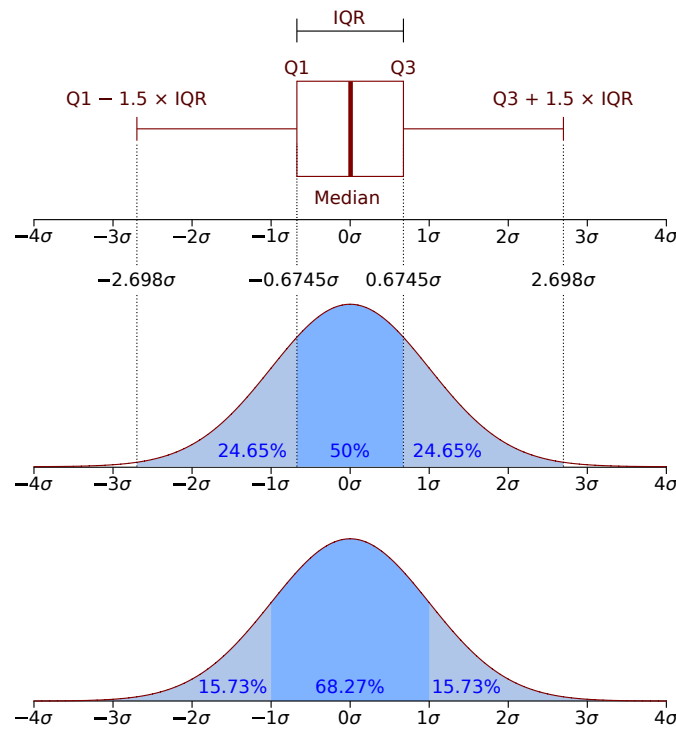


Figure 1.6: Box plot and comparison to the Gaussian probability density function. Figure from https://en.wikipedia.org/wiki/Box_plot

We can compute n_j by having A contain just one element, that is $A = \{v_j\}$. The condition $\xi \in A$ then simplifies to $\xi = a$ and the number of occurrences n_j is given by

$$n_j = \sum_{i=1}^n \mathbb{1}_{\{v_j\}}(x_i). \quad (1.45)$$

Note that $\sum_{i=1}^k n_j = n$ if $\{x_1, \dots, x_n\} \subseteq \{v_1, \dots, v_k\}$. When comparing populations with different number of samples, it is more useful to show *relative frequencies* $f_j = \frac{n_j}{n}$ instead of number of observations n_j .

Like the mode, the bar plot is always applicable and meaningful even if the data attributes are just categorical and unordered.

1.2.2 Box Plot

The box plot or box-and-whiskers-plot compactly shows the quartiles using a box and lines (typically vertically arranged). The box is drawn from Q_1 to Q_3 , illustrating the IQR with a crossing line at the median. From this IQR box, two lines (the “whiskers”) extend to indicate the spread of the data. The length of the lines is typically a multiple of the IQR, with a factor of 1.5 being a common default value, see Figure 1.6 for an example. Note, however, that different software packages may implement different defaults. Observations that fall beyond the limits set by the whiskers are considered outliers (see Section 1.3.2) and sometimes drawn individually as separate points. Since the box plot is based on robust measures, it is itself very robust.

1.2.3 Scatter Plot

The scatter plot is one of the most common and useful plots to visualise the distribution of two random variables. It shows the observations as symbols (typically dots, circles, triangles, ...) in a 2-dimensional coordinate system. Each symbol represents one data point. Colouring the symbols or changing their size or shape enables visualisation of further dimensions or class labels. Figure 1.4 shows two examples of scatter plots.

1.2.4 Histogram

A histogram can be thought of as a continuous extension of the bar plot where the (not necessarily complete) range of the observations is divided into k non-overlapping and successive bins B_1, \dots, B_k . We then count the number of samples falling into each bin B_j using the same procedure as for the bar plot:

$$n_j = \sum_{i=1}^n \mathbb{1}_{B_j}(x_i) \quad (1.46)$$

Optionally, we can normalise the n_j by n to show the relative frequencies. The bins are often chosen to have equal bin width h . For a starting value L , the sets B_j are then given by:

$$\begin{aligned} B_1 &= [L, L + h) \\ B_2 &= [L + h, L + 2h) \\ &\vdots \\ B_{k-1} &= [L + (k-2)h, L + (k-1)h) \\ B_k &= [L + (k-1)h, L + kh]. \end{aligned}$$

The starting value L and the value $L + kh$ correspond to the lower and upper bound, respectively, of data visualised in the histogram. The starting value L , the bin size h and the number of bins are parameters that need to be chosen. For $L \leq \min(x_1, \dots, x_n)$ and k such that $L + kh \geq \max(x_1, \dots, x_n)$ the whole dataset is visualised. B_j is centred at $L + jh - h/2$ and the corresponding bar showing n_j is located at this point on the x-axis. Figure 1.7(a) shows an example. We can see that different starting values L may lead to differently looking histograms.

For two dimensions, the bins can be 2-dimensional: $B_{1,1} = [L, L + h) \times [L, L + h)$, $B_{1,2} = [L, L + h) \times [L + h, L + 2h)$, $B_{2,1} = [L + h, L + 2h) \times [L, L + h)$, ... The bars can then be represented as a 2-dimensional surface plot.

1.2.5 Kernel Density Plot

Oftentimes, we know that the distribution underlying the observations is continuous and smooth. By construction, a histogram uses a discrete number of bins and will generally show abrupt changes from one bin to the next. Histograms, therefore, will generally misrepresent distributions that are continuous and smooth.

To remedy this shortcoming, we can try to estimate the density as a continuous function and then plot our estimate. To estimate the probability density, we

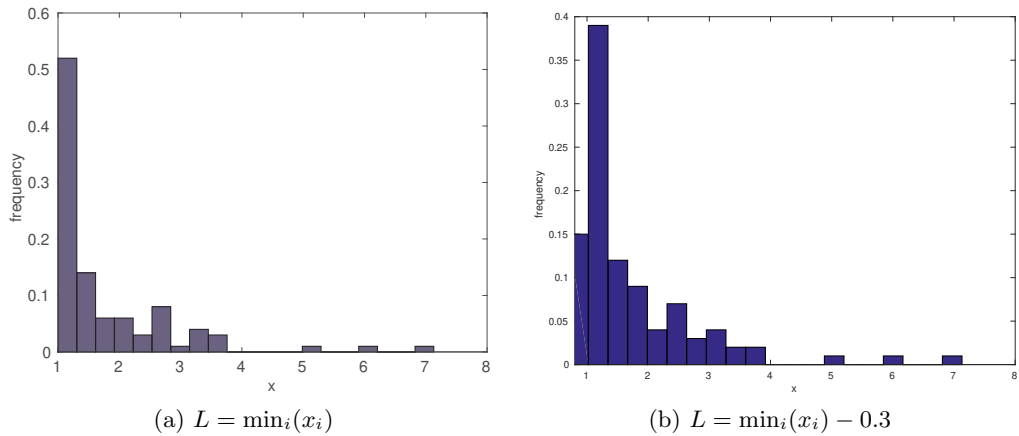


Figure 1.7: Describing data by histograms. The edges of the bars correspond to the edges of the bins used for the histogram. (a) and (b) show histograms with different starting values L .

can use a kernel density estimate:

$$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i) \quad (1.47)$$

where K_h is a non-negative function called the (scaled) *kernel*. To make sure that we obtain a density (which is normalised to one), K_h needs to integrate to one: $\int K_h(\xi) d\xi = 1$. For typical kernel functions, the parameter h is called the bandwidth and determines the smoothness of the estimate.

An example of K_h is the boxcar kernel

$$K_h(\xi) = \frac{1}{h} \mathbb{1}_{[-h/2, h/2]}(\xi), \quad (1.48)$$

which is a rectangular function located at zero. The bandwidth parameter determines the width of the rectangle and is inversely proportional to its height. The kernel density estimate in conjunction with the boxcar kernel counts the number of observations in the vicinity of each argument x and normalises that count with the bandwidth h .

When using the boxcar kernel, we stop counting observations as soon as they are further than $h/2$ away from x . Instead of this binary weight assignment, it is often more reasonable to assign to each data point a weight that decreases smoothly with the distance from the argument x . A popular choice is the Gaussian kernel

$$K_h(\xi) = \frac{1}{\sqrt{2\pi}h} \exp\left(-\frac{\xi^2}{2h^2}\right). \quad (1.49)$$

In these cases, the bandwidth parameter h is a free parameter that needs to be chosen. Other kernels are plotted in Figure 1.8. Figure 1.9 shows the densities estimated by rescaling the histogram, with the boxcar kernel, and with the Gaussian kernel.

Like the histogram, the kernel density plot can be readily extended to two dimensions using a surface plot and a 2-dimensional kernel and density estimator.

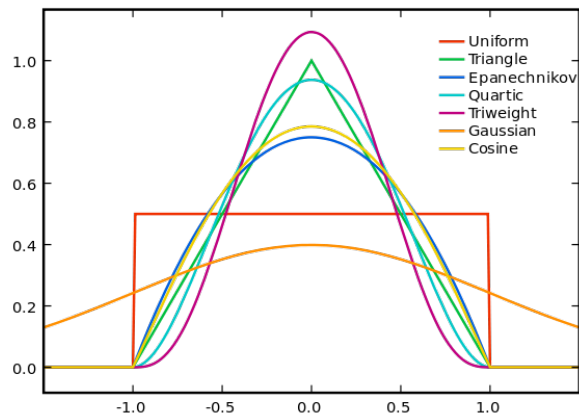
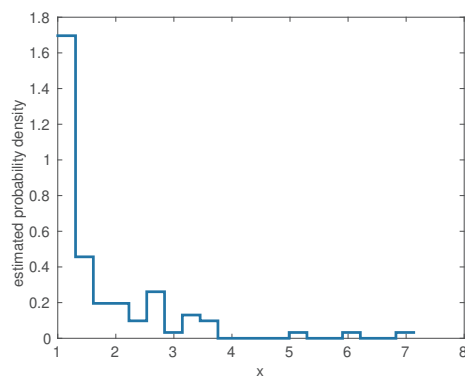
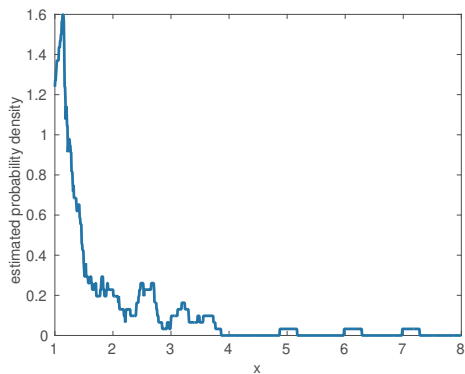


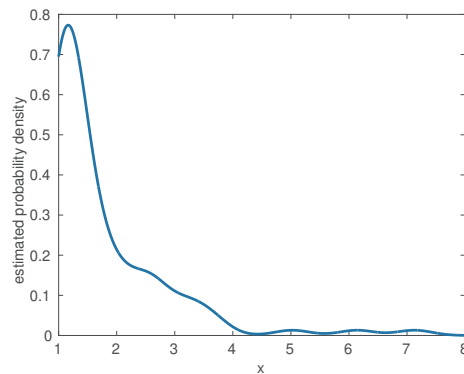
Figure 1.8: Kernels used in kernel density estimation. Figure from [https://en.wikipedia.org/wiki/Kernel_\(statistics\)](https://en.wikipedia.org/wiki/Kernel_(statistics))



(a) Scaled histogram



(b) Boxcar kernel



(c) Gaussian kernel

Figure 1.9: Estimating the probability density function from (a) the histogram or (b-c) via kernel density estimation.

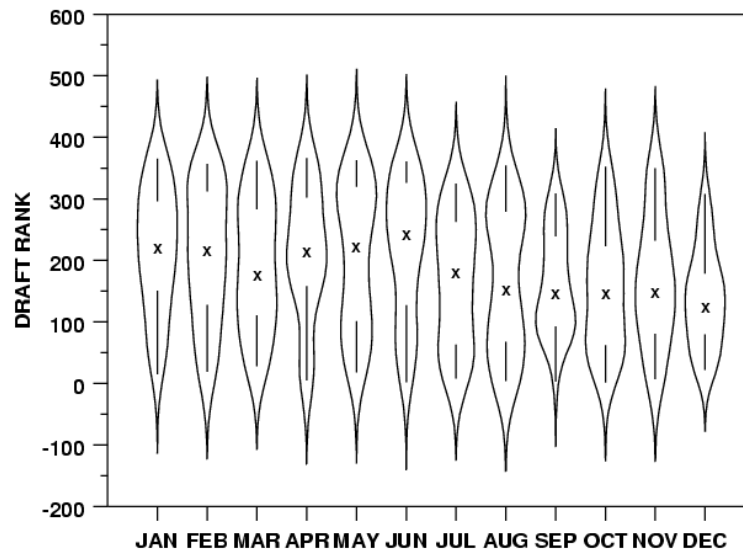


Figure 1.10: Example violin plot. Figure from https://en.wikipedia.org/wiki/violin_plot

1.2.6 Violin Plot

The Violin plot is a combination of a box plot and a rotated/reflected kernel density plot that is attached to both sides of the box plot. It describes robust (box plot) as well as non-robust but more informative (kernel density plot) aspects of the samples in a single plot. An example is plotted in Figure 1.10.

1.3 Data Pre-Processing

Oftentimes, it is useful to apply various transformations to the data to prepare it for further analysis. Typical pre-processing transformations include normalisation and outlier removal.

1.3.1 Standardisation

Data standardisation refers to data transformations that make variables comparable by dismissing certain characteristics. It classically refers to normalising the data to have zero (sample) mean and unit (sample) variance. We have seen this operation already in the definition of skewness and kurtosis. It can help to compare populations with different means and variances. It may, however, also refer to other kinds of transformations to make all variables comparable, for example, transforming the variables to be in the unit interval $[0, 1]$. Common further transformations that are being used are removing the average value of each single data vector, re-scaling the vector to unit norm, or computing the logarithm of its values. The transformations are often problem dependent.

Normalising the data to have zero sample mean and unit sample variance is a linear transformation. We will now express this transformation in terms of matrix operations.

Centring Matrix

We consider n observations $\mathbf{x}_1, \dots, \mathbf{x}_n$, where each observation is a d -dimensional vector. We can put these observations into a $d \times n$ matrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$. The sample mean now is a d -dimensional vector $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_d)^\top$ with elements

$$\bar{x}_i = \frac{1}{n} \sum_{j=1}^n (\mathbf{X})_{ij}, \quad (1.50)$$

where $(\mathbf{X})_{ij}$ denotes the (ij) -th element of \mathbf{X} . Subtracting this mean from each observation (i.e. each column vector in \mathbf{X}) results in new vectors

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}} \quad (1.51)$$

which have zero sample mean. These vectors give rise to another $d \times n$ matrix

$$\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n) = (\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}) = \mathbf{X} - \underbrace{(\bar{\mathbf{x}}, \dots, \bar{\mathbf{x}})}_{n \text{ times}} \quad (1.52)$$

We will now show that this sample mean removal can also be performed using a matrix multiplication

$$\tilde{\mathbf{X}} = \mathbf{X} \mathbf{C}_n, \quad (1.53)$$

where \mathbf{C}_n denotes the centring matrix

$$\mathbf{C}_n = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top, \quad (1.54)$$

\mathbf{I}_n is the $n \times n$ identity matrix and $\mathbf{1}_n = (1 \ 1 \ \dots \ 1)^\top$ is a vector of ones.

To see this, we write the sample mean as a matrix-vector product

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \frac{1}{n} \mathbf{X} \mathbf{1}_n \quad (1.55)$$

The matrix that contains the sample mean vector n -times $(\bar{\mathbf{x}}, \dots, \bar{\mathbf{x}})$ can be written as an outer product

$$\bar{\mathbf{x}} \mathbf{1}_n^\top = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_d \end{pmatrix} \underbrace{(1 \ 1 \ \dots \ 1)}_{n \text{ times}} = \begin{pmatrix} \bar{x}_1 & \bar{x}_1 & \dots & \bar{x}_1 \\ \bar{x}_2 & \bar{x}_2 & \dots & \bar{x}_2 \\ \vdots & & & \\ \bar{x}_d & \bar{x}_d & \dots & \bar{x}_d \end{pmatrix} = \underbrace{(\bar{\mathbf{x}}, \dots, \bar{\mathbf{x}})}_{n \text{ times}}. \quad (1.56)$$

Together with (1.55), we can represent this matrix as

$$(\bar{\mathbf{x}}, \dots, \bar{\mathbf{x}}) = \mathbf{X} \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \quad (1.57)$$

Plugging this into (1.52), we can show (1.53):

$$\tilde{\mathbf{X}} = \mathbf{X} - (\bar{\mathbf{x}}, \dots, \bar{\mathbf{x}}) \quad (1.58)$$

$$= \mathbf{X} - \mathbf{X} \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \quad (1.59)$$

$$= \mathbf{X} \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \quad (1.60)$$

$$= \mathbf{X} \mathbf{C}_n. \quad (1.61)$$

As we have seen, multiplying with the centring matrix from the right removes the sample mean of each row. Removing the sample mean again does not change anything. For this reason, the centring matrix is idempotent, that is $\mathbf{C}_n \mathbf{C}_n = \mathbf{C}_n$.

Multiplying the centring matrix from the *left* removes the sample mean of each *column*. For a vector $\mathbf{v} = (v_1 \ v_2 \ \dots \ v_n)^\top$ it holds

$$\mathbf{C}_n \mathbf{v} = \begin{pmatrix} v_1 - \bar{v} \\ v_2 - \bar{v} \\ \vdots \\ v_n - \bar{v} \end{pmatrix} \quad (1.62)$$

The centring matrix is thus a projection matrix that projects vectors on the space orthogonal to $\mathbf{1}_n$.

The centring matrix is not computationally useful for actually removing the mean, but it is useful in analytical calculations as we will see in later chapters.

Scaling to Unit Variance

We can use the centring matrix to express the sample covariance matrix (1.27) in the following way:

$$\text{cov}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top \quad (1.63)$$

$$= \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top \quad (1.64)$$

$$= \frac{1}{n} \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top \quad (1.65)$$

$$= \frac{1}{n} \mathbf{X} \mathbf{C}_n \mathbf{X}^\top \quad (1.66)$$

Recall that we can transform the sample covariance matrix to the sample correlation matrix using (1.41). The corresponding data transformation follows from (1.40) and is the multiplication from the left with the matrix that contains the inverse of the sample standard deviations on the diagonal and zeros off-diagonal, that is

$$\mathbf{z}_i = \text{diag} \left(\frac{1}{\text{std}(\mathbf{x})} \right) \tilde{\mathbf{x}}_i. \quad (1.67)$$

where $\tilde{\mathbf{x}}_i$ denotes the centred data. The elements of the \mathbf{z}_i now have zero sample mean and unit sample variance. Note, however, that the elements may still be correlated.

1.3.2 Outlier Detection and Removal

An outlier is an observation that seems unusual compared to others. This is a vague definition, which reflects the various possible causes for outliers.

An outlier can be due to an error in the data acquisition stage, for instance because a measurement device did not work properly. An observation may, however, also appear unusual because it does not conform to the current assumptions

that are made about the data. In the former case, we may omit the corrupted observations from the analysis, while in the latter case, the observations contain valuable information that should not be discarded.

Some outliers can be spotted by the methods above for describing univariate or bivariate data. If there is a strong difference between the mean and median, for example, the cause may be an outlier. Scatter plots, quantiles, histograms and in particular violin plots further enable one to spot outliers.

An additional way to detect outliers is to use *Tukey's fences* which define an interval for usual observations. The interval essentially corresponds to what is shown with the whiskers in box plots. Indeed, Tukey's fences are given by

$$[Q_1 - k(Q_3 - Q_1), Q_3 + k(Q_3 - Q_1)] = [Q_1 - k\text{IQR}(x), Q_3 + k\text{IQR}(x)] \quad (1.68)$$

for $k \geq 0$, and most commonly, $k = 1.5$. Observations outside of this interval are often labelled as outliers, in particular in box plots.

References

- [1] T.-H. Kim and H. White. "On more robust estimation of skewness and kurtosis". In: *Finance Research Letters* 1.1 (2004), pp. 56–73.

Chapter 2

Principal Component Analysis

This chapter presents several equivalent views on principal component analysis (PCA). The three main themes are finding directions in the data space along which the data are maximally variable, finding lower-dimensional yet accurate representations of the data, and formulating a probabilistic model of PCA. We assume that the data have been centred, i.e. that the sample mean has been subtracted from the data points \mathbf{x}_i , and that the corresponding random vector \mathbf{x} has zero mean.

2.1 PCA by Variance Maximisation

We first explain how to find the principal component direction sequentially and then formulate the equivalent simultaneous maximisation problem.

2.1.1 First Principal Component Direction

The first principal component direction is the unit vector \mathbf{w}_1 for which the projected data $\mathbf{w}_1^\top \mathbf{x}_i$ are maximally variable, where variability is measured by the sample variance. Equivalently, we can work with the random vector \mathbf{x} and look for the direction \mathbf{w}_1 for which the variance of $z_1 = \mathbf{w}_1^\top \mathbf{x}$ is maximal.

The variance $\text{Var}[z_1]$ can be expressed in terms of the covariance matrix Σ of \mathbf{x} ,

$$\text{Var}[z_1] = \text{Var}[\mathbf{w}_1^\top \mathbf{x}] = \mathbf{w}_1^\top \Sigma \mathbf{w}_1, \quad (2.1)$$

which follows from (1.40) with $\mathbf{A} = \mathbf{w}_1^\top$. The first principal component direction is thus the solution to the following optimisation problem:

$$\begin{aligned} & \underset{\mathbf{w}_1}{\text{maximise}} && \mathbf{w}_1^\top \Sigma \mathbf{w}_1 \\ & \text{subject to} && \|\mathbf{w}_1\| = 1 \end{aligned} \quad (2.2)$$

The eigenvalue decomposition of Σ allows us to find a solution in closed form. Let

$$\Sigma = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top, \quad (2.3)$$

where \mathbf{U} is an orthogonal matrix and where $\mathbf{\Lambda}$ is diagonal with eigenvalues $\lambda_i \geq 0$ (see Section A.8). We further assume that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$. As the columns

\mathbf{u}_i of \mathbf{U} form an orthogonal basis, we can express \mathbf{w}_1 as

$$\mathbf{w}_1 = \sum_{i=1}^d a_i \mathbf{u}_i = \mathbf{U} \mathbf{a}, \quad (2.4)$$

where $\mathbf{a} = (a_1, \dots, a_d)^\top$. The quadratic form $\mathbf{w}_1^\top \boldsymbol{\Sigma} \mathbf{w}_1$ can thus be written as

$$\mathbf{w}_1^\top \boldsymbol{\Sigma} \mathbf{w}_1 = \mathbf{a}^\top \underbrace{\mathbf{U}^\top \mathbf{U}}_{\mathbf{I}_d} \boldsymbol{\Lambda} \underbrace{\mathbf{U}^\top \mathbf{U}}_{\mathbf{I}_d} \mathbf{a} = \mathbf{a}^\top \boldsymbol{\Lambda} \mathbf{a} = \sum_{i=1}^d a_i^2 \lambda_i, \quad (2.5)$$

and the unit norm constraint on \mathbf{w}_1 becomes

$$\|\mathbf{w}_1\|^2 = \mathbf{w}_1^\top \mathbf{w}_1 = \mathbf{a}^\top \mathbf{U}^\top \mathbf{U} \mathbf{a} = \mathbf{a}^\top \mathbf{a} = \sum_{i=1}^d a_i^2 \stackrel{!}{=} 1 \quad (2.6)$$

An equivalent formulation of the optimisation problem in (2.2) is thus

$$\begin{aligned} & \underset{a_1, \dots, a_d}{\text{maximise}} && \sum_{i=1}^d a_i^2 \lambda_i \\ & \text{subject to} && \sum_{i=1}^d a_i^2 = 1 \end{aligned} \quad (2.7)$$

As $\lambda_1 \geq \lambda_i, i = 2, \dots, d$, setting a_1 to one and the remaining a_i to zero is a solution to the optimisation problem. This is the unique solution if λ_1 is the largest eigenvalue. But if, for example, $\lambda_1 = \lambda_2$, the solution is not unique any more: any a_1 and a_2 with $a_1^2 + a_2^2 = 1$ satisfy the constraint and yield the same objective. Assuming from now on that $\lambda_1 > \lambda_i, i = 2, \dots, d$, the unique \mathbf{w}_1 that solves the optimisation problem in (2.2) is

$$\mathbf{w}_1 = \mathbf{U} \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{u}_1. \quad (2.8)$$

The corresponding value of the objective $\mathbf{w}_1^\top \boldsymbol{\Sigma} \mathbf{w}_1$ is λ_1 .

The first principal component direction \mathbf{w}_1 is thus given by the eigenvector of the covariance matrix of \mathbf{x} that has the largest eigenvalue. The random variable $z_1 = \mathbf{w}_1^\top \mathbf{x}$ is called the first principal component of \mathbf{x} .

The variance of z_1 is equal to λ_1 —the largest eigenvalue of $\boldsymbol{\Sigma}$ and the maximal value of the objective in (2.2). We say that λ_1 is the variance of \mathbf{x} explained by the first principal component (direction). Since \mathbf{x} is assumed centred, the expected value of z_1 is zero,

$$\mathbb{E}[z_1] = \mathbb{E}[\mathbf{w}_1^\top \mathbf{x}] = \mathbf{w}_1^\top \underbrace{\mathbb{E}[\mathbf{x}]}_0 = 0. \quad (2.9)$$

In practice, we work with the centred data points \mathbf{x}_i . The projections $\mathbf{w}_1^\top \mathbf{x}_i, i = 1, \dots, n$ are often called the first principal components too, but also, more precisely, the first principal component scores. Collecting the centred data points into the $d \times n$ data matrix \mathbf{X} ,

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n), \quad (2.10)$$

the (row) vector \mathbf{z}_1^\top with all first principal component scores is given by $\mathbf{w}_1^\top \mathbf{X}$.

2.1.2 Subsequent Principal Component Directions

Given \mathbf{w}_1 , the next principal component direction \mathbf{w}_2 is chosen so that it maximises the variance of the projection $\mathbf{w}_2^\top \mathbf{x}$ and so that it reveals something “new” in the data, i.e. something that \mathbf{w}_1 has not uncovered. This puts a constraint on \mathbf{w}_2 , and in PCA, the constraint is implemented by requiring that \mathbf{w}_2 is orthogonal to \mathbf{w}_1 .

The second principal component direction is hence defined as the solution to the optimisation problem:

$$\begin{aligned} & \underset{\mathbf{w}_2}{\text{maximise}} && \mathbf{w}_2^\top \boldsymbol{\Sigma} \mathbf{w}_2 \\ & \text{subject to} && \|\mathbf{w}_2\| = 1 \\ & && \mathbf{w}_2^\top \mathbf{w}_1 = 0 \end{aligned} \tag{2.11}$$

As before, we decompose $\boldsymbol{\Sigma}$ as $\mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top$ and write \mathbf{w}_2 as $\mathbf{w}_2 = \mathbf{U}\mathbf{b}$. Since $\mathbf{w}_1 = \mathbf{u}_1$, $\mathbf{w}_2^\top \mathbf{w}_1$ equals

$$\mathbf{b}^\top \mathbf{U}^\top \mathbf{u}_1 = \mathbf{b}^\top \begin{pmatrix} 1 \\ 0 \\ \vdots \end{pmatrix} = b_1 \tag{2.12}$$

and the constraint $\mathbf{w}_2^\top \mathbf{w}_1 = 0$ becomes the constraint $b_1 = 0$. The optimisation problem in (2.11) can thus be equally expressed as:

$$\begin{aligned} & \underset{b_1, \dots, b_d}{\text{maximise}} && \sum_{i=1}^d b_i^2 \lambda_i \\ & \text{subject to} && \sum_{i=1}^d b_i^2 = 1 \\ & && b_1 = 0 \end{aligned} \tag{2.13}$$

We can insert the constraint $b_1 = 0$ directly into the other equations to obtain

$$\begin{aligned} & \underset{b_2, \dots, b_d}{\text{maximise}} && \sum_{i=2}^d b_i^2 \lambda_i \\ & \text{subject to} && \sum_{i=2}^d b_i^2 = 1 \end{aligned} \tag{2.14}$$

The optimisation problem is structurally the same as in (2.7); now we just optimise over b_2, \dots, b_d . As $\lambda_2 \geq \lambda_i$, $i = 3, \dots, d$, an optimal vector \mathbf{b} is $(0, 1, 0, \dots, 0)^\top$ and hence

$$\mathbf{w}_2 = \mathbf{U} \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{u}_2 \tag{2.15}$$

is a solution to the optimisation problem. Furthermore, the value of $\mathbf{w}_2^\top \boldsymbol{\Sigma} \mathbf{w}_2$ is λ_2 . As discussed for the first principal component, this solution is unique if $\lambda_2 > \lambda_i$, $i = 3, \dots, d$, which we here assume to be the case.

The second principal component direction \mathbf{w}_2 is thus given by the eigenvector of the covariance matrix of \mathbf{x} that has the second largest eigenvalue. Analogue to the first principal component z_1 , the second principal component is $z_2 = \mathbf{w}_2^\top \mathbf{x}$ with mean $\mathbb{E}[z_2] = 0$ and variance $\text{Var}[z_2] = \lambda_2$. The principal components are uncorrelated:

$$\mathbb{E}[z_1 z_2] = \mathbb{E}[\mathbf{w}_1^\top \mathbf{x} \mathbf{w}_2^\top \mathbf{x}] \tag{2.16}$$

$$= \mathbb{E}[\mathbf{w}_1^\top \mathbf{x} \mathbf{x}^\top \mathbf{w}_2] = \mathbf{w}_1^\top \mathbb{E}[\mathbf{x} \mathbf{x}^\top] \mathbf{w}_2 \tag{2.17}$$

$$= \mathbf{w}_1^\top \boldsymbol{\Sigma} \mathbf{w}_2 \tag{2.18}$$

$$= \mathbf{u}_1^\top \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top \mathbf{u}_2 \tag{2.19}$$

$$= (1 \ 0 \ 0 \ \dots) \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_d \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{pmatrix} \tag{2.20}$$

$$= (\lambda_1 \ 0 \ 0 \ \dots) \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{pmatrix} \tag{2.21}$$

$$= 0. \tag{2.22}$$

The procedure that we used to obtain \mathbf{w}_2 given \mathbf{w}_1 can be iterated to obtain further principal component directions. Assume that we have already computed $\mathbf{w}_1, \dots, \mathbf{w}_{m-1}$. The m -th principal component direction \mathbf{w}_m is then defined as the solution to:

$$\begin{aligned} & \underset{\mathbf{w}_m}{\text{maximise}} && \mathbf{w}_m^\top \boldsymbol{\Sigma} \mathbf{w}_m \\ & \text{subject to} && \|\mathbf{w}_m\| = 1 \\ & && \mathbf{w}_m^\top \mathbf{w}_i = 0 \quad i = 1, \dots, m-1 \end{aligned} \tag{2.23}$$

Arguing as before, the m -th principal component direction \mathbf{w}_m is given by eigenvector \mathbf{u}_m that corresponds to the m -th largest eigenvalue of $\boldsymbol{\Sigma}$ (assuming that there are no ties with other eigenvalues). The random variable $z_m = \mathbf{w}_m^\top \mathbf{x}$ is called the m -th principal component, its variance $\text{Var}[z_m]$ is λ_m , it is of zero mean (because \mathbf{x} is zero mean), and all principal components are uncorrelated. The $\mathbf{w}_m^\top \mathbf{x}_i$, $i = 1, \dots, n$, are the m -th principal component scores.

The total variance of the z_m , $m = 1, \dots, k$, is said to be the variance explained by the k principal components. It equals

$$\sum_{m=1}^k \text{Var}[z_m] = \sum_{m=1}^k \lambda_m. \tag{2.24}$$

The variance explained by the k principal components is often reported relative to the sum of the variances of the random variables x_i that make up the random

vector \mathbf{x} . The resulting number is the “fraction of variance explained”. With (1.34), the total variance of \mathbf{x} is

$$\sum_{i=1}^d \text{Var}[x_i] = \sum_{i=1}^d \lambda_i, \quad (2.25)$$

so that

$$\text{fraction of variance explained} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i}. \quad (2.26)$$

The fraction of variance explained is a useful number to compute for assessing how much of the variability in the data is captured by the k principal components.

2.1.3 Simultaneous Variance Maximisation

In the previous section, the k principal component directions $\mathbf{w}_1, \dots, \mathbf{w}_k$ were determined in a sequential manner, each time maximising the variance of each projection. Instead of the sequential approach, we can also determine all directions concurrently by maximising the total variance of all projections, i.e. by solving the optimisation problem:

$$\begin{aligned} & \underset{\mathbf{w}_1, \dots, \mathbf{w}_k}{\text{maximise}} && \sum_{i=1}^k \mathbf{w}_i^\top \boldsymbol{\Sigma} \mathbf{w}_i \\ & \text{subject to} && \|\mathbf{w}_i\| = 1 \quad i = 1, \dots, k \\ & && \mathbf{w}_i^\top \mathbf{w}_j = 0 \quad i \neq j \end{aligned} \quad (2.27)$$

It turns out that the optimal $\mathbf{w}_1, \dots, \mathbf{w}_k$ from the sequential approach, i.e. the eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_k$, are also solving the joint optimisation problem in (2.27), so that the maximal variance of all projections is $\sum_{i=1}^k \lambda_k$. This result may be intuitively understandable, but there is a subtle technical point: The sequential approach corresponds to solving the optimisation problem in (2.27) in a greedy manner, and greedy algorithms are generally not guaranteed to yield the optimal solution. Nevertheless, one can show that simultaneous and sequential variance maximisation yield the same solution (B.1, optional reading).

2.2 PCA by Minimisation of Approximation Error

A set of k orthonormal vectors $\mathbf{w}_1, \dots, \mathbf{w}_k$ of dimension d spans a k -dimensional subspace of \mathbb{R}^d denoted by $\text{span}(\mathbf{w}_1, \dots, \mathbf{w}_k)$, see Section A.5. Moreover, the matrix \mathbf{P} ,

$$\mathbf{P} = \sum_{i=1}^k \mathbf{w}_i \mathbf{w}_i^\top = \mathbf{W}_k \mathbf{W}_k^\top, \quad \mathbf{W}_k = (\mathbf{w}_1, \dots, \mathbf{w}_k), \quad (2.28)$$

projects any vector onto said subspace. This means that we can decompose our data points \mathbf{x}_i into elements $\hat{\mathbf{x}}_i = \mathbf{P} \mathbf{x}_i$ that belong to $\text{span}(\mathbf{w}_1, \dots, \mathbf{w}_k)$ and “residual” vectors orthogonal to it (see Figure 2.1 and Section A.6). The projections $\hat{\mathbf{x}}_i$ are lower dimensional approximations of the \mathbf{x}_i that can be represented

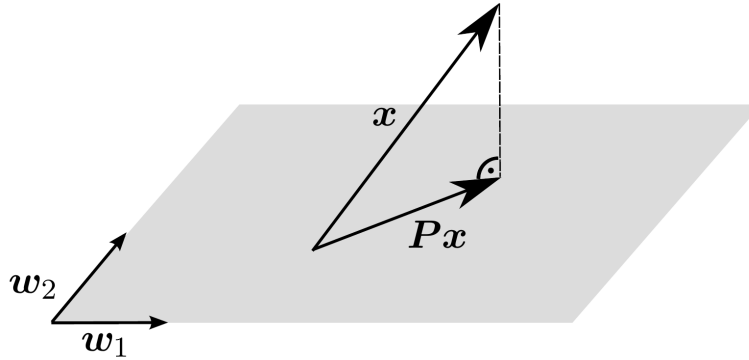


Figure 2.1: Orthogonal projection of \mathbf{x} onto the subspace spanned by the two orthonormal vectors \mathbf{w}_1 and \mathbf{w}_2 . The projection $\mathbf{P}\mathbf{x}$ can be written as a linear combination of \mathbf{w}_1 and \mathbf{w}_2 , and the residual $\mathbf{x} - \mathbf{P}\mathbf{x}$ is orthogonal to both vectors.

by the k coordinates $\mathbf{w}_1^\top \mathbf{x}_i, \dots, \mathbf{w}_k^\top \mathbf{x}_i$. Equivalently, the random vector \mathbf{x} can be approximated by $\hat{\mathbf{x}} = \mathbf{P}\mathbf{x} = \sum_{i=1}^k \mathbf{w}_i \mathbf{w}_i^\top \mathbf{x}$.

We now ask which subspace yields the approximations with the smallest error on average? Or equivalently, which subspace yields the smallest expected approximation error? The question can be formulated as the optimisation problem:

$$\begin{aligned} & \underset{\mathbf{w}_1, \dots, \mathbf{w}_k}{\text{minimise}} && \mathbb{E} \left\| \mathbf{x} - \sum_{i=1}^k \mathbf{w}_i \mathbf{w}_i^\top \mathbf{x} \right\|^2 \\ & \text{subject to} && \|\mathbf{w}_i\| = 1 && i = 1, \dots, k \\ & && \mathbf{w}_i^\top \mathbf{w}_j = 0 && i \neq j \end{aligned} \quad (2.29)$$

One can show that the optimisation problem is equivalent to the optimisation problem in (2.27), so that the optimal \mathbf{w}_i are the first k eigenvectors \mathbf{u}_i of the covariance matrix of \mathbf{x} , where “first k eigenvectors” means the eigenvectors with the k largest eigenvalues (B.2, optional reading). For this to make sense, it is assumed that the k -th eigenvalue is larger than the $(k + 1)$ -th eigenvalue.

In other words, the optimal k -dimensional subspace is spanned by $\mathbf{u}_1, \dots, \mathbf{u}_k$, the optimal projection matrix is $\mathbf{P} = \mathbf{U}_k \mathbf{U}_k^\top$, and the optimal lower dimensional representation of \mathbf{x} is $\hat{\mathbf{x}} = \mathbf{P}\mathbf{x}$. Since

$$\hat{\mathbf{x}} = \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x} = \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^\top \mathbf{x} = \sum_{i=1}^k \mathbf{u}_i z_i \quad (2.30)$$

the i -th principal component $z_i = \mathbf{u}_i^\top \mathbf{x}$ is the i -th coordinate of $\hat{\mathbf{x}}$ when represented in $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$.

We now compute the approximation error obtained by the optimal solution $\hat{\mathbf{x}}$. We have

$$\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \mathbb{E}[\mathbf{x}^\top \mathbf{x}] - 2\mathbb{E}[\mathbf{x}^\top \hat{\mathbf{x}}] + \mathbb{E}[\hat{\mathbf{x}}^\top \hat{\mathbf{x}}] \quad (2.31)$$

From Section 2.1.2, and (1.34), we know that $\mathbb{E}[\mathbf{x}^\top \mathbf{x}] = \mathbb{E}\|\mathbf{x}\|^2 = \sum_{i=1}^d \lambda_i$. Moreover, the \mathbf{u}_i are orthonormal vectors, so that the squared norm of $\hat{\mathbf{x}}$ equals

$\sum_{i=1}^k z_i^2$. Since we have assumed that the random variables have zero mean, previous results yield $\mathbb{E}[z_i^2] = \text{Var}[z_i] = \lambda_i$. With the expression for $\hat{\mathbf{x}}$ in (2.30), we further have

$$\mathbb{E}[\mathbf{x}^\top \hat{\mathbf{x}}] = \mathbb{E}[\mathbf{x}^\top \sum_{i=1}^k \mathbf{u}_i z_i] \tag{2.32}$$

$$= \sum_{i=1}^k \mathbb{E}[\mathbf{x}^\top \mathbf{u}_i z_i] \tag{2.33}$$

$$= \sum_{i=1}^k \mathbb{E}[z_i^2]. \tag{2.34}$$

The smallest expected approximation error when orthogonally projecting \mathbf{x} onto a k dimensional subspace thus equals

$$\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \sum_{i=1}^d \lambda_i - 2 \sum_{i=1}^k \lambda_i + \sum_{i=1}^k \lambda_i \tag{2.35}$$

$$= \sum_{i=k+1}^d \lambda_i, \tag{2.36}$$

which is the sum of the eigenvalues whose eigenvectors were omitted from the optimal subspace. Computing the relative approximation error highlights the connection between minimising approximation error and maximising the variance explained by the principal components,

$$\frac{\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|^2}{\mathbb{E} \|\mathbf{x}\|^2} = 1 - \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} = 1 - \text{fraction of variance explained.} \tag{2.37}$$

The fraction of variance explained by a principal component (direction) thus equals the relative reduction in approximation error that is achieved by including it into the subspace.

2.3 PCA by Low Rank Matrix Approximation

This section uses the theory of low rank matrix approximation to provide a complementary view on the PCA principles of variance maximisation and minimisation of approximation error.

2.3.1 Approximating the Data Matrix

We will here see that the principal component directions and scores together yield the best low rank approximation of the data matrix, and that the PC directions and scores can be computed by a singular value decomposition (SVD).

Let \mathbf{X} be the $d \times n$ data matrix that contains the centred data points \mathbf{x}_i in its columns,

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n). \tag{2.38}$$

We can express \mathbf{X} via its singular value decomposition as

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\top. \quad (2.39)$$

The $d \times d$ matrix \mathbf{U} and the $n \times n$ matrix \mathbf{V} are orthonormal with the vectors $\mathbf{u}_i \in \mathbb{R}^d$ and $\mathbf{v}_i \in \mathbb{R}^n$ in their columns. The \mathbf{u}_i are called the left singular vectors while the \mathbf{v}_i are called the right singular vectors. The matrix \mathbf{S} is $d \times n$ and zero everywhere but in the first r diagonal elements,

$$\mathbf{S} = \begin{pmatrix} s_1 & & & & & \\ & \ddots & & & & \\ & & s_r & & & \\ & & & & \mathbf{0} & \\ & & & & & \\ \mathbf{0} & & & & & \mathbf{0} \end{pmatrix}. \quad (2.40)$$

The s_i are the singular values. They are positive and assumed ordered from large to small. The number $r \leq \min(d, n)$ is the rank of \mathbf{X} . The matrix \mathbf{X} can further be written as

$$\mathbf{X} = \sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^\top. \quad (2.41)$$

Section A.7 provides further background on the SVD.

Assume we would like to approximate \mathbf{X} by a matrix $\hat{\mathbf{X}}$ of rank $k < r$. Judging the accuracy of the approximation by the sum of squared differences in the individual matrix elements, we can determine $\hat{\mathbf{X}}$ by solving the optimisation problem

$$\begin{aligned} & \underset{\mathbf{M}}{\text{minimise}} && \sum_{ij} ((\mathbf{X})_{ij} - (\mathbf{M})_{ij})^2 \\ & \text{subject to} && \text{rank}(\mathbf{M}) = k \end{aligned} \quad (2.42)$$

The sum of squared differences $\sum_{ij} ((\mathbf{X})_{ij} - (\mathbf{M})_{ij})^2$ is called the Frobenius norm between \mathbf{X} and \mathbf{M} and typically denoted by $\|\mathbf{X} - \mathbf{M}\|_F$.

It is known from linear algebra that the optimal low rank approximation is given by $\hat{\mathbf{X}}$,

$$\hat{\mathbf{X}} = \sum_{i=1}^k s_i \mathbf{u}_i \mathbf{v}_i^\top, \quad (2.43)$$

and that the corresponding approximation error is

$$\|\mathbf{X} - \hat{\mathbf{X}}\|_F = \sum_{i=k+1}^r s_i^2, \quad (2.44)$$

see (A.63) and (A.65) in Section A.10. The solution to the optimisation problem is thus rather simple: We just keep the first k terms in (2.41).

How does this relate to principal component analysis? It turns out that

- the left singular vectors \mathbf{u}_i are the eigenvectors of the (estimated) covariance matrix and hence equal to the principal component directions,

- the squared singular values s_i^2 are related to the eigenvalues λ_i of the covariance matrix by

$$\lambda_i = \frac{s_i^2}{n}, \quad (2.45)$$

- and that the principal component scores $\mathbf{z}_i^\top = \mathbf{u}_i^\top \mathbf{X}$ for principal component direction i are equal to the i -th right singular vector after scaling,

$$\mathbf{z}_i^\top = s_i \mathbf{v}_i^\top. \quad (2.46)$$

We can thus write the approximate data matrix $\hat{\mathbf{X}}$ in (2.43) as

$$\hat{\mathbf{X}} = \sum_{i=1}^k \mathbf{u}_i \mathbf{z}_i^\top, \quad (2.47)$$

which underlines how the k principal component directions \mathbf{u}_i and corresponding principal component scores \mathbf{z}_i together approximately represent the data \mathbf{X} .

The stated connections can be seen as follows: As we assume that the data points are centred, an estimate of the covariance matrix is given by the sample covariance matrix,

$$\mathbf{\Sigma} \approx \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \frac{1}{n} \mathbf{X} \mathbf{X}^\top. \quad (2.48)$$

Using $\mathbf{\Sigma}$ to denote both the covariance and the sample covariance matrix, we have

$$\mathbf{\Sigma} = \frac{1}{n} \mathbf{U} \mathbf{S} \mathbf{V}^\top (\mathbf{V} \mathbf{S}^\top \mathbf{U}^\top) = \mathbf{U} \left(\frac{1}{n} \mathbf{S} \mathbf{S}^\top \right) \mathbf{U}^\top \quad (2.49)$$

so that the eigenvectors of $\mathbf{\Sigma}$ are the left singular vectors \mathbf{u}_i of the data matrix \mathbf{X} with eigenvalues λ_i as in (2.45). The i -th principal component scores were defined as the projections $\mathbf{w}_i^\top \mathbf{x}_j$ of the data points \mathbf{x}_j onto the i -th principal component direction \mathbf{w}_i . Collecting all scores into the $1 \times n$ row vector \mathbf{z}_i^\top , we have

$$\mathbf{z}_i^\top = \mathbf{w}_i^\top \mathbf{X} = \mathbf{u}_i^\top \mathbf{X} = \mathbf{u}_i^\top \mathbf{U} \mathbf{S} \mathbf{V}^\top = \mathbf{u}_i^\top \sum_{j=1}^r \mathbf{u}_j s_j \mathbf{v}_j^\top = s_i \mathbf{v}_i^\top. \quad (2.50)$$

which means that the i -th principal component scores are given by the i -th right singular vector when multiplied with its singular value s_i as claimed in (2.46).

2.3.2 Approximating the Sample Covariance Matrix

Approximating the data matrix with a matrix of lower rank yielded the first k principal component directions and scores. We here show that the first principal component directions can also be obtained by a low rank approximation of the sample covariance matrix. This provides a complementary view to PCA, in that finding directions in the data space with maximal variance is also maximally preserving the variance structure of the original data. This approach does, however, not directly yield principal component scores.

The optimisation problem that we aim to solve is:

$$\begin{aligned} & \underset{\mathbf{M}}{\text{minimise}} && \|\boldsymbol{\Sigma} - \mathbf{M}\|_F \\ & \text{subject to} && \text{rank}(\mathbf{M}) = k \\ & && \mathbf{M}^\top = \mathbf{M} \end{aligned} \quad (2.51)$$

Much like the first k components of the SVD are solving the optimisation problem in (2.42), results from linear algebra tell us that the optimal low rank approximation of $\boldsymbol{\Sigma}$ is given by the first k components of its eigenvalue decomposition $\mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top$, i.e. by $\sum_{i=1}^k \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$, see (A.66) in Section A.10.

2.3.3 Approximating the Gram Matrix

The Gram matrix is defined as the $n \times n$ matrix \mathbf{G} ,

$$\mathbf{G} = \mathbf{X}^\top \mathbf{X}. \quad (2.52)$$

Its (ij) -th element $(\mathbf{G})_{ij}$ is the inner product between \mathbf{x}_i and \mathbf{x}_j . The matrix is positive semidefinite, i.e. its eigenvalues are non-negative. It is here shown that the first principal component scores provide an optimal low rank approximation of \mathbf{G} . Hence, finding coordinates that minimise the average approximation error of the data points \mathbf{x}_i is also maximally preserving the inner product structure between them.

With the singular value decomposition of \mathbf{X} in (2.39), the Gram matrix has the following eigenvalue decomposition

$$\mathbf{G} = (\mathbf{USV}^\top)^\top (\mathbf{USV}^\top) = (\mathbf{VS}^\top \mathbf{U}^\top) (\mathbf{USV}^\top) = \mathbf{VS}^\top \mathbf{SV}^\top = \mathbf{V}\tilde{\boldsymbol{\Lambda}}\mathbf{V}^\top \quad (2.53)$$

i.e. its eigenvectors are the right-singular vectors \mathbf{v}_i of \mathbf{X} , and the diagonal matrix $\tilde{\boldsymbol{\Lambda}} = \mathbf{S}^\top \mathbf{S}$ contains the eigenvalues s_i^2 ordered from large to small.

Like for the sample covariance matrix, we can determine the best rank k approximation of the Gram matrix \mathbf{G} . It is given by $\hat{\mathbf{G}}$,

$$\hat{\mathbf{G}} = \sum_{i=1}^k \mathbf{v}_i s_i^2 \mathbf{v}_i^\top. \quad (2.54)$$

In (2.46), we have seen that $\mathbf{z}_i = s_i \mathbf{v}_i$ is the column vector with all i -th principal component scores. We thus have

$$\hat{\mathbf{G}} = \sum_{i=1}^k \mathbf{z}_i \mathbf{z}_i^\top, \quad (2.55)$$

which shows that the k principal scores together maximally preserve the inner product structure of the data.

Denote by $\tilde{\boldsymbol{\Lambda}}_k$ the diagonal matrix with the top k eigenvalues of \mathbf{G} , and by \mathbf{V}_k ,

$$\mathbf{V}_k = (\mathbf{v}_1, \dots, \mathbf{v}_k) \quad (2.56)$$

the matrix with the corresponding eigenvectors. The $k \times n$ matrix with the principal component scores as its rows is then

$$\mathbf{Z} = \sqrt{\tilde{\Lambda}_k} \mathbf{V}_k^\top. \quad (2.57)$$

We have the square root because the singular values s_i are the square root of the eigenvalues of \mathbf{G} . Hence, we can compute the principal component scores directly from the Gram matrix of the centred data, without first computing the principal component directions. This can be done without knowing \mathbf{X} as long as \mathbf{G} is available.

2.4 Probabilistic PCA

In this section, we formulate a parametric probabilistic model of the data such that its maximum likelihood estimate corresponds to PCA. This description connects PCA to the powerful probabilistic modelling approach to data analysis.

2.4.1 Probabilistic Model

In Section 2.2, we saw that PCA can be understood as error minimisation in a subspace approximation. Building on this view, we will now define a k -dimensional latent variable \mathbf{z} corresponding to a variable living in the principal component subspace (Tipping and Bishop, 1999). We assume that the elements of \mathbf{z} are statistically independent and standard normal distributed, i.e.

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}) \quad (2.58)$$

where

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d |\det(\boldsymbol{\Sigma})|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (2.59)$$

denotes the multivariate normal distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Here, $\det(\boldsymbol{\Sigma})$ denotes the determinant of $\boldsymbol{\Sigma}$.

We further assume that our observable data are generated in the following way:

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}, \quad (2.60)$$

where \mathbf{W} denotes a $d \times k$ matrix, $\boldsymbol{\mu}$ denotes a d -dimensional constant mean vector and $\boldsymbol{\epsilon}$ denotes a d -dimensional zero-mean Gaussian-distributed noise variable with covariance $\sigma^2 \mathbf{I}$ and σ^2 denoting a scalar. The quantities \mathbf{W} , $\boldsymbol{\mu}$ and σ^2 are the parameters of the model. The generative process of this model is illustrated in Figure 2.2, left and centre panels.

2.4.2 Joint, Conditional and Observation Distributions

For constant \mathbf{z} , the term $\mathbf{W}\mathbf{z} + \boldsymbol{\mu}$ is constant which means that our data \mathbf{x} conditioned on the latent variable \mathbf{z} are also multivariate normal distributed:

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}). \quad (2.61)$$

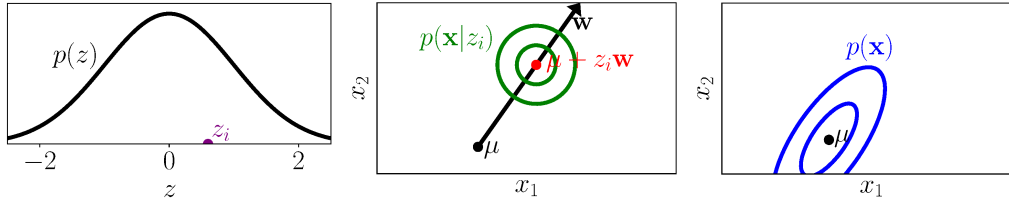


Figure 2.2: Illustration of the model of probabilistic PCA for a 2-dimensional data space and a 1-dimensional latent space. Left: Distribution of latent variable. Centre: An observation \mathbf{x} is drawn by first drawing a latent variable z_i and then adding independent Gaussian noise with mean $\boldsymbol{\mu} + z_i \mathbf{w}$ and covariance $\sigma^2 \mathbf{I}$ (green contour lines). Right: Contour lines of the density of the marginal distribution $p(\mathbf{x})$. Adapted from (Bishop, 2006, Figure 12.9).

Now, consider the joint distribution $p(\mathbf{z}, \mathbf{x}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$:

$$p(\mathbf{z}, \mathbf{x}) = \frac{1}{\text{const}} \exp \left(-\frac{1}{2} \left[(\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu})^\top \left(\frac{1}{\sigma^2} \mathbf{I} \right) (\mathbf{x} - \mathbf{W}\mathbf{z} - \boldsymbol{\mu}) + z^\top z \right] \right), \quad (2.62)$$

where “const” denotes terms that are independent of \mathbf{x} and \mathbf{z} . This is again a multivariate normal distribution over \mathbf{x} and \mathbf{z} . We will now determine the mean and covariance matrix of this joint distribution.

For a multivariate normal distribution, the term in the exponential function generally has the form

$$-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2} \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} + \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \text{const}, \quad (2.63)$$

where here “const” denotes terms that are independent of \mathbf{x} . The second order terms of the form $\mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \mathbf{x}$ contain the inverse of the covariance matrix and the linear terms $\mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$ contain the mean. Thus, if we encounter a term of the form

$$-\frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{x}^\top \boldsymbol{\xi} + \text{const} \quad (2.64)$$

then, by matching terms, we can obtain the covariance matrix by taking the inverse of the second order coefficient matrix

$$\boldsymbol{\Sigma} = \mathbf{A}^{-1} \quad (2.65)$$

and we can obtain the mean by multiplying the linear coefficient vector with the covariance matrix

$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \boldsymbol{\xi}. \quad (2.66)$$

To obtain the covariance matrix of the joint distribution $p(\mathbf{z}, \mathbf{x})$, we treat the tuple (\mathbf{z}, \mathbf{x}) as a single vector variable and consider the second order terms in the exponential of (2.62):

$$\begin{aligned} & -\frac{1}{2} \left[\mathbf{z}^\top \left(\mathbf{I} + \mathbf{W}^\top \frac{1}{\sigma^2} \mathbf{W} \right) \mathbf{z} + \mathbf{x}^\top \frac{1}{\sigma^2} \mathbf{x} - \mathbf{x}^\top \frac{1}{\sigma^2} \mathbf{W} \mathbf{z} - \mathbf{z}^\top \mathbf{W}^\top \frac{1}{\sigma^2} \mathbf{x} \right] \\ & = -\frac{1}{2} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^\top \begin{pmatrix} \mathbf{I} + \mathbf{W}^\top \frac{1}{\sigma^2} \mathbf{W} & -\mathbf{W}^\top \frac{1}{\sigma^2} \\ -\frac{1}{\sigma^2} \mathbf{W} & \frac{1}{\sigma^2} \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix} \quad (2.67) \end{aligned}$$

By means of block matrix inversion (A.22), we obtain the covariance matrix of $p(\mathbf{z}, \mathbf{x})$:

$$\text{Cov} \left[\begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix} \right] = \begin{pmatrix} \mathbf{I} + \mathbf{W}^\top \frac{1}{\sigma^2} \mathbf{W} & -\mathbf{W}^\top \frac{1}{\sigma^2} \\ -\frac{1}{\sigma^2} \mathbf{W} & \frac{1}{\sigma^2} \mathbf{I} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{I} & \mathbf{W}^\top \\ \mathbf{W} & \mathbf{W} \mathbf{W}^\top + \sigma^2 \mathbf{I} \end{pmatrix}. \quad (2.68)$$

To find the mean of the joint distribution $p(\mathbf{z}, \mathbf{x})$, we consider the linear terms in the exponential of (2.62):

$$-\mathbf{z}^\top \mathbf{W}^\top \frac{1}{\sigma^2} \boldsymbol{\mu} + \mathbf{x}^\top \frac{1}{\sigma^2} \boldsymbol{\mu} = \begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix}^\top \begin{pmatrix} -\mathbf{W}^\top \frac{1}{\sigma^2} \boldsymbol{\mu} \\ \frac{1}{\sigma^2} \boldsymbol{\mu} \end{pmatrix}. \quad (2.69)$$

With (2.66), we obtain the mean as

$$\mathbb{E} \left[\begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix} \right] = \begin{pmatrix} \mathbf{I} & \mathbf{W}^\top \\ \mathbf{W} & \sigma^2 \mathbf{I} + \mathbf{W} \mathbf{W}^\top \end{pmatrix} \begin{pmatrix} -\mathbf{W}^\top \frac{1}{\sigma^2} \boldsymbol{\mu} \\ \frac{1}{\sigma^2} \boldsymbol{\mu} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \boldsymbol{\mu} \end{pmatrix}. \quad (2.70)$$

Having this representation, we can immediately obtain the distribution $p(\mathbf{x})$ of the observations from (2.68) (lower right block matrix) and (2.70) (lower partition):

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{W} \mathbf{W}^\top + \sigma^2 \mathbf{I}). \quad (2.71)$$

This distribution is illustrated in Figure 2.2 right panel.

2.4.3 Maximum Likelihood

We now return to the case where we are given actual observations in the form of a $d \times n$ centred data matrix \mathbf{X} .

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n). \quad (2.72)$$

For given data that have not been centred, one can show that the (unique) maximum likelihood solution of the mean is the sample mean. Like earlier in this chapter, we now assume that the sample mean has been subtracted from the data points, and that the resulting random vector \mathbf{x} has zero mean, i.e. $\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} = \mathbf{0}$.

The probabilistic PCA model has parameters \mathbf{W} and σ^2 that we need to infer. Following (2.71), the log-likelihood has the form

$$\log p(\mathbf{X} | \mathbf{W}, \sigma^2) = \sum_{i=1}^n \log p(\mathbf{x}_i | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) \quad (2.73)$$

$$= \sum_{i=1}^n \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (2.74)$$

with $\boldsymbol{\mu} = \mathbf{0}$ and

$$\boldsymbol{\Sigma} = \mathbf{W} \mathbf{W}^\top + \sigma^2 \mathbf{I} \quad (2.75)$$

denoting the covariance matrix of \mathbf{x} under the model. To further develop the log-likelihood in (2.74), note that a quadratic form $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ can be expressed as follows

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \text{trace} \left[\mathbf{x}^\top \mathbf{A} \mathbf{x} \right] \quad (2.76)$$

$$\stackrel{\text{(A.8)}}{=} \text{trace} \left[\mathbf{A} \mathbf{x} \mathbf{x}^\top \right], \quad (2.77)$$

where the first line holds because $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is a scalar. Due to the linearity of the trace operator, we thus have

$$\sum_{i=1}^n \mathbf{x}_i^\top \mathbf{A} \mathbf{x}_i = n \operatorname{trace} \left[\mathbf{A} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right] \quad (2.78)$$

$$= n \operatorname{trace} \left[\mathbf{A} \hat{\Sigma} \right] \quad (2.79)$$

where

$$\hat{\Sigma} = \operatorname{cov}(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \frac{1}{n} \mathbf{X} \mathbf{X}^\top \quad (2.80)$$

denotes the sample covariance matrix. With this relationship in hand, we can write the log-likelihood in (2.74) compactly as follows

$$\log p(\mathbf{X} | \mathbf{W}, \sigma^2) = \sum_{i=1}^n \log \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \Sigma) \quad (2.81)$$

$$\stackrel{(2.59)}{=} \sum_{i=1}^n \left[-\frac{1}{2} d \log(2\pi) - \frac{1}{2} \log(|\det(\Sigma)|) - \frac{1}{2} \mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i \right] \quad (2.82)$$

$$= -\frac{n}{2} \left[d \log(2\pi) + \log(|\det(\Sigma)|) - \frac{1}{2} \sum_{i=1}^n \mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i \right] \quad (2.83)$$

$$\stackrel{(2.79)}{=} -\frac{n}{2} \left[d \log(2\pi) + \log(|\det(\Sigma)|) + \operatorname{trace}(\Sigma^{-1} \hat{\Sigma}) \right]. \quad (2.84)$$

To derive the maximum likelihood solutions \mathbf{W}_{ML} and σ_{ML}^2 , one can maximise (2.84) with respect to \mathbf{W} and σ^2 . It was shown by Tipping and Bishop, 1999 that the solution for \mathbf{W} is given by

$$\mathbf{W}_{\text{ML}} = \mathbf{U}_k (\mathbf{\Lambda}_k - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}, \quad (2.85)$$

where \mathbf{U}_k are the k principal eigenvectors of $\hat{\Sigma}$ with corresponding eigenvalues in the $k \times k$ diagonal matrix $\mathbf{\Lambda}_k$ and \mathbf{R} is an arbitrary $k \times k$ orthogonal matrix. \mathbf{R} can be interpreted as a rotation or reflection in the latent space and indicates that the solution is not unique. For instance, we can set \mathbf{R} to \mathbf{I} yielding $\mathbf{W}_{\text{ML}} = \mathbf{U}_k (\mathbf{\Lambda}_k - \sigma^2 \mathbf{I})^{1/2}$. For σ^2 the solution is

$$\sigma_{\text{ML}}^2 = \frac{1}{d-k} \sum_{i=k+1}^d \lambda_i, \quad (2.86)$$

where $\lambda_{k+1}, \dots, \lambda_d$ denote the smallest eigenvalues of $\hat{\Sigma}$. σ_{ML}^2 therefore represents the average lost variance per residual dimension. The derivation of this result is lengthy (see Tipping and Bishop, 1999) and will not be covered in this lecture. Briefly, one can consider different classes of \mathbf{W} for which $\partial \log p(\mathbf{X} | \mathbf{W}, \sigma^2) / \partial \mathbf{W} = \mathbf{0}$. Here, it helps to express \mathbf{W} in terms of its singular value decomposition.

Practically, to calculate the maximum likelihood solution, we can first compute the eigenvalue decomposition of $\hat{\Sigma}$, then compute σ_{ML}^2 using (2.86) and compute \mathbf{W}_{ML} using (2.85) with $\sigma^2 = \sigma_{\text{ML}}^2$.

2.4.4 Relation to PCA

PCA maps the observed data to principal component scores: $\mathbf{z}_i = \mathbf{U}_k^\top \mathbf{x}_i$. Probabilistic PCA, on the other hand, maps the latent space to the data space (2.60). In the probabilistic PCA framework, the closest thing to the PCA mapping is the posterior distribution $p(\mathbf{z}|\mathbf{x}_i)$ which represents a whole distribution in the latent space (instead of being just a single vector \mathbf{z}_i like in the case of PCA).

To find the posterior distribution, we fix \mathbf{x} in the equations of the joint distribution (2.62). This, again, yields a multivariate normal distribution. The second order coefficients for variable \mathbf{z} and constant \mathbf{x} are given by the upper left block matrix on the right hand side of (2.67). The covariance matrix of $p(\mathbf{z}|\mathbf{x})$ is therefore given by the inverse of that matrix:

$$\text{Cov}[\mathbf{z}|\mathbf{x}] = \left(\mathbf{I} + \mathbf{W}^\top \frac{1}{\sigma^2} \mathbf{W} \right)^{-1} = \sigma^2 (\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I})^{-1} = \sigma^2 \mathbf{M}^{-1}, \quad (2.87)$$

where

$$\mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}. \quad (2.88)$$

The linear coefficients for variable \mathbf{z} and constant \mathbf{x} are given by

$$\frac{1}{2} \mathbf{W}^\top \frac{1}{\sigma^2} \mathbf{x} + \frac{1}{2} \left(\mathbf{x}^\top \frac{1}{\sigma^2} \mathbf{W} \right)^\top - \mathbf{W}^\top \frac{1}{\sigma^2} \boldsymbol{\mu} = \mathbf{W}^\top \frac{1}{\sigma^2} (\mathbf{x} - \boldsymbol{\mu}). \quad (2.89)$$

With (2.66), we obtain the mean of $p(\mathbf{z}|\mathbf{x})$ as

$$\mathbb{E}[\mathbf{z}|\mathbf{x}] = \sigma^2 \mathbf{M}^{-1} \mathbf{W}^\top \frac{1}{\sigma^2} (\mathbf{x} - \boldsymbol{\mu}) = \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x} - \boldsymbol{\mu}) \quad (2.90)$$

and therefore

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x} - \boldsymbol{\mu}), \sigma^2 \mathbf{M}^{-1}). \quad (2.91)$$

The mean $\mathbb{E}[\mathbf{z}|\mathbf{x}]$ features a pseudo-inverse of \mathbf{W} , namely $\mathbf{W}^\dagger = \mathbf{M}^{-1} \mathbf{W}^\top$. For $\sigma = 0$, note that $\mathbf{W}^\dagger \mathbf{W} \mathbf{z} = (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{W} \mathbf{z} = \mathbf{z}$ if the $d \times k$ matrix \mathbf{W} has rank k , so that $\mathbf{W}^\top \mathbf{W}$ is invertible. If \mathbf{x} has been generated according to (2.60) with $\sigma = 0$, so that $\mathbf{x} - \boldsymbol{\mu} = \mathbf{W} \mathbf{z}^*$ where \mathbf{z}^* denotes some fixed value of \mathbf{z} , we thus obtain $\mathbb{E}[\mathbf{z}|\mathbf{x}] = \mathbf{z}^*$. This means that we can recover \mathbf{z} from \mathbf{x} if the noise in the generative process is zero and if the $d \times k$ matrix \mathbf{W} that was used to \mathbf{x} is known. If $\sigma > 0$, we can think that \mathbf{M}^{-1} is a regularised inverse of $\mathbf{W}^\top \mathbf{W}$, which exists even if \mathbf{W} does not have rank k , which makes

$$\mathbf{W}^\dagger = (\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^\top \quad (2.92)$$

a generalised pseudo-inverse of \mathbf{W} .

More realistically, we do not know \mathbf{W} that was used to generate \mathbf{x} . We now discuss the relation of the posterior distribution to the PCA mapping when using the maximum likelihood estimate \mathbf{W}_{ML} for \mathbf{W} . Recall that in PCA, the projection onto the subspace spanned by the first k principal components is given by $\hat{\mathbf{x}} = \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}$, see (2.30). Assuming that the data are centred, the corresponding projection in the probabilistic PCA framework is

$$\hat{\mathbf{x}} = \mathbf{W}_{\text{ML}} \mathbb{E}[\mathbf{z}|\mathbf{x}] \quad (2.93)$$

$$= \mathbf{W}_{\text{ML}} \mathbf{M}_{\text{ML}}^{-1} \mathbf{W}_{\text{ML}}^\top \mathbf{x}, \quad (2.94)$$

where $\mathbf{M}_{\text{ML}} = \mathbf{W}_{\text{ML}}^\top \mathbf{W}_{\text{ML}} + \sigma^2 \mathbf{I}$. For $\sigma^2 = 0$, we recover the PCA projection:

$$\mathbf{W}_{\text{ML}} \rightarrow \mathbf{U}_k \mathbf{\Lambda}_k^{1/2} \mathbf{R} \quad (2.95)$$

$$\mathbf{M}_{\text{ML}} \rightarrow \mathbf{R}^\top \mathbf{\Lambda}_k \mathbf{R} \quad (2.96)$$

$$\mathbf{W}_{\text{ML}} \mathbf{M}_{\text{ML}}^{-1} \mathbf{W}_{\text{ML}}^\top \mathbf{x} \rightarrow \mathbf{U}_k \mathbf{\Lambda}_k^{1/2} \mathbf{R} \mathbf{R}^\top \mathbf{\Lambda}_k^{-1} \mathbf{R} \mathbf{R}^\top \mathbf{\Lambda}_k^{1/2} \mathbf{U}_k^\top \mathbf{x} \quad (2.97)$$

$$= \mathbf{U}_k \mathbf{\Lambda}_k^{1/2} \mathbf{\Lambda}_k^{-1} \mathbf{\Lambda}_k^{1/2} \mathbf{U}_k^\top \mathbf{x} \quad (2.98)$$

$$= \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}, \quad (2.99)$$

where we have used that $\mathbf{U}_k^\top \mathbf{U}_k$ is the $k \times k$ identity matrix and that $\mathbf{R}^\top = \mathbf{R}^{-1}$ for rotation matrices. Thus, PCA can be seen as a special case of probabilistic PCA when the noise variance σ^2 is negligibly small.

References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006. ISBN: 0387310738.
- [2] Michael E Tipping and Christopher M Bishop. "Probabilistic principal component analysis". In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 61.3 (1999), pp. 611–622.

Chapter 3

Dimensionality Reduction

Dimensionality reduction is about representing the data in a lower dimensional space in such a way that certain properties of the data are preserved as much as possible. Dimensionality reduction can be used to visualise high-dimensional data if the plane is chosen as the lower dimensional space. Taking linear dimensionality reduction by principal component analysis (PCA) as starting point, several nonlinear dimensionality reduction methods are presented.

3.1 Linear Dimensionality Reduction

We can represent d -dimensional data by their first k principal components, or more precisely, their first k principal component scores. The principal components can be computed when the data are given in form of data vectors, and, importantly, also when given in form of inner products or distances between them.

3.1.1 From Data Points

Denote the uncentred data by $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ and the corresponding data matrix by $\tilde{\mathbf{X}}$,

$$\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n). \quad (3.1)$$

We first centre the data and form the matrix \mathbf{X} ,

$$\mathbf{X} = \tilde{\mathbf{X}}\mathbf{C}_n \quad \mathbf{C}_n = \mathbf{I}_n - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top, \quad (3.2)$$

where \mathbf{C}_n is the centring matrix from (1.53). Depending on the application, we may want to further process the data, e.g. by some form of standardisation that was introduced in Section 1.3.

We can now compute the principal components via an eigenvalue decomposition of the covariance matrix Σ ,

$$\Sigma = \frac{1}{n}\mathbf{X}\mathbf{X}^\top. \quad (3.3)$$

Denoting the matrix with the top k eigenvectors \mathbf{u}_i by \mathbf{U}_k ,

$$\mathbf{U}_k = (\mathbf{u}_1, \dots, \mathbf{u}_k), \quad (3.4)$$

the matrix with the principal component (PC) scores is

$$\mathbf{Z} = \mathbf{U}_k^\top \mathbf{X}. \quad (3.5)$$

While \mathbf{X} is $d \times n$, \mathbf{Z} is $k \times n$. The column vectors of \mathbf{Z} have dimension $k \leq d$ and form a lower dimensional representation of the data.

In dimensionality reduction, we are mostly interested in the PC scores, rather than the PC directions. We can thus bypass the computation of the PC directions and compute the PC scores directly from the Gram matrix \mathbf{G} introduced in (2.52),

$$\mathbf{G} = \mathbf{X}^\top \mathbf{X}. \quad (3.6)$$

With (2.57), the $k \times n$ matrix \mathbf{Z} in (3.5) equals

$$\mathbf{Z} = \sqrt{\tilde{\Lambda}_k} \mathbf{V}_k^\top, \quad (3.7)$$

where the diagonal $k \times k$ matrix $\tilde{\Lambda}_k$ contains the top k eigenvalues of \mathbf{G} ordered from large to small, and the $n \times k$ matrix \mathbf{V}_k contains the corresponding eigenvectors.

3.1.2 From Inner Products

The elements $(\mathbf{G})_{ij}$ of the Gram matrix $\mathbf{G} = \mathbf{X}^\top \mathbf{X}$ are the inner products between the centred data points \mathbf{x}_i and \mathbf{x}_j ,

$$(\mathbf{G})_{ij} = \mathbf{x}_i^\top \mathbf{x}_j = (\tilde{\mathbf{x}}_i - \boldsymbol{\mu})^\top (\tilde{\mathbf{x}}_j - \boldsymbol{\mu}). \quad (3.8)$$

Since we can compute the principal component scores (but not the directions) by an eigenvalue decomposition of the Gram matrix, we can do dimensionality reduction without actually having seen the data points \mathbf{x}_i . Knowing their inner products is enough.

But what should we do if we are only given the inner products between the original data points and not between the centred data points? That is, what should we do if we are only given the matrix $\tilde{\mathbf{G}}$,

$$\tilde{\mathbf{G}} = \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}, \quad (3.9)$$

and not \mathbf{G} ?

It turns out that we can compute \mathbf{G} from $\tilde{\mathbf{G}}$. With $\mathbf{X} = \tilde{\mathbf{X}}\mathbf{C}_n$, where \mathbf{C}_n is the centring matrix, we have

$$\mathbf{G} = \mathbf{X}^\top \mathbf{X} = \mathbf{C}_n^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{C}_n = \mathbf{C}_n \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{C}_n = \mathbf{C}_n \tilde{\mathbf{G}} \mathbf{C}_n, \quad (3.10)$$

where we have used that \mathbf{C}_n is a symmetric matrix. This operation is called double centring: Multiplying $\tilde{\mathbf{G}}$ with \mathbf{C}_n from the right makes all rows have zero average while multiplying it from the left makes all columns have a zero average.

Since inner products can be used to measure the similarity between data points, matrices like $\tilde{\mathbf{G}}$ and \mathbf{G} are sometimes called similarity matrices. We can thus say that we can do dimensionality reduction by PCA given a similarity matrix (with inner products) only.

3.1.3 From Distances

We here show that we can exactly recover the PC scores if we are only given the squared distances δ_{ij}^2 between the data points $\tilde{\mathbf{x}}_i$ and $\tilde{\mathbf{x}}_j$,

$$\delta_{ij}^2 = \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|^2 = (\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j)^\top (\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j). \quad (3.11)$$

The matrix $\mathbf{\Delta}$ with elements δ_{ij}^2 is called a distance matrix. Note that matrices with non-squared δ_{ij} are also called distance matrices. There is some ambiguity in the terminology.

The trick is to recover the Gram matrix \mathbf{G} in (2.52) from the distance matrix $\mathbf{\Delta}$: First, we note that the δ_{ij}^2 equal the squared distances between the centred data points \mathbf{x}_i ,

$$\delta_{ij}^2 = \|(\tilde{\mathbf{x}}_i - \boldsymbol{\mu}) - (\tilde{\mathbf{x}}_j - \boldsymbol{\mu})\|^2 = \|\mathbf{x}_i - \mathbf{x}_j\|^2 = (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j). \quad (3.12)$$

Multiplying out yields

$$\delta_{ij}^2 = \|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2 - 2\mathbf{x}_i^\top \mathbf{x}_j. \quad (3.13)$$

Importantly the first term $\|\mathbf{x}_i\|^2$ is constant along row i of the matrix $\mathbf{\Delta}$. We can thus eliminate it by multiplying $\mathbf{\Delta}$ with the centring matrix \mathbf{C}_n from the right. This is because

$$(\mathbf{\Delta}\mathbf{C}_n)_{ij} = (\mathbf{\Delta})_{ij} - \frac{1}{n} \sum_{j=1}^n (\mathbf{\Delta})_{ij}, \quad (3.14)$$

as, for example, in (1.61). In more detail, let us compute

$$\frac{1}{n} \sum_{j=1}^n (\mathbf{\Delta})_{ij} = \frac{1}{n} \sum_{j=1}^n \delta_{ij}^2 \quad (3.15)$$

$$= \|\mathbf{x}_i\|^2 + \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j\|^2 - 2\frac{1}{n} \sum_{j=1}^n \mathbf{x}_i^\top \mathbf{x}_j \quad (3.16)$$

which equals

$$\frac{1}{n} \sum_{j=1}^n (\mathbf{\Delta})_{ij} = \|\mathbf{x}_i\|^2 + \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j\|^2 - 2\mathbf{x}_i^\top \underbrace{\left(\frac{1}{n} \sum_{j=1}^n \mathbf{x}_j\right)}_0 \quad (3.17)$$

$$= \|\mathbf{x}_i\|^2 + \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j\|^2, \quad (3.18)$$

because the \mathbf{x}_j are centred and hence $\sum_j \mathbf{x}_j = 0$. We thus find that

$$(\mathbf{\Delta}\mathbf{C}_n)_{ij} = \|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2 - 2\mathbf{x}_i^\top \mathbf{x}_j - \|\mathbf{x}_i\|^2 - \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j\|^2 \quad (3.19)$$

$$= \|\mathbf{x}_j\|^2 - 2\mathbf{x}_i^\top \mathbf{x}_j - \frac{1}{n} \sum_{j=1}^n \|\mathbf{x}_j\|^2. \quad (3.20)$$

Now, the terms $\|\mathbf{x}_j\|^2$ and $1/n \sum_{j=1}^n \|\mathbf{x}_j\|^2$ are constant along column j of the matrix $\Delta \mathbf{C}_n$. We can thus eliminate them by multiplying $\Delta \mathbf{C}_n$ with \mathbf{C}_n from the left. Calculations as above show that

$$(\mathbf{C}_n \Delta \mathbf{C}_n)_{ij} = (\Delta \mathbf{C}_n)_{ij} - \frac{1}{n} \sum_i (\Delta \mathbf{C}_n)_{ij} = -2\mathbf{x}_i^\top \mathbf{x}_j. \quad (3.21)$$

We thus have $\mathbf{C}_n \Delta \mathbf{C}_n = -2\mathbf{G}$, and hence obtain the desired result,

$$\mathbf{G} = -\frac{1}{2} \mathbf{C}_n \Delta \mathbf{C}_n. \quad (3.22)$$

In the previous section, we double centred the similarity matrix $\tilde{\mathbf{G}}$ to obtain \mathbf{G} . Here, we double centre the distance matrix Δ , and swap the signs to convert distances to similarities. Once \mathbf{G} is available, we compute the principal component scores as before via an eigenvalue decomposition, see the previous section or Equations (2.53) and (2.57).

3.1.4 Example

Figure 3.1(a) shows data that can be well represented by one principal component. The data vary mostly along the diagonal and projecting them onto the first principal component (red line) captures most of the variability. In Figure 3.1(b) we colour-code each data point by its principal component score. The scores are the coordinates of the data with respect to the basis given by the principal component direction. We can see that there is a good correspondence between the value of the scores and the location of the data points. The scores change smoothly from large to small as we move along the diagonal: they faithfully represent the data and capture their structure well.

The data in Figure 3.2, on the other hand, are not as well represented by the first principal component. The first principal component direction captures the direction of maximal variance but the data are treated as a cloud of points and the scores roughly indicate the location of the data along the y -axis but not their position on the circle. The principal component scores do not capture the circular structure of the data.

Why is PCA doing better for data as in Figure 3.1 than for data as in Figure 3.2? This can be understood by considering that PCA projects the data onto a lower-dimensional subspace (Section 2.2). Subspaces are closed under addition and multiplication, which means that any point on a line going through two points from the subspace is also included in the subspace (see e.g. Section A.5). For the data in Figure 3.2, however, there are gaps of empty space between two data points that are unlikely to be filled even if we had more data. Such kind of data are said to lie on a manifold, and lines between two points on a manifold may not be part of the manifold (see, for example, Chapter 16 of Izenman (2008) or Chapter 1 of Lee and Verleysen (2007)). If the data are part of a subspace, it is reasonable to judge the distance between two points by the length of the straight line connecting them, like in PCA, but if the data are on a manifold, straight lines are a poor measure of their distance. Similarly, the linear projections that are used to compute the principal component scores do not take the manifold structure of the data into account.

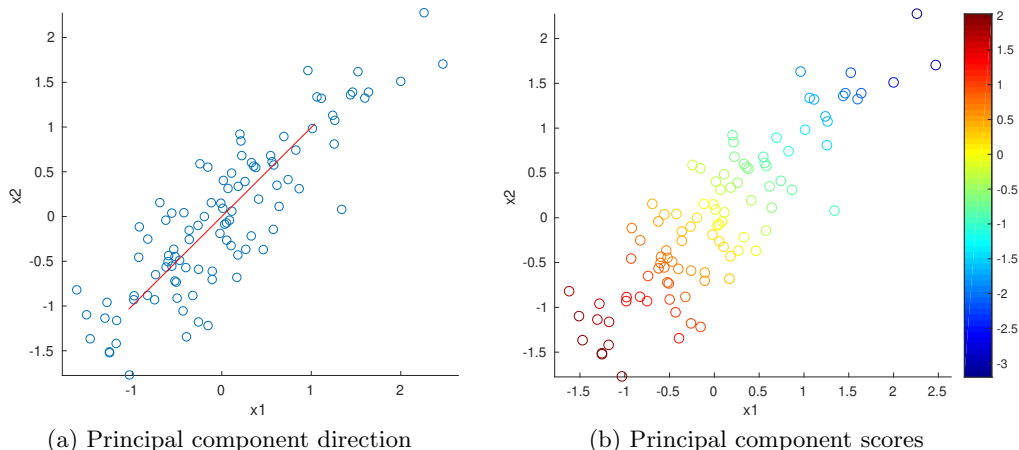


Figure 3.1: Dimensionality reduction by principal component analysis. (a) The red line shows the direction of the first PC direction. (b) The colours indicate the value of the PC score assigned to each data point.

3.2 Dimensionality Reduction by Kernel PCA

Principal component analysis uses linear projections to compute the lower dimensional representation of the data. We here discuss kernel PCA where the projections are typically nonlinear.

3.2.1 Idea

The principal components represent the data so that the variance is maximally preserved. Assume that we expand the dimensionality of the \mathbf{x}_i by transforming them to features $\phi(\mathbf{x}_i)$, which might be learned using neural networks. For the sake of argument, consider features of the form

$$\phi(\mathbf{x}) = (x_1, \dots, x_d, x_1x_2, \dots, x_1x_d, \dots, x_dx_d)^\top, \quad (3.23)$$

where $\mathbf{x} = (x_1, \dots, x_d)^\top$. The much higher dimensionality of the $\phi_i = \phi(\mathbf{x}_i)$ does not matter as long as we only compute k principal components from them.

Importantly, the k principal components maximally preserve the variance of the ϕ_i that contains much more information about the data than the variance of the \mathbf{x}_i . The covariance matrix for the particular $\phi(\mathbf{x})$ above, for example, contains terms like $\mathbb{E}(x_1x_2x_3^2)$ that measure non-linear correlations between the different dimensions of the data. Similarly, the principal components best approximate the ϕ_i which is much harder than approximating the \mathbf{x}_i , so that the principal components computed from the ϕ_i must capture more information about the data than the components computed from the \mathbf{x}_i .

Hence, we can power up PCA dimensionality reduction by choosing a transformation ϕ that maps the data points \mathbf{x}_i to $\phi_i = \phi(\mathbf{x}_i)$, and then computing the principal component scores from the new “data matrix” Φ ,

$$\Phi = (\phi_1, \dots, \phi_n), \quad (3.24)$$

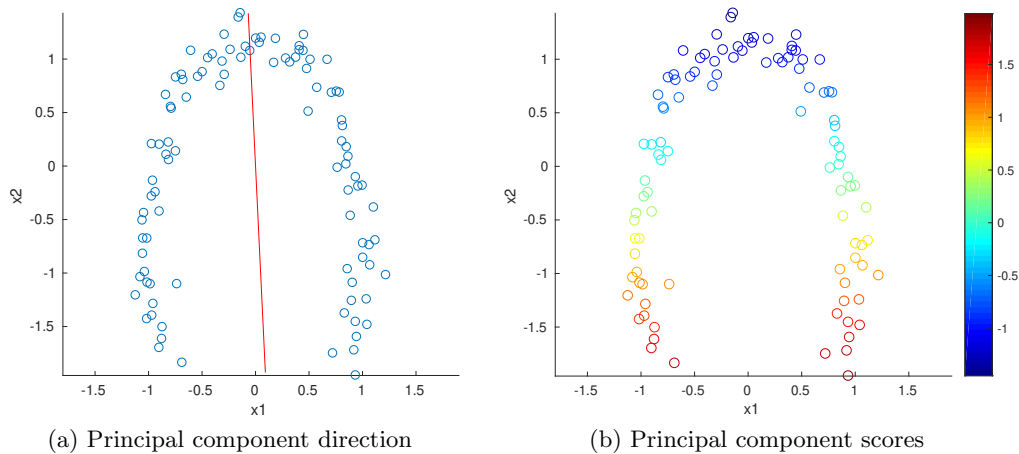


Figure 3.2: Dimensionality reduction by principal component analysis. (a) The red line shows the direction of the first PC direction. (b) The colours indicate the value of the PC score assigned to each data point.

rather than from the original data matrix \mathbf{X} . We call this approach dimensionality reduction by nonlinear PCA. (Note that “nonlinear PCA” sometimes refers to other kinds of methods too.)

3.2.2 Kernel Trick

Since we can compute the principal components scores from the Gram matrix of Φ , we actually do not need to know the individual ϕ_i , but only the inner products $\phi_i^\top \phi_j = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$.

The theory of reproducing kernel Hilbert spaces tells us that for some functions ϕ , the inner product can be computed as

$$\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j), \quad (3.25)$$

where $k(\mathbf{x}, \mathbf{x}')$ is called the kernel function (see, e.g. Schölkopf and Smola, 2002). This means that for some functions ϕ , we actually do not need to know the transformed data points ϕ_i to compute the inner product between them, it is enough to know the kernel $k(\mathbf{x}, \mathbf{x}')$. This is called the “kernel trick” and can be used to compute the (uncentred) Gram matrix of Φ as

$$(\tilde{\mathbf{G}})_{ij} = \phi_i^\top \phi_j = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j). \quad (3.26)$$

Performing PCA via a Gram matrix defined by kernels as above is called kernel PCA and has been introduced by Schölkopf, Smola, and Müller (1997).

Examples of kernels are the polynomial and Gaussian kernels,

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^a, \quad k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right), \quad (3.27)$$

where the exponent a and width-parameter σ^2 are hyperparameters that need to be chosen by the user. We see that the two kernels only require the inner

products or distances between the data points \mathbf{x}_i so that kernel PCA can also be used if that information is available only.

After specification of $\tilde{\mathbf{G}}$, we proceed exactly as in Section 3.1:

- Double centre $\tilde{\mathbf{G}}$ to compute

$$\mathbf{G} = \mathbf{C}_n \tilde{\mathbf{G}} \mathbf{C}_n. \quad (3.28)$$

- Compute the matrix \mathbf{Z} with the (kernel) PC scores by an eigenvalue decomposition of \mathbf{G} ,

$$\mathbf{Z} = \sqrt{\tilde{\mathbf{\Lambda}}_k} \mathbf{V}_k^\top, \quad (3.29)$$

where, as before, the diagonal $k \times k$ matrix $\tilde{\mathbf{\Lambda}}_k$ contains the top k eigenvalues of \mathbf{G} ordered from large to small, and the $n \times k$ matrix \mathbf{V}_k contains the corresponding eigenvectors.

3.2.3 Example

Let us reconsider the circularly structured data of Figure 3.2 and use nonlinear PCA to compute a one-dimensional representation. We map the data points \mathbf{x}_i to features ϕ_i using the function $\phi(\mathbf{x}) = \phi(x_1, x_2)$,

$$\phi(x_1, x_2) = (x_1, x_2, \sqrt{x_1^2 + x_2^2}, \text{atan}(x_2, x_1))^\top, \quad (3.30)$$

where x_1, x_2 are the first and second element of the vector \mathbf{x} . The last two elements in the vector $\phi(\mathbf{x})$ are the polar coordinates of \mathbf{x} , which should be helpful given the circular structure of the data. Figure 3.3 visualises the first principal component scores computed from the transformed data matrix Φ . We can see that the lower dimensional representation by the first PC scores is reflecting the circular structure of the data. But there is a discontinuity in the scores around the point $(-1, 0)$, and the scores still ignore that a piece of the circle is missing: data points on the lower left are assigned similar values as data points on the lower right.

Figure 3.4 visualises the one-dimensional representation achieved by kernel PCA with the Gaussian kernel in (3.27). The hyperparameter σ^2 was determined from the quantiles of all distances between all (different) data points. The results do not seem better than the results with ordinary PCA in Figure 3.2.

Centring was the only preprocessing for the results above. Let us next also scale each variable to unit variance before dimensionality reduction by kernel PCA (see Section 1.3.1 on data standardisation). Figure 3.5 shows that kernel PCA on the standardised data yields a mostly meaningful one-dimensional representation for a wide range of different tuning parameters σ^2 . The kernel PC scores change smoothly as we move on the data manifold. But the representation does ignore the gap between the lower-left and lower-right branch, so that points on the lower-left of the manifold (where $x_2 < -1.5$ and $x_1 \approx -1$) are considered closer to points on the lower-right of the manifold than points further up on the left branch. This may be considered a drawback of the representation.

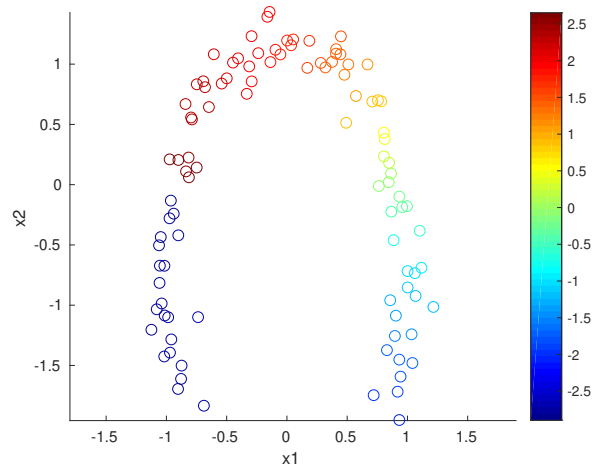


Figure 3.3: Dimensionality reduction by nonlinear PCA. Visualisation as in Figure 3.2(b).

3.3 Multidimensional Scaling

Multidimensional scaling (MDS) is an umbrella term for several methods that operate on dissimilarities δ_{ij} . Euclidean distances are examples of dissimilarities but dissimilarities are more general in that they can be any kind of measure of difference between two data items. The goal of MDS is to find a configuration of points in the plane, or the Euclidean space, so that their distances well represent the original dissimilarities.

3.3.1 Metric MDS

In metric MDS, the numerical values of the dissimilarities are assumed to carry information. This is in contrast to nonmetric MDS below where only the rank-order of the dissimilarities matters.

Denote the pairwise dissimilarities between n data points by δ_{ij} . A basic version of metric MDS consists in finding n points $\mathbf{z}_i \in \mathbb{R}^k$ that solve:

$$\underset{\mathbf{z}_1, \dots, \mathbf{z}_n}{\text{minimise}} \sum_{i < j} w_{ij} (\|\mathbf{z}_i - \mathbf{z}_j\| - \delta_{ij})^2, \quad (3.31)$$

where $\|\mathbf{z}_i - \mathbf{z}_j\|$ is the Euclidean distance between \mathbf{z}_i and \mathbf{z}_j ,

$$\|\mathbf{z}_i - \mathbf{z}_j\| = \sqrt{(\mathbf{z}_i - \mathbf{z}_j)^\top (\mathbf{z}_i - \mathbf{z}_j)}. \quad (3.32)$$

The $w_{ij} \geq 0$ are some weights specified by the user. The dimensionality k is typically set to two so that the data can be visualised on the plane. More complex versions of metric MDS exist where the dissimilarities δ_{ij} enter into the equation only after transformation with some monotonic function that is learned as well, see e.g. (Izenman, 2008, Section 13.7) or (Borg and Groenen, 2005, Chapter 9). The optimisation problem is typically solved by gradient descent.

For $w_{ij} = 1/\delta_{ij}$, the solution for the optimisation problem in (3.31) is called the Sammon nonlinear mapping. This choice of weights emphasises the faithful representation of small dissimilarities.

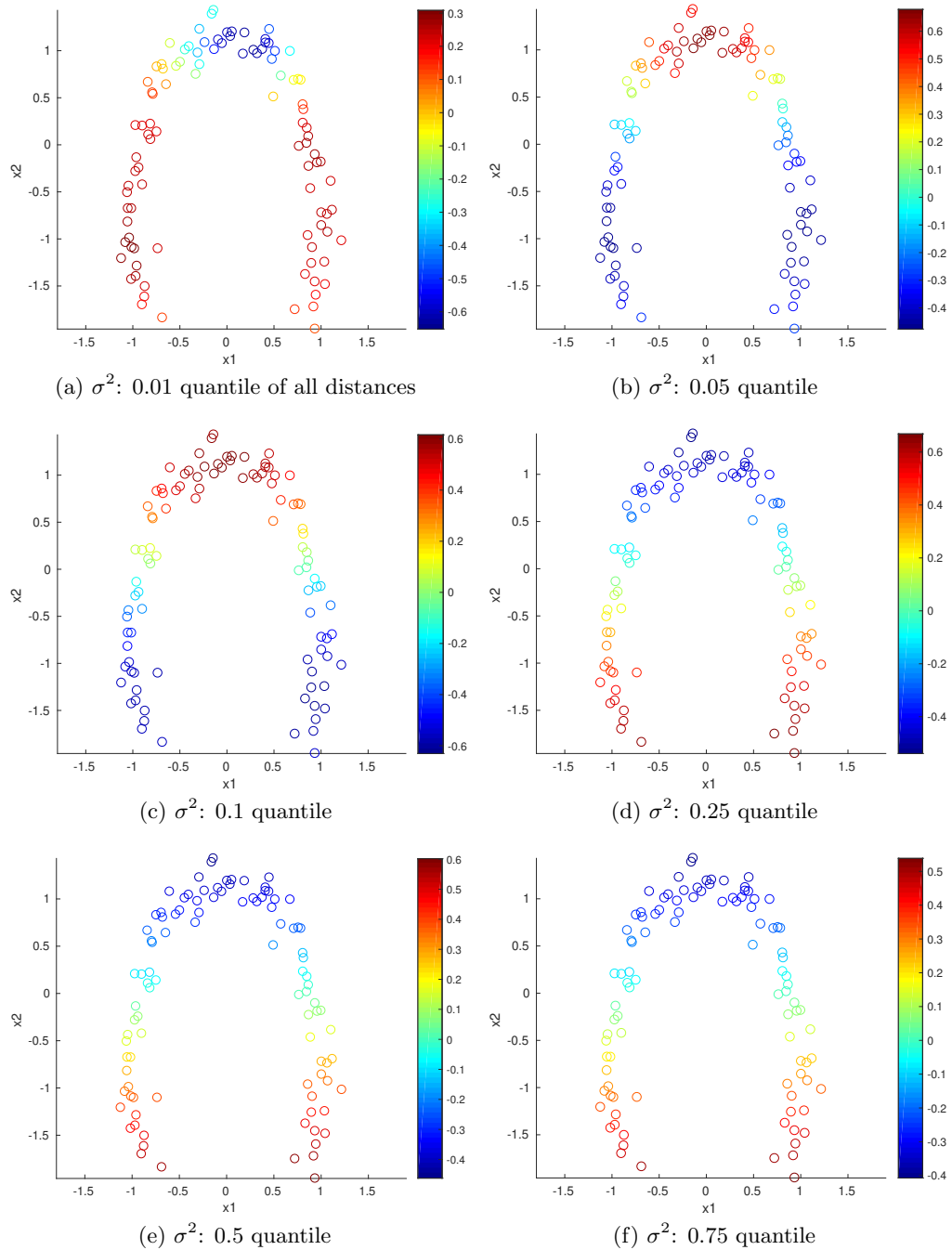


Figure 3.4: Dimensionality reduction by kernel principal component analysis. The Gaussian kernel was used where σ^2 was determined by the quantiles of the distances. The colours indicate the value of the (kernel) principal component score assigned to a data point. The sign of the scores in each panel is arbitrary.

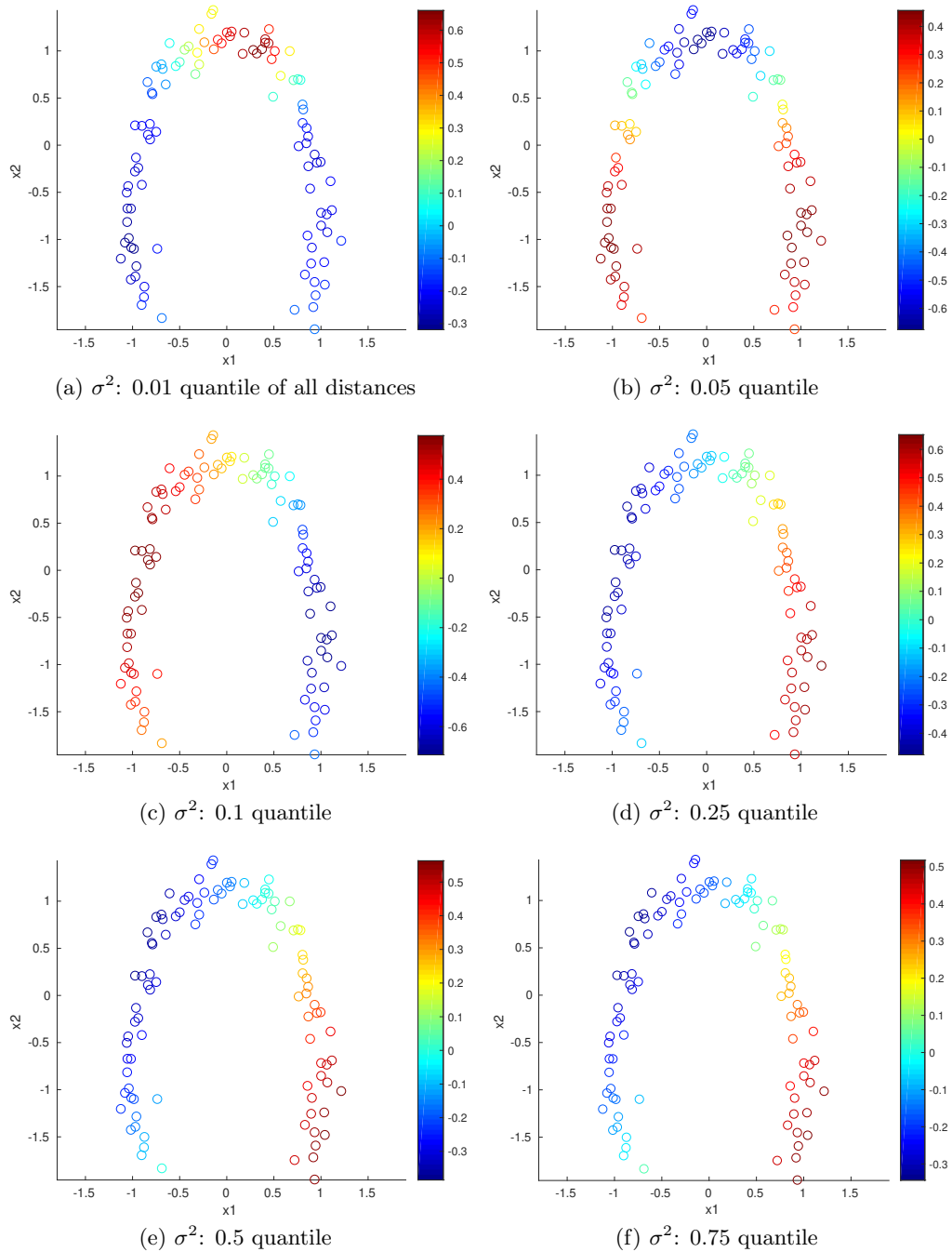


Figure 3.5: Dimensionality reduction by kernel principal component analysis on standardised data. The setup and visualisation is otherwise as in Figure 3.4.

3.3.2 Nonmetric MDS

In nonmetric MDS, only the relation between the δ_{ij} is assumed to matter, i.e. whether $\delta_{12} \geq \delta_{13}$ or $\delta_{12} \leq \delta_{13}$, and not the actual values of the dissimilarities. Such data are obtained, for example, when people are asked to rate the dissimilarity on a scale from 0 (“identical”) to 5 (“very different”).

Since the actual values of δ_{ij} do not matter, in nonmetric MDS, the optimisation problem in (3.31) is modified to

$$\underset{\mathbf{z}_1, \dots, \mathbf{z}_n, f}{\text{minimise}} \sum_{i < j} w_{ij} (\|\mathbf{z}_i - \mathbf{z}_j\| - f(\delta_{ij}))^2, \quad (3.33)$$

where f is a monotonic (non-decreasing) function that converts the dissimilarities to distances. The optimisation problem is typically solved by iterating between optimisation with respect to the \mathbf{z}_i and optimisation with respect to f , which can be done by regression (for further information, see, e.g. Izenman, 2008, Section 13.9).

3.3.3 Classical MDS

Classical MDS is also called classical scaling. It operates on the same kind of data as in metric scaling, that is, the actual numerical values of the dissimilarities are assumed to matter.

Classical scaling posits that the dissimilarities δ_{ij} are (squared) Euclidean distances between some unknown, hypothetical vectors of unknown dimensionality. Identifying the dissimilarity matrix $\mathbf{\Delta}$ that is formed by the δ_{ij} with a distance matrix between the unknown vectors brings us back to the setting from Section 3.1.3, and we can use the developed theory to determine the lower dimensional $\mathbf{z}_i \in \mathbb{R}^k, i = 1 \dots n$:

1. Compute the hypothetical Gram matrix \mathbf{G}' of the unknown centred data points,

$$\mathbf{G}' = -\frac{1}{2} \mathbf{C}_n \mathbf{\Delta} \mathbf{C}_n, \quad \mathbf{C}_n = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top, \quad (3.34)$$

as in (3.22). (The $'$ should emphasise that \mathbf{G}' is a hypothetical Gram matrix, it does not denote the transpose of the matrix.)

2. Compute the top k eigenvalues σ_k^2 and corresponding eigenvectors $\mathbf{v}_k \in \mathbb{R}^n$ of \mathbf{G} , and form the matrices $\tilde{\mathbf{\Lambda}}_k = \text{diag}(\sigma_1^2, \dots, \sigma_k^2)$ and $\mathbf{V}_k = (\mathbf{v}_1, \dots, \mathbf{v}_k)$.
3. The $k \times n$ matrix \mathbf{Z} with the \mathbf{z}_i as its columns,

$$\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n), \quad (3.35)$$

is then given by $\mathbf{Z} = \sqrt{\tilde{\mathbf{\Lambda}}_k} \mathbf{V}_k^\top$, as in (2.57).

Classical MDS can thus turn any dissimilarity matrix $\mathbf{\Delta}$ into a configuration of lower-dimensional vectors \mathbf{z}_i that represent the dissimilarities. It also has the nice property that it produces nested solutions because the classical MDS solution for $k' < k$ is directly given by the first k' coordinates of the k -dimensional \mathbf{z}_i .

There is one subtle technical caveat: The matrix Δ is symmetric but not necessarily positive semidefinite. This is because we only pretended that Δ corresponds to Euclidean distances in some unknown space, but this may only hold approximately. Since Δ is not necessarily positive semidefinite, some of its eigenvalues may be negative so that taking square roots as above in the third step would not yield meaningful representations. The simple fix is to choose k small enough that all eigenvalues contained in $\tilde{\Lambda}_k$ are indeed positive.

For matrices Δ that are not positive semidefinite, eigenvectors corresponding to negative eigenvalues are thus excluded. We can think of this operation as a way to approximate Δ by a positive semidefinite matrix. It turns out that the simple operation of excluding directions with negative eigenvalues is actually the optimal positive semidefinite approximation of Δ with respect to the Frobenius norm (see Section A.10.3). Further results from linear algebra show that the lower dimensional representation $Z = \sqrt{\tilde{\Lambda}_k} V_k^\top$ yields the best low rank approximation of G' with respect to the Frobenius norm (see Section A.10.4). That is, Z is the solution to

$$\begin{aligned} \underset{M}{\text{minimise}} \quad & \|(-\frac{1}{2}C_n\Delta C_n) - M^\top M\|_F \\ \text{subject to} \quad & \text{rank}(M^\top M) = k \end{aligned} \quad (3.36)$$

This is a different optimisation problem than the one in (3.31) for metric MDS, and the solution returned by classical and metric MDS are generally not the same.

3.3.4 Example

Let us apply the Sammon nonlinear mapping to the circularly shaped data in Figure 3.2, reducing the dimension from two to one. The Sammon mapping is the solution of the optimisation problem in (3.31), which can have multiple local minima. The algorithm was run multiple times and Figure 3.6 shows the two different solutions typically obtained. The solution in Figure 3.6(a) basically corresponds to the PCA-solution. The solution figure (b), however, shows that the learned one-dimensional representation, i.e. the points z_1, \dots, z_n in (3.31), do take the manifold structure of the data into account.

The solution in 3.6(a) assigns the same low dimensional coordinate to the point on the left and right branch of the manifold. Their distance is thus practically zero even though in the original data space, their distance is rather large. The solution (b) assigns different values to the points on the left and the right half of the circle, so that their distances in the lower dimensional space better matches their distances in the original space.

Figure 3.7 plots the distances in the original space against the distances in the lower dimensional space. The phenomenon described above is well visible in that the solution from 3.6(a) has distances equal to zero even though the original distances are around two.

3.4 Isomap

Classical MDS is used as part of the isometric feature mapping (Isomap) algorithm (Tenenbaum, Silva, and Langford, 2000) where the dissimilarity δ_{ij} between

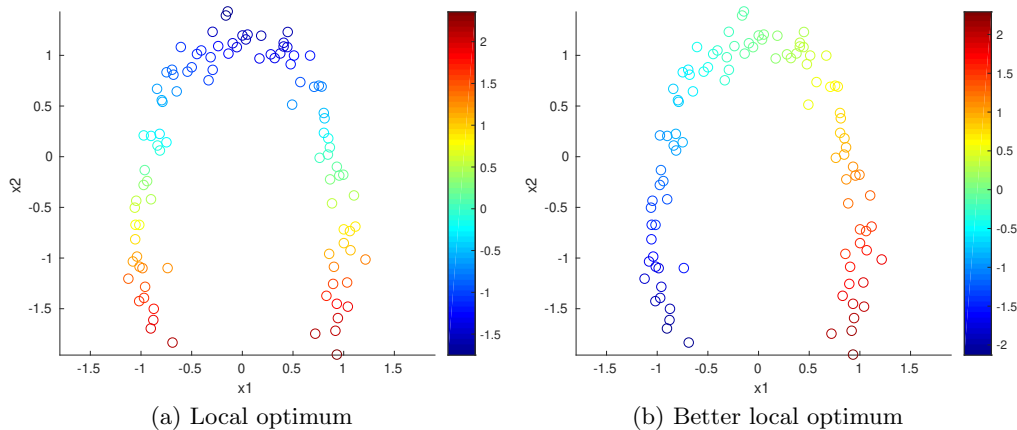


Figure 3.6: Dimensionality reduction by the Sammon nonlinear mapping. The method is prone to local optima. The solution in (b) has a smaller cost than (a). The colours indicate the value of the one-dimensional coordinate assigned to a data point.

two data points $\mathbf{x}_i \in \mathbb{R}^d$ and $\mathbf{x}_j \in \mathbb{R}^d$ is measured by the shortest distance between them when only allowed to travel on the data manifold from one neighbouring data point to the next. This is called a geodesic distance. The neighbourhood of a data point can be taken to be the m -nearest neighbours or also all points that are within a certain (Euclidean) distance. The set of neighbourhoods defines a graph on which one is allowed to move. For further information on Isomap, see the original algorithm or the books by Izenman (2008, Section 16.6) and Lee and Verleysen (2007, Section 4.3).

Figure 3.8 shows the graphs for the circularly shaped data in Figure 3.2. We see that for a neighbourhood that is specified by 5 nearest neighbours, the graph has two unconnected components. In this case, Isomap is often applied to each component separately.

Figure 3.9 visualises the one-dimensional coordinates z_1, \dots, z_n that are obtained by applying classical MDS on the geodesic distances. They well represent the circular structure when the learned graph is connected.

3.5 UMAP

Uniform Manifold Approximation and Projection (UMAP) is a dimensionality reduction method by McInnes, Healy, and Melville (2018) that has become popular due to its good performance on large datasets.

Like Isomap, UMAP is a nearest neighbour based graph dimensionality reduction method. For Isomap, a nearest neighbour graph is used to transform the original data to a dissimilarity matrix Δ . For UMAP, separate (weighted) graphs are calculated for the original data $\tilde{\mathbf{X}}$ and for the low dimensional data representation \mathbf{M} . The representation \mathbf{M} is then varied to minimise a graph distance to the original data.

For given points $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$, the graph is constructed in the following

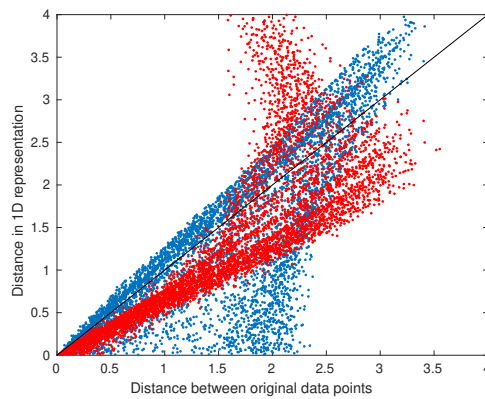


Figure 3.7: Dimensionality reduction by the Sammon nonlinear mapping. Comparison of the distances in the original and lower dimensional space. The blue points correspond to the solution in Figure 3.6(a); the red points to the solution in Figure 3.6(b)

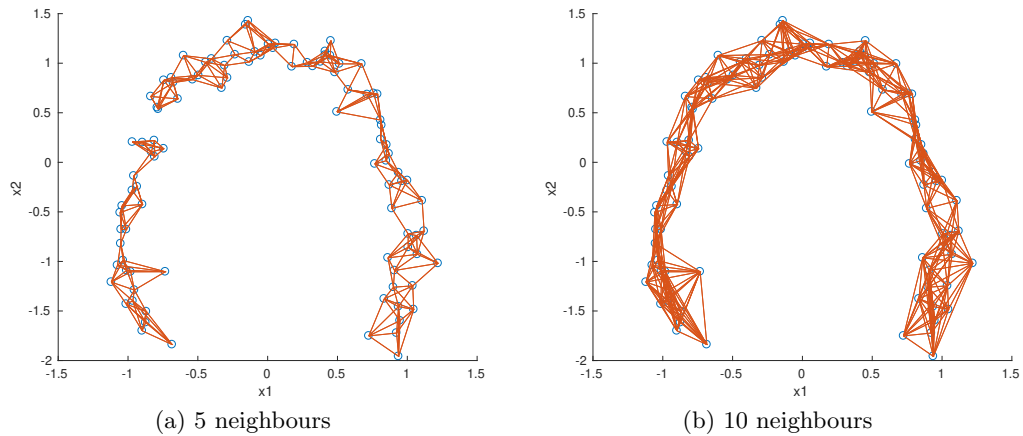


Figure 3.8: Dimensionality reduction by Isomap. Comparison of graphs constructed from different neighbourhoods.

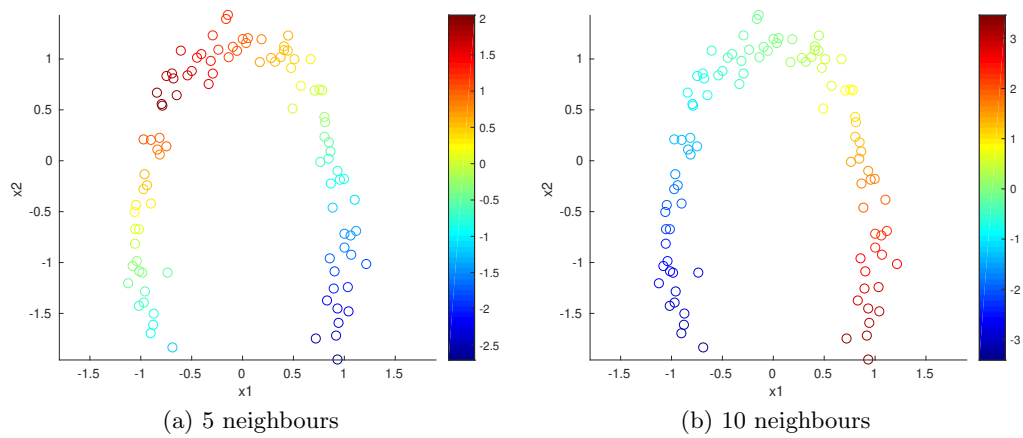


Figure 3.9: Dimensionality reduction with Isomap. The colours indicate the value of the one-dimensional coordinate assigned to each data point.

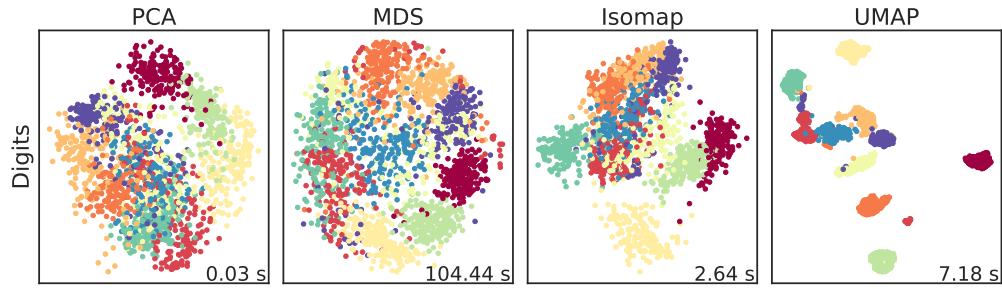


Figure 3.10: Comparison of dimensionality reduction with PCA, Metric MDS, Isomap and UMAP on a dataset of handwritten digits from 0 to 9. Colours indicate the digit classes. Adapted from <https://umap-learn.readthedocs.io/>.

way. First, generate the m nearest neighbour graph for \mathbf{Y} . Then, for each \mathbf{y}_i , compute weights for each outgoing edge:

$$w_i^{(\mathbf{Y})}(\mathbf{y}_i, \mathbf{y}_j) = \exp\left(-\frac{\|\mathbf{y}_i - \mathbf{y}_j\| - \rho_i}{\sigma_i}\right), \quad (3.37)$$

where ρ_i denotes the distance to the nearest neighbour and σ_i is a measure of the size of the neighbourhood around \mathbf{y}_i (the “diameter”, see the original paper on how it is computed). The metrics used to compute the distances may depend on the nature of the data; their choice is a hyperparameter of the method. The weight can be interpreted as the probability that the edge exists. Since ρ_i is the distance to the nearest neighbour of \mathbf{y}_i , the weight to the nearest neighbour is one, and this ensures that each point connects to at least one other point.

Note that we obtain two weights per edge which can be asymmetric, i.e. $w_i^{(\mathbf{Y})}(\mathbf{y}_i, \mathbf{y}_j) \neq w_j^{(\mathbf{Y})}(\mathbf{y}_j, \mathbf{y}_i)$. To symmetrise, UMAP sets

$$w^{(\mathbf{Y})}(\mathbf{y}_i, \mathbf{y}_j) = w_i^{(\mathbf{Y})}(\mathbf{y}_i, \mathbf{y}_j) + w_j^{(\mathbf{Y})}(\mathbf{y}_j, \mathbf{y}_i) - w_i^{(\mathbf{Y})}(\mathbf{y}_i, \mathbf{y}_j)w_j^{(\mathbf{Y})}(\mathbf{y}_j, \mathbf{y}_i) \quad (3.38)$$

which can be interpreted as the probability that at least one of the edges exists. We thus obtain a weight matrix $\mathbf{W}^{(\mathbf{Y})}$ with $(\mathbf{W}^{(\mathbf{Y})})_{ij} = w^{(\mathbf{Y})}(\mathbf{y}_i, \mathbf{y}_j)$ representing the graph in a probabilistic manner.

Using the above approach, we can calculate $\mathbf{W}^{(\tilde{\mathbf{X}})}$ for the observed data $\tilde{\mathbf{X}}$ and $\mathbf{W}^{(\mathbf{M})}$ for a candidate lower dimensional representation \mathbf{M} . To quantify the distance between $\mathbf{W}^{(\tilde{\mathbf{X}})}$ and $\mathbf{W}^{(\mathbf{M})}$, we compute the average Kullback-Leiber divergence between the weights (probabilities) $\mathbf{W}^{(\tilde{\mathbf{X}})}$ and $\mathbf{W}^{(\mathbf{M})}$

$$C(\mathbf{W}^{(\tilde{\mathbf{X}})}, \mathbf{W}^{(\mathbf{M})}) = \sum_{i < j} w_{ij}^{(\tilde{\mathbf{X}})} \log\left(\frac{w_{ij}^{(\tilde{\mathbf{X}})}}{w_{ij}^{(\mathbf{M})}}\right) + (1 - w_{ij}^{(\tilde{\mathbf{X}})}) \log\left(\frac{(1 - w_{ij}^{(\tilde{\mathbf{X}})})}{(1 - w_{ij}^{(\mathbf{M})})}\right). \quad (3.39)$$

The lower dimensional representation \mathbf{Z} is then the solution to

$$\underset{\mathbf{M}}{\text{minimise}} C(\mathbf{W}^{(\tilde{\mathbf{X}})}, \mathbf{W}^{(\mathbf{M})}). \quad (3.40)$$

Details on how the optimisation problem is solved efficiently can be found in the original publication (McInnes, Healy, and Melville, 2018, Sections 2, 3).

References

- [1] I. Borg and P.J.F. Groenen. *Modern Multidimensional Scaling: Theory and Applications*. Springer, 2005.
- [2] A.J. Izenman. *Modern Multivariate Statistical Techniques: Regression, Classification, and Manifold Learning*. Springer, 2008.
- [3] J.A. Lee and M. Verleysen. *Nonlinear Dimensionality Reduction*. Springer, 2007.
- [4] Leland McInnes, John Healy, and James Melville. “Umap: Uniform manifold approximation and projection for dimension reduction”. In: *arXiv preprint arXiv:1802.03426* (2018).
- [5] B. Schölkopf, A. Smola, and K.-R. Müller. “Kernel principal component analysis”. In: *Artificial Neural Networks ICANN’97*. Springer Berlin Heidelberg, 1997, pp. 583–588.
- [6] B. Schölkopf and A.J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization and Beyond*. MIT Press, 2002.
- [7] J. B. Tenenbaum, V. de Silva, and J. C. Langford. “A Global Geometric Framework for Nonlinear Dimensionality Reduction”. In: *Science* 290.5500 (2000), pp. 2319–2323.

Chapter 4

Predictive Modelling and Generalisation

Regression and classification are typical examples of predictive modelling. The general goal in predictive modelling of data is to identify a relationship between some predictor (input) and some target (output) variables that enables one to accurately predict the values of the target variables for some newly observed values of the predictor variables. This chapter is about evaluating the performance of prediction models and methods, and about selecting among competing alternatives.

4.1 Prediction and Training Loss

We here introduce the key notions of prediction and training loss.

4.1.1 Prediction Loss

Let us denote the predictor variables by \mathbf{x} and let us assume that we are only interested in a single target variable y . In regression, y is real-valued while in classification, y is the class label, e.g. minus one and one. Both \mathbf{x} and y are considered random variables that have a joint probability density function $p(\mathbf{x}, y)$. For any fixed value of \mathbf{x} , the target variable y thus follows the conditional distribution $p(y|\mathbf{x})$. Both the joint pdf and the conditional pdf are unknown.

From a probabilistic perspective, the goal of predictive modelling is to estimate the conditional distribution $p(y|\mathbf{x})$ from observed data. In many cases, however, we need to report a single estimated value of y rather than a whole distribution. That is, we are looking for a prediction function $h(\mathbf{x})$ that provides an estimate $\hat{y} = h(\mathbf{x})$ for any value of \mathbf{x} .

Making a prediction \hat{y} may incur a loss $\mathcal{L}(\hat{y}, y)$ so that certain prediction functions are better than others. Due to the stochasticity of the predictors and the target, the quality of a prediction function h is measured via the expected value of $\mathcal{L}(\hat{y}, y)$,

$$\mathcal{J}(h) = \mathbb{E}_{\hat{y}, y} [\mathcal{L}(\hat{y}, y)] = \mathbb{E}_{\mathbf{x}, y} [\mathcal{L}(h(\mathbf{x}), y)], \quad (4.1)$$

which is called the prediction loss. The term $\mathbb{E}_{\mathbf{x}, y}$ means expectation with respect to the joint distribution of \mathbf{x} and y , i.e. $p(\mathbf{x}, y)$.

The goal of predictive modelling can be formulated as the optimisation problem

$$\underset{h}{\text{minimise}} \mathcal{J}(h). \quad (4.2)$$

While concise, the formulation hides some fundamental issues: First, we generally cannot compute the expectation over (\mathbf{x}, y) analytically. Secondly, the loss function \mathcal{L} may not be easy to evaluate—it could, for example, be given by user ratings that indicate the quality of a prediction \hat{y} . And thirdly, minimising the prediction loss with respect to a function is generally difficult.

4.1.2 Training Loss

The objective in (4.2) can typically not be computed and the optimisation problem not be solved exactly. We make a number of approximations to obtain a computable loss function for which optimisation is, at least in principle, feasible.

If n samples (\mathbf{x}_i, y_i) are available that are each independently drawn from $p(\mathbf{x}, y)$,

$$(\mathbf{x}_i, y_i) \stackrel{iid}{\sim} p(\mathbf{x}, y), \quad (4.3)$$

the expectation in the definition of the prediction loss can be approximated by a sample average,

$$\mathcal{J}(h) \approx \frac{1}{n} \sum_{i=1}^n \mathcal{L}(h(\mathbf{x}_i), y_i). \quad (4.4)$$

The samples (\mathbf{x}_i, y_i) are called the training data $\mathcal{D}^{\text{train}}$,

$$\mathcal{D}^{\text{train}} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}. \quad (4.5)$$

In the sample-average approximation, we assumed that training data are available that come from the same distribution $p(\mathbf{x}, y)$ as the data for which we would like to perform predictions. In many cases, however, this assumption is violated and the training data come from a different distribution. This can lead to inaccurate predictions, so that care should be taken that, if possible, at least parts of the training data are representative of the conditions for which the prediction function will be ultimately used.

Instead of minimising the (approximate) prediction loss with respect to any function h , we typically search for h inside model families that are parametrised by some parameters $\boldsymbol{\theta}$, so that $h(\mathbf{x}) = h_{\boldsymbol{\lambda}}(\mathbf{x}; \boldsymbol{\theta})$, where $\boldsymbol{\lambda}$ is a vector of hyperparameters indicating the model family and some tuning parameters associated with it. The hyperparameters could for example indicate whether we use regression trees or neural networks. And when using neural networks, they could additionally indicate the number of hidden units, whereas $\boldsymbol{\theta}$ would correspond to the weights in the network.

The number of parameters $\boldsymbol{\theta}$ may be rather large so that gradient information is needed in the optimisation. But some loss functions \mathcal{L} , like for example classification error, are not differentiable so that gradient descent is not possible. Other loss functions \mathcal{L} may be expensive to evaluate, like for example when based on user ratings. For practical reasons, we may thus prefer to determine $\boldsymbol{\theta}$ by minimising a proxy loss function L rather than the loss function \mathcal{L} that we are really interested in.

In summary, instead of working with $\mathcal{J}(h)$ we work with the training loss function $J_{\lambda}(\boldsymbol{\theta})$,

$$J_{\lambda}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n L(h_{\lambda}(\mathbf{x}_i; \boldsymbol{\theta}), y_i). \quad (4.6)$$

Minimisation of $J_{\lambda}(\boldsymbol{\theta})$ is typically done by minimising the loss function with respect to $\boldsymbol{\theta}$ separately for fixed values of the hyperparameters. We then obtain a set of prediction function $\hat{h}_{\lambda}(\mathbf{x})$ indexed by λ ,

$$\hat{h}_{\lambda}(\mathbf{x}) = h_{\lambda}(\mathbf{x}; \hat{\boldsymbol{\theta}}_{\lambda}), \quad \hat{\boldsymbol{\theta}}_{\lambda} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J_{\lambda}(\boldsymbol{\theta}). \quad (4.7)$$

Determining $\hat{h}_{\lambda}(\mathbf{x})$ from training data is called model estimation. The associated minimal value of the training loss function is the training loss J_{λ}^* ,

$$J_{\lambda}^* = \min_{\boldsymbol{\theta}} J_{\lambda}(\boldsymbol{\theta}). \quad (4.8)$$

The training loss function $J_{\lambda}(\boldsymbol{\theta})$, the prediction function $\hat{h}_{\lambda}(\mathbf{x})$, and the corresponding training loss J_{λ}^* all depend on the training data $\mathcal{D}^{\text{train}}$. Different training data sets will result in different loss functions, different prediction functions, and different training losses. This means that they are all random quantities whose stochasticity is induced by the variability of the training data.

Minimising $J_{\lambda}(\boldsymbol{\theta})$ for several λ yields a set of prediction functions $\hat{h}_{\lambda}(\mathbf{x})$. Choosing from them the prediction function $\hat{h}(\mathbf{x})$ that is actually used for making predictions is done by a process called hyperparameter selection. If the hyperparameter indicates the model family, the process is called model selection. We will see below that choosing the hyperparameters that yield the smallest training loss is generally a bad idea because the corresponding prediction function tends to be highly specific to the particular training data used and thus may perform poorly when making predictions for new (unseen) values of \mathbf{x} .

4.1.3 Example

Let us illustrate the above concepts on a simple example where the joint distribution of the prediction and target variable is given by

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right), \quad (4.9)$$

$$p(y|x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - g(x))^2\right), \quad g(x) = \frac{1}{4}x + \frac{3}{4}x^2 + x^3. \quad (4.10)$$

The function $g(x)$ is the conditional mean $\mathbb{E}(y|x)$. It minimises the expected square loss, i.e. (4.1) for

$$\mathcal{L}(\hat{y}, y) = (\hat{y} - y)^2. \quad (4.11)$$

We assume that we have a training data set $\mathcal{D}^{\text{train}}$ with n data points (x_i, y_i) . Figure 4.1(a) shows $g(x)$ and an example training set.

Training Loss for Linear Regression

Let us first work with a linear prediction model so that

$$h_1(x; \boldsymbol{\theta}) = \theta_0 + \theta_1 x, \quad \boldsymbol{\theta} = (\theta_0, \theta_1)^\top, \quad (4.12)$$

where θ_0 is the intercept and θ_1 the slope parameter. When using the quadratic loss, the training loss function is

$$J_1(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2. \quad (4.13)$$

For any value of θ_1 , the optimal value of the constant θ_0 is

$$\hat{\theta}_0 = \bar{y} - \theta_1 \bar{x}, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad \bar{x} = \sum_{i=1}^n \frac{1}{n} x_i, \quad (4.14)$$

so that θ_1 is the only unknown when working with centred data. The training loss function becomes

$$J_1(\theta_1) = \frac{1}{n} \sum_{i=1}^n ((y_i - \bar{y}) - \theta_1(x_i - \bar{x}))^2. \quad (4.15)$$

Minimising $J_1(\theta_1)$ yields $\hat{\theta}_1$,

$$\hat{\theta}_1 = \underset{\theta_1}{\operatorname{argmin}} J_1(\theta_1), \quad (4.16)$$

and the estimated prediction model $\hat{h}_1(x)$ thus equals

$$\hat{h}_1(x) = \hat{\theta}_0 + \hat{\theta}_1 x = \bar{y} + \hat{\theta}_1(x - \bar{x}). \quad (4.17)$$

Figure 4.2(a) shows the training loss function $J_1(\theta_1)$ and the estimated regression function $\hat{h}_1(x)$.

The training loss function $J_1(\boldsymbol{\theta})$ in (4.13) varies as the training data vary. The training loss function $J_1(\boldsymbol{\theta})$ is a random quantity. Its minimiser inherits the randomness, and the minimal training loss

$$J_1^* = \min_{\boldsymbol{w}} J_1(\boldsymbol{w}) \quad (4.18)$$

is a random variable too. The randomness is due to the variability of the training data $\mathcal{D}^{\text{train}}$. Random quantities have a probability distribution, and Figure 4.3 visualises the distribution of the training loss function and the distribution of its minima J_1^* . The probability density function of J_1^* was estimated from relative frequencies.

Training loss for Polynomial Regression

Instead of working with a linear prediction model, let us now consider more general prediction models of the form

$$h_\lambda(x; \boldsymbol{\theta}) = \sum_{k=0}^{\lambda} \theta_k x^k, \quad \boldsymbol{\theta} = (\theta_0, \dots, \theta_\lambda)^\top. \quad (4.19)$$

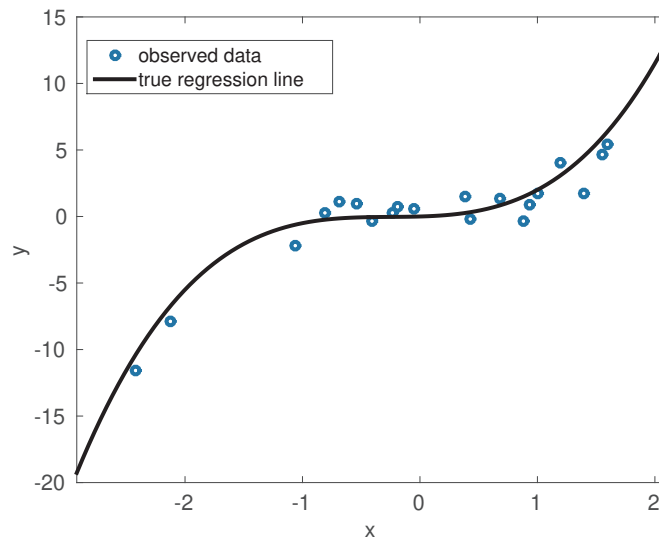


Figure 4.1: Example nonlinear regression (prediction) problem. The true regression curve is shown in black and the example training data in blue. The size of the training data is $n = 20$.

The functions $h_\lambda(x; \boldsymbol{\theta})$ are polynomials of degree λ . The $h_\lambda(x; \boldsymbol{\theta})$ correspond to a set of prediction models: We have one prediction model for each value of λ . Its complexity and number of free parameters increases with increasing values of λ . For $\lambda = 0$, the prediction model is a constant, and for $\lambda = 1$ we obtain the linear model used above.

We can estimate the prediction models by minimising the average square loss as before,

$$J_\lambda(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left(\sum_{k=0}^{\lambda} \theta_k x_i^k - y_i \right)^2, \quad (4.20)$$

Minimising $J_\lambda(\boldsymbol{\theta})$ yields the prediction functions \hat{h}_λ with training loss J_λ^* .

Figure 4.4(a) shows the estimated probability density function (pdf) of the training loss J_λ^* . The pdf for the polynomial of degree one is the same as in Figure 4.3(b). We see that the training loss tends to become smaller if the complexity of the model increases.

Another view is provided in Figure 4.4(b). The figure shows the training loss J_λ^* as a function of the degree of the polynomials. We see that both the variance and the mean of the training loss becomes smaller with increasing complexity of the model. The same holds for more general models than the polynomial one used here. Indeed, it is generally possible to increase the complexity of the model to the point where the minimal training loss becomes zero (see Section 4.2.2).

4.2 Generalisation Performance

The training loss can be made smaller by using more complex models. But we are ultimately interested in the prediction rather than in the training loss. In other

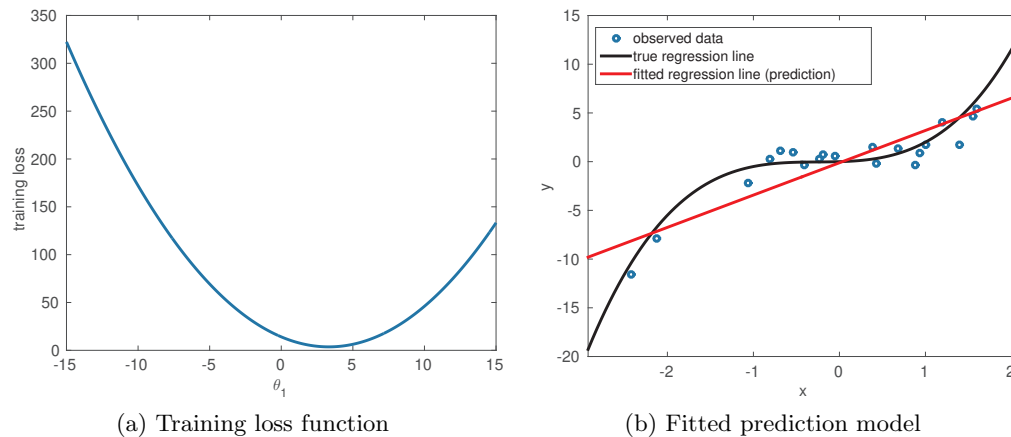


Figure 4.2: The true regression curve is shown in black and the estimated regression curve in red.

words, we are interested in how well we perform on unseen data after learning. This is called generalisation performance.

4.2.1 Generalisation for Prediction Functions and Algorithms

The training loss function in (4.6) was used as a proxy of the prediction loss $\mathcal{J}(h)$ in (4.1) that we are really interested in minimising. We say that a prediction function \hat{h} generalises well if its prediction loss $\mathcal{J}(\hat{h})$ is small,

$$\mathcal{J}(\hat{h}) = \mathbb{E}_{\mathbf{x},y} [\mathcal{L}(\hat{h}(\mathbf{x}), y)]. \quad (4.21)$$

The prediction loss $\mathcal{J}(\hat{h})$ is called the generalisation loss or the test loss of \hat{h} . This is because $\mathcal{J}(\hat{h})$ measures whether the performance of \hat{h} generalises from the training data $\mathcal{D}^{\text{train}}$ and training loss function L to new “test” data $(\mathbf{x}, y) \sim p(\mathbf{x}, y)$ and the prediction loss function \mathcal{L} . As argued above, $\mathcal{J}(\hat{h})$ can generally not be computed. But unlike before, we here do not need to solve an optimisation problem. We only need to evaluate \mathcal{J} at \hat{h} , which is considerably easier. It amounts to estimating the expected value of $\mathcal{L}(\hat{h}(\mathbf{x}), y)$, which can be done with hold-out data (see Section 4.3.1).

Since the prediction function \hat{h} depends on the training data $\mathcal{D}^{\text{train}}$, the prediction loss $\mathcal{J}(\hat{h})$ depends on the training data too. The prediction loss $\mathcal{J}(\hat{h})$ is thus a random variable whose stochasticity is induced by the variability of the training sets. We will now see that its expected value $\bar{\mathcal{J}}$ can be used to measure the generalisation performance of prediction algorithms.

Let us denote the prediction algorithm that is used to turn training data $\mathcal{D}^{\text{train}}$ into a prediction function \hat{h} by \mathcal{A} so that

$$\hat{h} = \mathcal{A}(\mathcal{D}^{\text{train}}). \quad (4.22)$$

The algorithm \mathcal{A} subsumes all operations needed to turn training data into a prediction function, including for example the minimisation of the loss function or the selection of hyperparameters. Think of it as a piece of code that takes

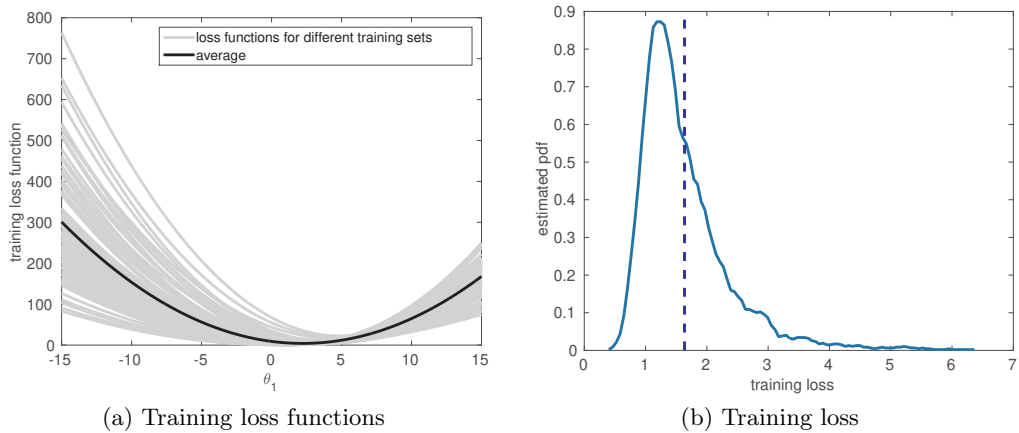


Figure 4.3: The loss functions and their minima are random quantities. The figures illustrate their distribution. (a) Loss functions for different training sets. (b) Distribution of the square root of the training loss J_1^* for different training sets. The dashed vertical line indicates the mean of the estimated distribution.

training data as input and returns a prediction function \hat{h} as output. We can then write the expected prediction loss $\bar{\mathcal{J}}$ as a function of \mathcal{A}

$$\bar{\mathcal{J}}(\mathcal{A}) = \mathbb{E}_{\mathcal{D}^{\text{train}}} [\mathcal{J}(\hat{h})] = \mathbb{E}_{\mathcal{D}^{\text{train}}} [\mathcal{J}(\mathcal{A}(\mathcal{D}^{\text{train}}))]. \quad (4.23)$$

While $\mathcal{J}(\hat{h})$ in (4.21) measures the performance of a specific \hat{h} , $\bar{\mathcal{J}}(\mathcal{A})$ measures the performance of the process, or algorithm, that is used to obtain \hat{h} from the training data. The purpose of the two performance measures in (4.21) and (4.23) is thus different: $\mathcal{J}(\hat{h})$ can be used to compare different prediction functions while $\bar{\mathcal{J}}(\mathcal{A})$ can be used to compare different prediction algorithms.

If we consider algorithms \mathcal{A}_λ that operate with different (fixed) hyperparameters λ , we can use $\bar{\mathcal{J}}(\mathcal{A}_\lambda)$ to compare and select among them. Like $\mathcal{J}(\hat{h})$, however, the expected prediction loss $\bar{\mathcal{J}}(\mathcal{A})$ can typically not be computed in closed form and needs to be estimated, for which cross-validation can be used (see Section 4.3.1).

4.2.2 Overfitting and Underfitting

Let us consider the training and (expected) prediction loss of the prediction functions $\hat{h}_\lambda(\mathbf{x})$ in (4.7) for different models. By using a model with n free parameters $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)^\top$, we can make the training loss always equal to zero. Indeed, if

$$h_{\text{flexible}}(\mathbf{x}; \boldsymbol{\theta}) = \begin{cases} \theta_i & \text{if } \mathbf{x} = \mathbf{x}_i \\ 0 & \text{otherwise} \end{cases} \quad (4.24)$$

we can set $\hat{\theta}_i = y_i$ and the training loss is zero (assuming that $L(y_i, y_i) = 0$). Unless \mathbf{x} and y can only take discrete values that are all included in the training data, the (expected) prediction loss of $\hat{h}_{\text{flexible}}$ will be large. The prediction function is overfitting the training data. More generally, a model has been overfitted

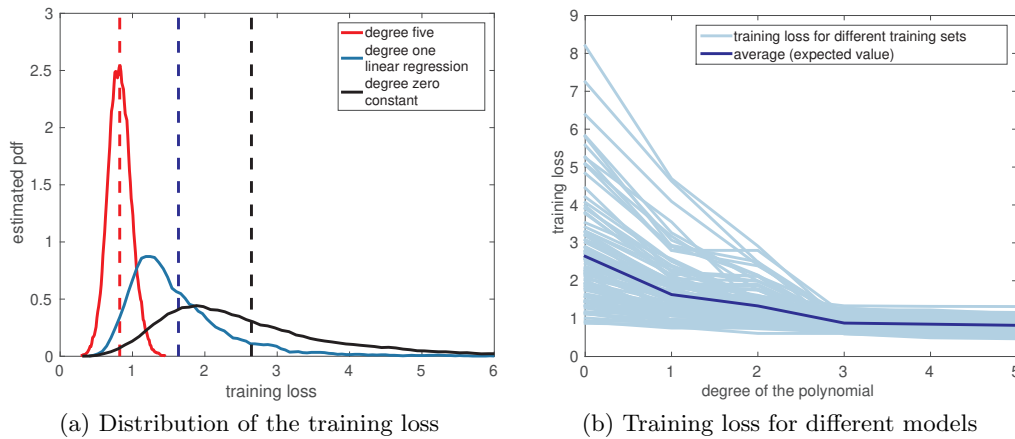


Figure 4.4: Distribution of the training loss for different degrees of the polynomial prediction model. The complexity of the model increases with the degree of the polynomial.

to the training data if reducing its complexity reduces the (expected) prediction loss.

On the other hand, a prediction model

$$h_{\text{rigid}}(\mathbf{x}; \theta) = \theta, \quad (4.25)$$

which always takes on a constant value, will have a training loss that is rather large. Unless the response variable y does indeed not depend on the predictors \mathbf{x} , the (expected) prediction loss will be large, too, and could be decreased by choosing a more flexible model that better captures the relationship between \mathbf{x} and y . Prediction functions like $h_{\text{rigid}}(\mathbf{x}; \theta)$ are said to underfit the training data.

The problem of over- and underfitting can be addressed by model selection and by means of regularisation. In regularisation, we work with flexible models but augment the training loss function $J_{\lambda}(\theta)$, which measures the quality of the prediction, with an additional term that penalises flexibility of the prediction function. For training, we thus solve the optimisation problem

$$\underset{\theta}{\text{minimise}} J_{\lambda}(\theta) + \lambda_{\text{reg}} R(\theta), \quad (4.26)$$

where $R(\theta)$ is the penalty term on the parameters of $h_{\lambda}(\mathbf{x}; \theta)$ and λ_{reg} indicates the strength of the regularisation. Typical penalty terms are

$$R(\theta) = \sum_i \theta_i^2 \quad (L_2 \text{ or Tikhonov regularisation}) \quad (4.27)$$

$$R(\theta) = \sum_i |\theta_i| \quad (L_1 \text{ regularisation}) \quad (4.28)$$

but also terms that penalises rapidly varying functions. The amount of regularisation depends on λ_{reg} . We can consider it to be another hyperparameter that we can select in order to maximise generalisation performance.

4.2.3 Example

We continue the example of polynomial regression to illustrate how the generalisation performance depends on the model complexity and the size of the training data.

Generalisation Performance and Model Complexity

Figure 4.5(a) shows the training and prediction loss of the fitted polynomial regression model \hat{h}_λ as a function of the degree of the polynomial (model complexity) λ . We can see that the prediction loss and training loss are generally not the same, i.e.

$$\mathcal{J}(\hat{h}_\lambda) \neq J_\lambda^*. \quad (4.29)$$

In the figure, the prediction loss is smallest for $\lambda = 4$, and while a five-degree polynomial has the smallest training loss, it has the largest prediction loss. Such a mismatch between training and prediction performance is due to overfitting. The estimated model \hat{h}_λ is highly tuned to the specific training data $\mathcal{D}^{\text{train}}$ and does not reflect the general relationship between the predictor and the target variable. In contrast, we see that increasing the complexity of the degree-zero or degree-one polynomial will decrease the prediction loss. That is, these models are underfitting the training data.

While Figure 4.5(a) depicts the training and prediction loss for a particular training set, Figure 4.5(b) shows their distribution over different training data sets. We can see that the variability of the prediction loss increases with the flexibility of the model. This is due to overfitting because the estimated model then depends strongly on the particularities of each training set that are bound to vary when the training data change. Underfitting, in contrast, leads to a small variability of the prediction loss because the fitted model captures comparably few properties of the training data.

The red solid line in Figure 4.5(b) shows the expected (average) prediction loss $\bar{\mathcal{J}}$ in (4.23) as a function of λ . While a model of degree $\lambda = 4$ performed best for the particular training data used in (a), models of degree $\lambda = 3$ yield the best performance on average. We see that there is here a difference between the generalisation performance of a specific fitted model and the generalisation performance of a model-family across different training sets, which reflects the general difference between $\mathcal{J}(\hat{h}_\lambda)$ and $\bar{\mathcal{J}}(\mathcal{A}_\lambda)$ discussed in Section 4.2.1.

Generalisation Performance and the Size of the Training Data

The results so far were obtained for training sets of size $n = 20$. We saw that flexible models tended to overfit the training data, so that there was stark difference between training and prediction performance. Here, we illustrate how the size of the training data influences the generalisation performance.

Figure 4.6 shows the expected training and prediction loss as a function of the size n of the training data for polynomial models of different degree. We can generally see that the training and prediction loss approach each other as the sample size increases. Note that they may generally not reach the same limit as n increases because the training and prediction loss functions L and \mathcal{L} , for example, may not be the same.

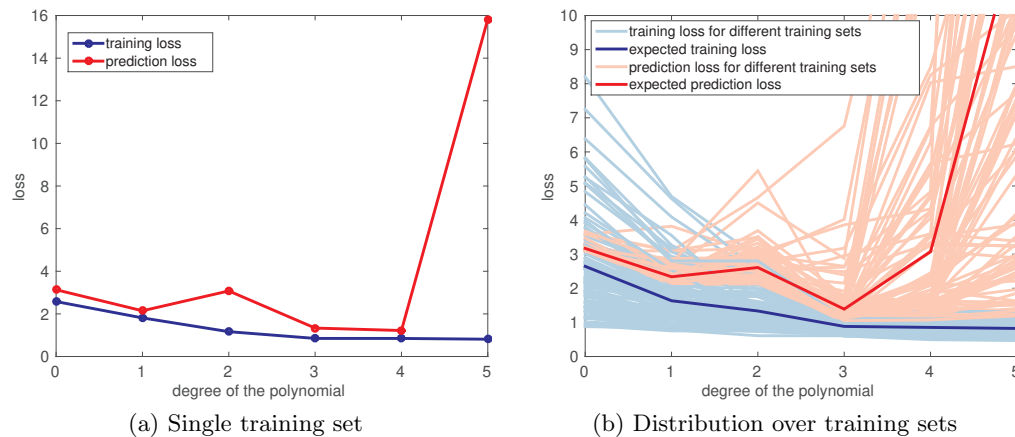


Figure 4.5: Training versus prediction performance of different prediction models.

Figure 4.6(a) shows that increasing the model complexity decreases the prediction loss for the models of degree zero and one. Moreover, their prediction loss does not decrease below a certain level even if the size of the training data increases. Both phenomena are a sign of underfitting.

Figure 4.6(b) shows the average training and prediction loss for the polynomial model of degree five. In this example, the large difference between training and prediction loss for small sample sizes is due to overfitting. As the size of the training data increases, however, the gap between the two losses becomes smaller, which means that the amount of overfitting decreases. Generally, the gap should not be used as an indicator for the amount of overfitting, though. The gap can also arise due to insufficient training data for fitting the right model. This would not be overfitting because reducing model complexity would not reduce prediction loss.

Comparing Figure 4.6(a) and (b) shows us further that even for large samples, on average, the model of degree five does here not achieve a smaller prediction loss than the model of degree three. Hence, for this problem, there is no advantage in using a more complex model than the model of degree three. In general, we can use model selection to choose among candidate models, or regularisation to avoid overfitting flexible models on small training data. Both model selection and choosing the right amount of regularisation correspond to hyperparameter selection.

4.3 Estimating the Generalisation Performance

We typically need to estimate the generalisation performance twice: Once for hyperparameter selection, and once for final performance evaluation. We first discuss two methods for estimating the generalisation performance and then apply them to the two aforementioned tasks.

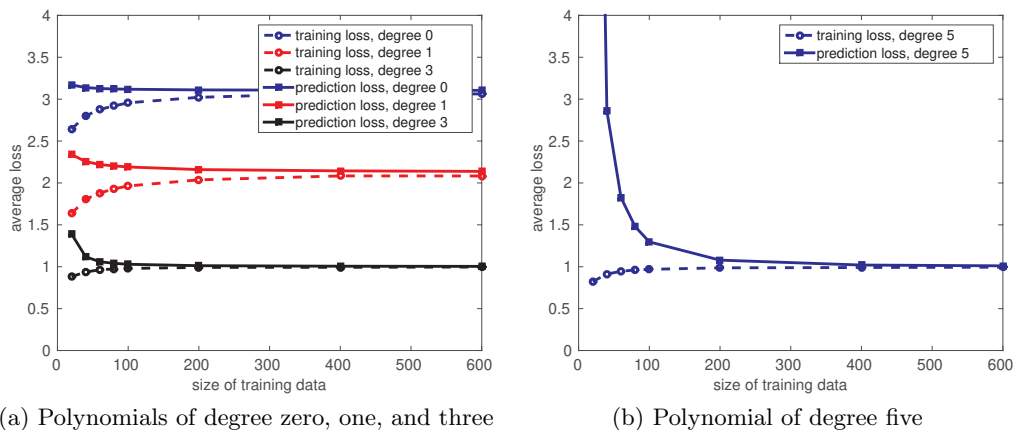


Figure 4.6: Average training versus average prediction performance for different sizes of the training data.

4.3.1 Methods for Estimating the Generalisation Performance

The hold-out and the cross-validation approach to estimate the generalisation performance are presented.

Hold-out Approach

Assume that the prediction function \hat{h} has been obtained using training data $\mathcal{D}^{\text{train}}$, i.e.

$$\hat{h} = \mathcal{A}(\mathcal{D}^{\text{train}}). \quad (4.30)$$

If another data set $\tilde{\mathcal{D}}$ is available with \tilde{n} samples $(\tilde{\mathbf{x}}_i, \tilde{y}_i) \sim p(\mathbf{x}, y)$ that are statistically independent from the samples in $\mathcal{D}^{\text{train}}$, we can use $\tilde{\mathcal{D}}$ to estimate the prediction loss $\mathcal{J}(\hat{h})$ via a sample average

$$\hat{\mathcal{J}}(\hat{h}; \tilde{\mathcal{D}}) = \frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} \mathcal{L}(\hat{h}(\tilde{\mathbf{x}}_i), \tilde{y}_i). \quad (4.31)$$

Depending on the context, $\tilde{\mathcal{D}}$ is called a test or a validation set.

We are typically given the union of the two data sets $\mathcal{D}^{\text{train}}$ and $\tilde{\mathcal{D}}$, and it is up to us how to split them into the two sets. Common split ratios are $n/\tilde{n} = 60/40$, $70/30$, or $80/20$. If the number of (hyper) parameters is large, it is better to increase the ratio so that more data are available for training.

While the splitting is often done randomly, particularly in classification, it is important that the different values of the target variable (e.g. the class labels) represented in a balanced way in both $\mathcal{D}^{\text{train}}$ and $\tilde{\mathcal{D}}$. Stratification methods can be used so that e.g. the classes are present in the same proportions in both $\mathcal{D}^{\text{train}}$ and $\tilde{\mathcal{D}}$.

The value of the estimated prediction loss in (4.31) may vary strongly for different hold-out data sets $\tilde{\mathcal{D}}$ unless \tilde{n} is large. This is often seen as a drawback of the hold-out approach. Figure 4.7 illustrates the variability that can be introduced by randomly splitting a data set into a training set $\mathcal{D}^{\text{train}}$ and test set $\tilde{\mathcal{D}}$. Cross-validation is often used to avoid such issues.

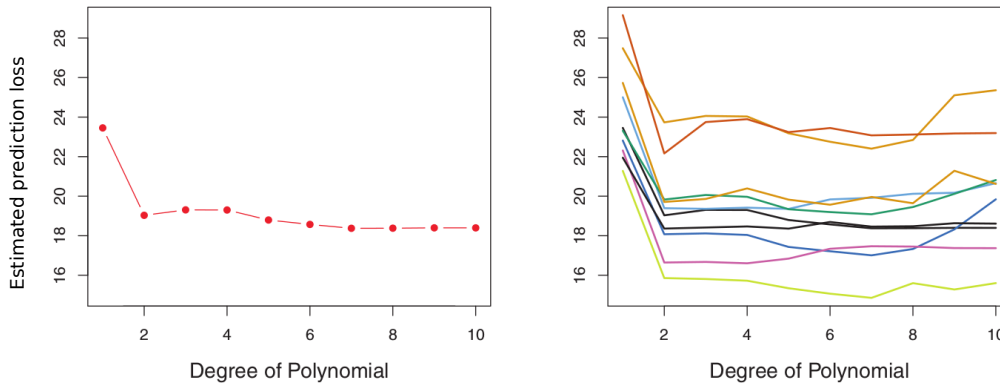


Figure 4.7: Possible variability in the estimated prediction loss. (a) The estimated prediction loss for a classification problem with a polynomial prediction model. (b) Variability induced by random splitting of the available data (392 data points) into training set and test set (here: of equal size). Each curve shows a different realisation of the random variable $\hat{\mathcal{J}}(\hat{h}; \tilde{\mathcal{D}})$. Adapted from (James, Witten, and Hastie, 2016, Figure 5.2).

Cross-validation

Cross-validation consists in randomly dividing the data that are available for training into K (roughly) equally-sized subsets (folds) $\mathcal{D}_1, \dots, \mathcal{D}_K$ without overlap. For the same reasons as in the hold-out approach, we may want to use here stratification. From the folds, we construct K pairs of training data sets $\mathcal{D}_k^{\text{train}}$ and hold-out (validation) sets $\mathcal{D}_k^{\text{val}}$,

$$\mathcal{D}_k^{\text{train}} = \bigcup_{i \neq k} \mathcal{D}_i, \quad \mathcal{D}_k^{\text{val}} = \mathcal{D}_k, \quad (4.32)$$

as illustrated in Figure 4.8. The K training sets are used to obtain K prediction functions \hat{h}_k ,

$$\hat{h}_k = \mathcal{A}(\mathcal{D}_k^{\text{train}}), \quad (4.33)$$

whose performance $\hat{\mathcal{J}}_k$ is evaluated on the data $\mathcal{D}_k^{\text{val}}$ that was held-out during training,

$$\hat{\mathcal{J}}_k = \hat{\mathcal{J}}(\hat{h}_k; \mathcal{D}_k^{\text{val}}). \quad (4.34)$$

The performance $\hat{\mathcal{J}}(\hat{h}_k; \mathcal{D}_k^{\text{val}})$ is computed via (4.31). We are essentially repeating the hold-out approach K times, each time with different data. The cross-validation (cv) score CV is then the average of all $\hat{\mathcal{J}}_k$,

$$\text{CV} = \frac{1}{K} \sum_{i=1}^K \hat{\mathcal{J}}_k. \quad (4.35)$$

The cv score is sometimes used as an improved version of $\hat{\mathcal{J}}$ in (4.31). But it is actually rather an estimate of the expected prediction performance $\bar{\mathcal{J}}(\mathcal{A})$ in (4.23). The cv score does indeed not depend on a prediction function but on the prediction algorithm \mathcal{A} . This can be more clearly seen when writing

$$\hat{\mathcal{J}}_k = \hat{\mathcal{J}}(\hat{h}_k; \mathcal{D}_k^{\text{val}}) = \hat{\mathcal{J}}\left(\mathcal{A}(\mathcal{D}_k^{\text{train}}); \mathcal{D}_k^{\text{val}}\right) \quad (4.36)$$


 Figure 4.8: Sketch of K -fold cross-validation for $K = 5$.

so that

$$\text{CV} = \frac{1}{K} \sum_{i=1}^K \hat{\mathcal{J}}_i = \frac{1}{K} \sum_{i=1}^K \hat{\mathcal{J}} \left(\mathcal{A}(\mathcal{D}_k^{\text{train}}); \mathcal{D}_k^{\text{val}} \right), \quad (4.37)$$

which does depend on the algorithm \mathcal{A} but not on a particular prediction function \hat{h} . The cv score is thus an estimate of $\bar{\mathcal{J}}(\mathcal{A})$ that we denote by $\hat{\bar{\mathcal{J}}}(\mathcal{A})$,

$$\hat{\bar{\mathcal{J}}}(\mathcal{A}) = \text{CV}. \quad (4.38)$$

The cross-validation score CV, and hence the estimate of $\bar{\mathcal{J}}$, depends on the particular assignment of the data points into the K folds, so that the score is a random variable. One can assess its distribution by performing cross-validation several times but this tends to be a computationally expensive procedure.

Alternatively, if K is not too large, e.g. $K = 5$, one can assess the variability of the cv-score by estimating its variance as

$$\text{Var}(\text{CV}) \approx \frac{1}{K} \text{Var}(\hat{\mathcal{J}}), \quad \text{Var}(\hat{\mathcal{J}}) \approx \frac{1}{K} \sum_{k=1}^K (\hat{\mathcal{J}}_k - \text{CV})^2. \quad (4.39)$$

We have here approximations because the formulae assume statistical independence of the $\hat{\mathcal{J}}_k$, which is not the case as they were all computed from the same data. The square root of $\text{Var}(\text{CV})$ is called the standard error of the cv-score.

The value of K is a tuning parameter. A typical choice is $K = 5$, so that the training sets $\mathcal{D}_k^{\text{train}}$ consist of 4/5 of all data points available and the validation sets $\mathcal{D}_k^{\text{val}}$ of 1/5 of them. If the validation sets consist of one data point only, the method is called leave-one-out cross-validation (LOOCV). While generally very expensive, for some problems, the computation can be done quickly. For a further discussion of the choice of K , see e.g. Section 7.10 in the textbook by Hastie, Tibshirani, and Friedman (2009).

4.3.2 Hyperparameter Selection and Performance Evaluation

We consider a scenario where we have several prediction models $h_{\lambda}(\mathbf{x}; \boldsymbol{\theta})$ that we can possibly use for solving our prediction task, and that we need to select among them. An algorithm that depends on the hyperparameters $\boldsymbol{\lambda}$ will be denoted by \mathcal{A}_{λ} . Two approaches to hyperparameter selection and performance evaluation of the final prediction function are presented: The first uses hold-out data to select the hyperparameters and hold-out data for performance evaluation while the second uses cross-validation for hyperparameter selection and hold-out data for performance evaluation.

Two Times Hold-out

This approach to hyperparameter selection and performance evaluation proceeds as follows:

1. From all the data \mathcal{D} that are available to us, we split off some test data $\mathcal{D}^{\text{test}}$ to estimate the performance of our final prediction function \hat{h} . The test data will never be touched until the final performance evaluation. A typical size of the test data is 20% of \mathcal{D} .
2. We split the remaining data into a training set $\mathcal{D}^{\text{train}}$ and a validation set \mathcal{D}^{val} , using, for example, again the 80/20 ratio ($\mathcal{D}^{\text{train}}$ contains 80% of the data that remain after the initial splitting while \mathcal{D}^{val} contains 20% of them).

3. Running an algorithm with hyperparameters $\boldsymbol{\lambda}$ on $\mathcal{D}^{\text{train}}$ returns a set of functions

$$\hat{h}_{\lambda} = \mathcal{A}_{\lambda}(\mathcal{D}^{\text{train}}) \quad (4.40)$$

indexed by the hyperparameters $\boldsymbol{\lambda}$.

4. We evaluate the performance of \hat{h}_{λ} on \mathcal{D}^{val} by computing the estimated prediction loss $\text{PL}(\boldsymbol{\lambda})$

$$\text{PL}(\boldsymbol{\lambda}) = \hat{\mathcal{J}}(\hat{h}_{\lambda}; \mathcal{D}^{\text{val}}), \quad (4.41)$$

where $\hat{\mathcal{J}}$ is defined in (4.31). We choose $\boldsymbol{\lambda}$ by minimising $\text{PL}(\boldsymbol{\lambda})$,

$$\hat{\boldsymbol{\lambda}} = \underset{\boldsymbol{\lambda}}{\text{argmin}} \text{PL}(\boldsymbol{\lambda}). \quad (4.42)$$

5. Using $\hat{\boldsymbol{\lambda}}$, we re-estimate the parameters $\boldsymbol{\theta}$ on the union of the training and validation data $\mathcal{D}^{\text{train}} \cup \mathcal{D}^{\text{val}}$. By using more data, we can estimate the prediction model more accurately. Denote the resulting prediction function by \hat{h} ,

$$\hat{h} = \mathcal{A}_{\hat{\boldsymbol{\lambda}}}(\mathcal{D}^{\text{train}} \cup \mathcal{D}^{\text{val}}). \quad (4.43)$$

6. We take the test data $\mathcal{D}^{\text{test}}$ out of the vault to compute an estimate $\hat{\mathcal{J}}$ of the prediction loss of \hat{h} ,

$$\hat{\mathcal{J}} = \hat{\mathcal{J}}(\hat{h}; \mathcal{D}^{\text{test}}), \quad (4.44)$$

using (4.31).

7. We re-estimate \hat{h} using all data available,

$$\hat{h}(\mathbf{x}) = \mathcal{A}_{\hat{\lambda}}(\mathcal{D}), \quad (4.45)$$

which provides us with the final prediction function \hat{h} . An estimate of its generalisation performance is given by $\hat{\mathcal{J}}$ in (4.44).

In some cases the re-estimation needs to be skipped because of computational reasons. Optimisation over the hyperparameters λ is typically not possible by gradient descent. Grid search can be used if the number of hyperparameters is small. Alternative methods are random search where different values of the hyperparameters are randomly tried out (Bergstra and Bengio, 2012), or Bayesian optimisation where the functional relationship between the hyperparameters and the prediction loss is modelled via (Gaussian process) regression, which is used to guide the optimisation (e.g. Snoek, Larochelle, and Adams, 2012).

Cross-validation and Hold-out

In this approach, we choose the hyperparameters by cross-validation and estimate the prediction performance by a hold-out test set. In more detail, we proceed as follows:

1. As above, from all the data \mathcal{D} that are available to us, we split off some test data $\mathcal{D}^{\text{test}}$ to estimate the performance our final prediction function \hat{h} . The test data will never be touched until the final performance evaluation. A typical size of the test data is 20% of \mathcal{D} .
2. We use the remaining data, call it $\mathcal{D}^{\text{train}}$, to compute the cv-score CV as a function of the hyperparameters. The cv-score is an estimate $\hat{\mathcal{J}}$ of the expected prediction loss $\bar{\mathcal{J}}$, see (4.38). Let us denote it by $\text{EPL}(\lambda)$,

$$\text{EPL}(\lambda) = \text{CV} = \hat{\mathcal{J}}(\mathcal{A}_{\lambda}). \quad (4.46)$$

3. We choose $\hat{\lambda}$ by minimising $\text{EPL}(\lambda)$. Since the cv-score is an estimate with standard-deviation $\sqrt{\text{Var}(\text{CV})}$, an alternative method is to choose the hyperparameters so that they result in the simplest model while still having a cv-score that is within one standard deviation of the minimal cv-score.
4. Using $\hat{\lambda}$, we re-estimate the parameters θ from $\mathcal{D}^{\text{train}}$. Denote the resulting prediction function by \hat{h} ,

$$\hat{h} = \mathcal{A}_{\hat{\lambda}}(\mathcal{D}^{\text{train}}). \quad (4.47)$$

5. We take the test data $\mathcal{D}^{\text{test}}$ out of the vault to compute an estimate $\hat{\mathcal{J}}$ of the prediction loss of \hat{h} ,

$$\hat{\mathcal{J}} = \hat{\mathcal{J}}(\hat{h}; \mathcal{D}^{\text{test}}), \quad (4.48)$$

using (4.31).

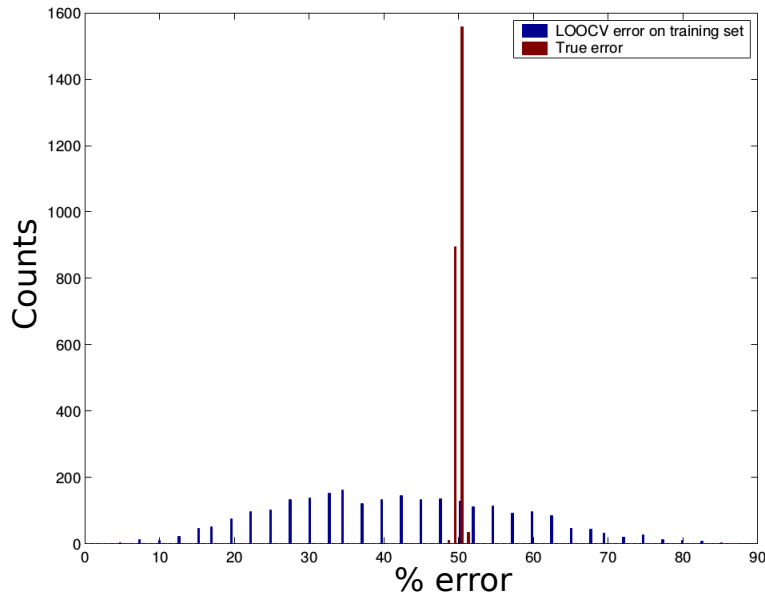


Figure 4.9: Distribution of the minimal cv-score (blue) and true prediction losses (prediction errors, red) for several artificially generated classification problems where the true prediction error is 0.5. Sample size was 40 and classification was done by support-vector machines. The figure is from Varma and Simon (2006), Figure 2.

6. We re-estimate \hat{h} using all data available,

$$\hat{h} = \mathcal{A}_{\hat{\lambda}}(\mathcal{D}), \quad (4.49)$$

which provides us with the final prediction function \hat{h} . An estimate of its generalisation performance is given by $\hat{\mathcal{J}}$ in (4.48).

In some cases the re-estimation needs to be skipped because of computational reasons. Minimisation of the cv-score can typically not be done by gradient descent. As before, gradient-free minimisation methods such as grid search, random search, or Bayesian optimisation can be used.

Like a training loss, the minimal cv-score is typically an optimistic estimate of the prediction loss because the hyperparameters are chosen such that the cv-score is minimised. The prediction loss tends to be underestimated as illustrated in Figure 4.9. That is why we need the hold-out test data $\mathcal{D}^{\text{test}}$ to determine the generalisation performance.

4.4 Loss Functions in Predictive Modelling

This section provides a brief overview of loss functions that are widely used in regression and classification.

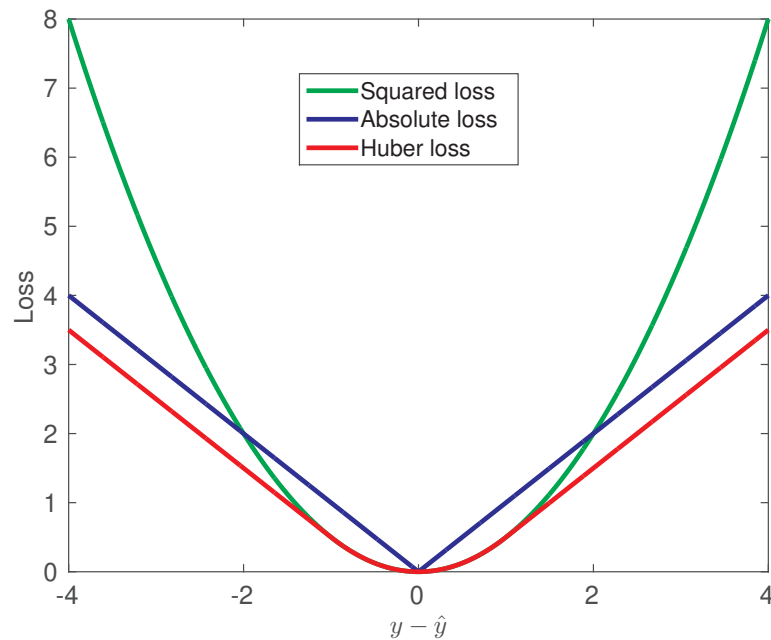


Figure 4.10: Loss functions that are often used in regression.

4.4.1 Loss Functions in Regression

Typical loss functions L in regression are

$$L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2 \quad (\text{square loss}) \quad (4.50)$$

$$L(\hat{y}, y) = |\hat{y} - y| \quad (\text{absolute loss}) \quad (4.51)$$

$$L(\hat{y}, y) = \begin{cases} \frac{1}{2}(\hat{y} - y)^2 & \text{if } |\hat{y} - y| < \delta \\ \delta|\hat{y} - y| - \frac{1}{2}\delta^2 & \text{otherwise} \end{cases} \quad (\text{Huber loss}) \quad (4.52)$$

Figure 4.10 shows plots of the different loss functions. The absolute loss is more robust than the square loss since it does not grow as quickly, but it is not differentiable when the residual $\hat{y} - y$ is zero. The Huber loss combines the good properties of the square and the absolute loss.

4.4.2 Loss Functions in Classification

We distinguish between loss functions that are differentiable with respect to parameters of the classifier and those that are not.

Non-differentiable Loss Functions

We assume here that y and \hat{y} can take K different values, for instance $\{1, \dots, K\}$. This corresponds to classification with K different classes. The loss function

$L(\hat{y}, y)$ can then be represented as a $K \times K$ matrix \mathbf{L} ,

$$\mathbf{L} = \begin{pmatrix} L(1,1) & L(1,2) & \cdots & L(1,K) \\ L(2,1) & L(2,2) & \cdots & L(2,K) \\ \vdots & \vdots & & \vdots \\ L(K,1) & L(K,2) & \cdots & L(K,K) \end{pmatrix}. \quad (4.53)$$

The diagonal elements $L(i, i)$ are zero as they correspond to correct predictions. The off-diagonal elements $L(i, j)$ are positive; they correspond to the loss incurred when predicting i instead of j . Since \hat{y} takes on discrete values, we cannot compute derivatives with respect to parameters θ that might govern the classifier.

If $L(i, j) = 1$ for $i \neq j$ and zero otherwise, the loss is said to be the zero-one loss. Its expectation $\mathcal{J}(h)$ equals

$$\mathcal{J}(h) = \mathbb{E}_{\mathbf{x}, y} L(h(\mathbf{x}), y) \quad (4.54)$$

$$= \mathbb{E}_{\hat{y}, y} L(\hat{y}, y) \quad (4.55)$$

$$= \sum_{i,j} L(i, j) p(i, j) \quad (4.56)$$

$$= \sum_{i \neq j} p(i, j) \quad (4.57)$$

$$= \mathbb{P}(y \neq \hat{y}), \quad (4.58)$$

which is the misclassification or error rate. The term $p(i, j) = \mathbb{P}(\hat{y} = i, y = j)$ denotes the joint probability of (\hat{y}, y) . The joint probability of (\hat{y}, y) is induced by the joint probability of (\mathbf{x}, y) and the prediction function h . The $p(i, j)$ for $i \neq j$ indicate the probabilities that h wrongly predicts i if the true class is j . We generally want h to be such that these probabilities are small.

If there are only two classes, for example $\{-1, 1\}$, the random variables (\hat{y}, y) can take four possible values and the predictions are typically called “true positive”, “false negative”, “false positive”, or “true negative”, see Table 4.1. The possible conditional probabilities are:

$$\text{true-positive rate of } h: \quad \mathbb{P}(\hat{y} = 1|y = 1) \quad (4.59)$$

$$\text{true-negative rate of } h: \quad \mathbb{P}(\hat{y} = -1|y = -1) \quad (4.60)$$

$$\text{false-positive rate of } h \quad \mathbb{P}(\hat{y} = 1|y = -1) = 1 - \text{true-negative rate} \quad (4.61)$$

$$\text{false-negative rate of } h: \quad \mathbb{P}(\hat{y} = -1|y = 1) = 1 - \text{true-positive rate} \quad (4.62)$$

The probabilities all depend on h since $\hat{y} = h(\mathbf{x})$. The true-positive rate is also called sensitivity, hit rate, or recall. Another name for the true-negative rate is specificity. The false-positive rate is the probability that h wrongly declares a “1”. It is also called the type 1 error. The false-negative rate is the probability that h wrongly declares a “-1”. It is also called the type 2 error. While the true-positive and true-negative rates highlight the benefits of h , measuring what it gets right, the false-positive and false-negative rates highlight the costs associated with using h , measuring its errors.

The loss function $L(\hat{y}, y)$ can be defined such that $\mathcal{J}(h)$ penalises false-positive and false-negative rates. If we let

$$\mathbf{L} = \begin{pmatrix} 0 & \frac{1}{\mathbb{P}(y=1)} \\ \frac{1}{\mathbb{P}(y=-1)} & 0 \end{pmatrix} \quad (4.63)$$

\hat{y}	y	meaning	probability	shorthand notation
1	1	true positive	$\mathbb{P}(\hat{y} = 1, y = 1)$	$p(1, 1)$
1	-1	false positive	$\mathbb{P}(\hat{y} = 1, y = -1)$	$p(1, -1)$
-1	1	false negative	$\mathbb{P}(\hat{y} = -1, y = 1)$	$p(-1, 1)$
-1	-1	true negative	$\mathbb{P}(\hat{y} = -1, y = -1)$	$p(-1, -1)$

Table 4.1: Possible events and their probabilities in binary classification.

the expected loss equals the sum of the false-positive and the false-negative rate:

$$\mathcal{J}(h) = \mathbb{E}_{\mathbf{x},y} L(h(\mathbf{x}), y) \tag{4.64}$$

$$= \mathbb{E}_{\hat{y},y} L(\hat{y}, y) \tag{4.65}$$

$$= \sum_{i,j} L(i, j)p(i, j) \tag{4.66}$$

$$= \frac{p(1, -1)}{\mathbb{P}(y = -1)} + \frac{p(-1, 1)}{\mathbb{P}(y = 1)} \tag{4.67}$$

$$= \mathbb{P}(\hat{y} = 1|y = -1) + \mathbb{P}(\hat{y} = -1|y = 1) \tag{4.68}$$

Such a cost function can be advantageous over the misclassification rate if there is, for instance, an imbalance between the probabilities for $y = 1$ and $y = -1$.

Minimising either the false-positive or the false-negative rate alone is not a very meaningful strategy. For example, if we solely minimised the false-positive rate, the trivial classifier $h(\mathbf{x}) = \hat{y} = -1$ would be the optimal solution (you can't have false-positives if there are no positives). However, for this classifier the true-positive rate would be zero.

There is generally a trade-off between true-positive and false-positive rates. This trade-off can be visualised by plotting the false-positive rate, or “cost” of h versus the true-positive rate, or “benefit” of h , see Figure 4.11. Such a plot is said to visualise the classifier in the “ROC space”, where ROC stands for “receiver operating characteristic”.

For classifiers or models with a hyperparameter, the performance of the classifier in the ROC space traces out a curve as the value of the hyperparameter is changed. The curve can be used for hyperparameter selection because classifiers that are located closest to the upper-left corner have the best trade-off between true-positive and false-positive rate (maximal benefit at the smallest cost). Classifiers that are located on a line parallel to the diagonal trade a better true-positive rate against a larger false-positive rate. We may consider such classifiers to be equivalent and the choice of working with one rather than the other is problem dependent. The area under the curve in the ROC space can be used to compare two classifiers or models irrespective of the value of the hyperparameter.

Differentiable Loss Functions in Classification

For simplicity, we consider here binary classification only. Let us assume that $\hat{y} \in \{-1, 1\}$ is given by

$$\hat{y}(\mathbf{x}) = \text{sign}(h(\mathbf{x})) \tag{4.69}$$

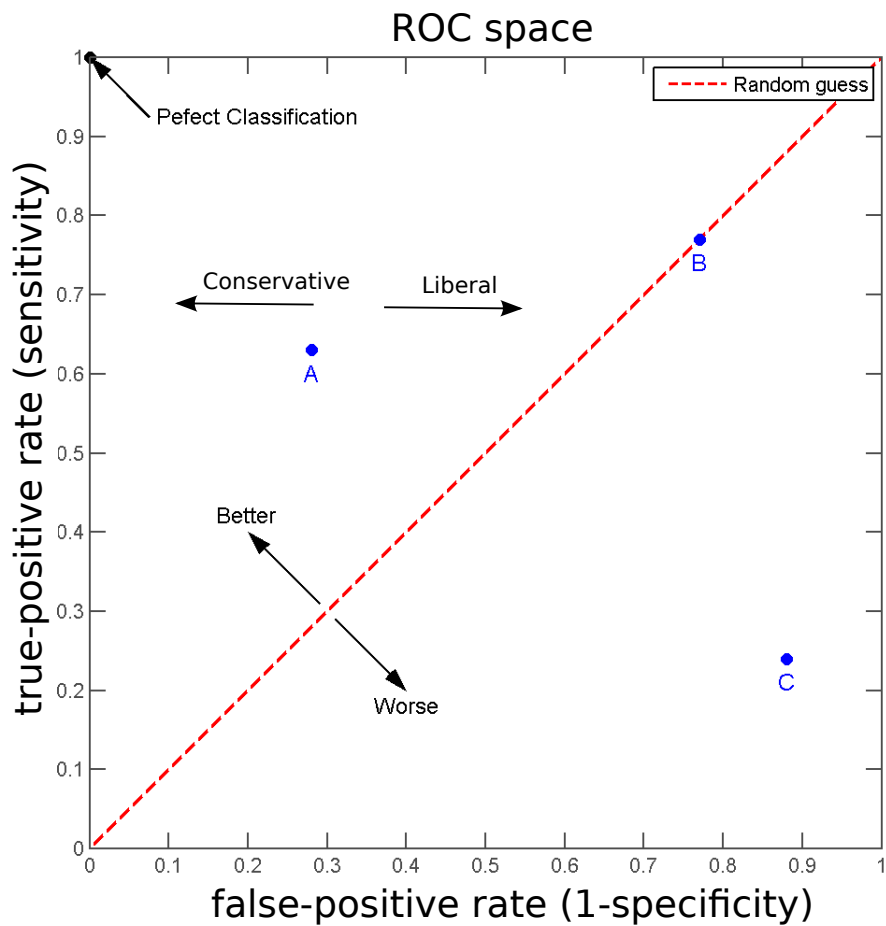


Figure 4.11: Plotting the false-positive rate (“cost”) of a classifier versus its true-positive rate (“benefit”). Classifier B is obtained by setting $\hat{y} = 1$ with probability 0.8 irrespective of the data. Classifier A takes advantage of the data and its benefit outweighs its cost while classifier C incurs a larger cost than benefit. Adapted from https://en.wikipedia.org/wiki/Receiver_operating_characteristic

where $h(\mathbf{x})$ is real-valued.

An input \mathbf{x} gets correctly classified if $h(\mathbf{x})$ takes positive values for $y = 1$ and negative values for $y = -1$. That is,

$$\text{correct classification of } \mathbf{x} \iff yh(\mathbf{x}) > 0.$$

The quantity $yh(\mathbf{x})$ is called the margin and it plays a similar role as the residual $y - h(\mathbf{x})$ in regression. The zero-one loss introduced above can be obtained by operating on the margin rather than on \hat{y} . Indeed, the zero-one loss is obtained for

$$L(h(\mathbf{x}), y) = \begin{cases} 1 & \text{if } yh(\mathbf{x}) < 0 \\ 0 & \text{otherwise.} \end{cases} \quad (4.70)$$

Several loss functions operate on the margin $yh(\mathbf{x})$. Typical ones are:

$$L(h(\mathbf{x}), y) = (h(\mathbf{x}) - y)^2 = (1 - yh(\mathbf{x}))^2 \quad (\text{square loss}) \quad (4.71)$$

$$L(h(\mathbf{x}), y) = \log(1 + \exp(-yh(\mathbf{x}))) \quad (\text{log loss}) \quad (4.72)$$

$$L(h(\mathbf{x}), y) = \exp(-yh(\mathbf{x})) \quad (\text{exponential loss}) \quad (4.73)$$

$$L(h(\mathbf{x}), y) = \max(0, 1 - yh(\mathbf{x})) \quad (\text{hinge loss}) \quad (4.74)$$

$$L(h(\mathbf{x}), y) = \max(0, 1 - yh(\mathbf{x}))^2 \quad (\text{square hinge loss}) \quad (4.75)$$

$$L(h(\mathbf{x}), y) = \begin{cases} -4yh(\mathbf{x}) & \text{if } yh(\mathbf{x}) < -1 \\ \max(0, 1 - yh(\mathbf{x}))^2 & \text{otherwise} \end{cases} \quad (\text{Huberised square hinge loss})$$

(Hastie, Tibshirani, and Friedman, 2009, Section 10.6 and Table 12.1). The different loss functions are visualised in Figure 4.12. Unlike the standard hinge loss, the square hinge loss is differentiable everywhere. The remaining loss functions are differentiable with respect to h , so that a smoothly parametrised model $h(\mathbf{x}; \boldsymbol{\theta})$ can be optimised by gradient-based optimisation methods. The different loss functions can be considered to approximate the zero-one loss. Most of them assign a loss to small positive margins, thus encouraging more confident decisions about the label. The square loss function is both sensitive to outliers and penalises large (positive) margins, which can be seen as a key disadvantage of the loss function.

Minimising the log-loss over a sample of n data points (\mathbf{x}_i, y_i) , drawn from $p(\mathbf{x}, y)$, is equivalent to maximising the log-likelihood in logistic regression. In logistic regression, we model the conditional probabilities of $y|\mathbf{x}$ as

$$\mathbb{P}(y = 1|\mathbf{x}; h) = \frac{1}{1 + \exp(-h(\mathbf{x}))} \quad \mathbb{P}(y = -1|\mathbf{x}; h) = \frac{1}{1 + \exp(h(\mathbf{x}))} \quad (4.76)$$

and estimate h by maximising the log-likelihood

$$\ell(h) = \sum_{\mathbf{x}_i: y_i=1} \log \mathbb{P}(y_i = 1|\mathbf{x}_i; h) + \sum_{\mathbf{x}_i: y_i=-1} \log \mathbb{P}(y_i = -1|\mathbf{x}_i; h) \quad (4.77)$$

$$= - \sum_{\mathbf{x}_i: y_i=1} \log(1 + \exp(-h(\mathbf{x}_i))) - \sum_{\mathbf{x}_i: y_i=-1} \log(1 + \exp(h(\mathbf{x}_i))) \quad (4.78)$$

$$= - \sum_{\mathbf{x}_i} \log(1 + \exp(-y_i h(\mathbf{x}_i))). \quad (4.79)$$

We can see that $\ell(h)$ is n times the negated sample average of the log loss.

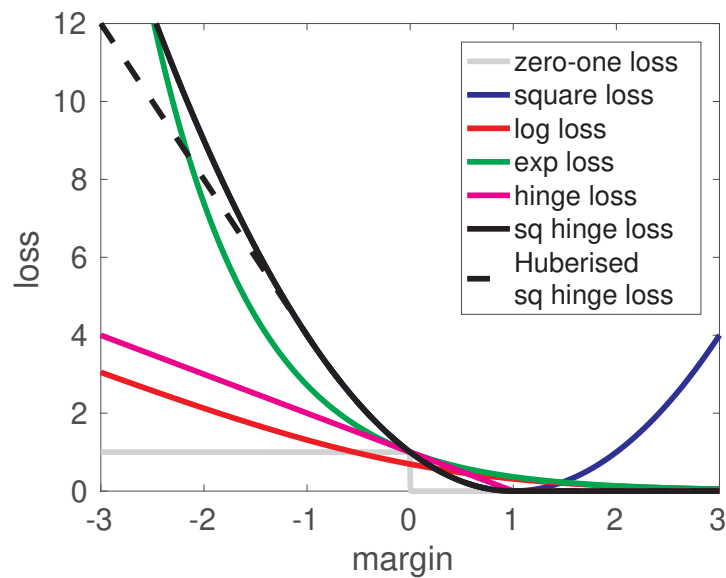


Figure 4.12: Loss functions that are often used in classification.

References

- [1] J. Bergstra and Y. Bengio. “Random Search for Hyper-Parameter Optimization”. In: *Journal of Machine Learning Research* 13 (2012), pp. 281–305.
- [2] T. Hastie, R. Tibshirani, and J.H. Friedman. *The Elements of Statistical Learning*. Springer, 2009.
- [3] G. James, D. Witten, and T. Hastie. *An Introduction to Statistical Learning: with Applications in R*. 6th ed. Springer, 2016.
- [4] J. Snoek, H. Larochelle, and R.P. Adams. “Practical Bayesian Optimization of Machine Learning Algorithms”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2012.
- [5] S. Varma and R. Simon. “Bias in error estimation when using cross-validation for model selection”. In: *BMC Bioinformatics* 7.91 (2006).

Appendix A

Linear Algebra

The material in this chapter is mostly a refresher of some basic results from linear algebra. But it also contains some proofs of results that may be harder to find. The proofs are not examinable.

A.1 Matrices

A $m \times n$ matrix \mathbf{A} is a $m \times n$ array of numbers arranged into m rows and n columns. The element at row i and column j is denoted by a_{ij} so that

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & & & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix}. \quad (\text{A.1})$$

We will sometimes use the indexing notation $(\mathbf{A})_{ij}$ to refer to element a_{ij} . The transpose \mathbf{A}^\top of a matrix \mathbf{A} is the matrix where the entries of \mathbf{A} are mirrored at the diagonal, i.e. $(\mathbf{A}^\top)_{ij} = (\mathbf{A})_{ji}$. If $\mathbf{A}^\top = \mathbf{A}$, the matrix is said to be symmetric.

Multiplying a matrix with a scalar produces a matrix where each element is scaled by said scalar, for example

$$\alpha \mathbf{A} = \begin{pmatrix} \alpha a_{11} & \alpha a_{12} & \dots & \alpha a_{1n} \\ \vdots & & & \vdots \\ \alpha a_{m1} & \alpha a_{m2} & \dots & \alpha a_{mn} \end{pmatrix} \quad (\text{A.2})$$

Two matrices of the same size can be added together by adding their corresponding elements, for example

$$\mathbf{A} + \mathbf{B} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & & & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} + \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ \vdots & & & \vdots \\ b_{m1} & b_{m2} & \dots & b_{mn} \end{pmatrix} \quad (\text{A.3})$$

$$= \begin{pmatrix} a_{11} + b_{11} & a_{12} + b_{12} & \dots & a_{1n} + b_{1n} \\ \vdots & & & \vdots \\ a_{m1} + b_{m1} & a_{m2} + b_{m2} & \dots & a_{mn} + b_{mn} \end{pmatrix} \quad (\text{A.4})$$

If matrix \mathbf{A} has size $m \times n$ and matrix \mathbf{B} size $n \times p$, the two matrices can be multiplied together. The result is a $m \times p$ matrix $\mathbf{C} = \mathbf{AB}$ whose elements $(\mathbf{C})_{ij} = c_{ij}$ are given by

$$(\mathbf{C})_{ij} = \sum_{k=1}^n (\mathbf{A})_{ik} (\mathbf{B})_{kj}, \quad \text{or, equivalently,} \quad c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}. \quad (\text{A.5})$$

The equations mean that to compute the (ij) -th element of \mathbf{C} , we multiply the elements of the i -th row of \mathbf{A} with the elements of the j -th column of \mathbf{B} and sum them all up.

The trace of a $m \times m$ matrix \mathbf{A} is the sum of its diagonal elements,

$$\text{trace}(\mathbf{A}) = \sum_{i=1}^m a_{ii}. \quad (\text{A.6})$$

The trace of \mathbf{AB} equals the trace of \mathbf{BA} : Let \mathbf{A} be $m \times n$ and \mathbf{B} $n \times m$. We then have

$$\text{trace}(\mathbf{AB}) = \sum_{i=1}^m (\mathbf{AB})_{ii} = \sum_{i=1}^m \left(\sum_{j=1}^n a_{ij} b_{ji} \right) = \sum_{j=1}^n \sum_{i=1}^m b_{ji} a_{ij}, \quad (\text{A.7})$$

which equals $\sum_{j=1}^n (\mathbf{BA})_{jj}$ and hence

$$\text{trace}(\mathbf{AB}) = \text{trace}(\mathbf{BA}) \quad (\text{A.8})$$

as claimed.

A.2 Vectors

A n -dimensional vector \mathbf{v} can be seen as $n \times 1$ matrix. We denote its i -th element by v_i or sometimes also by $(\mathbf{v})_i$. By default, \mathbf{v} is a column vector, i.e.

$$\mathbf{v} = \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix}. \quad (\text{A.9})$$

Its transpose \mathbf{v}^\top is the row vector (v_1, \dots, v_n) . Like matrices, vectors can be scaled, added or multiplied together. The product between a $1 \times n$ (row) vector \mathbf{r} and a $n \times 1$ (column) vector \mathbf{v} is with (A.5) a number equal to

$$\mathbf{r}\mathbf{v} = \sum_{i=1}^n r_i v_i. \quad (\text{A.10})$$

The inner product or scalar product $\mathbf{u}^\top \mathbf{v}$ between two n dimensional vectors \mathbf{u} and \mathbf{v} is

$$\mathbf{u}^\top \mathbf{v} = \sum_{i=1}^n u_i v_i, \quad (\text{A.11})$$

that is, the vector \mathbf{u} is first transposed to be row vector after which (A.5) is applied. Importantly, it does not matter whether \mathbf{u} or \mathbf{v} is transposed, i.e.

$$\mathbf{u}^\top \mathbf{v} = \mathbf{v}^\top \mathbf{u}. \quad (\text{A.12})$$

The outer product $\mathbf{u}\mathbf{v}^\top$ between a m dimensional vector \mathbf{u} and a n dimensional vector \mathbf{v} is a $m \times n$ matrix

$$\mathbf{u}\mathbf{v}^\top = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_m \end{pmatrix} (v_1 \ v_2 \ \dots \ v_n) = \begin{pmatrix} u_1v_1 & u_1v_2 & \dots & u_1v_n \\ u_2v_1 & u_2v_2 & \dots & u_2v_n \\ \vdots & \vdots & & \vdots \\ u_mv_1 & u_mv_2 & \dots & u_mv_n \end{pmatrix}. \quad (\text{A.13})$$

It can be seen that the (i, j) -th element of the matrix is equal to $u_i v_j$ in line with (A.5).

Equation (A.5) also tells us that the product between a $m \times n$ matrix \mathbf{A} and n -dimensional vector \mathbf{u} equals a m -dimensional vector \mathbf{v} with elements v_i ,

$$v_i = \sum_{j=1}^n a_{ij} u_j \quad i = 1, \dots, m. \quad (\text{A.14})$$

A.3 Matrix Operations as Operations on Column Vectors

It is often helpful to consider a $m \times n$ matrix \mathbf{A} as a collection of n column vectors \mathbf{a}_j of dimension m that are arranged next to each other,

$$\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n). \quad (\text{A.15})$$

Note that the i -th element of the j -th column of \mathbf{A} is $(\mathbf{A})_{ij} = (\mathbf{a}_j)_i$.

A.3.1 Matrix-vector Products

By computing the i -th element, we see that $\mathbf{v} = \mathbf{A}\mathbf{u}$ can be written as weighted combination of the column vectors \mathbf{a}_j ,

$$\mathbf{A}\mathbf{u} = \sum_{j=1}^n \mathbf{a}_j u_j = \underbrace{\begin{pmatrix} a_{11} \\ \vdots \\ a_{m1} \end{pmatrix}}_{\mathbf{a}_1} u_1 + \dots + \underbrace{\begin{pmatrix} a_{1j} \\ \vdots \\ a_{mj} \end{pmatrix}}_{\mathbf{a}_j} u_j + \dots + \underbrace{\begin{pmatrix} a_{1n} \\ \vdots \\ a_{mn} \end{pmatrix}}_{\mathbf{a}_n} u_n, \quad (\text{A.16})$$

The equation shows that for vectors \mathbf{u} that are zero everywhere but in slot k , $\mathbf{A}\mathbf{u} = \mathbf{a}_k u_k$, which means that we can “pick” column k of \mathbf{A} by multiplying the matrix with the k unit vector.

From (A.5), we know that

$$(\mathbf{A}\mathbf{B}^\top)_{ij} = \sum_{k=1}^n (\mathbf{A})_{ik}(\mathbf{B}^\top)_{kj} = \sum_{k=1}^n (\mathbf{A})_{ik}(\mathbf{B})_{jk} \quad (\text{A.25})$$

We now show that $\mathbf{A}\mathbf{B}^\top$ can also be written as sum of outer products between the column vectors of \mathbf{A} and \mathbf{B} ,

$$\mathbf{A}\mathbf{B}^\top = \sum_{k=1}^n \mathbf{a}_k \mathbf{b}_k^\top. \quad (\text{A.26})$$

This identity can be verified by computing the (i, j) -th element of the matrix on the right-hand-side:

$$\left(\sum_{k=1}^n \mathbf{a}_k \mathbf{b}_k^\top \right)_{i,j} = \sum_{k=1}^n (\mathbf{a}_k \mathbf{b}_k^\top)_{i,j} \quad (\text{A.27})$$

$$= \sum_{k=1}^n (\mathbf{a}_k)_i (\mathbf{b}_k)_j. \quad (\text{A.28})$$

Since $(\mathbf{a}_k)_i$ is the i -th element of the k -th column of \mathbf{A} , we have $(\mathbf{a}_k)_i = (\mathbf{A})_{ik}$. For the same reason, $(\mathbf{b}_k)_j = (\mathbf{B})_{jk}$, so that $(\mathbf{a}_k)_i (\mathbf{b}_k)_j = (\mathbf{A})_{ik}(\mathbf{B})_{jk}$ and

$$\left(\sum_{k=1}^n \mathbf{a}_k \mathbf{b}_k^\top \right)_{i,j} = \sum_{k=1}^n (\mathbf{A})_{ik}(\mathbf{B})_{jk}, \quad (\text{A.29})$$

which equals (A.25) and thus proves the identity in (A.26).

A.4 Orthogonal Basis

Two vectors $\mathbf{u}_1 \in \mathbb{R}^n$ and $\mathbf{u}_2 \in \mathbb{R}^n$ are said to be orthogonal if their inner product (scalar product) $\mathbf{u}_1^\top \mathbf{u}_2$ is zero. If additionally the vectors are of unit norm, i.e.

$$\|\mathbf{u}_i\| = \sqrt{\mathbf{u}_i^\top \mathbf{u}_i} = 1, \quad i = 1, 2, \quad (\text{A.30})$$

the vectors are said to be orthonormal. A set of n orthonormal vectors $\mathbf{u}_i \in \mathbb{R}^n$ forms an orthogonal basis of \mathbb{R}^n . This means that any vector $\mathbf{x} \in \mathbb{R}^n$ can be written as a weighted combinations of the $\mathbf{u}_1, \dots, \mathbf{u}_n$,

$$\mathbf{x} = \sum_{i=1}^n c_i \mathbf{u}_i. \quad (\text{A.31})$$

The weights c_i are the coordinates of \mathbf{x} with respect to the basis. Due to the orthogonality of the \mathbf{u}_i , the coordinates c_i can be computed via an inner product between the \mathbf{u}_i and \mathbf{x} ,

$$c_i = \mathbf{u}_i^\top \mathbf{x}, \quad i = 1, \dots, n, \quad (\text{A.32})$$

We can form a matrix \mathbf{U} by putting all the orthonormal basis vectors next to each other as the columns of the matrix,

$$\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n). \quad (\text{A.33})$$

The matrix U is said to be an orthogonal matrix. Since the vectors \mathbf{u}_i have unit norm and are orthogonal to each other, we have that $U^\top U = I_n$ where I_n is the n -dimensional identity matrix.

Collecting all coordinates c_i into the vector $\mathbf{c} = (c_1, \dots, c_n)^\top$, we have with (A.32)

$$\mathbf{c} = U^\top \mathbf{x}. \quad (\text{A.34})$$

With (A.16), we can similarly write (A.31) more compactly as

$$\mathbf{x} = U\mathbf{c}. \quad (\text{A.35})$$

It follows that $\mathbf{x} = UU^\top \mathbf{x}$, from where we see that not only $U^\top U = I_n$ but also $UU^\top = I_n$ for orthogonal matrices U .

A.5 Subspaces

An orthogonal basis $\mathbf{u}_1, \dots, \mathbf{u}_n$ enables us to represent any vector $\mathbf{x} \in \mathbb{R}^n$ as a weighted combination of the vectors. If we do not have n orthonormal vectors but only k of them, e.g. $\mathbf{u}_1, \dots, \mathbf{u}_k$, we cannot represent all n -dimensional vectors but only those vectors $\mathbf{z} \in \mathbb{R}^n$ that can be written as

$$\mathbf{z} = \sum_{i=1}^k a_i \mathbf{u}_i, \quad a_i \in \mathbb{R}. \quad (\text{A.36})$$

This set of vectors is said to be spanned by the $\mathbf{u}_1, \dots, \mathbf{u}_k$ and denoted by $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$. In other words,

$$\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k) = \left\{ \mathbf{z} \in \mathbb{R}^n : \mathbf{z} = \sum_{i=1}^k a_i \mathbf{u}_i \right\}. \quad (\text{A.37})$$

If $\mathbf{z}_1 \in \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ and $\mathbf{z}_2 \in \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$, i.e. if

$$\mathbf{z}_1 = \sum_{i=1}^k a_i \mathbf{u}_i, \quad \mathbf{z}_2 = \sum_{i=1}^k b_i \mathbf{u}_i, \quad (\text{A.38})$$

their weighted sum $\alpha \mathbf{z}_1 + \beta \mathbf{z}_2$ equals

$$\alpha \mathbf{z}_1 + \beta \mathbf{z}_2 = \sum_{i=1}^k \alpha a_i \mathbf{u}_i + \sum_{i=1}^k \beta b_i \mathbf{u}_i = \sum_{i=1}^k (\alpha a_i + \beta b_i) \mathbf{u}_i \quad (\text{A.39})$$

and thus belongs to $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ as well. This means that the span is closed under addition and scalar multiplication, which makes it a subspace of \mathbb{R}^n . Since any vector \mathbf{z} of $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ can be expressed using k coordinates only, namely the $\mathbf{u}_i^\top \mathbf{z}, i = 1, \dots, k$, $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ is a k -dimensional subspace of \mathbb{R}^n .

We now show that any vector $\mathbf{x} \in \mathbb{R}^n$ can be split into a part \mathbf{x}_\parallel that belongs to $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ and a part \mathbf{x}_\perp that belongs to $\text{span}(\mathbf{u}_{k+1}, \dots, \mathbf{u}_n)$, the span of the remaining basis vectors $\mathbf{u}_{k+1}, \dots, \mathbf{u}_n$. Since

$$\mathbf{x} = \sum_{j=1}^n \mathbf{u}_j c_j = \sum_{j=1}^k \mathbf{u}_j c_j + \sum_{j=k+1}^n \mathbf{u}_j c_j \quad (\text{A.40})$$

we have that

$$\mathbf{x} = \mathbf{x}_{\parallel} + \mathbf{x}_{\perp}, \quad \mathbf{x}_{\parallel} = \sum_{j=1}^k \mathbf{u}_j c_j, \quad \mathbf{x}_{\perp} = \sum_{j=k+1}^n \mathbf{u}_j c_j. \quad (\text{A.41})$$

As \mathbf{x}_{\parallel} is a weighted sum of the $\mathbf{u}_1, \dots, \mathbf{u}_k$, and \mathbf{x}_{\perp} a weighted sum of the $\mathbf{u}_{k+1}, \dots, \mathbf{u}_n$, the vectors \mathbf{x}_{\parallel} and \mathbf{x}_{\perp} are orthogonal to each other. The subspace $\text{span}(\mathbf{u}_{k+1}, \dots, \mathbf{u}_n)$ is said to be orthogonal to $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$ and is thus also denoted by $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)^{\perp}$.

A.6 Orthogonal Projections

Let us collect the k vectors \mathbf{u}_k into the $n \times k$ matrix \mathbf{U}_k ,

$$\mathbf{U}_k = (\mathbf{u}_1, \dots, \mathbf{u}_k). \quad (\text{A.42})$$

Since the \mathbf{u}_k are orthonormal, $\mathbf{U}_k^{\top} \mathbf{U}_k = \mathbf{I}_k$, but, unlike for orthogonal matrices, $\mathbf{U}_k \mathbf{U}_k^{\top}$ is not the identity matrix. We next show that $\mathbf{U}_k \mathbf{U}_k^{\top} \mathbf{x}$ equals the part \mathbf{x}_{\parallel} of \mathbf{x} that belongs to the k -dimensional subspace $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$.

This can be most easily seen by writing $\mathbf{U}_k \mathbf{U}_k^{\top}$ as a sum of elementary matrices $\mathbf{u}_i \mathbf{u}_i^{\top}$,

$$\mathbf{U}_k \mathbf{U}_k^{\top} = \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^{\top}, \quad (\text{A.43})$$

which we can do according to (A.26). Applying $\mathbf{U}_k \mathbf{U}_k^{\top}$ on a vector \mathbf{x} thus gives

$$\mathbf{U}_k \mathbf{U}_k^{\top} \mathbf{x} \stackrel{(\text{A.43})}{=} \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^{\top} \mathbf{x} \quad (\text{A.44})$$

$$\stackrel{(\text{A.31})}{=} \sum_{i=1}^k \mathbf{u}_i \mathbf{u}_i^{\top} \sum_{j=1}^n \mathbf{u}_j c_j \quad (\text{A.45})$$

$$= \sum_{i=1}^k \mathbf{u}_i \sum_{j=1}^n \mathbf{u}_i^{\top} \mathbf{u}_j c_j \quad (\text{A.46})$$

$$= \sum_{i=1}^k \mathbf{u}_i c_i \quad (\text{A.47})$$

$$\stackrel{(\text{A.41})}{=} \mathbf{x}_{\parallel}, \quad (\text{A.48})$$

where we have used that $\mathbf{u}_i^{\top} \mathbf{u}_j$ equals zero unless $j = i$. The mapping of \mathbf{x} to $\mathbf{U}_k \mathbf{U}_k^{\top} \mathbf{x} = \mathbf{x}_{\parallel}$ is called the orthogonal projection of \mathbf{x} onto $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)$. It follows that $(\mathbf{I}_d - \mathbf{U}_k \mathbf{U}_k^{\top}) \mathbf{x}$ equals \mathbf{x}_{\perp} , and that the matrix $(\mathbf{I}_d - \mathbf{U}_k \mathbf{U}_k^{\top})$ is the orthogonal projection of \mathbf{x} onto $\text{span}(\mathbf{u}_1, \dots, \mathbf{u}_k)^{\perp}$.

A.7 Singular Value Decomposition

The singular value decomposition (SVD) of a $m \times n$ matrix \mathbf{A} is the factorisation of the matrix into the product $\mathbf{U} \mathbf{S} \mathbf{V}^{\top}$,

$$\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^{\top}, \quad (\text{A.49})$$

or negative. Matrix \mathbf{U} is orthogonal with orthonormal column vectors \mathbf{u}_i . As for the SVD, the vectors \mathbf{u}_i for which $(\mathbf{\Lambda})_{ii} = 0$ can actually be ignored so that

$$\mathbf{A} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top = \mathbf{U}_r \mathbf{\Lambda}_r \mathbf{U}_r^\top, \quad (\text{A.57})$$

where

$$\mathbf{U}_r = (\mathbf{u}_1, \dots, \mathbf{u}_r), \quad \mathbf{\Lambda}_r = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_r \end{pmatrix} \quad (\text{A.58})$$

The vectors \mathbf{u}_i are called the eigenvectors and the λ_i the eigenvalues. It follows from (A.57) that

$$\mathbf{A} \mathbf{u}_k = \lambda_k \mathbf{u}_k, \quad (\text{A.59})$$

i.e. the matrix \mathbf{A} only scales the vectors \mathbf{u}_i by their corresponding eigenvalue λ_i .

A.9 Positive Semi-definite and Definite Matrices

A symmetric $m \times m$ matrix is called positive semi-definite if all m eigenvalues are non-negative and positive definite if they are all positive. A positive definite matrix has full rank, $r = m$, and the eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_m$ form an orthogonal basis of \mathbb{R}^m .

If a matrix \mathbf{M} has the singular value decomposition $\mathbf{M} = \mathbf{U}_r \mathbf{S}_r \mathbf{V}_r^\top$ as in (A.54), the eigenvalue decomposition of $\mathbf{M} \mathbf{M}^\top$ is

$$\mathbf{M} \mathbf{M}^\top = \mathbf{U}_r \mathbf{S}_r \underbrace{\mathbf{V}_r^\top \mathbf{V}_r}_{\mathbf{I}_r} \mathbf{S}_r \mathbf{U}_r^\top = \mathbf{U}_r \mathbf{S}_r^2 \mathbf{U}_r^\top, \quad (\text{A.60})$$

on the other hand, the eigenvalue decomposition of $\mathbf{M}^\top \mathbf{M}$ is

$$\mathbf{M}^\top \mathbf{M} = \mathbf{V}_r \mathbf{S}_r \underbrace{\mathbf{U}_r^\top \mathbf{U}_r}_{\mathbf{I}_r} \mathbf{S}_r \mathbf{V}_r^\top = \mathbf{V}_r \mathbf{S}_r^2 \mathbf{V}_r^\top, \quad (\text{A.61})$$

where in both cases \mathbf{S}_r^2 refers to the diagonal matrix with elements s_i^2 . Both $\mathbf{M}^\top \mathbf{M}$ and $\mathbf{M} \mathbf{M}^\top$ have the s_i^2 as eigenvalues. We see that the eigenvalues are non-negative so that $\mathbf{M}^\top \mathbf{M}$ and $\mathbf{M} \mathbf{M}^\top$ are positive semi-definite matrices.

A.10 Matrix Approximations

A.10.1 Low Rank Approximation of General Matrices

The singular value decomposition allows us to decompose a $m \times n$ matrix \mathbf{A} of rank r as

$$\mathbf{A} = \sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^\top = \mathbf{U}_r \mathbf{S}_r \mathbf{V}_r^\top, \quad (\text{A.62})$$

see (A.54). The r singular values $s_i > 0$ are decreasing. Intuitively, the “later” rank-one matrices $\mathbf{u}_i \mathbf{v}_i^\top$ with smaller singular values contribute less to \mathbf{A} than

the “earlier” rank-one matrices with larger singular values. In fact the best approximation $\hat{\mathbf{A}}$ of the matrix \mathbf{A} by a matrix $\tilde{\mathbf{A}}$ of rank $k < r$ is given by the first k terms of the expansion above,

$$\hat{\mathbf{A}} = \sum_{i=1}^k s_i \mathbf{u}_i \mathbf{v}_i^\top. \quad (\text{A.63})$$

This result is unique if and only if $s_k > s_{k+1}$. The result is obtained when the quality of the approximation is measured by the Frobenius norm

$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F = \sum_{ij} ((\mathbf{A})_{ij} - (\tilde{\mathbf{A}})_{ij})^2 \quad (\text{A.64})$$

but also for other matrix norms (e.g. the spectral norm). For the Frobenius norm, the error when approximating \mathbf{A} with $\hat{\mathbf{A}}$ is the sum of the squares of the remaining singular values

$$\sum_{ij} ((\mathbf{A})_{ij} - (\hat{\mathbf{A}})_{ij})^2 = \sum_{i=k+1}^r s_i^2. \quad (\text{A.65})$$

This result is known as the Eckart–Young–Mirsky theorem and a proof can be found in e.g. (Gentle, 2007, Section 3.10) or (Björck, 2015, Theorem 2.2.11).

A.10.2 Low rank Approximation of Positive Semi-definite Matrices

For positive semi-definite matrices, the above approximation based on the singular value decomposition carries over: The best approximation $\hat{\mathbf{A}}$ of a positive semi-definite matrix \mathbf{A} of rank r by a matrix $\tilde{\mathbf{A}}$ of rank $k < r$ is

$$\hat{\mathbf{A}} = \sum_{i=1}^k \lambda_i \mathbf{u}_i \mathbf{u}_i^\top. \quad (\text{A.66})$$

The smallest approximation error for the Frobenius norm is

$$\|\mathbf{A} - \hat{\mathbf{A}}\|_F = \sum_{ij=1}^m ((\mathbf{A})_{ij} - (\hat{\mathbf{A}})_{ij})^2 = \sum_{i=k+1}^r \lambda_i^2, \quad (\text{A.67})$$

so that $\|\mathbf{A} - \tilde{\mathbf{A}}\|_F \geq \sum_{i=k+1}^r \lambda_i^2$ for other candidates $\tilde{\mathbf{A}}$.

A.10.3 Approximating Symmetric Matrices by Positive Semi-definite Matrices

A rank r symmetric matrix \mathbf{A} that is not positive definite has the eigenvalue decomposition

$$\mathbf{A} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top, \quad (\text{A.68})$$

where some λ_i are negative. Let us assume that there are $p \geq 1$ positive eigenvalues and that $\lambda_1 \geq \dots \geq \lambda_p > 0 > \lambda_{p+1} \geq \dots \geq \lambda_r$. We would like to determine

the positive semi-definite matrix closest to \mathbf{A} . Measuring closeness by the Frobenius norm, a result by Higham (1988) shows that the closest matrix $\hat{\mathbf{A}}$ is obtained by retaining the terms with positive eigenvalues only,

$$\hat{\mathbf{A}} = \sum_{i=1}^p \lambda_i \mathbf{u}_i \mathbf{u}_i^\top = \sum_{i=1}^r \max(\lambda_i, 0) \mathbf{u}_i \mathbf{u}_i^\top. \quad (\text{A.69})$$

The approximation error is

$$\|\mathbf{A} - \hat{\mathbf{A}}\|_F = \sum_{i=p+1}^r \lambda_i^2, \quad (\text{A.70})$$

and matrix $\hat{\mathbf{A}}$ has rank p .

Following (Higham, 1988), the proof exploits that the Frobenius norm is invariant under rotations, i.e. $\|\mathbf{A}\|_F = \|\mathbf{A}\mathbf{U}\|_F = \|\mathbf{U}\mathbf{A}\|_F$ for any orthogonal matrix \mathbf{U} . Let $\tilde{\mathbf{A}}$ be a positive semi-definite matrix. We then have

$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F = \|\mathbf{U}_r \mathbf{\Lambda}_r \mathbf{U}_r^\top - \tilde{\mathbf{A}}\|_F \quad (\text{A.71})$$

$$= \|\mathbf{U}_r^\top \mathbf{U}_r \mathbf{\Lambda}_r \mathbf{U}_r^\top \mathbf{U}_r - \mathbf{U}_r^\top \tilde{\mathbf{A}} \mathbf{U}_r\|_F \quad (\text{A.72})$$

$$= \|\mathbf{\Lambda}_r - \underbrace{\mathbf{U}_r^\top \tilde{\mathbf{A}} \mathbf{U}_r}_{\mathbf{B}}\|_F \quad (\text{A.73})$$

$$= \sum_{i=1}^r (\lambda_i - b_{ii})^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^r b_{ij}^2 \quad (\text{A.74})$$

where b_{ij} are the elements of the matrix $\mathbf{B} = \mathbf{U}_r^\top \tilde{\mathbf{A}} \mathbf{U}_r$. Because the $b_{ij}^2 \geq 0$, we have

$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F \geq \sum_{i=1}^r (\lambda_i - b_{ii})^2 \quad (\text{A.75})$$

$$= \sum_{i=1}^p \underbrace{(\lambda_i - b_{ii})^2}_{\geq 0} + \sum_{i=p+1}^r (\lambda_i - b_{ii})^2 \quad (\text{A.76})$$

$$\geq \sum_{i=p+1}^r (\lambda_i - b_{ii})^2 \quad (\text{A.77})$$

Since $b_{ii} \geq 0$ as $\tilde{\mathbf{A}}$ is restricted to be positive semi-definite and $\lambda_i < 0$ for $i > p$, we have in the equation above that $\lambda_i - b_{ii} \leq \lambda_i < 0$ and thus $(\lambda_i - b_{ii})^2 \geq \lambda_i^2$. We thus obtain the following lower bound for $\|\mathbf{A} - \tilde{\mathbf{A}}\|_F$:

$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F \geq \sum_{i=p+1}^r \lambda_i^2 \quad (\text{A.78})$$

A diagonal matrix \mathbf{B} with elements $b_i = \max(\lambda_i, 0)$ achieves the lower bound. The result in (A.69) now follows from $\tilde{\mathbf{A}} = \mathbf{U}_r \mathbf{B} \mathbf{U}_r^\top$.

A.10.4 Low Rank Approximation of Symmetric Matrices by Positive Semi-definite Matrices

As before let the symmetric matrix \mathbf{A} of rank r have p positive eigenvalues,

$$\mathbf{A} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top, \quad (\text{A.79})$$

where $\lambda_1 \geq \dots \geq \lambda_p > 0 > \lambda_{p+1} \geq \dots \geq \lambda_r$. Combining (A.69) with (A.66) we show here that the best positive semi-definite approximation of rank $k < p$ is

$$\hat{\mathbf{A}} = \sum_{i=1}^k \lambda_i \mathbf{u}_i \mathbf{u}_i^\top, \quad (\text{A.80})$$

and that the smallest approximation error is

$$\|\mathbf{A} - \hat{\mathbf{A}}\|_F = \sum_{i=k+1}^r \lambda_i^2. \quad (\text{A.81})$$

Let $\tilde{\mathbf{A}}$ be a positive semi-definite matrix of rank $k < p$. As for the proof of (A.66), we write

$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F = \|\mathbf{U}_r \mathbf{\Lambda}_r \mathbf{U}_r^\top - \tilde{\mathbf{A}}\|_F \quad (\text{A.82})$$

$$= \|\mathbf{U}_r^\top \mathbf{U}_r \mathbf{\Lambda}_r \mathbf{U}_r^\top \mathbf{U}_r - \mathbf{U}_r^\top \tilde{\mathbf{A}} \mathbf{U}_r\|_F \quad (\text{A.83})$$

$$= \|\mathbf{\Lambda}_r - \underbrace{\mathbf{U}_r^\top \tilde{\mathbf{A}} \mathbf{U}_r}_{\mathbf{B}}\|_F \quad (\text{A.84})$$

$$= \sum_{i=1}^r (\lambda_i - b_{ii})^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^r b_{ij}^2 \quad (\text{A.85})$$

where $b_{ij} = \mathbf{u}_i^\top \tilde{\mathbf{A}} \mathbf{u}_j$ are the elements of the matrix $\mathbf{B} = \mathbf{U}_r^\top \tilde{\mathbf{A}} \mathbf{U}_r$. Because the $b_{ij}^2 \geq 0$, we have

$$\sum_{\substack{i,j=1 \\ i \neq j}}^r b_{ij}^2 \geq \sum_{\substack{i,j=1 \\ i \neq j}}^p b_{ij}^2 \quad (\text{A.86})$$

and hence

$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F \geq \sum_{i=1}^r (\lambda_i - b_{ii})^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^p b_{ij}^2 \quad (\text{A.87})$$

$$= \sum_{i=1}^p (\lambda_i - b_{ii})^2 + \sum_{\substack{i,j=1 \\ i \neq j}}^p b_{ij}^2 + \sum_{i=p+1}^r (\lambda_i - b_{ii})^2 \quad (\text{A.88})$$

$$= \|\mathbf{\Lambda}_p - \mathbf{U}_p^\top \tilde{\mathbf{A}} \mathbf{U}_p\|_F + \sum_{i=p+1}^r (\lambda_i - b_{ii})^2 \quad (\text{A.89})$$

As $\tilde{\mathbf{A}}$ is restricted to be positive semi-definite $b_{ii} \geq 0$, and since $\lambda_i < 0$ for $i > p$, we have in the equation above that $\lambda_i - b_{ii} \leq \lambda_i < 0$ and thus $(\lambda_i - b_{ii})^2 \geq \lambda_i^2$. Hence:

$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F \geq \|\mathbf{\Lambda}_p - \mathbf{U}_p^\top \tilde{\mathbf{A}} \mathbf{U}_p\|_F + \sum_{i=p+1}^r \lambda_i^2 \quad (\text{A.90})$$

The matrix $\mathbf{\Lambda}_p$ is a positive definite $p \times p$ matrix, while the matrix $\mathbf{U}_p^\top \tilde{\mathbf{A}} \mathbf{U}_p$ is a $p \times p$ matrix of rank k . The smallest approximation error of a positive definite matrix by a matrix of lower rank is with (A.67) equal to $\sum_{i=k+1}^p \lambda_i^2$. We can thus bound $\|\mathbf{A} - \tilde{\mathbf{A}}\|_F$ from below by $\sum_{i=k+1}^r \lambda_i^2$,

$$\|\mathbf{A} - \tilde{\mathbf{A}}\|_F \geq \sum_{i=k+1}^r \lambda_i^2. \quad (\text{A.91})$$

The matrix $\hat{\mathbf{A}}$ in (A.80) achieves the lower bound which completes the proof.

References

- [1] Å Björck. *Numerical Methods in Matrix Computations*. Springer, 2015.
- [2] J.E. Gentle. *Matrix Algebra: Theory, Computations, and Applications in Statistics*. Springer, 2007.
- [3] N.J. Higham. “Computing a nearest symmetric positive semidefinite matrix”. In: *Linear Algebra and its Applications* 103 (1988), pp. 103–118.

Appendix B

Proofs Related to PCA

In this chapter, we present two additional proofs on equivalence of PCA formulations. These proofs are optional reading.

B.1 Sequential Maximisation Yields Simultaneous Maximisation

A proof that simultaneous and sequential variance maximisation yield the same solution is given below. As in the sequential approach, we work in the orthogonal basis provided by the eigenvectors of Σ , i.e.

$$\mathbf{w}_i = U \mathbf{a}_i, \quad (\text{B.1})$$

so that we can write the optimisation problem as

$$\begin{aligned} & \underset{\mathbf{a}_1, \dots, \mathbf{a}_k}{\text{maximise}} && \sum_{i=1}^k \mathbf{a}_i^\top \Lambda \mathbf{a}_i \\ & \text{subject to} && \|\mathbf{a}_i\| = 1 \quad i = 1, \dots, k \\ & && \mathbf{a}_i^\top \mathbf{a}_j = 0 \quad i \neq j \end{aligned} \quad (\text{B.2})$$

We see that the k vectors \mathbf{a}_i are required to be orthonormal. They can be extended by orthonormal vectors $\mathbf{a}_{k+1}, \dots, \mathbf{a}_d$ so that the matrix

$$\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_k, \mathbf{a}_{k+1}, \dots, \mathbf{a}_d) \quad (\text{B.3})$$

is orthogonal and thus satisfies $\mathbf{A}\mathbf{A}^\top = \mathbf{I}_d$. This means that the row vectors of \mathbf{A} have norm one,

$$\sum_{j=1}^d (\mathbf{A})_{ij}^2 = 1, \quad (\text{B.4})$$

and thus that

$$\sum_{j=1}^k (\mathbf{A})_{ij}^2 \leq 1. \quad (\text{B.5})$$

Below, we will denote $\sum_{j=1}^k (\mathbf{A})_{ij}^2$ by b_i . Note that $\sum_{i=1}^d b_i = k$ since the column vectors of \mathbf{A} have unit norm.

Since $\mathbf{\Lambda}$ is a diagonal matrix, the objective in (B.2) can be written as

$$\sum_{j=1}^k \mathbf{a}_j^\top \mathbf{\Lambda} \mathbf{a}_j = \sum_{j=1}^k \sum_{i=1}^d (\mathbf{a}_j)_i^2 \lambda_i = \sum_{j=1}^k \sum_{i=1}^d (\mathbf{A})_{ij}^2 \lambda_i. \quad (\text{B.6})$$

We now show that $\sum_{i=1}^k \lambda_i$ is the maximal sum that can be obtained by any set of k orthogonal vectors \mathbf{a}_i . This proves our claim about the solution of the optimisation problem in (2.27). We start with re-writing $\sum_{j=1}^k \mathbf{a}_j^\top \mathbf{\Lambda} \mathbf{a}_j$ as

$$\sum_{j=1}^k \mathbf{a}_j^\top \mathbf{\Lambda} \mathbf{a}_j = \sum_{j=1}^k \sum_{i=1}^d (\mathbf{A})_{ij}^2 \lambda_i \quad (\text{B.7})$$

$$= \sum_{i=1}^d \underbrace{\sum_{j=1}^k (\mathbf{A})_{ij}^2}_{b_i} \lambda_i \quad (\text{B.8})$$

$$= \sum_{i=1}^d b_i \lambda_i \quad (\text{B.9})$$

$$= \sum_{i=1}^k b_i \lambda_i + \sum_{i=k+1}^d b_i \lambda_i \quad (\text{B.10})$$

For $i > k$, $\lambda_i \leq \lambda_k$, as we assume that the eigenvalues are ordered from large to small. We thus obtain an upper bound for $\sum_{j=1}^k \mathbf{a}_j^\top \mathbf{\Lambda} \mathbf{a}_j$,

$$\sum_{j=1}^k \mathbf{a}_j^\top \mathbf{\Lambda} \mathbf{a}_j = \sum_{i=1}^k b_i \lambda_i + \sum_{i=k+1}^d b_i \lambda_i \quad (\text{B.11})$$

$$\leq \sum_{i=1}^k b_i \lambda_i + \lambda_k \sum_{i=k+1}^d b_i. \quad (\text{B.12})$$

We now write $\sum_{i=k+1}^d b_i = \sum_{i=1}^d b_i - \sum_{i=1}^k b_i$ and use that $\sum_{i=1}^d b_i = k$, so that

$$\sum_{i=k+1}^d b_i = k - \sum_{i=1}^k b_i \quad (\text{B.13})$$

and hence

$$\sum_{j=1}^k \mathbf{a}_j^\top \mathbf{\Lambda} \mathbf{a}_j \leq \sum_{i=1}^k b_i \lambda_i + k \lambda_k - \sum_{i=1}^k b_i \lambda_k \quad (\text{B.14})$$

$$= \sum_{i=1}^k b_i (\lambda_i - \lambda_k) + k \lambda_k. \quad (\text{B.15})$$

Since $\lambda_i - \lambda_k \geq 0$ for $i \leq k$ and $0 \leq b_i \leq 1$ we have $b_i(\lambda_i - \lambda_k) \leq (\lambda_i - \lambda_k)$ so that

$$\sum_{j=1}^k \mathbf{a}_j^\top \Lambda \mathbf{a}_j \leq \sum_{i=1}^k (\lambda_i - \lambda_k) + k\lambda_k \quad (\text{B.16})$$

$$= \sum_{i=1}^k \lambda_i - \sum_{i=1}^k \lambda_k + k\lambda_k \quad (\text{B.17})$$

$$= \sum_{i=1}^k \lambda_i - k\lambda_k + k\lambda_k, \quad (\text{B.18})$$

from where the desired result follows:

$$\sum_{j=1}^k \mathbf{a}_j^\top \Lambda \mathbf{a}_j \leq \sum_{i=1}^k \lambda_i. \quad (\text{B.19})$$

The upper bound is achieved if \mathbf{a}_j is the j -th unit vector, i.e. if $\mathbf{a}_1 = (1, 0, \dots)^\top$, $\mathbf{a}_2 = (0, 1, 0, \dots)^\top, \dots$, that is, if $\mathbf{w}_i = \mathbf{u}_j$. They are the unique solution if there are not ties in the first eigenvalues, i.e. if $\lambda_1 > \dots > \lambda_k > \lambda_{k+1}$.

B.2 Equivalence to PCA by Variance Maximisation

To prove the equivalence of (2.27) and (2.29), we first write $\sum_{i=1}^k \mathbf{w}_i \mathbf{w}_i^\top \mathbf{x}$ more compactly as $\mathbf{W}_k \mathbf{W}_k^\top \mathbf{x}$ and expand the norm of the approximation error,

$$\|\mathbf{x} - \mathbf{W}_k \mathbf{W}_k^\top \mathbf{x}\|^2 = (\mathbf{x} - \mathbf{W}_k \mathbf{W}_k^\top \mathbf{x})^\top (\mathbf{x} - \mathbf{W}_k \mathbf{W}_k^\top \mathbf{x}) \quad (\text{B.20})$$

$$= \mathbf{x}^\top \mathbf{x} - 2\mathbf{x}^\top \mathbf{W}_k \mathbf{W}_k^\top \mathbf{x} + \mathbf{x}^\top \underbrace{\mathbf{W}_k \mathbf{W}_k^\top \mathbf{W}_k \mathbf{W}_k^\top}_{\mathbf{I}_k} \mathbf{x} \quad (\text{B.21})$$

$$= \mathbf{x}^\top \mathbf{x} - \mathbf{x}^\top \mathbf{W}_k \mathbf{W}_k^\top \mathbf{x} \quad (\text{B.22})$$

Using again that $\mathbf{W}_k \mathbf{W}_k^\top = \sum_{i=1}^k \mathbf{w}_i \mathbf{w}_i^\top$, we obtain

$$\|\mathbf{x} - \mathbf{W}_k \mathbf{W}_k^\top \mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x} - \mathbf{x}^\top \left(\sum_{i=1}^k \mathbf{w}_i \mathbf{w}_i^\top \right) \mathbf{x} \quad (\text{B.23})$$

$$= \mathbf{x}^\top \mathbf{x} - \sum_{i=1}^k (\mathbf{x}^\top \mathbf{w}_i) (\mathbf{w}_i^\top \mathbf{x}) \quad (\text{B.24})$$

$$= \mathbf{x}^\top \mathbf{x} - \sum_{i=1}^k \mathbf{w}_i^\top \mathbf{x} \mathbf{x}^\top \mathbf{w}_i \quad (\text{B.25})$$

and the expected approximation error is

$$\mathbb{E} \|\mathbf{x} - \mathbf{W}_k \mathbf{W}_k^\top \mathbf{x}\|^2 = \mathbb{E}[\mathbf{x}^\top \mathbf{x}] - \mathbb{E} \left[\sum_{i=1}^k \mathbf{w}_i^\top \mathbf{x} \mathbf{x}^\top \mathbf{w}_i \right] \quad (\text{B.26})$$

$$= \mathbb{E}[\mathbf{x}^\top \mathbf{x}] - \sum_{i=1}^k \mathbf{w}_i^\top \mathbb{E}[\mathbf{x} \mathbf{x}^\top] \mathbf{w}_i \quad (\text{B.27})$$

due to the linearity of the expectation. As we assume that the expected value $\mathbb{E}[\mathbf{x}]$ is zero, due to the centring, we have $\mathbf{\Sigma} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ and

$$\mathbb{E} \|\mathbf{x} - \mathbf{W}_k \mathbf{W}_k^\top \mathbf{x}\|^2 = \mathbb{E}[\mathbf{x}^\top \mathbf{x}] - \sum_{i=1}^k \mathbf{w}_i^\top \mathbf{\Sigma} \mathbf{w}_i. \quad (\text{B.28})$$

Since $\mathbb{E}[\mathbf{x}^\top \mathbf{x}]$ is a constant, minimising the expected approximation error is equivalent to maximising $\sum_{i=1}^k \mathbf{w}_i^\top \mathbf{\Sigma} \mathbf{w}_i$, which is the total variance of the projections $\mathbf{w}_i^\top \mathbf{x}$ and the objective in (2.27). The constraints in (2.29) and (2.27) are also the same so that the two optimisation problems are equivalent.