

Supervised Machine Learning

Andreas Lindholm, Niklas Wahlström,
Fredrik Lindsten, Thomas B. Schön

Draft version: August 13, 2020

This material will be published by Cambridge University Press. This pre-publication version is free to view and download for personal use only. Not for re-distribution, re-sale or use in derivative works. © The authors, 2020.

Feedback and exercise problems: <http://smlbook.org>

Contents

1	Introduction	5
1.1	For whom is this book intended?	5
1.2	How to proceed throughout the book	6
2	Supervised learning: a first approach	7
2.1	The supervised learning problem	7
2.2	A distance-based method: k -NN	12
2.3	A rule-based method: Decision trees	17
3	Basic parametric models for regression and classification	25
3.1	Linear regression	25
3.2	Classification and logistic regression	32
3.3	Polynomial regression and regularization	38
3.4	Nonlinear regression and generalized linear models	40
3.A	Derivation of the normal equations	41
4	Understanding, evaluating and improving the performance	43
4.1	Expected new data error E_{new} : performance in production	43
4.2	Estimating E_{new}	45
4.3	The training error-generalization gap decomposition of E_{new}	49
4.4	The bias-variance decomposition of E_{new}	54
4.5	Evaluation for imbalanced and asymmetric classification problems	60
5	Learning parametric models	63
5.1	Loss functions	63
5.2	Regularization	70
5.3	Parameter optimization	73
5.4	Optimization with large datasets	82
6	Neural networks and deep learning	85
6.1	Neural networks	85
6.2	Convolutional neural networks	92
6.3	Training a neural network	95
6.4	Perspective and further reading	98
7	Ensemble methods: Bagging and boosting	101
7.1	Bagging	101
7.2	Random forests	107
7.3	Boosting and AdaBoost	109
7.4	Gradient boosting	115
8	Nonlinear input transformations and kernels	117
8.1	Creating features by nonlinear input transformations	117
8.2	Kernel ridge regression	119
8.3	Support vector regression	122
8.4	Kernel theory	125

Contents

8.5 Support vector classification	129
8.A The representer theorem	131
8.B Derivation of support vector classification	131
9 The Bayesian approach and Gaussian processes	133
9.1 The Bayesian idea	133
9.2 Bayesian linear regression	135
9.3 The Gaussian process	140
9.4 Practical usage of the Gaussian process	149
9.A The multivariate Gaussian distribution	153
10 User aspects of machine learning	155
10.1 Defining the machine learning problem	155
10.2 Improving a machine learning model	157
10.3 What if we cannot collect more data?	161
10.4 Practical data issues	164
10.5 Can I trust my machine learning model?	166
10.6 Further reading	167
11 Generative models and learning from unlabeled data	169
11.1 The Gaussian mixture model and the LDA & QDA classifiers	169
11.2 The Gaussian mixture model when some or all labels are missing	174
11.3 More unsupervised methods: k -means and PCA	180

1 Introduction

1.1 For whom is this book intended?

This book is a result of an introductory machine learning course that we started to teach at Uppsala University in 2017. The course was taken mainly by engineering students who already had a fairly solid mathematical background as well as some previous programming experience, but had not taken any previous course on machine learning. We did keep that audience in mind also when we wrote this book. In particular we assume that the reader has the basics of calculus, linear algebra and probability theory fresh in mind, but not necessarily any previous familiarity with machine learning.

Of course, everyone is welcome to read the book, either in parts or in its entirety, no matter their background. We hope that all readers will find something they like.

1.2 How to proceed throughout the book

This book is written such that it can be read from start to end. There are, however, multiple possible paths through the book that are more selective depending on the interest of the reader. Figure 1.1 illustrates the major dependencies between the chapters. In particular are the most fundamental topics discussed in Chapter 2, 3 and 4, and we do recommend the reader to read those chapters before proceeding to the later chapters that contains technically more advanced topics (Chapter 5-9). Chapter 10 focuses on some of the more practical aspects on designing a successful machine learning solution and has a less technical nature than the other chapters, whereas Chapter 11 gives some directions for how to go beyond the supervised setting of machine learning.

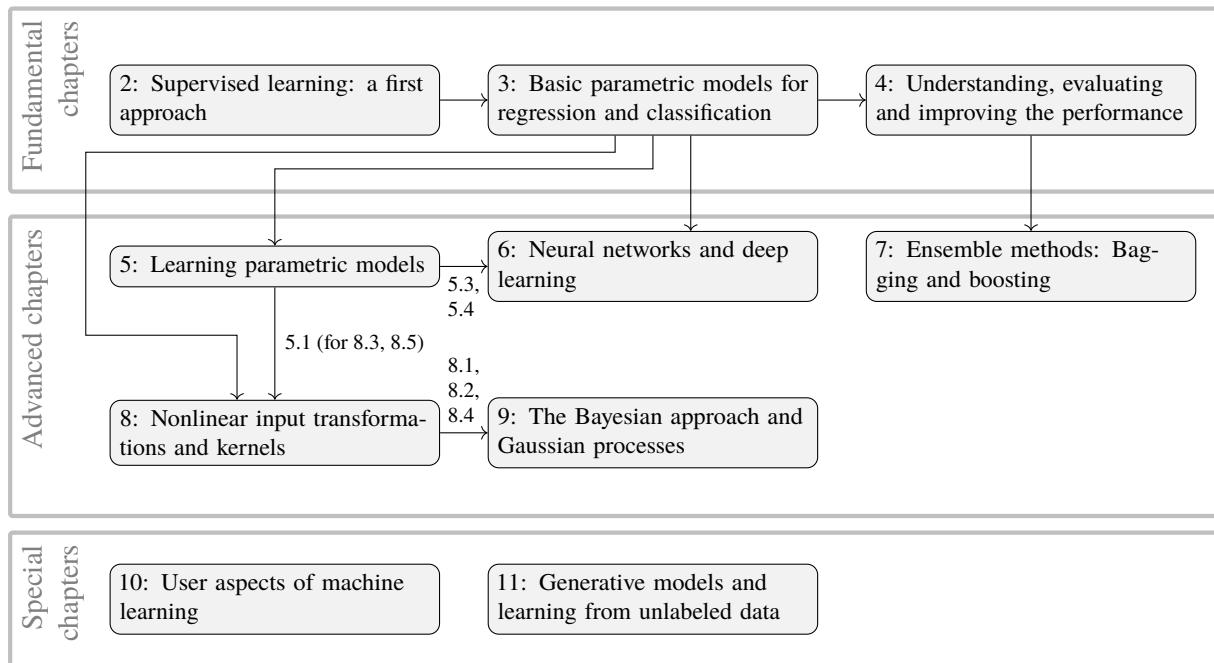


Figure 1.1: The structure of this book, illustrated by blocks (chapters) and arrows (recommended order to read the chapters). We do recommend everyone to read (or at least skim) the fundamental chapters Chapter 2, 3 and 4 first. The path through the technically more advanced Chapter 5-9, however, can be chosen to match the particular interest of the reader. For Chapter 10 and 11 we do recommend the reader to, at least, have read the fundamental chapters first.

2 Supervised learning: a first approach

In this chapter we will introduce the supervised machine learning problem, as well as two basic machine learning methods for solving it. In short the problem amounts to use data for learning the relationship between an input and an output variable. The methods we will introduce are called k -NN and decision trees. These two methods are relatively simple but useful on their own, and therefore a good place to start. Understanding the inner workings, advantages and shortcomings of these basic methods also lays a good foundation for the more advanced methods that are to come in later chapters.

2.1 The supervised learning problem

The supervised machine learning problem is to use *training data* to *learn* (or, equivalently, *train*) a mathematical model that relates an *output*¹ variable y to an *input*² variable \mathbf{x} . That mathematical model can thereafter be used for *predicting* the output y for a new, previously unseen, *test data* for which only \mathbf{x} is known. The beauty of machine learning is that it is quite arbitrary what these variables represent, and it varies from application to application. We saw several examples in Chapter 1 of this problem: Inferring the type of a skin lesion (y) from an image of the skin (\mathbf{x}), deducing the type of road user (y) from a lidar measurement (\mathbf{x}) or predicting the effect on carbon dioxide emissions (y) from a political intervention (\mathbf{x}).

Learning from labeled data

In most interesting supervised machine learning applications, the relationship between input \mathbf{x} and output y is difficult to describe explicitly. It may be too cumbersome or complicated to fully unravel from application domain knowledge, or even unknown. The problem can therefore usually not be solved by writing a traditional computer program that takes \mathbf{x} as input and returns y as output from a set of rules. The supervised machine learning approach is instead to learn the relationship between \mathbf{x} and y from data, which contains examples of observed pairs of input and output values. In other words, supervised machine learning amounts to learning from examples.

The data used for learning is called *training data*, and it has to consist of several input-output samples (\mathbf{x}_i, y_i) , in total n of them. We will compactly write the training data as $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$. Each sample in the training data provides a snapshot of how y depends on \mathbf{x} , and the goal in supervised machine learning is to squeeze as much information as possible out of \mathcal{T} . In this book we will only consider problems where the individual data points are assumed to be (probabilistically) *independent*. This excludes, for example, applications in time series analysis where it is of interest to model the correlation between \mathbf{x}_i and \mathbf{x}_{i+1} .

The fact that the training data contains not only input values \mathbf{x}_i , but also output values y_i , is the reason for the term “supervised” machine learning. We may say that each input \mathbf{x}_i is accompanied with a label y_i , or simply that we have *labeled data*. For some applications, it is only a matter of jointly recording \mathbf{x} and y . In other applications, the output y has to be created by labeling the training data inputs \mathbf{x} by a domain expert. For instance, to construct a training data set for the skin lesion application, a dermatologist needs to look at all training data inputs (images) \mathbf{x}_i and label them by assigning to the variable y_i the type of lesion that is seen in the image. The entire learning process is thus “supervised” by the domain expert.

We use a vector boldface notation \mathbf{x} to denote the input, since we assume it to be a p -dimensional vector, $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_p]^T$. Each element of the input vector \mathbf{x} represent some information that is considered to

¹The output is commonly also called response, regressand, label, explained variable, predicted variable or dependent variable.

²The input is commonly also called feature, attribute, predictor, regressor, covariate, explanatory variable, controlled variable and independent variable.

be relevant, for example the outdoor temperature or the unemployment rate. In an application where the input is a black and white image \mathbf{x} can be all pixel values in the image, so $p = h \times w$, where h and w denote the height and width of the input image.³ The output y , on the other hand, is often of low dimension and throughout most of this book we will assume that it is a scalar value. The *type* of the output value, numerical or categorical, turns out to be important and is used to distinguish between two subtypes of the supervised machine learning problems: *regression* and *classification*. We will discuss this next.

Numerical and categorical variables

The variables contained in our data (input as well as output) can be of two different types, *numerical* or *categorical*. A numerical variable has a natural ordering. We can say that one instance of a numerical variable is larger or smaller than another instance of the same variable. A numerical variable could for instance be represented by a continuous real number, but it could also be discrete, such as an integer. Categorical variables, on the other hand, are always discrete and importantly they lack a natural ordering. In this book we assume that any categorical variable only can take a finite number of different values. A few examples are given in Table 2.1 below.

Table 2.1: Examples of numerical and categorical variables.

Variable type	Example	Handled as
Number (continuous)	32.23 km/h, 12.50 km/h, 42.85 km/h	Numerical
Number (discrete) with natural ordering	0 children, 1 child, 2 children	Numerical
Number (discrete) without natural ordering	1 = Sweden, 2 = Denmark, 3 = Norway	Categorical
Text string	Hello, Goodbye, Welcome	Categorical

The distinction between numerical and categorical is sometimes somewhat arbitrary. We could for instance argue that having no children is something qualitatively different than having children, and use the categorical variable “children: yes/no”, instead of the numerical “0, 1 or 2 children”. In other words, it is sometimes left to the machine learning engineer to decide whether a certain variable is considered to be numerical or categorical.

The notion of categorical vs. numerical applies to both the output variable y and to the p elements x_j of the input vector $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_p]^T$. All p input variables do not have to be of the same type. It is perfectly fine (and common in practice) to have a mix of categorical and numerical inputs.

Classification and regression

We distinguish between different supervised machine learning problems by the type of the output y .

Regression means that the output is numerical, and *classification* means that the output is categorical.

The reason for this distinction is that the regression and classification problems have somewhat different properties, and require different methods to be solved.

Note that the p input variables $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_p]^T$ can be either numerical or categorical for both regression and classification problems. It is only the type of the output that determines whether a problem is a regression or a classification problem. A method which solves a classification problems is often called a *classifier*, and when a classifier predicts the wrong output to a data point it is called a *misclassification*.

For classification the output is categorical and can therefore only take values in a finite set. We use M to denote the number of elements in the set of possible output values. It can, for instance, be `{false, true}` ($M = 2$) or `{Sweden, Norway, Finland, Denmark}` ($M = 4$). We will refer to these elements as *classes* or

³For image-based problems it is often more convenient to represent the input as a matrix of size $h \times w$, than as a vector of length $p = hw$, but the dimension is nevertheless the same. We will get back to this in Chapter 6 when discussing the convolutional neural network, a model structure tailored to image-type inputs.

labels. The number of classes M is assumed to be known in the classification problem. To prepare for a concise mathematical notation, we use integers $1, 2, \dots, M$ to denote the output classes if $M > 2$. The ordering of the integers is arbitrary, and does *not* imply any ordering of the classes. When there are only $M = 2$ classes, we have the important special case of *binary* classification. In binary classification we use the labels -1 and 1 (instead of 1 and 2). Occasionally we will also use the equivalent terms *negative* and *positive* class. The only reason for using a different convention for binary classification is that it gives a more compact mathematical notation for some of the methods, and carries no deeper meaning. Let us now have a look at a classification and a regression problem, which we both will return to later in this book.

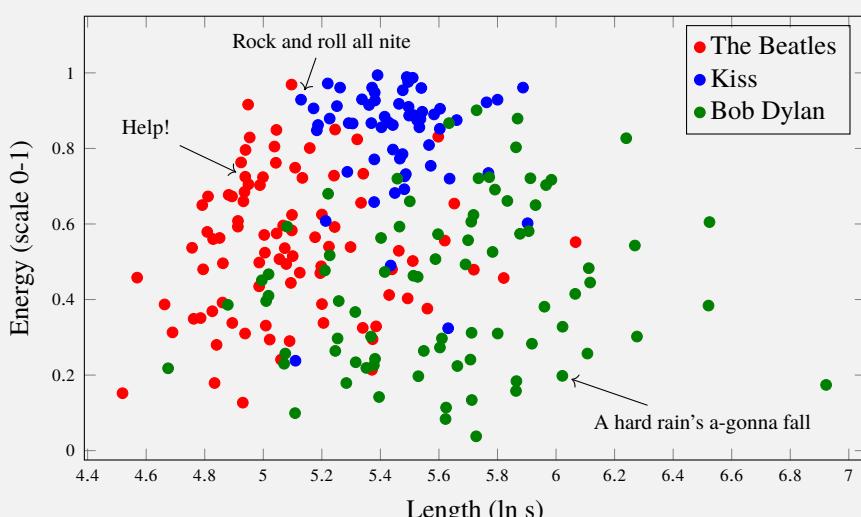
Example 2.1: Classifying songs

Say that we want to build a “song categorizer” app, where the user records a song and the app answers by telling if the song has the artistic style of either the Beatles, Kiss or Bob Dylan. At the heart of this fictitious app there has to be a machinery that takes an audio recording as an input and returns an artist name (the Beatles, Kiss or Bob Dylan).

If we first collect some recordings with songs from the three groups/artists (where we know which artist is behind each song; a labeled dataset), we could use supervised machine learning to *learn* the characteristics of their different styles and therefrom *predict* the artist of the new user-provided song. In the supervised machine learning terminology, the artist name (the Beatles, Kiss or Bob Dylan) is the output y . In this problem y is categorical, and we are hence facing a classification problem.

One of the important design choices for a machine learning engineer is a detailed specification of what the input \mathbf{x} really is. It would in principle be possible to consider the raw audio information as input, but that would give a very high-dimensional \mathbf{x} which (unless an audio-specific machine learning method is used) most likely would require an unrealistically large amount of training data in order to be successful (we will discuss this aspect in detail in Chapter 4). A better option could therefore be to define some summary statistics (features) of audio recordings, and use those as input \mathbf{x} instead. As input we could for example use the length of the audio recording and the amount of energy in it. The length of a recording is easy to measure, and since some songs are rather long it seems natural to take the logarithm of it. The energy is a bit more trickier (the exact definition may even be ambiguous), but we can leave that to the audio experts and re-use a piece of software that they have written for this purpose^a without bothering too much about its inner workings. As long as this piece of software returns a number for any recording that is sent into it (and always returns the same number for the same recording), we can use it as an input to a machine learning method.

Below we have plotted a dataset with 230 songs from the three artists. Each song is represented by a dot, where the horizontal axis is the logarithm of its length (measured in seconds) and the vertical axis the energy (on a scale 0-1). When we later return to this example, and apply different supervised machine learning methods to it, this data will be the training data.

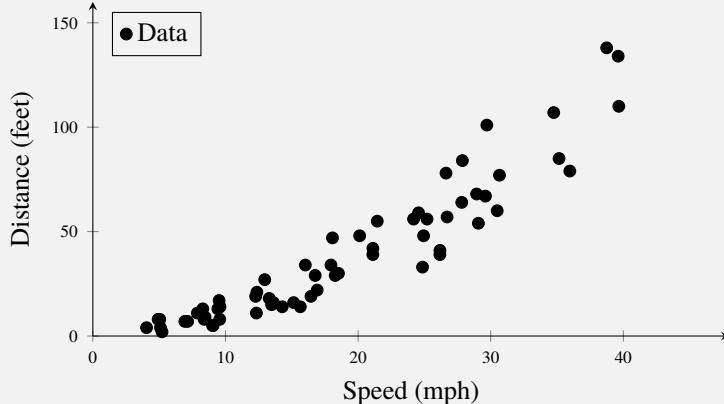


^aWe use <http://api.spotify.com/> here.

Example 2.2: Car stopping distances

Ezekiel and Fox (1959) presents a dataset with 62 observations of the distance needed for various cars at different initial speeds to break to a full stop.^a The dataset has the following two variables:

- Speed: The speed of the car when the break signal is given.
- Distance: The distance traveled after the signal is given until the car has reached a full stop.



To make a supervised machine learning problem out of it, we interpret Speed as the input variable x , and Distance as the output variable y . Since y is numerical, this is a regression problem. We then ask ourselves what the stopping distance would be if the initial speed would be, for example, at 33 mph and 45 mph respectively (two speeds at which no data has been recorded). Another way to frame this question is to ask for the prediction $\hat{y}(x_*)$ for $x_* = 33$ and $x_* = 45$.

^aThe data is somewhat dated, so the conclusions are perhaps not applicable to modern cars.

Time to reflect 2.1: For each application example discussed in Chapter 1, can you determine the input x , the output y , and whether it is a regression or classification problem?

Generalizing beyond training data

There are two primary reasons for why it is of interest to design methods for learning mathematical models of input–output relationships from training data.

- (i) To *reason about and explore* how input and output variables are connected. An often encountered task in sciences such as medicine and sociology is to determine whether a correlation between a pair of variables exists or not (“does sea food increase the life expectancy?”). Such questions can be addressed by learning a mathematical model and carefully reason about the chance that the learned relationships between input x and output y are due only to random effects in the data, or if there appear to be some substance to the found relationships.
- (ii) To *predict* the output value y_* for some new, previously unseen input x_* . By first learning a mathematical model of the input–output relationship from the training data, we can make a prediction by inserting the test input x_* into the model and use the output from the model $\hat{y}(x_*)$ as a prediction of the true (but unknown) output y_* associated with x_* . The hat $\hat{\cdot}$ indicates that the prediction is an estimate (that is, an educated guess) of the true output.

These two objectives are sometimes used to roughly distinguish between classical statistics (i) and machine learning (ii), even though it is not completely true (predictive modeling is a topic of classical statistics and explainable models is a topic within machine learning). The primary focus in this book, however, is on learning models for making predictions, objective (ii) above, which is the foundation of supervised machine learning. Our overall goal is to obtain as accurate predictions $\hat{y}(x_*)$ as possible

(measured in some appropriate way) for a wide range of possible test inputs \mathbf{x}_* . We say that the we are interested in models that *generalize* well beyond the training data.

A model that generalizes well for the music example above would be able to correctly tell the artist of a new song which was not in the training data (assuming that the artist of the new song is one of the three that was present in the training data, of course). The ability to generalize to new data is a key concept of machine learning. It is not difficult to construct models that give very accurate predictions if they are only evaluated on the training data (we will see an example in the next section). If the model is not able to generalize, meaning that the predictions are poor when the model is applied to new test data points, then the model is of little use in practice for making predictions. If that is the case we say that the model is *overfitting* to the training data. We will illustrate the issue of overfitting for a specific machine learning model in the next section and in Chapter 4 we will return to this concept using a more general and mathematical approach.

2.2 A distance-based method: k -NN

It is now high time to encounter the first actual machine learning method. We will start with the relatively simple k -NN method, which can be used for both regression and classification. Remember that the setting is that we have access to training data $\{\mathbf{x}_i, y_i\}_{i=1}^n$, which consists of n input samples \mathbf{x}_i with corresponding output samples y_i . From this we want to construct a prediction $\hat{y}(\mathbf{x}_\star)$ for what we believe the output y_\star would be for a new \mathbf{x}_\star , which we have not seen previously.

The k -nearest neighbors method

Most methods for supervised machine learning builds on the intuition that *if the test data point \mathbf{x}_\star is close to training data sample \mathbf{x}_i , then the prediction $\hat{y}(\mathbf{x}_\star)$ should be close to y_i* . This is a general idea, but one simple way to implement it in practice is the following way: First, compute the distance⁴ between the test input and all training inputs, $\|\mathbf{x}_i - \mathbf{x}_\star\|$ for $i = 1, \dots, n$. Second, we find the data point \mathbf{x}_j with the *shortest distance* to \mathbf{x}_\star and let the prediction be $\hat{y}(\mathbf{x}_\star) = y_j$.

This simple prediction method is referred to as the 1-nearest neighbor method. It is not very complicated, but for most machine learning applications of interest it is too simplistic. In practice we can rarely say *for certain* what the output value y will be, but we should mathematically describe y as random variable. Consequently, it makes sense to consider the data as *noisy*, that is, being affected by random noise. In this perspective, the shortcoming of 1-nearest neighbor is that the prediction relies only on *one* sample from the training data, which makes it quite “erratic” and sensitive to noisy training data.

To improve the 1-nearest neighbor method we can extend it to make use of the k nearest neighbors instead. Formally we define the set $R_\star = \{i : \mathbf{x}_i \text{ is one of the } k \text{ training data points closest to } \mathbf{x}_\star\}$ and aggregate the information from the k outputs y_j for $j \in R_\star$ to make the prediction. For regression problems we take the average of all y_j for $j \in R_\star$, and a majority vote⁵ for classification problems. We illustrate the k -nearest neighbors (k -NN) method by Example 2.3 and summarize by Algorithm 1.

Methods that explicitly use the training data when making predictions are referred to as *non-parametric*. This is in contrast with *parametric* methods, where the prediction is given by some function (a model), governed by a fixed number of parameters. For parametric methods the training data is used to *learn* the parameters in an initial training phase, but once the model has been learned, the training data can be discarded since it is not used explicitly when making predictions. We will introduce some parametric models in Chapter 3.

Algorithm 1: k -nearest neighbor, k -NN

Data: Training data $\{\mathbf{x}_i, y_i\}_{i=1}^n$ and test input \mathbf{x}_\star

Result: Predicted test output $\hat{y}(\mathbf{x}_\star)$

- 1 Compute the distances $\|\mathbf{x}_i - \mathbf{x}_\star\|$ for all training data points $i = 1, \dots, n$
- 2 Let $R_\star = \{i : \mathbf{x}_i \text{ is one of the } k \text{ data points closest to } \mathbf{x}_\star\}$
- 3 Compute the prediction $\hat{y}(\mathbf{x}_\star)$ as

$$\hat{y}(\mathbf{x}_\star) = \begin{cases} \text{Average}\{y_j : j \in R_\star\} & (\text{Regression problems}) \\ \text{MajorityVote}\{y_j : j \in R_\star\} & (\text{Classification problems}) \end{cases}$$

⁴Other distance functions can also be used, and will be discussed in Chapter 8. Categorical input variable can be handled as we will discuss in Section 3.1.

⁵Ties can be handled in different ways, for instance by a coin-flip, or by reporting the actual vote count to the end user who gets to decide what to do with it.

Example 2.3: Predicting colors with k -NN

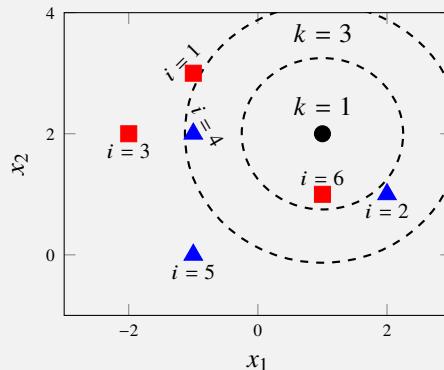
We consider a binary classification problem ($M = 2$). We are given a training data set with $n = 6$ observations of $p = 2$ input variables x_1, x_2 and one categorical output y , the color Red or Blue,

i	x_1	x_2	y
1	-1	3	Red
2	2	1	Blue
3	-2	2	Red
4	-1	2	Blue
5	-1	0	Blue
6	1	1	Red

and we are interested in predicting the output for $\mathbf{x}_\star = [1 \ 2]^\top$. For this purpose we will explore two different k -NN classifiers, one using $k = 1$ and one using $k = 3$.

First, we compute the Euclidian distance $\|\mathbf{x}_i - \mathbf{x}_\star\|$ between each training data point \mathbf{x}_i and the test data point \mathbf{x}_\star , and then sort them in ascending order.

i	$\ \mathbf{x}_i - \mathbf{x}_\star\ $	y_i
6	$\sqrt{1}$	Red
2	$\sqrt{2}$	Blue
4	$\sqrt{4}$	Blue
1	$\sqrt{5}$	Red
5	$\sqrt{8}$	Blue
3	$\sqrt{9}$	Red



Since the closest training data point to \mathbf{x}_\star is the data point $i = 6$ (Red), it means that for k -NN with $k = 1$, we get the prediction $\hat{y}(\mathbf{x}_\star) = \text{Red}$. For $k = 3$, the 3 nearest neighbors are $i = 6$ (Red), $i = 2$ (Blue), and $i = 4$ (Blue). Taking a majority vote among these three training data points, Blue wins with 2 votes against 1, so our prediction becomes $\hat{y}(\mathbf{x}_\star) = \text{Blue}$.

This is also illustrated by the figure above where the training data points \mathbf{x}_i are represented with red squares and the blue triangles depending on which class they belong to. The test data point \mathbf{x}_\star is represented with a black filled circle. For $k = 1$ the closest training data point is identified by the inner circle and for $k = 3$ the three closest points are identified by the outer circle.

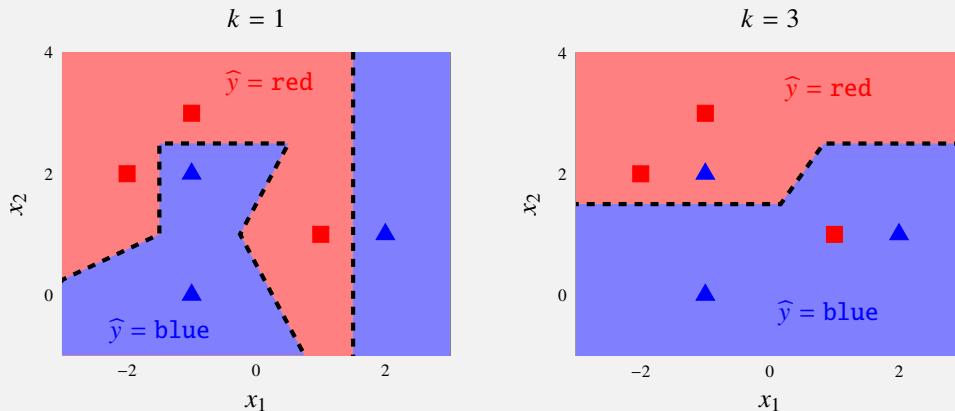
Decision boundaries for a classifier

In Example 2.3 we only computed a prediction for one single test data point \mathbf{x}_* . That prediction might indeed be the ultimate goal of the application, but in order to visualize and better understand a classifier we can also study its *decision boundary*, which illustrates the prediction for all possible test inputs. We introduce the decision boundary using Example 2.4. It is a general concept for classifiers, not only k -NN, but it is only possible to visualize when the dimension of \mathbf{x} is $p = 2$.

Example 2.4: Decision boundaries for the color example

In Example 2.3 we computed the prediction for $\mathbf{x}_* = [1 \ 2]^\top$. If we would shift that test point by one step to the left at $\mathbf{x}_{*}^{\text{alt}} = [0 \ 2]^\top$ the three closest training data points would still include $i = 6$ and $i = 4$ but now $i = 2$ is exchanged for $i = 1$. For $k = 3$ this would give two votes for Red and one vote for Blue and we would therefore predict $\hat{y} = \text{Red}$. In between these two test data points \mathbf{x}_* and $\mathbf{x}_{*}^{\text{alt}}$ at $[0.5 \ 2]^\top$ it is equally far to $i = 1$ as to $i = 2$ and it is undecided if the 3-NN classifier should predict Red or Blue. (In practice this is most often not a problem, since the test data points rarely end up exactly at the decision boundary. If they do, this can be handled by a coin-flip.) For all classifiers we always end up with such points in the input space where the class prediction abruptly changes from one class to another. These points are said to be on the *decision boundary* of the classifier.

Continuing in a similar way, changing the location of the test input across the entire input space and recording the class prediction, we can compute the full decision boundaries for Example 2.3. We plot the decision boundaries using dashed lines below, both for $k = 1$ and $k = 3$.



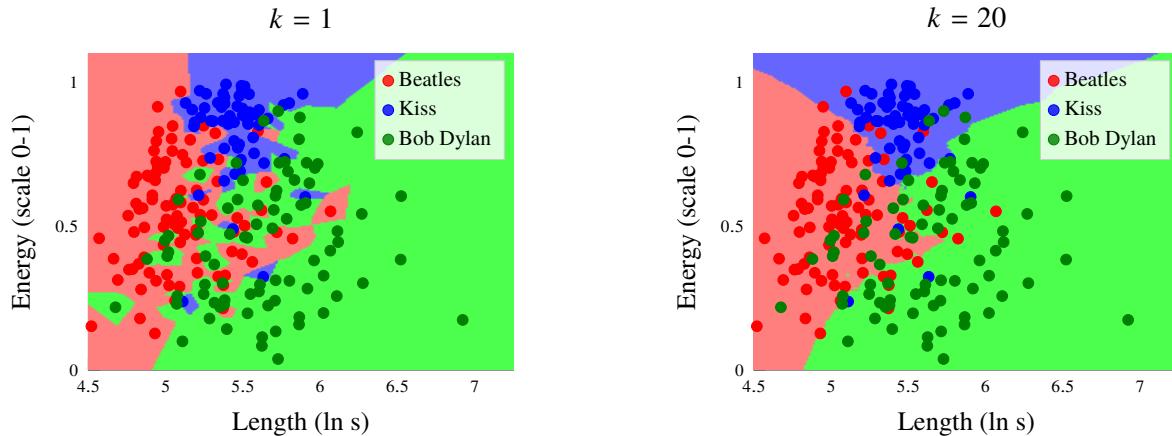
These figures show the decision boundaries, that is, the points in input space where the class prediction changes. This type of figure gives a concise summary of a classifier, but requires of course that the problem has a 2-dimensional input in order to make such a plot. As we can see, the decision boundaries of k -NN are not linear. In the terminology we will introduce later, k -NN is thereby a nonlinear classifier.

Choosing k

The user has to decide on which k to use in k -NN and this decision has a big impact on the final predictions. To understand the impact of k , we study how the decision boundary changes as k changes in Figure 2.1, where k -NN is applied to the music classification Example 2.1 and car stopping distance Example 2.2, both with $k = 1$ as well as $k = 20$.

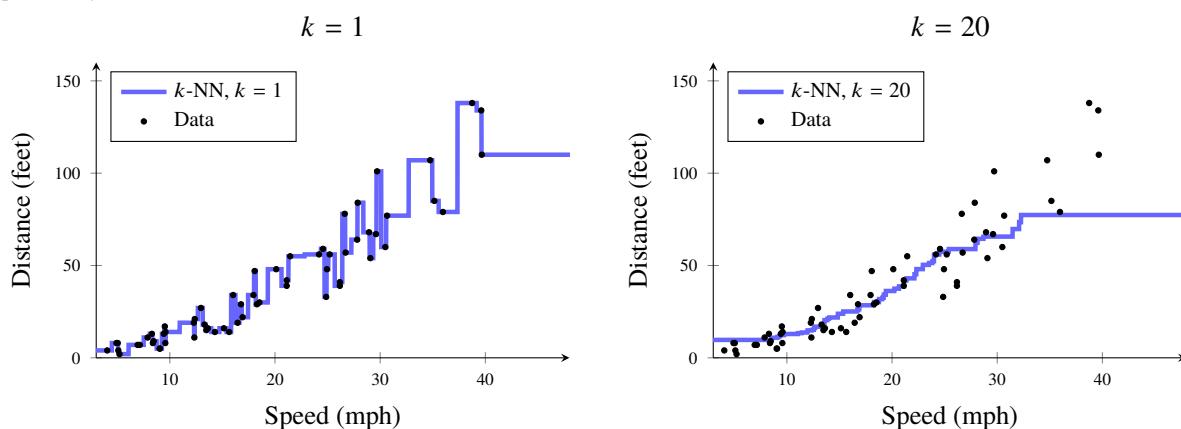
With $k = 1$ all training data points will, by construction, be correctly predicted and the model is adapted to the exact \mathbf{x} and y values of the training data. In the classification problem there are for instance small green (Bob Dylan) regions within the red (the Beatles) area that are most likely not relevant when it comes to accurately predict the artist of a new song. In order to make good predictions it would probably be better to instead predict red (the Beatles) for a new song in the entire middle-left region since the vast majority of training data points in that area are red. For the regression problem $k = 1$ gives a quite shaky

behavior, and also for this problem it is intuitively clear that this does not describe an interesting effect, but rather noise.



(a) Decision boundaries for the music classification problem using $k = 1$. This is a typical example of overfitting, meaning that the model has adapted too much to the training data so that it does not generalize well to new previously unseen data.

(b) The music classification problem again, now using $k = 20$. A higher value of k gives a more smooth behavior which, hopefully, predicts the artist of new songs more accurately.



(c) The black dots are the car stopping distance data, and the blue line shows the prediction for k -NN with $k = 1$ for any x . As for the classification problem above, k -NN with $k = 1$ overfits to the training data.

(d) The car stopping distance , this time with $k = 20$. Except for the boundary effect at the right, this seems like a much more useful model which captures the interesting effects of the data and ignores the noise.

Figure 2.1: k -NN applied to the music classification Example 2.1 (a and b) and the car stopping distance Example 2.2 (c and d). For both problems k -NN is applied with $k = 1$ and well as $k = 20$.

The drawbacks of using $k = 1$ is not specific to these two examples. In most real world problems there is a certain amount of randomness involved in collecting the training data (or at least some decisions which we know little about, and may consider as being random). In the music example the $n = 230$ songs were somehow selected from all songs ever recorded from these artists, and for the car stopping distance there appears to also be a certain amount of random effects also in y . Thus, by using $k = 1$ and thereby adapting very closely to the training data, the predictions will depend not only on the interesting patterns in the problem, but also on the (more or less) random effects that have shaped the training data. Typically we are not interested in capturing these effects, and we refer to this as *overfitting*.

With the k -NN classifier we can mitigate overfit by increasing the region of the neighborhood used to compute the prediction, that is, increase the parameter k . With, for example, $k = 20$ the predictions are no longer based on only the closest neighbor, but instead a majority vote among the 20 closest neighbors. As a consequence all training data samples are no longer perfectly classified, but some of the songs end up in the wrong region in Figure 2.1b. The predictions are however less adapted to the peculiarities of the

training data and thereby less overfitted, and Figure 2.1b as well as Figure 2.1d are indeed less “noisy” than Figure 2.1a and Figure 2.1c. Selecting k is thus a trade-off between flexibility and rigidity. Indeed, selecting k too big will lead to a meaningless classifier, so there must exist a sweet spot for some moderate k (possibly 20, but could be less or more) where the classifier generalizes the best. Unfortunately there is no general answer for which k this happens, but it is different for different problems. In the music classification it seems reasonable that $k = 20$ will predict new test data points better than $k = 1$, but there might very well be an even better choice of k . For the car stopping problem the behavior is also more reasonable for $k = 20$ than $k = 1$, except for the boundary effect for large x where k -NN is unable to capture the trend in the data as x increases (simply because the 20 nearest neighbors are the same for all test points x_\star around and above 35). A systematic way of choosing a good value of k is to use cross-validation, which we will discuss in Chapter 4.

Time to reflect 2.2: The prediction $\hat{y}(x_\star)$ obtained using the k -NN method is a piece-wise constant function of the input x_\star . For a classification problem this is natural, since the output is categorical (see, e.g., Figure 2.1 where the colored regions correspond to areas of the input space where the prediction is constant according to the color of that region). However, also for regression does k -NN have piece-wise constant predictions. Why?

Input normalization

A final important practical aspect when using k -NN is the importance of normalization of the input data. Since k -NN is based on the Euclidean distances between points, it is important that it is a relevant measure of how close two data points are. Imagine a training data set with $p = 2$ input variables $\mathbf{x} = [x_1 \ x_2]^\top$ where all values of x_1 are in the range [100, 1100] and the values for x_2 are in the much smaller range [0, 1]. It could for example be that x_1 and x_2 are measured in different units. The Euclidean distance between a test point \mathbf{x}_\star and a training data point \mathbf{x}_i is $\|\mathbf{x}_i - \mathbf{x}_\star\| = \sqrt{(x_{i1} - x_{\star 1})^2 + (x_{i2} - x_{\star 2})^2}$, and this expression will be totally dominated by the first term $(x_{i1} - x_{\star 1})^2$ whereas the second term $(x_{i2} - x_{\star 2})^2$ will almost not matter in practice, due to the different magnitude of x_1 and x_2 . That is, the different ranges leads to x_1 being considered much more important than x_2 by k -NN.

To avoid this undesired effect we have to scale the input variables. One option, in the mentioned example, could be to subtract 100 from x_1 and thereafter divide it by 1000 and create $x_{i1}^{\text{new}} = \frac{x_{i1} - 100}{1000}$ such that x_1^{new} and x_2 both are in the range [-1, 1]. More generally, this normalization procedure for the input data can be written as

$$x_{ij}^{\text{new}} = \frac{x_{ij} - \min_\ell(x_{\ell j})}{\max_\ell(x_{\ell j}) - \min_\ell(x_{\ell j})}, \quad \text{for all } j = 1, \dots, p, \quad i = 1, \dots, n. \quad (2.1)$$

Another common way of normalizing (sometimes called standardizing) is by using the mean and standard deviation in the training data:

$$x_{ij}^{\text{new}} = \frac{x_{ij} - \bar{x}_j}{\sigma_j}, \quad \forall j = 1, \dots, p, \quad i = 1, \dots, n, \quad (2.2)$$

where \bar{x}_j and σ_j are the mean and standard deviation for each input variable, respectively.

It is crucial for k -NN to apply some type of input normalization (that was indeed done in Figure 2.1), but it is a good practice to apply also when using other methods, at least for numerical stability. It is however important to compute the scaling factors ($\min_\ell(x_{\ell j})$, \bar{x}_j , etc) using training data only and apply that scaling also to future test data points. Failing this, for example by performing normalization before setting test data aside (which we will discuss more in Chapter 4), might lead to wrong conclusions on how well the method will perform in predicting future not yet seen data points.

2.3 A rule-based method: Decision trees

The k -NN method results in a prediction $\hat{y}(\mathbf{x}_\star)$ that is a piece-wise constant function of the input \mathbf{x}_\star . That is, the method partitions the input space into disjoint regions and each region is associated with a certain (constant) prediction. For k -NN, these regions are given implicitly by the k -neighborhood of each possible test input. An alternative approach, that we will study in this section, is to come up with a set of *rules* that define the regions explicitly. For instance, considering the data in Figure 2.1, a simple set of high-level rules for constructing a classifier would be: inputs on the right side are classified as green, in the lower left corner as blue, and in the upper left corner as red. We will now see how such rules can be learned systematically from the training data.

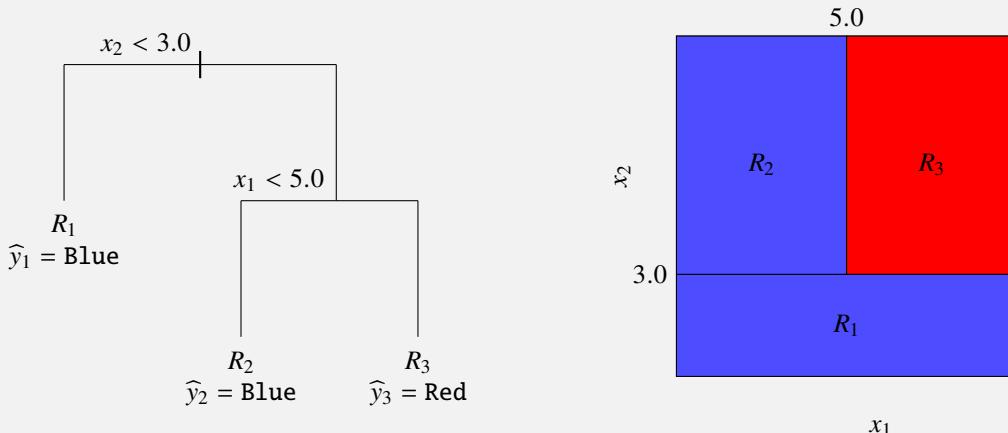
Predicting using a decision tree

The rule-based models that we consider here are referred to as *decision trees*. The reason is that the rules used to define the model can be organized in a graph structure referred to as a binary tree. The decision tree effectively divides the input space into multiple disjoint regions and in each region a constant value is used for the prediction $\hat{y}(\mathbf{x}_\star)$. We illustrate this with an example.

Example 2.5: Predicting colors with a decision tree

We consider a classification problem with two numerical input variables $\mathbf{x} = [x_1 \ x_2]^\top$ and one categorical output y , the color Red or Blue. For now we do not consider any training data or how to actually learn the tree, but only how an already existing decision tree can be used to predict $\hat{y}(\mathbf{x}_\star)$.

The rules defining the model are organized in the graph below (left figure) which is referred to as a binary tree. To use this tree to predict a label for the test input $\mathbf{x}_\star = [x_{\star 1} \ x_{\star 2}]^\top$ we start at the top, referred to as the *root node* of the tree (in the metaphor the tree is growing upside down, with the root at the top and the leaves at the bottom). If the condition stated at the root is true, i.e. if $x_{\star 2} < 3.0$, then we proceed down the left *branch*, otherwise along the right *branch*. If we reach a new *internal node* of the tree, we check the rule associated with that node and pick the left or the right branch accordingly. We continue and work our way down until we reach the end of a branch. Each such final branch corresponds to a constant prediction \hat{y}_m , in this case one of the two classes Red or Blue.



A classification tree. At each internal node a rule on the form $x_j < s_k$ indicates the left branch coming from that split and the right branch then consequently corresponds to $x_j \geq s_k$. This tree has two internal nodes (including the root) and three leaf nodes.

A region partition, where each region corresponds to a leaf node in the tree. Each border between regions correspond to a split in the tree. Each region is colored with the prediction corresponding to that region, and the boundary between red and blue is therefore the decision boundary.

The decision tree partitions the input space into axis-aligned “boxes”, as shown in the right panel above. By increasing the depth of the tree (the number of steps from the root to the leaves), the partitioning can be made finer and finer and thereby describing more complicated functions of the input variable.

A pseudo code for predicting a test input with the tree above would look like

```

if x_2 < 3.0 then
    return Blue
else
    if x_1 < 5.0 then
        return Blue
    else
        return Red
    end
end

```

As an example if we have $\mathbf{x}_\star = [2.5 \ 3.5]^\top$, in the first split we would take the right branch since $x_{\star 2} = 3.5 \geq 3.0$ and in the second split we would take the left branch since $x_{\star 1} = 2.5 < 5.0$. The prediction for this test point would be $\hat{y}(\mathbf{x}_\star) = \text{Blue}$.

To set the terminology, the endpoint of each branch R_1 , R_2 and R_3 in Example 2.5 are called *leaf nodes* and the internal splits, $x_2 < 3.0$ and $x_1 < 5.0$ are known as *internal nodes*. The lines that connect the nodes are referred to as *branches*. The tree is referred to as *binary* since each internal node splits into exactly two branches.

With more than two input variables it is difficult to illustrate the partitioning of the inputs space into regions (right figure in the example), but the tree representation can still be used in precisely the same way. Each internal node corresponds to a rule where one of the p input variables x_j , $j = 1, \dots, p$, is compared with a threshold s . If $x_j < s$ we continue along the left branch and if $x_j \geq s$ we continue along the right branch.

The constant predictions that we associate with the leaf nodes can be either categorical, as in the example above, or numerical and decision trees can thus be used to address both classification and regression problems. In fact, the models that we describe here also go by the name Classification And Regression Trees (CART). Example 2.5 illustrates how a decision tree is used for making a prediction. We will now turn to the question of how the tree can be learned from training data.

Learning a regression tree

We will now discuss how to learn (or, equivalently, train) a decision tree for a regression problem. Training a decision tree for classification is conceptually similar and explained in the next section.

As discussed above, the prediction $\hat{y}(\mathbf{x}_\star)$ corresponding to a classification or regression tree is a piece-wise constant function of the input \mathbf{x}_\star . We can write this mathematically as,

$$\hat{y}(\mathbf{x}_\star) = \sum_{\ell=1}^L \hat{y}_\ell \mathbb{I}\{\mathbf{x}_\star \in R_\ell\}, \quad (2.3)$$

where L is the total number of regions (leaf nodes) in the tree, R_ℓ is the ℓ th region, and \hat{y}_ℓ is the constant prediction for the ℓ th region. Note that in the regression setting \hat{y}_ℓ is a numerical variable, and we will consider it to be a real number for simplicity. In the equation above we have used the indicator function, $\mathbb{I}\{\mathbf{x} \in R_\ell\} = 1$ if $\mathbf{x} \in R_\ell$ and $\mathbb{I}\{\mathbf{x} \in R_\ell\} = 0$ if $\mathbf{x} \notin R_\ell$ otherwise.

Learning the tree from data corresponds to finding suitable values for the parameters defining the function (2.3), namely the regions R_ℓ and the constant predictions \hat{y}_ℓ , $\ell = 1, \dots, L$, as well as the total size of the tree L . If we start by assuming that the shape of the tree, the partition $(L, \{R_\ell\}_{\ell=1}^L)$ is known, then we can compute the constants $\{\hat{y}_\ell\}_{\ell=1}^L$ in a natural way, simply as the average of the training data points falling in each region:

$$\hat{y}_\ell = \text{Average}\{y_i : \mathbf{x}_i \in R_\ell\}$$

It remains however to find the shape of the tree, the regions R_ℓ , which requires a bit more work. The basic idea is of course to select the regions so that the tree fits the training data. This means that the output predictions from the tree should match the output values in the training data. Unfortunately, even when

restricting ourselves to seemingly simple regions such as the “boxes” obtained from a decision tree, finding the tree (a collection of splitting rules) that optimally partitions the input space to fit the training data as well as possible turns out to be computationally infeasible. The problem is that there is a combinatorial explosion in the number of ways in which we can partition the input space. Searching through all possible binary trees is not possible in practice unless the tree size is so small that it is practically useless.

To handle this situation we use a *greedy* algorithm known as *recursive binary splitting* for learning the tree. The word “recursive” means that we will determine the splitting rules one after the other, starting with the first split at the root and then build the tree from top to bottom. The word greedy means that tree is constructed one split at a time, without having the complete tree “in mind”. That is, when determining the splitting rule at the root node, the objective is to obtain a model that explains the training data as well as possible after *a single split*, without taking into consideration that additional splits may be added before arriving at the final model. When we have decided on the first split of the input space (corresponding to the root node of the tree), this split is kept fixed and we continue in a similar way for the two resulting half-spaces (corresponding to the two branches of the tree), etc.

To see in detail how one step of this algorithm works, consider the setting when we are about to do our very first split at the root of the tree. Hence, we want to select one of the p input variables x_1, \dots, x_p and a corresponding cutpoint s which divides the input space into two half-spaces,

$$R_1(j, s) = \{\mathbf{x} \mid x_j < s\} \quad \text{and} \quad R_2(j, s) = \{\mathbf{x} \mid x_j \geq s\}. \quad (2.4)$$

Note that the regions depend on the index j of the splitting variable as well as the value of the cutpoint s , which is why we write them as functions of j and s . This is the case also for the predictions associated with the two regions,

$$\hat{y}_1(j, s) = \text{Average}\{y_i : \mathbf{x}_i \in R_1(j, s)\} \quad \text{and} \quad \hat{y}_2(j, s) = \text{Average}\{y_i : \mathbf{x}_i \in R_2(j, s)\}$$

since the sums in these expression range over different data points depending on the regions.

For each training data point (\mathbf{x}_i, y_i) we can compute a *prediction error* by first determining which region the data point falls in, and then computing the difference between y_i and the constant prediction associated with that region. Doing this for all training data points the sum of squared errors can be written as

$$\sum_{i:\mathbf{x}_i \in R_1(j, s)} (y_i - \hat{y}_1(j, s))^2 + \sum_{i:\mathbf{x}_i \in R_2(j, s)} (y_i - \hat{y}_2(j, s))^2. \quad (2.5)$$

The square is added to ensure that the expression above is non-negative. The squared error is a common *loss function* used for measuring the closeness of a model’s prediction and the training data, but other loss functions can also be used. We will discuss the choice of loss function in more detail in later chapters.

To find the optimal split we select the values for j and s that minimize the squared error (2.5). This minimization problem can be solved easily by looping through all possible values for $j = 1, \dots, p$. For each j we can scan through the finite number of possible splits, and pick the pair (j, s) for which the expression above is minimized. As pointed out above, when we have found the optimal split at the root node, this splitting rule is fixed. We then continue in the same way for the left and right branch independently. Each branch (corresponding to a half-space) is split again by minimizing the squared prediction error over all training data points correspond to that branch.

In principle, we can continue in this way until there is only a single training data point in each of the regions, i.e. $L = n$. Such a *full grown tree* will result in predictions that exactly match the training data points, and the resulting model is quite similar to k -NN with $k = 1$. As pointed out above, this will typically result in a too erratic model that has overfitted to (possibly noisy) training data. To mitigate this issue, it is common to stop the growth of the tree at an earlier stage, for instance by adding a constraint on the minimum number of training data points associated with each leaf node. Forcing the model to have more training data points in each leaf will result in an averaging effect, similarly to increasing the value of k in the k -NN method. Note that using such a stopping criteria means that the value of L is not set manually, but determined adaptively based on the result of the learning procedure. The full algorithm is

summarized in Algorithm 2 and Algorithm 3. Note that Algorithm 2 includes a recursive call, where we in each recursion grow one branch of the tree one step further. Also the Algorithm 2 can be formulated as a recursion, but here we have chosen an easier presentation.

Algorithm 2: Learning a decision tree using recursive binary splitting

Data: Training data $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$
Result: Decision tree with regions $R_1, \dots, R_L \in \mathcal{R}$

- 1 Initialize the tree with R being the whole input space
- 2 **return** $\text{Split}(R, \mathcal{T})$
- 3 **Function** $\text{Split}(R, \mathcal{T})$:
- 4 **if** splitting criteria fulfilled **then**
- 5 **return** R
- 6 **else**
- 7 Go through all possible splits $x_j < s$ for all input variables $j = 1, \dots, p$.
- 8 Pick the pair (j, s) that minimizes (2.5)/(2.6) for regression/classification problems.
- 9 Split region R into R_1 and R_2 according to (2.4).
- 10 Split data \mathcal{T} into \mathcal{T}_1 and \mathcal{T}_2 accordingly.
- 11 **return** $\text{Split}(R_1, \mathcal{T}_1), \text{Split}(R_2, \mathcal{T}_2)$
- 12 **end**
- 13 **end**

Algorithm 3: Making predictions from a decision tree

Data: Decision tree with regions R_1, \dots, R_L , training data $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ and test data point \mathbf{x}_*

Result: Predicted test output $\hat{y}(\mathbf{x}_*)$

- 1 Find the region R_ℓ which \mathbf{x}_* belongs to.
- 2 Compute the prediction $\hat{y}(\mathbf{x}_*)$ as

$$\hat{y}(\mathbf{x}_*) = \begin{cases} \text{Average}\{y_i : \mathbf{x}_i \in R_\ell\} & \text{(Regression problems)} \\ \text{MajorityVote}\{y_i : \mathbf{x}_i \in R_\ell\} & \text{(Classification problems)} \end{cases}$$

Classification trees

Trees can also be used for classification. We use the same procedure of recursive binary splitting but with two main differences.

Firstly, we then use a majority vote instead of an average to compute the prediction associated with each region,

$$\hat{y}_\ell = \text{MajorityVote}\{y_i : \mathbf{x}_i \in R_\ell\}.$$

Secondly, when learning the tree we need a different splitting criteria than the squared prediction error to take into account the fact that the output is categorical. To define these criteria we first introduce

$$\hat{\pi}_{\ell m}(j, s) = \frac{1}{n_\ell} \sum_{i: \mathbf{x}_i \in R_\ell(j, s)} \mathbb{I}\{y_i = m\}$$

to be the proportion of training observations in the ℓ th region that belong to the m th class.

Similar to (2.5), we now want to minimize

$$\min_{j, s} n_1 Q_1 + n_2 Q_2 \quad (2.6)$$

where n_ℓ is the number of data points in region R_ℓ and where Q_ℓ is the splitting criteria.

For Q_ℓ , we have a few options. One simple alternative is the *misclassification rate*

$$Q_\ell = 1 - \max_m \hat{\pi}_{\ell m}, \quad (2.7)$$

which is simply the proportion of data points in region R_ℓ which do not belong to the most common class. Another common splitting criteria is the *Gini index*

$$Q_\ell = \sum_{m=1}^M \hat{\pi}_{\ell m} (1 - \hat{\pi}_{\ell m}) \quad (2.8)$$

and the *entropy* criteria.

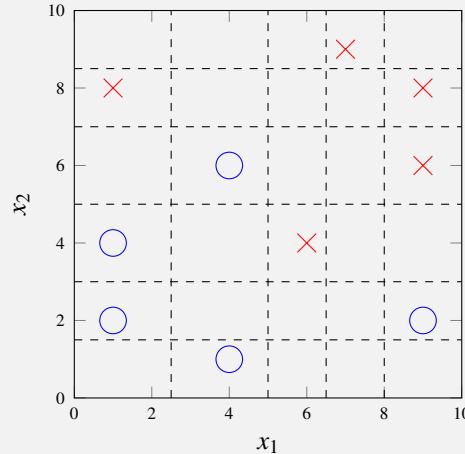
$$Q_\ell = - \sum_{m=1}^M \hat{\pi}_{\ell m} \ln \hat{\pi}_{\ell m}. \quad (2.9)$$

In Example 2.5 we illustrate how to construct a classification tree using recursive binary splitting and with the entropy as the slitting criteria.

Example 2.6: Learning a classification tree (continuation of Example 2.5)

We consider the same setup as in Example 2.5 with the following dataset

x_1	x_2	y
9.0	2.0	Blue
1.0	4.0	Blue
4.0	6.0	Blue
4.0	1.0	Blue
1.0	2.0	Blue
1.0	8.0	Red
6.0	4.0	Red
7.0	9.0	Red
9.0	8.0	Red
9.0	6.0	Red



We want to learn a classification tree, by using the entropy criteria in (2.9) and growing the tree until there are no regions with more than five data points left.

First split: There are infinitely many possible splits we can make, but all splits which gives the same partition of the data points will be the same. Hence, in practice we only have nine different splits to consider in this dataset. The data and these splits (dashed lines) are visualized in the figure above.

We consider all nine splits in turn. We start with the split at $x_1 = 2.5$, which splits the input space into the two regions, $R_1 = x_1 < 2.5$ and $R_2 = x_1 \geq 2.5$. In region R_1 we have two blue data points and one red, in total $n_1 = 3$ data points. The proportion of the two classes in region R_1 will therefore be $\hat{\pi}_{1B} = 2/3$ and $\hat{\pi}_{1R} = 1/3$. The entropy is calculated as

$$Q_1 = -\hat{\pi}_{1B} \ln(\hat{\pi}_{1B}) - \hat{\pi}_{1R} \ln(\hat{\pi}_{1R}) = -\frac{2}{3} \ln\left(\frac{2}{3}\right) - \frac{1}{3} \ln\left(\frac{1}{3}\right) = 0.64.$$

In region R_2 we have $n_2 = 7$ data points with the proportions $\hat{\pi}_{2B} = 3/7$ and $\hat{\pi}_{2R} = 4/7$. The entropy for this regions will be

$$Q_2 = -\hat{\pi}_{2B} \ln(\hat{\pi}_{2B}) - \hat{\pi}_{2R} \ln(\hat{\pi}_{2R}) = -\frac{3}{7} \ln\left(\frac{3}{7}\right) - \frac{4}{7} \ln\left(\frac{4}{7}\right) = 0.68$$

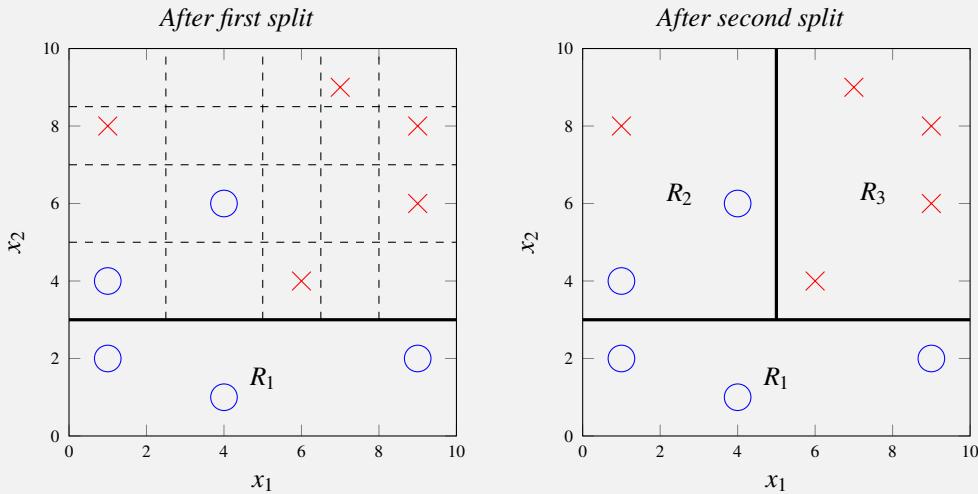
and inserted in 2.6 the total weighted entropy for this split becomes

$$n_1 Q_1 + n_2 Q_2 = 3 \cdot 0.64 + 7 \cdot 0.68 = 6.69.$$

We compute the cost for all other splits in the same manner, and summarize it in the table below.

Split (R_1)	n_1	$\hat{\pi}_{1B}$	$\hat{\pi}_{1R}$	Q_1	n_2	$\hat{\pi}_{2B}$	$\hat{\pi}_{2R}$	Q_2	$n_1 Q_1 + n_2 Q_2$
$x_1 < 2.5$	3	2/3	1/3	0.64	7	3/7	4/7	0.68	6.69
$x_1 < 5.0$	5	4/5	1/5	0.50	5	1/5	4/5	0.50	5.00
$x_1 < 6.5$	6	4/6	2/6	0.64	4	1/4	3/4	0.56	6.07
$x_1 < 8.0$	7	4/7	3/7	0.68	3	1/3	2/3	0.64	6.69
$x_2 < 1.5$	1	1/1	0/1	0.00	9	4/9	5/9	0.69	6.18
$x_2 < 3.0$	3	3/3	0/3	0.00	7	2/7	5/7	0.60	4.18
$x_2 < 5.0$	5	4/5	1/5	0.50	5	1/5	4/5	0.06	5.00
$x_2 < 7.0$	7	5/7	2/7	0.60	3	0/3	3/3	0.00	4.18
$x_2 < 8.5$	9	5/9	4/9	0.69	1	0/1	1/1	0.00	6.18

From the table we can read that the two splits at $x_2 < 3.0$ and $x_2 < 7.0$ are both equally good. We choose to continue with $x_2 < 3.0$.



Second split: We notice that only R_2 has more than five data points. Also there is no point splitting region R_1 further since it only contains data points from the same class. In the next step we therefore split the second region into two new regions R_2 and R_3 . All possible splits are displayed above to the left (dashed lines) and we compute their cost in the same manner as before.

Splits (R_1)	n_2	$\hat{\pi}_{2B}$	$\hat{\pi}_{2R}$	Q_2	n_3	$\hat{\pi}_{3B}$	$\hat{\pi}_{3R}$	Q_3	$n_2 Q_2 + n_3 Q_3$
$x_1 < 2.5$	2	1/2	1/2	0.69	5	1/5	4/5	0.50	3.89
$x_1 < 5.0$	3	2/3	1/3	0.63	4	0/4	4/4	0.00	1.91
$x_1 < 6.5$	4	2/4	2/4	0.69	3	0/3	3/3	0.00	2.77
$x_1 < 8.0$	5	2/5	3/5	0.67	2	0/2	2/2	0.00	3.37
$x_2 < 5.0$	2	1/2	1/2	0.69	5	1/5	4/5	0.50	3.88
$x_2 < 7.0$	4	2/4	2/4	0.69	3	0/3	3/3	0.00	2.77
$x_2 < 8.5$	6	2/6	4/6	0.64	1	0/1	1/1	0.00	3.82

The best split is the one at $x_1 < 5.0$ visualized above to the right. The final tree and partition were displayed in Example 2.5. None of the three regions has more than five data points. Therefore, we terminate the training.

If we want to use the tree for prediction, we predict blue if $\mathbf{x}_\star \in R_1$ or $\mathbf{x}_\star \in R_2$ since the blue training data points are in majority in each of these two regions. Similarly, we predict red if $\mathbf{x}_\star \in R_3$. This tree is also visualized in Example 2.5.

When choosing between the different splitting criteria mentioned above, the misclassification rate sounds like a reasonable choice since that is typically the criteria we want the final model to do well on⁶. However, one drawback is that it does not favor pure nodes. With pure nodes we mean nodes where most

⁶This is not always true, for example for imbalanced and asymmetric classification problems, see further in Section 4.5

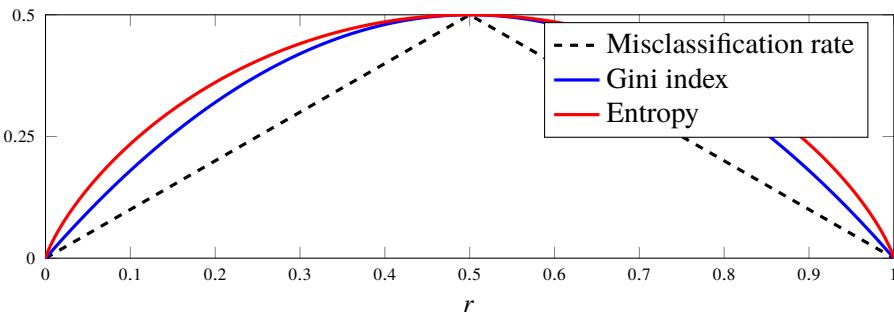


Figure 2.2: Three splitting criteria for classification trees as a function of the proportion of the first class $r = \pi_{\ell 1}$ in a certain region R_ℓ . The entropy criteria has been scaled such that it passes through $(0.5, 0.5)$.

of the data points belong to a certain class. It is usually an advantage to favor pure nodes in the greedy procedure that we use to grow the tree, since it can lead to a total of fewer splits. Both the entropy criteria and Gini index favors node purity more than misclassification rate does.

This advantage can also be illustrated in Example 2.6. Consider the first split in this example. If we would use the misclassification rate as the splitting criteria, both the split $x_2 < 5.0$ as well as the split $x_2 < 3.0$ would provide a total misclassification rate of 0.2. However, the split at $x_2 < 3.0$, which the entropy criteria favored, provides a pure node R_1 . If we would go with the split $x_2 < 5.0$ the misclassification after the second split would still be 0.2. If we would continue to grow the tree until no data points are misclassified we would need three splits if we used the entropy criteria whereas we would need five splits if we would use the misclassification criteria and started with the split at $x_2 < 5.0$.

To generalize this discussion, consider a problem with two classes where we denote the proportion of the first class as $\pi_{\ell 1} = r$ and hence the proportion of the second class as $\pi_{\ell 2} = 1 - r$, the three criteria can then in terms of r be expressed as

$$\begin{aligned}\text{Misclassification rate: } Q_\ell &= 1 - \max(r, 1 - r), \\ \text{Gini index: } Q_\ell &= 2r(1 - r), \\ \text{Entropy: } Q_\ell &= -r \ln r - (1 - r) \ln(1 - r).\end{aligned}$$

These functions are shown in Figure 2.2. All three criteria are similar in the sense that they provide zero loss if all data points belongs to either of the two classes, and maximum loss the data points are equally divided between the two classes. However, the Gini index and entropy have a higher loss for all other proportions. In other words, the gain of having a pure node (r close to 0 or 1) is higher for the Gini index and the entropy than for the misclassification rate. As a consequence, when doing a split using Gini index or the entropy, they both tend to favour (more than misclassification does) make one of the two nodes pure (or close to being pure) since that provides a smaller total loss.

How deep should a decision tree be?

The depth of a decision tree (the maximum number of splits in a certain branch) has a big impact on the final predictions. This impact is very similar the discussion of k in k -NN, and we again use the music classification Example 2.1 and car stopping distance Example 2.2 to study how the decision boundaries change depending on the depth of the trees. In Figure 2.3 the decision boundaries are illustrated for different trees. In Figure 2.3a and Figure 2.3c we have used fully grown trees In Figure 2.3b and Figure 2.3d the depth has been restricted to 4 and 3 splits in each branch, respectively.

Similar to choosing $k = 1$ in k -NN, for a fully grown tree all training data points will by construction be correctly predicted since each region only contains data points with the same output. As a result, for the music classification problem we get thin and small regions adapted to single training data points and for the car stopping distance problem we get an overly adaptive line passing exactly through the observations. Even though these trees give excellent performance on the training data, they are not likely to be the best

models for new yet unseen data. As we discussed previously in the context of k -NN, we refer to this as overfitting.

In decision trees we can mitigate overfitting by using trees that are not as deep. Consequently, we get fewer and larger regions and the decision boundaries are less adapted to the peculiarities of the training data. This is illustrated in Figure 2.3b and Figure 2.3d for the two problems. As for k in k -NN, the optimal size of the tree depends on many properties of the problem and it is a trade-off between flexibility and rigidity. Similar trade-offs have to be made in almost all models presented in this book and systematic strategies will be discussed in Chapter 4.

How can the user control the depth of the tree? Here we have different strategies. The most obvious strategy is to adjust the stopping criteria, i.e. the criteria that should be fulfilled for not proceeding with further splits in a certain node. As mentioned earlier, this criteria could be that we do not attempt further splits if there are less than a certain number of training data points the corresponding region, or as in the example above, stop splitting when we reach a certain depth. Another strategy of controlling the depth is to use *pruning*. In pruning we grow a deep and then in a second post-processing step prune it back to a smaller subtree. However, we will not discuss pruning further here.

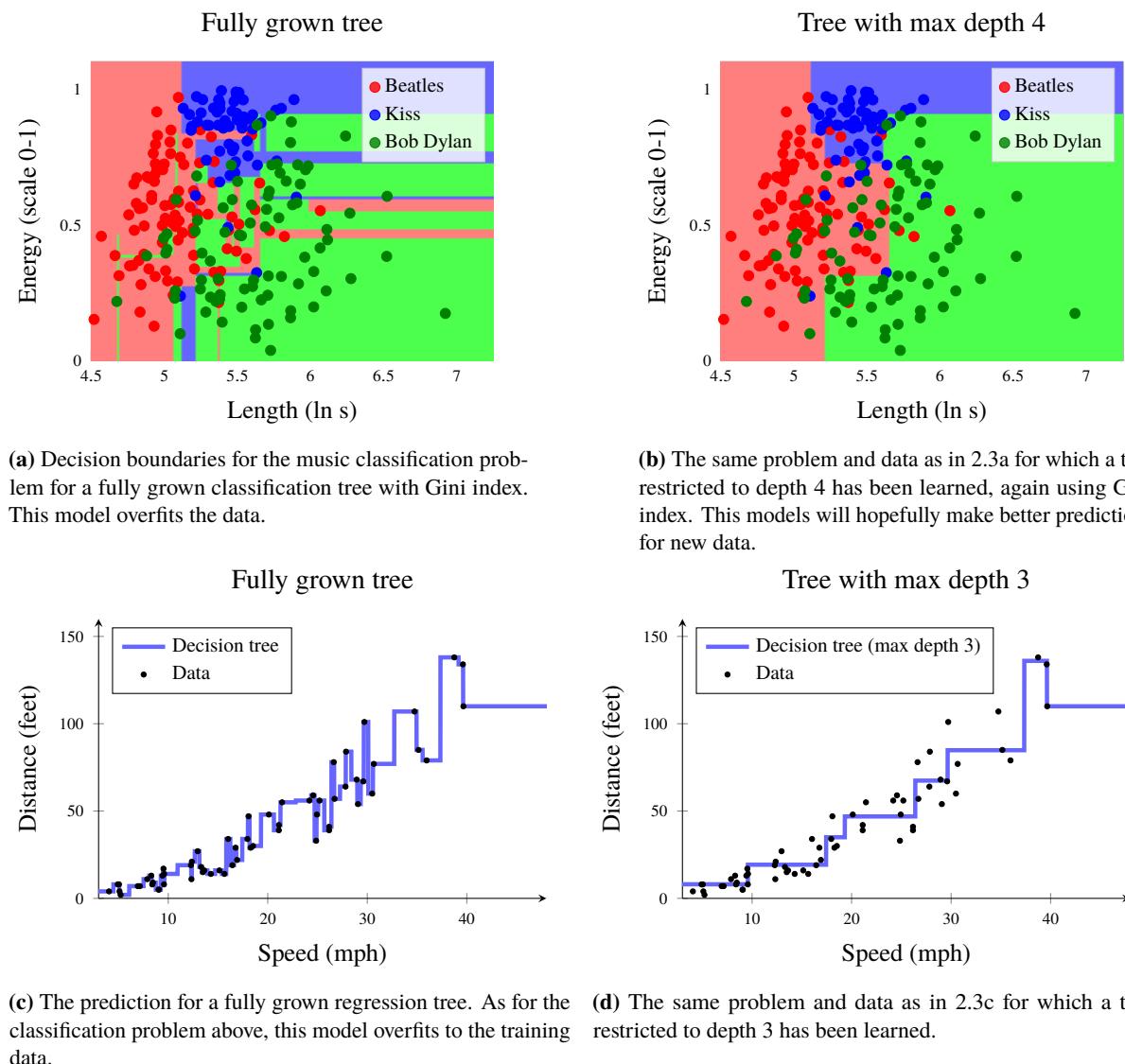


Figure 2.3: Decision trees applied to the music classification Example 2.1 (a and b) and the car stopping distance Example 2.2 (c and d).

3 Basic parametric models for regression and classification

In the previous chapter we saw the supervised machine learning problem, as well as two methods for solving it. In this chapter we introduce a systematic approach referred to as *parametric* modeling, by first looking at *linear regression* and *logistic regression* which are two such parametric models. The key point of a parametric model is that it contains some parameters θ , which are learned from training data. However, once the parameters are learned, the training data may be discarded, since the prediction only depends on θ .

3.1 Linear regression

Regression is one of the two fundamental tasks of supervised learning (the other one is classification). We will now introduce the *linear regression* model, which might (at least historically) be the most popular method for solving regression problems. Despite its relative simplicity, it is a surprisingly useful and is an important stepping stone for more advanced methods, such as deep learning, see Chapter 6.

As discussed in the previous chapter, regression amounts to learning the relationships between some input variables $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_p]^\top$ and a numerical output variable y . The inputs can be either categorical or numerical, but we will start by assuming that all p the inputs also are numerical and introduce categorical inputs later. In a more statistical framework, regression is about learning a model f

$$y = f(\mathbf{x}) + \varepsilon, \quad (3.1)$$

mapping the input to the output, where ε is an error term that describes everything about the input–output relationship that cannot be captured by the model. With a statistical perspective, we consider ε as random variable, referred to as a *noise*, that is independent of \mathbf{x} and has mean value of zero. As a running example of regression, we will use the car stopping distance regression problem introduced in the previous chapter as Example 2.2.

The linear regression model

The linear regression model assumes that the output variable y (a scalar) can be described as an affine¹ combination of the input variables x_1, x_2, \dots, x_p (each a scalar) plus a noise term ε ,

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_p x_p + \varepsilon. \quad (3.2)$$

We refer to the coefficients $\theta_0, \theta_1, \dots, \theta_p$ as the *parameters* of the model, and we sometimes refer to θ_0 specifically as the intercept (or offset) term. The noise term ε accounts for random errors in the data not captured by the model. The noise is assumed to have mean zero and to be independent of \mathbf{x} . The zero-mean assumption is nonrestrictive, since any (constant) non-zero mean can be incorporated in the offset term θ_0 .

To have a more compact notation, we introduce the parameter vector $\theta = [\theta_0 \ \theta_1 \ \dots \ \theta_p]^\top$ and extend the vector \mathbf{x} with a constant one in its first position, such that we can write the linear regression model (3.2)

¹An affine function is a linear function plus a constant offset.

compactly

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p + \varepsilon = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_p] \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{bmatrix} + \varepsilon = \boldsymbol{\theta}^\top \mathbf{x} + \varepsilon \quad (3.3)$$

This notation means that the symbol \mathbf{x} unfortunately is used both for the $p+1$ and the p -dimensional version of the input vector, with or without the constant one in the leading position. This is only a matter of book-keeping for handling the intercept term θ_0 and will be clear from context, and carries no deeper meaning.

The linear regression model is a *parametric* function of the form (3.3). The parameters $\boldsymbol{\theta}$ can take arbitrary values, and the actual values that we assign to them will control the input–output relationship described by the model. *Learning* of the model therefore amounts to finding suitable values for $\boldsymbol{\theta}$ based on observed training data. Before discussing how to do this, however, let us have a closer look at how the model can be used to make predictions once it has been learned.

Making predictions with the linear regression model

The goal in supervised machine learning is making predictions $\hat{y}(\mathbf{x}_*)$ for new, previously unseen, test inputs $\mathbf{x}_* = [1 \quad x_{*1} \quad x_{*2} \quad \dots \quad x_{*p}]^\top$. Let us assume that we have already learned some parameter values $\hat{\boldsymbol{\theta}}$ for the linear regression model (how this is done will be described next). We use the symbol $\hat{\cdot}$ to indicate that $\hat{\boldsymbol{\theta}}$ contains learned values of the otherwise unknown parameter vector $\boldsymbol{\theta}$. Since we assume that the noise term ε is random with zero mean and independent of all observed variables, it makes statistical sense to let $\varepsilon = 0$ in the prediction. That is, a prediction from the linear regression model takes the form

$$\hat{y}(\mathbf{x}_*) = \hat{\theta}_0 + \hat{\theta}_1 x_{*1} + \hat{\theta}_2 x_{*2} + \cdots + \hat{\theta}_p x_{*p} = \hat{\boldsymbol{\theta}}^\top \mathbf{x}_*. \quad (3.4)$$

This is in fact applicable for all regression models of the type (3.1), where the model used for prediction in general is $\hat{y}(\mathbf{x}_*) = f(\mathbf{x}_*)$. The noise term ε is often referred to as an *irreducible error* or an *aleatoric*² uncertainty in the prediction. We illustrate the predictions made by a linear regression model in Figure 3.1.

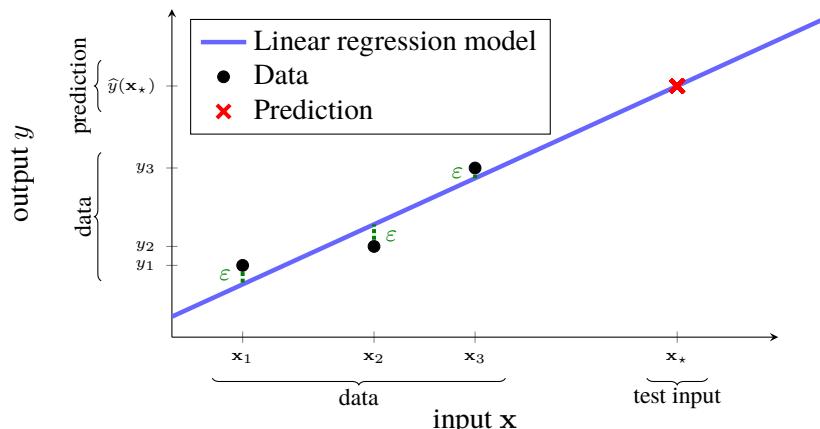


Figure 3.1: Linear regression with $p = 1$: The black dots represent $n = 3$ data samples, from which a linear regression model (blue line) is learned. The model does not fit the data perfectly, but there is a remaining error/noise ε (green). The model can be used to predict (red cross) the output $\hat{y}(\mathbf{x}_*)$ for a test input \mathbf{x}_* .

²From the Latin word *aleator*, meaning dice-player.

Learning linear regression from training data

Let us now discuss how to learn a linear regression model, that is learn θ , from training data $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$. We collect the training data, which consists of n samples \mathbf{x}_i of the input and corresponding n samples y_i of the output in the $n \times (p + 1)$ matrix \mathbf{X} and n -dimensional vector \mathbf{y} ,

$$\mathbf{X} = \begin{bmatrix} -\mathbf{x}_1^\top - \\ -\mathbf{x}_2^\top - \\ \vdots \\ -\mathbf{x}_n^\top - \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \text{where each } \mathbf{x}_i = \begin{bmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix}. \quad (3.5)$$

Example 3.1: Car stopping distances

We continue Example 2.2, and attempt to learn a liner regression model for it. We therefore form the matrices \mathbf{X} and \mathbf{y} . Since we only have one input and one output, both x_i and y_i are scalar. We get

$$\mathbf{X} = \begin{bmatrix} 1 & 4.0 \\ 1 & 4.9 \\ 1 & 5.0 \\ 1 & 5.1 \\ 1 & 5.2 \\ \vdots & \vdots \\ 1 & 39.6 \\ 1 & 39.7 \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} 4.0 \\ 8.0 \\ 8.0 \\ 4.0 \\ 2.0 \\ \vdots \\ 134.0 \\ 110.0 \end{bmatrix}. \quad (3.6)$$

Altogether we can use this vector ant matrix notation to describe the linear regression model for each and every training data point $i = 1, \dots, n$ \mathbf{X} in one equation as a matrix multiplication

$$\mathbf{y} = \mathbf{X}\theta + \epsilon, \quad (3.7)$$

where ϵ is a vector of errors/noise. Moreover we can also define a vector of predicted outputs for the training data $\hat{\mathbf{y}} = [\hat{y}(\mathbf{x}_1) \quad \hat{y}(\mathbf{x}_2) \quad \dots \quad \hat{y}(\mathbf{x}_n)]^\top$ which also allows a compact matrix formulation

$$\hat{\mathbf{y}} = \mathbf{X}\theta. \quad (3.8)$$

Note that whereas \mathbf{y} is a vector of recorded training data values, $\hat{\mathbf{y}}$ is a vector whose entries are functions of θ . Learning the unknown parameters θ amounts to finding values such that $\hat{\mathbf{y}}$ is similar to \mathbf{y} , that is, the model should fit the training data well. There are multiple ways to define what ‘similar’ or ‘well’ actually means, but it somehow amounts to find θ such that $\hat{\mathbf{y}} - \mathbf{y} = \epsilon$ is small. We will approach this by formulating a loss function, which will imply a meaning of ‘fitting the data well’. We will thereafter interpret the loss function from a statistical perspective, by understanding this as selecting the value of θ which makes the observed training data \mathbf{y} as likely as possible to have been observed—the so-called *maximum likelihood* solution. Later in Chapter 9 we will also introduce a conceptually different way to learn θ .

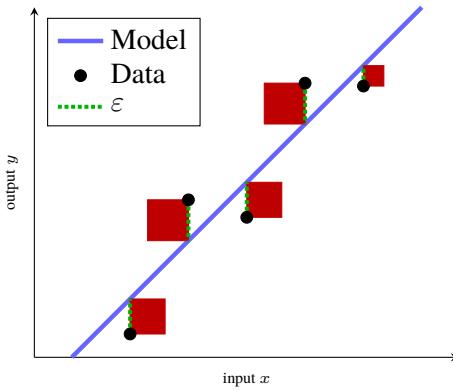


Figure 3.2: A graphical explanation of the squared error loss function: the goal is to choose the model (blue line) such that the sum of the squares (red) of each error ε (green dashed lines) is minimized. That is, the blue line is to be chosen so that the amount of red color is minimized. The black dots, the training data, are fixed. This motivates the name *least squares*.

Loss functions and cost functions

One principled way to define the learning problem is to introduce a *loss function* $L(\hat{y}, y)$ which measures how close the model's prediction \hat{y} is to the observed data y . If the model fits the data well, i.e. if $\hat{y} \approx y$, then the loss function should take a small value, and vice versa. Based on the chosen loss function we also define the *cost function* as the average loss over the training data. Learning a model then amounts to finding the parameter values that minimize the cost

$$\widehat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \underbrace{L(\hat{y}(\mathbf{x}_i; \boldsymbol{\theta}), y_i)}_{\text{cost function } J(\boldsymbol{\theta})}. \quad (3.9)$$

Note that each term in the expression above corresponds to evaluating the loss function for the prediction $\hat{y}(\mathbf{x}_i; \boldsymbol{\theta})$, given by (3.4), for the training point with index i and the true output value y_i at that point. To emphasize that the prediction depends on the parameters $\boldsymbol{\theta}$, we have included $\boldsymbol{\theta}$ as an argument to \hat{y} for clarity. The operator $\arg \min_{\boldsymbol{\theta}}$ means “the value of $\boldsymbol{\theta}$ for which the cost function attains its minimum”. The relationship between loss and cost functions (3.9) is general for all cost functions in this book.

Least squares and the normal equations

For regression, a commonly used loss function is the *squared error* loss

$$L(\hat{y}(\mathbf{x}; \boldsymbol{\theta}), y) = (\hat{y}(\mathbf{x}; \boldsymbol{\theta}) - y)^2. \quad (3.10)$$

This loss function attains 0 if $\hat{y}(\mathbf{x}; \boldsymbol{\theta}) = y$, and grows fast (quadratic) as the difference between y and the prediction $\hat{y}(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{x}$ increases. The corresponding cost function for the linear regression model (3.7) can be written with matrix multiplications as

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (\hat{y}(\mathbf{x}_i; \boldsymbol{\theta}) - y_i)^2 = \frac{1}{n} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 = \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \frac{1}{n} \|\boldsymbol{\epsilon}\|_2^2, \quad (3.11)$$

where $\|\cdot\|_2$ denotes the usual Euclidean vector norm, and $\|\cdot\|_2^2$ its square. Due to the square, this particular cost function is also commonly referred to as the *least squares* cost. It is illustrated in Figure 3.2. We will discuss other loss functions in Chapter 5.

When using the squared error loss for learning a linear regression model from \mathcal{T} , we thus need to solve the problem

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (\theta^\top \mathbf{x}_i - y_i)^2 = \arg \min_{\theta} \frac{1}{n} \|\mathbf{X}\theta - \mathbf{y}\|_2^2. \quad (3.12)$$

From a linear algebra point of view, this can be seen as the problem of finding the closest vector to \mathbf{y} (in an Euclidean sense) in the subspace of \mathbb{R}^n spanned by the columns of \mathbf{X} . The solution to this problem is the orthogonal projection of \mathbf{y} onto this subspace, and the corresponding $\hat{\theta}$ can be shown (see Section 3.A) to fulfill

$$\mathbf{X}^\top \mathbf{X} \hat{\theta} = \mathbf{X}^\top \mathbf{y}. \quad (3.13)$$

Equation (3.13) is often referred to as the *normal equations*, and gives the solution to the least squares problem (3.12). If $\mathbf{X}^\top \mathbf{X}$ is invertible, which is often the case, then $\hat{\theta}$ has the closed form expression

$$\hat{\theta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (3.14)$$

The fact that this closed-form solution exists is important, and is probably the reason for why the linear regression with squared error loss is so extremely commonly used in practice. Other loss functions lead to optimization problems that often lack closed-form solutions.

We now have everything in place for using linear regression, and we summarize by Algorithm 4 and Algorithm 5, as well as Example 3.2.

Time to reflect 3.1: What does it mean in practice if $\mathbf{X}^\top \mathbf{X}$ is not invertible?

Time to reflect 3.2: If the columns of \mathbf{X} are linearly independent and $p = n - 1$, \mathbf{X} spans the entire \mathbb{R}^n . That means a unique solution exists such that $\mathbf{y} = \mathbf{X}\theta$ exactly, i.e., the model fits the training data perfectly. If that is the case, (3.14) reduces to $\theta = \mathbf{X}^{-1} \mathbf{y}$, and the model fits the data perfectly. Why would that not be a desired property in practice?

Algorithm 4: Learning a linear regression model with squared error loss

Data: Training data $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$

Result: Learned parameter vector $\hat{\theta}$

- 1 Construct the matrix \mathbf{X} and vector \mathbf{y} according to (3.5).
 - 2 Compute $\hat{\theta}$ by solving (3.13).
-

Algorithm 5: Making predictions from a linear regression model

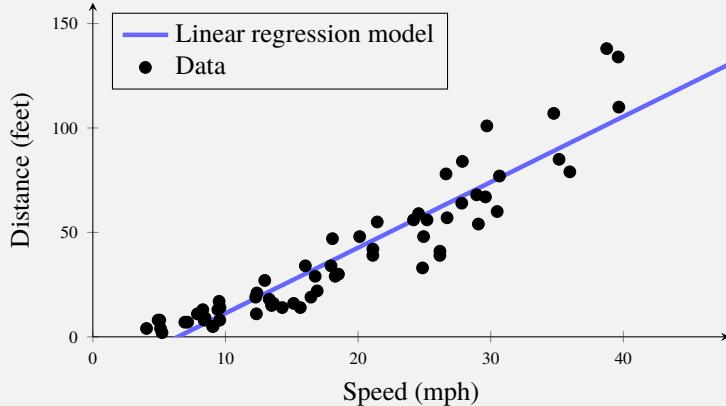
Data: Learned parameters $\hat{\theta}$ and test input \mathbf{x}_*

Result: Prediction $\hat{y}(\mathbf{x}_*)$

- 1 Compute $\hat{y}(\mathbf{x}_*) = \hat{\theta}^\top \mathbf{x}_*$.
-

Example 3.2: Car stopping distances

By inserting the matrices (3.6) from Example 3.1 into the normal equations (3.7), we obtain $\hat{\theta}_0 = -20.1$ and $\hat{\theta}_1 = 3.14$. If we plot the resulting model, it looks like this:



This can be compared to how k -NN and decision trees solves the same problem, Figure 2.1 and 2.3. Clearly the linear regression model behaves differently to k -NN and decision trees; linear regression does not share the ‘‘local’’ nature of k -NN and decision trees (only training data points close to \mathbf{x}_\star affects $\hat{y}(\mathbf{x}_\star)$), which is related to the fact that linear regression is a parametric model.

The maximum likelihood perspective

To get another perspective of the squared error loss we will now reinterpret the least squares method above as a *maximum likelihood* solution. The word ‘likelihood’ refers to the statistical concept of the likelihood function, and maximizing the likelihood function amounts to finding the value of θ that makes observing \mathbf{y} as likely as possible. That is, instead of (somewhat arbitrarily) selecting a loss function, we start with the problem

$$\hat{\theta} = \arg \max_{\theta} p(\mathbf{y} | \mathbf{X}; \theta). \quad (3.15)$$

Here $p(\mathbf{y} | \mathbf{X}; \theta)$ is the probability density of all observed outputs \mathbf{y} in the training data, given all inputs \mathbf{X} and parameters θ . This determines mathematically what ‘likely’ means, but we need to specify it in more detail. We do that by considering the noise term ε as a stochastic variable with a certain distribution. A common assumption is the noise terms are independent, each with a Gaussian distribution with zero mean and variance σ_ε^2 ,

$$\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2). \quad (3.16)$$

This implies that the n observed training data points are independent, and $p(\mathbf{y} | \mathbf{X}; \theta)$ factorizes as

$$p(\mathbf{y} | \mathbf{X}; \theta) = \prod_{i=1}^n p(y_i | \mathbf{x}_i, \theta). \quad (3.17)$$

Together with the linear regression model (3.3) we have

$$p(y_i | \mathbf{x}_i, \theta) = \mathcal{N}\left(y_i; \theta^\top \mathbf{x}_i, \sigma_\varepsilon^2\right) = \frac{1}{\sqrt{2\pi\sigma_\varepsilon^2}} \exp\left(-\frac{1}{2\sigma_\varepsilon^2} (\theta^\top \mathbf{x}_i - y_i)^2\right). \quad (3.18)$$

Recall that we want to maximize the likelihood w.r.t. θ . For numerical reasons, it is usually better to work with the logarithm of $p(\mathbf{y} | \mathbf{X}; \theta)$,

$$\ln p(\mathbf{y} | \mathbf{X}; \theta) = \sum_{i=1}^n \ln p(y_i | \mathbf{x}_i, \theta). \quad (3.19)$$

Since the logarithm is a monotonically increasing function, maximizing the log-likelihood (3.19) is equivalent to maximizing the likelihood itself. Putting (3.18) and (3.19) together, we get

$$\ln p(\mathbf{y} | \mathbf{X}; \boldsymbol{\theta}) = -\frac{n}{2} \ln(2\pi\sigma_\varepsilon^2) - \frac{1}{2\sigma_\varepsilon^2} \sum_{i=1}^n (\boldsymbol{\theta}^\top \mathbf{x}_i - y_i)^2. \quad (3.20)$$

Removing terms and factors independent of $\boldsymbol{\theta}$ does not change the maximizing argument, and we see that we can rewrite (3.15) as

$$\widehat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} p(\mathbf{y} | \mathbf{X}; \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} - \sum_{i=1}^n (\boldsymbol{\theta}^\top \mathbf{x}_i - y_i)^2 = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n \frac{1}{n} (\boldsymbol{\theta}^\top \mathbf{x}_i - y_i)^2. \quad (3.21)$$

This is indeed linear regression with the least squares cost (the cost function implied by the squared error loss function (3.10)). The key to arrive at this expression was the assumption of ε having a Gaussian distribution. Other assumptions on ε leads to other loss functions, as we will discuss more in Chapter 5.

Categorical input variables

The regression problem is characterized by a numerical output y , and inputs \mathbf{x} of arbitrary type. We have, however, only discussed the case of numerical inputs so far, but categorical inputs are perfectly possible as well.

Assume that we have a categorical input variable that only takes two different values. We refer to those two values as A and B. We can then create a *dummy variable* x as

$$x = \begin{cases} 0 & \text{if A,} \\ 1 & \text{if B,} \end{cases} \quad (3.22)$$

and use this variable in any supervised machine learning method as if it was numerical. For linear regression, this effectively gives us a model which looks like

$$y = \theta_0 + \theta_1 x + \varepsilon = \begin{cases} \theta_0 + \varepsilon & \text{if A,} \\ \theta_0 + \theta_1 + \varepsilon & \text{if B.} \end{cases} \quad (3.23)$$

The choice is somewhat arbitrary, since A and B of course can be switched. If the categorical variable takes more than two values, let us say A, B, C and D, we can make a so-called one-hot encoding by constructing a four-dimensional vector

$$\mathbf{x} = [x_A \ x_B \ x_C \ x_D]^\top \quad (3.24)$$

where $x_A = 1$ if A, $x_B = 1$ if B, and so on. That is, only one element of \mathbf{x} will be 1, the rest are 0. Again, this construction can be used for any supervised machine learning method, not only linear regression.

3.2 Classification and logistic regression

After presenting a parametric method for solving the regression problem, we now turn our attention to classification. As we will see, with a modification of the linear regression model we can apply it to the classification problem as well, however at the cost of not being able to use the convenient normal equations for learning the parameters. To perform Instead we have to resort to numerical optimization, which we will discuss later in Section 5.3.

A statistical view of the classification problem

Supervised machine learning amounts to predicting the output from the input. With a statistical perspective, classification amounts to predicting the conditional class probabilities

$$p(y = m | \mathbf{x}), \quad (3.25)$$

where y is the output (1, 2, ..., or M) and \mathbf{x} is the input.³ In words, $p(y = m | \mathbf{x})$ describes *the probability for class m given that we know the input \mathbf{x}* . Talking about $p(y | \mathbf{x})$ implies that we think about the class label y as a random variable. Why? Because we choose to model the real world, from where the data originates, as involving a certain amount of randomness (much like the random error ε in regression). Let us illustrate with an example:

Example 3.3: Describing voting behavior using probabilities

We want to construct a model that can predict voting preferences ($= y$, the categorical output) for different population groups ($= \mathbf{x}$, the input). However, we then have to face the fact that not everyone in a certain population group will vote for the same political party. We can therefore think of y as a random variable which follows a certain probability distribution. If we knew that the vote count in the group of 45 year old women ($= \mathbf{x}$) is 13% for the cerise party, 39% for the turquoise party and 48% for the purple party (here we have $M = 3$), we could describe it as

$$\begin{aligned} p(y = \text{cerise party} | \mathbf{x} = 45 \text{ year old women}) &= 0.13, \\ p(y = \text{turquoise party} | \mathbf{x} = 45 \text{ year old women}) &= 0.39, \\ p(y = \text{purple party} | \mathbf{x} = 45 \text{ year old women}) &= 0.48. \end{aligned}$$

In this way, the probabilities $p(y | \mathbf{x})$ describe the non-trivial fact that

- (a) all 45 year old women do not vote for the same party, but
- (b) the choice of party does not appear to be completely random among 45 year old women either; the purple party is the most popular, and the cerise party is the least popular.

Thus, it can be useful to have a classifier which predicts not only a class \hat{y} (one party), but a distribution over classes $p(y | \mathbf{x})$.

We now aim to construct a classifier which can not only predict classes, but also learn the class probabilities $p(y | \mathbf{x})$. More specifically, for binary classification problems ($M = 2$, and y is either 1 or -1), we learn a model $g(\mathbf{x})$ for which

$$p(y = 1 | \mathbf{x}) \text{ is modeled by } g(\mathbf{x}). \quad (3.26a)$$

By the laws of probabilities, it holds that $p(y = 1 | \mathbf{x}) + p(y = -1 | \mathbf{x}) = 1$, which gives that

$$p(y = -1 | \mathbf{x}) \text{ is modeled by } 1 - g(\mathbf{x}). \quad (3.26b)$$

Since $g(\mathbf{x})$ is a model for a probability, it is natural to require that $0 \leq g(\mathbf{x}) \leq 1$ for any \mathbf{x} . We will see how this constraint can be enforced below.

³We use the notation $p(y | \mathbf{x})$ to denote probability masses (y discrete) as well as probability densities (y continuous).

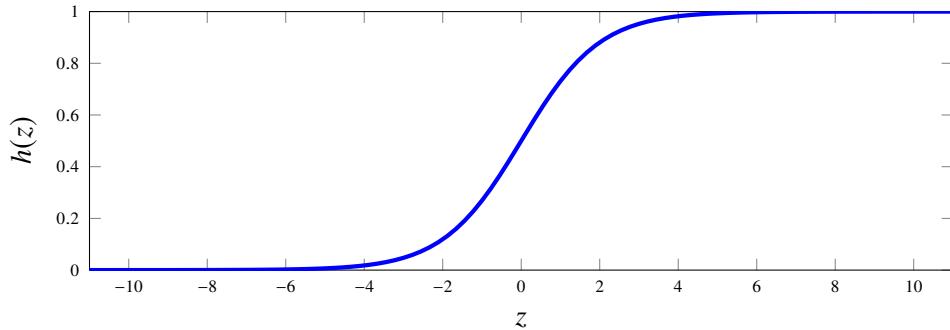


Figure 3.3: The logistic function $h(z) = \frac{e^z}{1+e^z}$.

For the multiclass problem, we instead let the classifier return a vector-valued function $\mathbf{g}(\mathbf{x})$, where

$$\begin{bmatrix} p(y = 1 | \mathbf{x}) \\ p(y = 2 | \mathbf{x}) \\ \vdots \\ p(y = M | \mathbf{x}) \end{bmatrix} \text{ is modeled by } \begin{bmatrix} g_1(\mathbf{x}) \\ g_2(\mathbf{x}) \\ \vdots \\ g_M(\mathbf{x}) \end{bmatrix} = \mathbf{g}(\mathbf{x}). \quad (3.27)$$

In words, each element $g_m(\mathbf{x})$ of $\mathbf{g}(\mathbf{x})$ corresponds to the conditional class probability $p(y = m | \mathbf{x})$. Since $\mathbf{g}(\mathbf{x})$ models a probability vector, we require that each element $g_m(\mathbf{x}) \geq 0$ and that $\|\mathbf{g}(\mathbf{x})\|_1 = \sum_{m=1}^M |g_m(\mathbf{x})| = 1$ for any \mathbf{x} .

The logistic regression model

We will now introduce the logistic regression model, which is one possible way of modeling conditional class probabilities. Logistic regression can be viewed as a modification of the linear regression model so that it fits the classification (rather than regression) problem.

Let us start with binary classification, that is, learning a function $g(\mathbf{x})$ that approximates the conditional probability of the positive class. The linear regression model, without the noise term, is given by

$$z = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p = \boldsymbol{\theta}^\top \mathbf{x}. \quad (3.28)$$

This is a mapping which takes \mathbf{x} and returns z , which in this context is called the *logit*. Note that z takes values on the entire real line, whereas we need a function which instead returns a value on the interval $[0, 1]$. The key idea of logistic regression is thus to ‘squeeze’ z from (3.28) to the interval $[0, 1]$ by using the logistic function $h(z) = \frac{e^z}{1+e^z}$, see Figure 3.3. This gives the function

$$g(\mathbf{x}) = \frac{e^{\boldsymbol{\theta}^\top \mathbf{x}}}{1 + e^{\boldsymbol{\theta}^\top \mathbf{x}}}, \quad (3.29a)$$

which is restricted to $[0, 1]$ and hence can be interpreted as a probability. (3.29a) is the logistic regression model for $p(y = 1 | \mathbf{x})$. Note that this implicitly also gives a model for $p(y = -1 | \mathbf{x})$,

$$1 - g(\mathbf{x}) = 1 - \frac{e^{\boldsymbol{\theta}^\top \mathbf{x}}}{1 + e^{\boldsymbol{\theta}^\top \mathbf{x}}} = \frac{1}{1 + e^{\boldsymbol{\theta}^\top \mathbf{x}}} = \frac{e^{-\boldsymbol{\theta}^\top \mathbf{x}}}{1 + e^{-\boldsymbol{\theta}^\top \mathbf{x}}}. \quad (3.29b)$$

These equations defines the logistic regression model, which in a nutshell is linear regression appended with the logistic function. The reason why there is no noise term ε in (3.28) when comparing to the linear regression model (3.3) is that the randomness in classification is statistically modeled by the class probability construction $p(y = m | \mathbf{x})$ instead of an additive noise ε .

As for linear regression, we have a model (3.29) which contains unknown parameters $\boldsymbol{\theta}$. Logistic regression is thereby also a parametric model, and we learn the parameters from training data.

Remark 3.1 Despite its name, logistic regression is a method for classification, not regression! The (somewhat confusing) name is due to historical reasons.

Learning the logistic regression model by maximum likelihood

By using the logistic function, we have transformed linear regression (a model for regression problems) into logistic regression (a model for classification problems). As it will turn out, the price to pay is that we will not be able to use the convenient normal equations for learning θ in logistic regression (as we could for linear regression if we used the squared error loss).

In order to derive a principled way of learning θ in (3.29) from training data $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$, we start with the maximum likelihood approach. From a maximum likelihood perspective, learning a classifier amounts to solving

$$\hat{\theta} = \arg \max_{\theta} p(\mathbf{y} | \mathbf{X}; \theta) = \arg \max_{\theta} \sum_{i=1}^n \ln p(y_i | \mathbf{x}_i; \theta), \quad (3.30)$$

where we, similarly to linear regression (3.19), assume that the training data points are independent and we consider the logarithm of the likelihood function for numerical reasons. We have also added θ explicitly to the notation to emphasize the dependence on the model parameters. Remember that our model of $p(y = 1 | \mathbf{x}; \theta)$ is $g(\mathbf{x}; \theta)$, which gives

$$\ln p(y_i | \mathbf{x}_i; \theta) = \begin{cases} \ln g(\mathbf{x}_i; \theta) & \text{if } y_i = 1, \\ \ln(1 - g(\mathbf{x}_i; \theta)) & \text{if } y_i = -1. \end{cases} \quad (3.31)$$

It is common to turn the maximization problem (3.30) into an equivalent minimization problem by using the negative log-likelihood as cost function, $J(\theta) = -\frac{1}{n} \sum \ln p(y_i | \mathbf{x}_i; \theta)$, that is

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{\begin{cases} -\ln g(\mathbf{x}_i; \theta) & \text{if } y_i = 1, \\ -\ln(1 - g(\mathbf{x}_i; \theta)) & \text{if } y_i = -1. \end{cases}}_{\text{Binary cross-entropy loss } L(g(\mathbf{x}_i; \theta), y_i)} \quad (3.32)$$

The loss function in the expression above is called the *cross-entropy loss*. It is not specific to logistic regression, but can be used for any binary classifier that predicts class probabilities $g(\mathbf{x}; \theta)$.

Considering specifically the logistic regression model we can write out the cost function (3.32) in more detail. In doing so, the particular choice of labeling $\{-1, 1\}$ turns out to be convenient since we for $y_i = 1$ we can write

$$g(\mathbf{x}_i; \theta) = \frac{e^{\theta^\top \mathbf{x}_i}}{1 + e^{\theta^\top \mathbf{x}_i}} = \frac{e^{y_i \theta^\top \mathbf{x}_i}}{1 + e^{y_i \theta^\top \mathbf{x}_i}}, \quad (3.33a)$$

and for $y_i = -1$

$$1 - g(\mathbf{x}_i; \theta) = \frac{e^{-\theta^\top \mathbf{x}_i}}{1 + e^{-\theta^\top \mathbf{x}_i}} = \frac{e^{y_i \theta^\top \mathbf{x}_i}}{1 + e^{y_i \theta^\top \mathbf{x}_i}}. \quad (3.33b)$$

Since we get the same expression in both cases, we write (3.32) compactly as

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n -\ln \frac{e^{y_i \theta^\top \mathbf{x}_i}}{1 + e^{y_i \theta^\top \mathbf{x}_i}} = \frac{1}{n} \sum_{i=1}^n -\ln \frac{1}{1 + e^{-y_i \theta^\top \mathbf{x}_i}} = \frac{1}{n} \sum_{i=1}^n \underbrace{\ln \left(1 + e^{-y_i \theta^\top \mathbf{x}_i}\right)}_{\text{Logistic loss } L(\mathbf{x}_i, y_i, \theta)}. \quad (3.34)$$

The loss function $L(\mathbf{x}, y_i, \theta)$ above, which is the cross-entropy loss specialized to logistic regression, is called the *logistic loss* (or sometimes binomial deviance). Learning a logistic regression model thus amounts to solving

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \ln \left(1 + e^{-y_i \theta^\top \mathbf{x}_i}\right). \quad (3.35)$$

Contrary to linear regression with squared error loss, (3.35) has no closed-form solution, so we have to use numerical optimization instead. We will come back to that topic in Section 5.3.

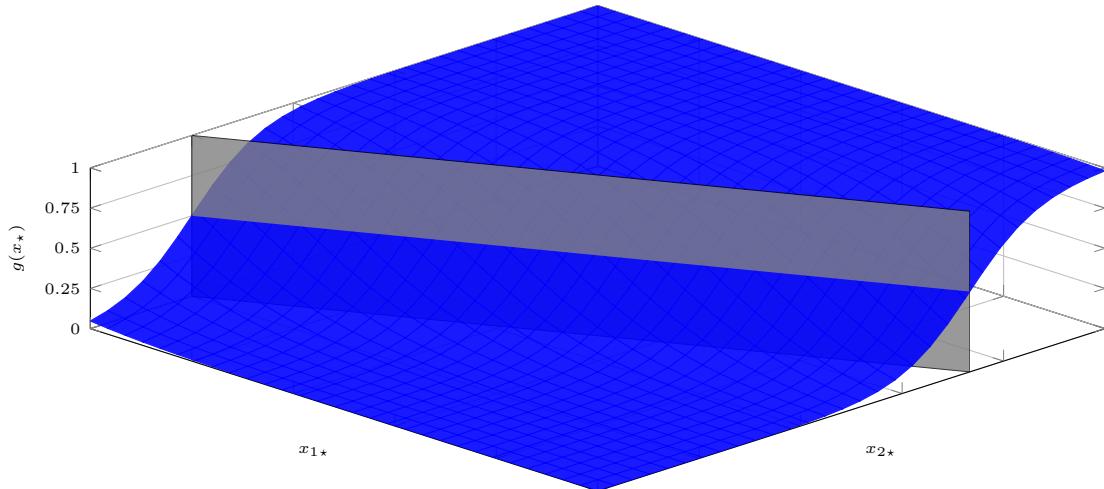


Figure 3.4: In binary classification ($y = -1$ or 1) logistic regression predicts $g(\mathbf{x}_\star)$ (\mathbf{x} is two-dimensional here), which is an attempt to determine $p(y = 1 | \mathbf{x}_\star)$. This implicitly also gives a prediction for $p(y = -1 | \mathbf{x}_\star)$ as $1 - g(\mathbf{x}_\star)$. To turn these probabilities into actual class predictions ($\hat{y}(\mathbf{x}_\star)$ is either -1 or 1), the class which is modeled to have the highest probability can be taken as the prediction, as in Equation (3.36). The point(s) where the prediction changes from one class to another is the decision boundary (gray plane).

Predictions and decision boundaries

So far, we have discussed logistic regression as a method for predicting class probabilities for a test input \mathbf{x}_\star by first learning $\boldsymbol{\theta}$ from training data and thereafter computing $g(\mathbf{x}_\star)$ (our model for $p(y = 1 | \mathbf{x}_\star)$ in binary classification) or $\mathbf{g}(\mathbf{x}_\star)$ (our model for $p(y = m | \mathbf{x}_\star)$ in multiclass classification). However, sometimes we want to make a ‘‘hard’’ prediction for the test input \mathbf{x}_\star , that is, predicting $\hat{y}(\mathbf{x}_\star) = -1$ or $\hat{y}(\mathbf{x}_\star) = 1$ in binary classification, just like with k -NN or decision trees. We then have to append logistic regression with a final step in which the predicted probabilities are turned into a class prediction. The most common approach is to *let the most probable class be $\hat{y}(\mathbf{x}_\star)$* . For binary classification, we can express this⁴ as

$$\hat{y}(\mathbf{x}_\star) = \begin{cases} 1 & \text{if } g(\mathbf{x}) > r \\ -1 & \text{if } g(\mathbf{x}) \leq r \end{cases}, \quad (3.36)$$

with decision threshold $r = 0.5$, which is illustrated in Figure 3.4. With binary logistic regression, this is equivalent to

$$\hat{y}(\mathbf{x}_\star) = \text{sign}(\boldsymbol{\theta}^\top \mathbf{x}_\star). \quad (3.37)$$

We now have everything in place for summarizing binary logistic regression in Algorithm 6.

Algorithm 6: Logistic regression for binary classification

Data: Training data $\{\mathbf{x}_i, y_i\}_{i=1}^n$ (with output classes $y = \{-1, 1\}$) and a test input \mathbf{x}_\star
Result: Predicted test output $\hat{y}(\mathbf{x}_\star)$

Learn

1 Compute $\hat{\boldsymbol{\theta}}$ by solving (3.35) numerically.

Predict

2 Compute $g(\mathbf{x}_\star)$ (3.29a).

3 If $g(\mathbf{x}_\star) > 0.5$, return $\hat{y}(\mathbf{x}_\star) = 1$, otherwise return $\hat{y}(\mathbf{x}_\star) = -1$.

In some applications, however, it can be beneficial to explore different thresholds than $r = 0.5$. It can be shown that if $g(\mathbf{x}) = p(y = 1 | \mathbf{x})$, that is, the model provides a correct description of the real-world

⁴It is arbitrary what happens if $g(\mathbf{x}) = 0.5$.

class probabilities, then the choice $r = 0.5$ will give the smallest possible number of misclassification on average. In other words, $r = 0.5$ minimizes the so-called *misclassification rate*. The misclassification rate is, however, not always the most important aspect of classifier. Many classification problems are asymmetric (the classes are of different importance to correctly predict) or imbalanced (the classes occur with very different frequency). In a medical diagnosis application, for example, it can be of higher importance not to falsely predict the negative class (that is, by mistake predict a sick patient being healthy) than falsely predict the positive class (by mistake predict a healthy patient as sick). For such a problem, minimizing the misclassification rate might lead not lead to the desired performance. Furthermore, the medical diagnosis problem could be imbalanced if the disorder is very rare, meaning that the vast majority of the data samples (patients) belong to the negative class. By only considering the misclassification rate in such a situation we implicitly value accurate predictions of the negative class higher than accurate predictions of the positive class, simply because the negative class is more common in the data. We will discuss how we can evaluate such situations more systematically in Section 4.5. In the end, however, the decision threshold r is a choice that the user has to make.

For the multiclass problem, taking the prediction as the most probable class amounts to solving

$$\hat{y}(\mathbf{x}_*) = \arg \max_m g_m(\mathbf{x}_*). \quad (3.38)$$

As for the binary case, it is possible to modify this when working with an asymmetric or imbalanced problem.

The decision boundary for binary classification can be computed by solving the equation

$$g(\mathbf{x}) = 1 - g(\mathbf{x}). \quad (3.39)$$

The solutions to this equation are points in the input space for which the two classes are predicted to be equally probable. These points therefore lie on the decision boundary. For binary logistic regression, this means

$$\frac{e^{\theta^\top \mathbf{x}}}{1 + e^{\theta^\top \mathbf{x}}} = \frac{1}{1 + e^{\theta^\top \mathbf{x}}} \Leftrightarrow e^{\theta^\top \mathbf{x}} = 1 \Leftrightarrow \theta^\top \mathbf{x} = 0. \quad (3.40)$$

The equation $\theta^\top \mathbf{x} = 0$ parameterizes a (linear) hyperplane. Hence, the decision boundaries in logistic regression always have the shape of a (linear) hyperplane. The same argument can be generalized to multiclass logistic regression, but the decision boundaries will then be given by a combination of $M - 1$ hyperplanes.

In general we distinguish between different types of classifiers by the shape of their decision boundaries.

A classifier whose decision boundaries are linear hyperplanes is a *linear classifier*.

All other classifiers are *nonlinear classifiers*. Logistic regression is an example of a linear classifier, whereas k -NN and decision trees are nonlinear classifiers.

Note that the term “linear” is used in a different sense for linear regression; linear regression is a model which is linear in its parameters, whereas a linear classifier is a model whose decision boundaries are linear.

Logistic regression for more than two classes

Logistic regression can be used also for the multiclass problem when there are more than two classes, $M > 2$. There are several ways of generalizing logistic regression to the multiclass problem. We will follow one path using the so-called softmax function which will be useful also later when introducing deep learning models in Chapter 6.

For the binary problem, we used the logistic function to design a model for $g(\mathbf{x})$ (a scalar-valued function representing $p(y = 1 | \mathbf{x})$). For the multiclass problem we instead have to design a vector-valued function $\mathbf{g}(\mathbf{x})$, whose elements should be non-negative and sum to one. For this purpose, we first use M instances of (3.28), each denoted z_m and each with a different set of parameters θ_m . We stack all z_m into a vector of logits $\mathbf{z} = [z_1 \ z_2 \ \dots \ z_M]^\top$ and use the *softmax function* as a vector-valued generalization of the logistic function,

$$\text{softmax}(\mathbf{z}) \triangleq \frac{1}{\sum_{m=1}^M e^{z_m}} \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ \vdots \\ e^{z_M} \end{bmatrix}. \quad (3.41)$$

Note that the argument \mathbf{z} to the softmax function is an M -dimensional vector, and that it also returns a vector of the same dimension. By construction, the output vector from the softmax function always sums to 1, and each element is always ≥ 0 . Similarly to how we combined linear regression and the logistic function for the binary classification problem (3.29), we have now combined linear regression and the softmax function to model the class probabilities,

$$\mathbf{g}(\mathbf{x}) = \text{softmax}(\mathbf{z}), \quad \text{where } \mathbf{z} = \begin{bmatrix} \theta_1^\top \mathbf{x}_i \\ \theta_2^\top \mathbf{x}_i \\ \vdots \\ \theta_M^\top \mathbf{x}_i \end{bmatrix}. \quad (3.42)$$

Equivalently, we can write out the individual class probabilities, i.e. the elements of the vector $\mathbf{g}(\mathbf{x})$, as

$$g_m(\mathbf{x}) = \frac{e^{\theta_m^\top \mathbf{x}_i}}{\sum_{j=1}^M e^{\theta_j^\top \mathbf{x}_i}}, \quad m = 1, \dots, M. \quad (3.43)$$

This is the multiclass logistic regression model. Note that this construction uses M parameter vectors $\theta_1, \dots, \theta_M$ (one for each class), meaning that the number of parameters to learn grows with M . As for binary logistic regression, we can learn those parameters using the maximum likelihood method. We use θ to denote *all* model parameters, $\theta = \{\theta_1, \dots, \theta_M\}$. Since $g_m(\mathbf{x}; \theta)$ is our model for $p(y_i = m | \mathbf{x}_i)$, the cost function for the cross-entropy (or negative log-likelihood) loss for the multiclass problem is

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{-\ln g_{y_i}(\mathbf{x}_i; \theta)}_{\substack{\text{Multi-class cross-entropy} \\ \text{loss } L(\mathbf{g}(\mathbf{x}_i; \theta), y_i)}}. \quad (3.44)$$

Note that we use the training data labels y_i as index variables to select the correct conditional probability for the loss function. That is, the i th term of the sum is the negative logarithm of the y_i th element of the vector $\mathbf{g}(\mathbf{x}_i; \theta)$. We illustrate the meaning of this further in example 3.4.

Inserting the model (3.43) into the loss function (3.44) gives the cost function to optimize when learning multiclass logistic regression (the softmax version),

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \left(-\theta_{y_i}^\top \mathbf{x}_i + \ln \sum_{j=1}^M e^{\theta_j^\top \mathbf{x}_i} \right). \quad (3.45)$$

Example 3.4: The cross-entropy loss for multiclass problems

Consider the following (very small) data set with $n = 6$ data samples, $p = 2$ input dimensions and $M = 3$ classes, which we want to use for learning a multiclass classifier:

$$\mathbf{X} = \begin{bmatrix} 0.20 & 0.86 \\ 0.41 & 0.18 \\ 0.96 & -1.84 \\ -0.25 & 1.57 \\ -0.82 & -1.53 \\ -0.31 & 0.58 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 2 \\ 3 \\ 1 \\ 2 \\ 1 \\ 3 \end{bmatrix}.$$

Multiclass logistic regression with softmax, or any other multiclass classifier which predicts conditional class probabilities, return a 3-dimensional probability vector $\mathbf{g}(\mathbf{x}; \theta)$ for any \mathbf{x} and θ . If we stack the logarithms of the transpose of all vectors $\mathbf{g}(\mathbf{x}_i; \theta)$ for $i = 1, \dots, 6$, we obtain the matrix

$$\mathbf{G} = \begin{bmatrix} \ln g_1(\mathbf{x}_1; \theta) & \ln g_2(\mathbf{x}_1; \theta) & \ln g_3(\mathbf{x}_1; \theta) \\ \ln g_1(\mathbf{x}_2; \theta) & \ln g_2(\mathbf{x}_2; \theta) & \ln g_3(\mathbf{x}_2; \theta) \\ \ln g_1(\mathbf{x}_3; \theta) & \ln g_2(\mathbf{x}_3; \theta) & \ln g_3(\mathbf{x}_3; \theta) \\ \ln g_1(\mathbf{x}_4; \theta) & \ln g_2(\mathbf{x}_4; \theta) & \ln g_3(\mathbf{x}_4; \theta) \\ \ln g_1(\mathbf{x}_5; \theta) & \ln g_2(\mathbf{x}_5; \theta) & \ln g_3(\mathbf{x}_5; \theta) \\ \ln g_1(\mathbf{x}_6; \theta) & \ln g_2(\mathbf{x}_6; \theta) & \ln g_3(\mathbf{x}_6; \theta) \end{bmatrix}.$$

Computing the multi-class cross-entropy cost (3.44) now simply amounts to compute the average of all circled elements, and multiply that with -1 . The element that we have circled in row i is given by the training label y_i . Training the model amounts to finding θ such that this average is maximized.

Time to reflect 3.3: Can you derive (3.32) as a special case of (3.44)?

Hint: think of the binary case as a special case of the multiclass case with $\mathbf{g}(\mathbf{x}) = \begin{bmatrix} g(\mathbf{x}) \\ 1 - g(\mathbf{x}) \end{bmatrix}$.

Time to reflect 3.4: The softmax-based logistics regression is actually over-parameterized, in the sense that we can construct an equivalent model with fewer parameters. That is often not a problem in practice, but compare the multiclass model (3.42) for the case $M = 2$ with binary logistic regression (3.29) and see if you can spot the over-parametrization!

3.3 Polynomial regression and regularization

In comparison to k -NN and decision trees in Chapter 2, linear and logistic regression might appear to be rigid and non-flexible models with their straight lines (such as Figure 3.1 and 3.4). However, both models are able to adapt to the training data well if the input dimension p is large, or the number of data samples n is small.

A common way of increasing the input dimension in linear and logistic regression, which we will discuss more thoroughly in Chapter 8, is to make a nonlinear transformation of the input. A simple nonlinear transformation is to replace a one-dimensional input x with itself raised to different powers, which makes the linear regression model a polynomial

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \dots + \varepsilon. \quad (3.46)$$

This is called *polynomial regression*, and can be applied also to logistic regression. Note that if we let $x_1 = x$, $x_2 = x^2$ and $x_3 = x^3$, this is still a linear model (3.2) with input $\mathbf{x} = [1 \ x \ x^2 \ x^3]$, but we have ‘lifted’ the input from being one-dimensional ($p = 1$) to three-dimensional ($p = 3$). Using nonlinear input transformations can be very useful in practice, but it effectively increases p and we may easily end up *overfitting* the model to the noise—rather than the interesting patterns—in the training data, as in the example below.

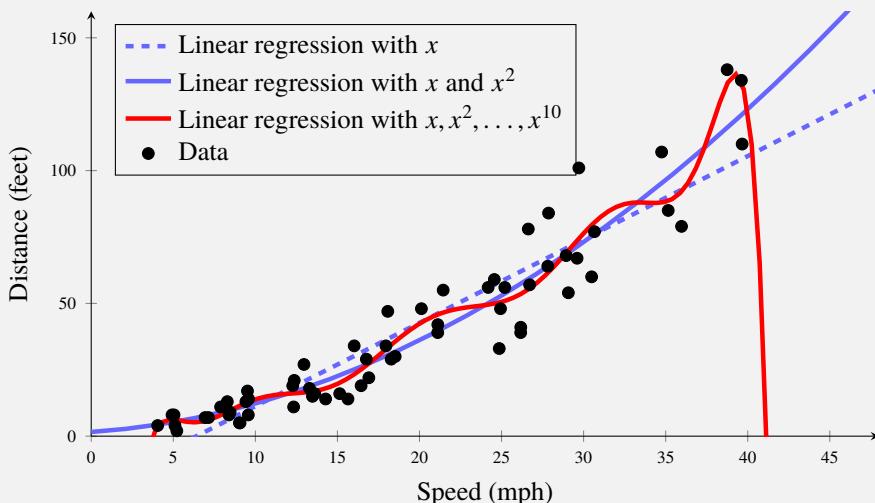
Example 3.5: Car stopping distances with polynomial regression

We return to Example 2.2, but this time we also add the squared speed as an input and thereby use a 2nd order polynomial in linear regression. This gives the new matrices (instead of Example 3.1)

$$\mathbf{X} = \begin{bmatrix} 1 & 4.0 & 16.0 \\ 1 & 4.9 & 24.0 \\ 1 & 5.0 & 25.0 \\ \vdots & \vdots & \vdots \\ 1 & 39.6 & 1568.2 \\ 1 & 39.7 & 1576.1 \end{bmatrix}, \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 4.0 \\ 8.0 \\ 8.0 \\ \vdots \\ 134.0 \\ 110.0 \end{bmatrix}, \quad (3.47)$$

and when we insert them into the normal equations (3.13), the new parameter estimates are $\hat{\theta}_0 = 1.58$, $\hat{\theta}_1 = 0.42$ and $\hat{\theta}_2 = 0.07$. (Note that also $\hat{\theta}_0$ and $\hat{\theta}_1$ change, compared to Example 3.2.)

In a completely analogous way we also learn a 10th order polynomial, and we illustrate them all below.



The second order polynomial (blue solid line) appears sensible, and using a 2nd order polynomial seems to give some advantage compared to plain linear regression (blue dashed line, from Example 3.2). However, using a 10th order polynomial (red solid line) seems to make the model less useful than even plain linear regression, and instead suffering from overfit. In conclusion it seems that there is some merit to the idea of polynomial regression, but it has to be applied carefully.

For linear (or logistic) regression with nonlinear input transformations, one could indeed mitigate overfit issues by very carefully selecting the nonlinear transformations. Such a careful selection can be quite hard to perform in practice, and overfit can occur also in other situations when p is large compared to n .

A useful approach in practice is therefore to use *regularization*. The idea of regularization can be described as ‘keeping the parameters $\hat{\theta}$ small unless the data really convinces us otherwise’, or alternatively ‘if a model with small parameter $\hat{\theta}$ values fits the data almost as well as a model with larger parameter values, the one with small parameter values should be preferred’. There are several ways to implement this idea mathematically, which leads to different regularization methods. We will give a more complete treatment of this in Section 5.2, and only discussed the so-called L^2 regularization for now. When paired with regularization, the idea of using nonlinear input transformations can be very powerful and enables a whole family of supervised machine learning methods that we will properly introduce and discuss in

Chapter 8.

To keep $\hat{\theta}$ small (in order to prevent overfit) an extra penalty term $\lambda \|\theta\|_2^2$ is added to the cost function when using L^2 regularization. The purpose of the penalty term is to prevent overfit, whereas the original cost function only rewards the fit to training data (which does not prevent overfit). The regularization parameter $\lambda \geq 0$ is therefore important to select wisely, in order to obtain the right amount of regularization. With $\lambda = 0$ the regularization has no effect, whereas $\lambda \rightarrow \infty$ will force all parameters $\hat{\theta}$ to 0. A common solution is to use cross-validation (Chapter 4) to select λ . Applied to linear regression with square error lost (3.12), it becomes⁵

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \|\mathbf{X}\theta - \mathbf{y}\|_2^2 + \lambda \|\theta\|_2^2. \quad (3.48)$$

It turns out that (3.48) has a closed-form solution, as a modified version of the normal equations, namely

$$(\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}_{p+1}) \hat{\theta} = \mathbf{X}^\top \mathbf{y}, \quad (3.49)$$

where \mathbf{I}_{p+1} is the identity matrix of size $(p+1) \times (p+1)$.

Regularization is not restricted to linear regression. The very same L^2 regularization idea can be applied to any method that involves optimizing a cost function, such as logistic regression

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n \ln \left(1 + \exp \left(-y_i \theta^\top \mathbf{x}_i \right) \right) + \lambda \|\theta\|_2^2. \quad (3.50)$$

It is common in practice to train logistic regression using (3.50) instead of (3.29). One reason is indeed to decrease possible overfit issues. Another reason is that for the non-regularized cost function (3.29), the optimal $\hat{\theta}$ is not finite if the training data is linearly separable (meaning there exists a linear decision boundary which separates the classes perfectly). In practice it means that the logistic regression learning diverges with some datasets, unless (3.50) (with $\lambda > 0$) is used instead of (3.29).

3.4 Nonlinear regression and generalized linear models

In this chapter we have so far introduced two basic parametric models for regression and classification, linear regression and logistic regression, respectively. The concept of parametric modeling is however much more general. Before leaving this chapter we will briefly discuss how the models introduced above can be generalized to describe more intricate input–output relationships. We will also dig deeper into a special case of this later in Chapter 6, where we discuss neural networks.

Generalized linear models

To be written.

Nonlinear parametric functions

Let us remind ourselves about the general regression model (3.1),

$$y = f_\theta(\mathbf{x}) + \varepsilon. \quad (3.51)$$

We have introduced an explicit dependence on the parameters θ in the notation to emphasize that f_θ is a model of the input–output relationship which depends on some parameters θ . To turn this model into a linear regression, that could be trained using least squares with a closed form solution, we made two assumptions in Section 3.1. First, the function f_θ was assumed to be linear in the model parameters, $f_\theta(\mathbf{x}) = \theta^\top \mathbf{x}$. Second, the noise term ε was assumed to be Gaussian, $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$. The latter assumption

⁵In practice, it can be wise to exclude θ_0 , the intercept, from the regularization.

is sometimes implicit, but as we saw above it makes the maximum likelihood formulation equivalent to least squares.

Both of these assumptions can be replaced with other assumptions. Replacing the Gaussian assumption will be discussed in Section 5.1, and we will here discuss what happens if we allow the function f_θ to be some arbitrary nonlinear function. We still want to learn the model from data, so f_θ has to depend on some parameters θ which we can fit to the training data, no matter if it is a linear function or not. Let us give an example of a nonlinear f_θ from biochemistry.

Example 3.6: Michaelis–Menten kinetics

An example of a relatively simple nonlinear parametric function is the Michaelis–Menten equation for modeling enzyme kinetics. The model is given by

$$y = \underbrace{\frac{\theta_1 x}{\theta_2 + x}}_{=f_\theta(x)} + \varepsilon$$

where y corresponds to a reaction rate and x a substrate concentration. The model is parameterized by the maximum reaction rate $\theta_1 > 0$ and the so called Michaelis constant of the enzyme $\theta_2 > 0$. Just like the parameters in linear regression, these parameters can be learned from input-output data $\{x_i, y_i\}_{i=1}^n$ by formulating a cost function and solving the resulting optimization problem.

Within biochemistry this model is often expressed as a deterministic relationship without the noise term ε , but we include it as an error term for consistency with our framework.

In the example above the parameters θ_1 and θ_2 have physical interpretations and are restricted to be positive. However, in machine learning we most often lack such physical interpretations of the parameters. The way we use a model is more of a “black box” which is adapted to fit the training data as well as possible. The archetype of such nonlinear black-box models are neural networks, which we will discuss in more detail in Chapter 6.

Nonlinear classification models can be constructed in a very similar way, as a generalization of the logistic regression model (3.42). In multiclass logistic regression we compute a vector of logits $\mathbf{z} = [z_1 z_2 \dots z_M]^\top$ where each element of the vector is given by a linear model $z_m = \theta_m^\top \mathbf{x}$. The class probabilities are then obtained by propagating the logit vector through the softmax function as $g_\theta(\mathbf{x}) = \text{softmax}(\mathbf{z})$. To turn this into a nonlinear classification model we can replace the expression for the logit vector with

$$\mathbf{z} = f_\theta(\mathbf{x})$$

where f_θ is some arbitrary function that maps \mathbf{x} to an M -dimensional real-valued vector. Propagating this logit vector through the softmax function, analogously to the logistic regression case, results in a nonlinear model for the conditional class probabilities, $g_\theta(\mathbf{x}) = \text{softmax}(f_\theta(\mathbf{x}))$. We will return to nonlinear classification models of this form in Chapter 6, where we use neural networks to construct the function f_θ .

3.A Derivation of the normal equations

The normal equations (3.13)

$$\mathbf{X}^\top \mathbf{X} \hat{\theta} = \mathbf{X}^\top \mathbf{y}.$$

can be derived from (3.12) (the scaling $\frac{1}{n}$ does not affect the minimizing argument)

$$\hat{\theta} = \arg \min_{\theta} \|\mathbf{X}\theta - \mathbf{y}\|_2^2,$$

in different ways. We will present one based on (matrix) calculus and one based on geometry and linear algebra.

No matter how (3.13) is derived, if $\mathbf{X}^\top \mathbf{X}$ is invertible, it (uniquely) gives

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y},$$

If $\mathbf{X}^\top \mathbf{X}$ is not invertible, then (3.13) has infinitely many solutions $\hat{\boldsymbol{\theta}}$, which all are equally good solutions to the problem (3.12).

A calculus approach

Let

$$V(\boldsymbol{\theta}) = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta}, \quad (3.52)$$

and differentiate $V(\boldsymbol{\theta})$ with respect to the vector $\boldsymbol{\theta}$,

$$\frac{\partial}{\partial \boldsymbol{\theta}} V(\boldsymbol{\theta}) = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\theta}. \quad (3.53)$$

Since $V(\boldsymbol{\theta})$ is a positive quadratic form, its minimum must be attained at $\frac{\partial}{\partial \boldsymbol{\theta}} V(\boldsymbol{\theta}) = 0$, which characterizes the solution $\hat{\boldsymbol{\theta}}$ as

$$\frac{\partial}{\partial \boldsymbol{\theta}} V(\hat{\boldsymbol{\theta}}) = 0 \Leftrightarrow -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} = 0 \Leftrightarrow \mathbf{X}^\top \mathbf{X}\hat{\boldsymbol{\theta}} = \mathbf{X}^\top \mathbf{y}, \quad (3.54)$$

i.e., the normal equations.

A linear algebra approach

Denote the $p+1$ columns of \mathbf{X} as $c_j, j = 1, \dots, p+1$. We first show that $\|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2$ is minimized if $\boldsymbol{\theta}$ is chosen such that $\mathbf{X}\boldsymbol{\theta}$ is the orthogonal projection of \mathbf{y} onto the (sub)space spanned by the columns c_j of \mathbf{X} , and then show that the orthogonal projection is found by the normal equations.

Let us decompose \mathbf{y} as $\mathbf{y}_\perp + \mathbf{y}_\parallel$, where \mathbf{y}_\perp is orthogonal to the (sub)space spanned by all columns c_i , and \mathbf{y}_\parallel is in the (sub)space spanned by all columns c_i . Since \mathbf{y}_\perp is orthogonal to both \mathbf{y}_\parallel and $\mathbf{X}\boldsymbol{\theta}$, it follows that

$$\|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \|\mathbf{X}\boldsymbol{\theta} - (\mathbf{y}_\perp + \mathbf{y}_\parallel)\|_2^2 = \|(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}_\parallel) - \mathbf{y}_\perp\|_2^2 \geq \|\mathbf{y}_\perp\|_2^2, \quad (3.55)$$

and the triangle inequality also gives us

$$\|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}_\perp - \mathbf{y}_\parallel\|_2^2 \leq \|\mathbf{y}_\perp\|_2^2 + \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}_\parallel\|_2^2. \quad (3.56)$$

This implies that if we choose $\boldsymbol{\theta}$ such that $\mathbf{X}\boldsymbol{\theta} = \mathbf{y}_\parallel$, the criterion $\|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2$ must have reached its minimum. Thus, our solution $\hat{\boldsymbol{\theta}}$ must be such that $\mathbf{X}\hat{\boldsymbol{\theta}} - \mathbf{y}$ is orthogonal to the (sub)space spanned by all columns c_i , i.e.,

$$(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}})^\top c_j = 0, j = 1, \dots, p+1 \quad (3.57)$$

(remember that two vectors \mathbf{u}, \mathbf{v} are, by definition, orthogonal if their scalar product, $\mathbf{u}^\top \mathbf{v}$, is 0.) Since the columns c_j together form the matrix \mathbf{X} , we can write this compactly as

$$(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}})^\top \mathbf{X} = 0, \quad (3.58)$$

where the right hand side is the $p+1$ -dimensional zero vector. This can equivalently be written as

$$\mathbf{X}^\top \mathbf{X}\hat{\boldsymbol{\theta}} = \mathbf{X}^\top \mathbf{y},$$

i.e., the normal equations.

4 Understanding, evaluating and improving the performance

We have so far encountered four different methods for supervised machine learning, and more are to come in later chapters. We always learn the models by adapting them to training data, and hope that the models thereby will give us good predictions also when faced with new, previously unseen, data. But can we really expect that to work? This may sound like a trivial question, but on a second thought it is perhaps not (so) obvious, and we will give it some attention in this chapter before we dive into more advanced methods in later chapters. By doing so we will unveil some interesting concepts, and also discover some practical tools for evaluating, choosing between and improving supervised machine learning methods.

4.1 Expected new data error E_{new} : performance in production

We start by introducing some concepts and notation. First, we define an error function $E(\hat{y}, y)$ which encodes the purpose of classification or regression. The error function compares a prediction $\hat{y}(\mathbf{x})$ to a measured data point y , and returns a small value (possibly zero) if $\hat{y}(\mathbf{x})$ is a good prediction of y , and a larger value otherwise. One could consider many different error functions, but our default choices are misclassification and squared error, respectively:

$$\text{Misclassification: } E(\hat{y}, y) \triangleq \begin{cases} 0 & \text{if } \hat{y} = y \\ 1 & \text{if } \hat{y} \neq y \end{cases} \quad (\text{classification}) \quad (4.1a)$$

$$\text{Squared error: } E(\hat{y}, y) \triangleq (\hat{y} - y)^2 \quad (\text{regression}) \quad (4.1b)$$

When we compute the average misclassification (4.1a), we usually refer to it as the *misclassification rate*. The misclassification rate is often a natural quantity to consider in classification, but for imbalanced or asymmetric problems other aspects might be more important, as we discuss in Section 4.5.

The error function $E(\hat{y}, y)$ has similarities to a loss function $L(\hat{y}, y)$. However, they are used differently: A loss function is used to *train* a model, whereas we use the error function to *analyze performance* of an already trained model. There are reasons for choosing $E(\hat{y}, y)$ and $L(\hat{y}, y)$ differently, which we will come back to soon.

In the end, supervised machine learning amounts to designing a method which performs well when faced with an endless stream of new, unseen data. Imagine for example all real-time recordings of street views that have to be processed by a vision system in a self-driving car once it is sold to a customer, or all incoming patients that have to be classified by a medical diagnosis system once it is implemented in clinical practice. The performance on fresh unseen data can in mathematical terms be understood as the average of the error function—how often the classifier is right, or how good the regression method predicts. To be able to mathematically describe the endless stream of new data, we introduce a *distribution over data* $p(\mathbf{x}, y)$. In most other chapters, we only consider the output y as a random variable whereas the input \mathbf{x} is considered fixed. In this chapter, however, we have to think of also the input \mathbf{x} as a random variable with a certain probability distribution. In any real-world machine learning scenario $p(\mathbf{x}, y)$ can be extremely complicated and practically impossible to write down. We will nevertheless use $p(\mathbf{x}, y)$ to *reason* about supervised machine learning methods, and the bare notion of $p(\mathbf{x}, y)$ (even though it is unknown in practice) will be helpful for that.

No matter which specific classification or regression method we consider, once it has been trained on training data $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$, it will return predictions $\hat{y}(\mathbf{x}_\star)$ for any new input \mathbf{x}_\star we give to it. We will

in this chapter write $\hat{y}(\mathbf{x}; \mathcal{T})$ to emphasize that the training data \mathcal{T} was used to train the model. Indeed, different training datasets will train the model differently and, consequently, give different predictions.

In the other chapters we mostly discuss how a model predicts one, or a few, test inputs \mathbf{x}_\star . Let us take that to the next level by integrating (averaging) the error function (4.1) over *all* possible test data points with respect to the distribution $p(\mathbf{x}, y)$. We refer to this as the *expected new data error*

$$E_{\text{new}} \triangleq \mathbb{E}_\star [E(\hat{y}(\mathbf{x}_\star; \mathcal{T}), y_\star)], \quad (4.2)$$

where the expectation \mathbb{E}_\star is the expectation over all possible test data points with respect to the distribution $(\mathbf{x}_\star, y_\star) \sim p(\mathbf{x}, y)$, that is,

$$\mathbb{E}_\star [E(\hat{y}(\mathbf{x}_\star; \mathcal{T}), y_\star)] = \int E(\hat{y}(\mathbf{x}_\star; \mathcal{T}), y_\star) p(\mathbf{x}_\star, y_\star) d\mathbf{x}_\star dy_\star. \quad (4.3)$$

Remember that the model (regardless if it is linear regression, a classification tree, an ensemble of trees, a neural network or something else) is trained on a given training dataset \mathcal{T} and represented by $\hat{y}(\cdot; \mathcal{T})$. What is happening in equation (4.2) is an averaging over possible test data points $(\mathbf{x}_\star, y_\star)$. Thus, E_{new} describes how well the model *generalizes* from the training data \mathcal{T} to new situations.

We also introduce the *training error*

$$E_{\text{train}} \triangleq \frac{1}{n} \sum_{i=1}^n E(\hat{y}(\mathbf{x}_i; \mathcal{T}), y_i), \quad (4.4)$$

where $\{\mathbf{x}_i, y_i\}_{i=1}^n$ is the training data \mathcal{T} . E_{train} simply describes how well a method performs on the training data on which it was trained, but gives no information on how well the method will perform for new unseen data points.¹

Time to reflect 4.1: What is E_{train} for k -NN with $k = 1$?

Whereas the training error E_{train} describes how well the method is able to “reproduce” the data from which it was learned, the expected new data error E_{new} tells us how well a method performs when we put it into production; what proportions of predictions a classifier will get right, and how well a regression method will predict in terms of average squared error. Or, in a more applied setting, what rate of false and missed detections of pedestrians we can expect a vision system in a self-driving car to make, or how big a proportion of all future patients a medical diagnosis system will get wrong.

The overall goal in supervised machine learning is to achieve as small E_{new} as possible.

This actually sheds some additional light upon the comment we made previously, that the loss function $L(\hat{y}, y)$ and the error function $E(\hat{y}, y)$ do not have to be the same. As we will discuss thoroughly in this chapter, a model which fits the training data well and consequently has a small E_{train} might still have a large E_{new} when faced with new unseen data. The best strategy to achieve as small E_{new} as possible is therefore not necessarily to minimize E_{train} . Besides the fact that the misclassification (4.1a) is unsuited as optimization objective (it is discontinuous and has derivative zero almost everywhere) it can also, depending on the method, be argued that E_{new} can be made smaller by a more clever choice of loss function. Such examples include gradient boosting and support vector machines. (Of course, this is only relevant for methods that are trained using a loss function. k -NN, for example, is not. The idea of evaluating using an error function applies however to all methods, no matter how it is trained.)

In practical cases we can, unfortunately, never compute E_{new} to assess how well we are doing. The reason is that $p(\mathbf{x}, y)$ —which we do not know in practice—is part of the definition of E_{new} . It is, however,

¹The term “risk function” is used in some literature both for loss and error functions, assuming they are chosen equally. In that terminology, E_{new} is referred to as “expected risk”, E_{train} as “empirical risk”, and the idea of minimizing the cost function as “empirical risk minimization”.

a too important construction to be abandoned, just because we cannot compute it. We will instead spend quite some effort both trying to *estimate* E_{new} (essentially by replacing the integral with a sum) as well as analyze how E_{new} behaves, to better understand how we can decrease it.

Remark 4.1 Note that E_{new} is a property of a trained model and a specific machine learning problem. That is, we cannot talk about “ E_{new} for logistic regression” in general, but instead we have to make more specific statements, like “ E_{new} for the handwritten digit recognition problem, with a logistic regression classifier trained with the MNIST data²”.

4.2 Estimating E_{new}

There are multiple reasons for a machine learning engineer to be interested in E_{new} , such as:

- judging if the performance is satisfying (whether E_{new} is small enough), or if more work should be put into the solution and/or more training data should be collected
- choosing between different methods
- choosing hyperparameters (such as k in k -NN, the regularization parameter in ridge regression or the number of hidden layers in deep learning) in order to minimize E_{new}
- reporting the expected performance to the customer

As discussed above, we can unfortunately not compute E_{new} in any practical situation. We will therefore explore some possibilities to *estimate* E_{new} , which will lead us to a very useful concept known as cross-validation.

$E_{\text{train}} \neq E_{\text{new}}$: We cannot estimate E_{new} from training data

We have both introduced the expected new data error, E_{new} , and the training error E_{train} . In contrast to E_{new} , we can always compute E_{train} .

We assume for now that \mathcal{T} consists of samples from $p(\mathbf{x}, y)$. This assumption means that the training data is collected under similar circumstances as the ones the learned model will be used under, which seems reasonable.

When an integral is hard to compute, it can be numerically approximated with a sum. Now, the question is if the integral in E_{new} can be well approximated by the sum in E_{train} , like

$$E_{\text{new}} = \int E(\hat{y}(\mathbf{x}; \mathcal{T}), y) p(\mathbf{x}, y) d\mathbf{x} dy \stackrel{??}{\approx} \frac{1}{n} \sum_{i=1}^n E(\hat{y}(\mathbf{x}_i; \mathcal{T}), y_i) = E_{\text{train}}. \quad (4.5)$$

Or, put differently: Can we expect a method to perform equally well when faced with new, previously unseen, data, as it did on the training data?

The answer is, unfortunately, **no**.

Time to reflect 4.2: Why can we not expect the performance on training data (E_{train}) to be a good approximation for how a method will perform on new, previously unseen data (E_{new}), even though the training data is drawn from the distribution $p(\mathbf{x}, y)$?

Equation (4.5) does *not* hold, and the reason is that the training data are not just any data points, but the predictions \hat{y} depends on them since they are used for training the model. We can therefore not expect (4.5) to hold. (Technically, the conditions for approximating an integral with a sum are not fulfilled since \hat{y} depends on \mathcal{T} .)

As we will discuss more thoroughly later, the average behavior of E_{train} and E_{new} is, in fact, typically $E_{\text{train}} < E_{\text{new}}$. That means that a method usually performs worse on new, unseen data, than on training data. *The performance on training data E_{train} is therefore not a good estimate of E_{new} .*

²<http://yann.lecun.com/exdb/mnist/>

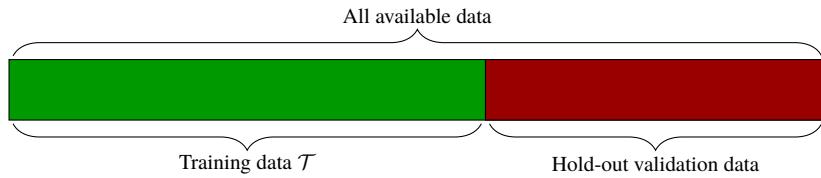


Figure 4.1: The hold-out validation dataset approach: If we split the available data in two sets and train the model on the training set, we can compute $E_{\text{hold-out}}$ using the hold-out validation set. $E_{\text{hold-out}}$ is an unbiased estimate of E_{new} , and the more data in the hold-out validation dataset the less variance (better estimate) in $E_{\text{hold-out}}$, but the less data left for training the model (larger E_{new}). The split here is only pictorial, in practice one should always split the data randomly.

$E_{\text{hold-out}} \approx E_{\text{new}}$: We can estimate E_{new} from hold-out validation data

We could not use the “approximate-the-integral-with-a-finite-sum” trick to estimate E_{new} by E_{train} , due to the fact that it effectively meant using the training data twice: first, to train the model (\hat{y} in (4.4)) and second, to evaluate the error function (the sum in (4.4)). A remedy is to set aside some *hold-out validation data* $\{\mathbf{x}_j, y_j\}_{j=1}^{n_v}$, which are not in \mathcal{T} used for training, and then use the hold-out validation data only for estimating the model performance as the *hold-out validation error*

$$E_{\text{hold-out}} \triangleq \frac{1}{j} \sum_{j=1}^{n_v} E(\hat{y}(\mathbf{x}_j; \mathcal{T}), y_j). \quad (4.6)$$

In this way, not all data will be used for training, but some data points (the hold-out validation data) will be saved and used only for computing $E_{\text{hold-out}}$. This procedure is a simple version of *cross-validation*, and is illustrated by Figure 4.1.

Be aware! *If you are splitting your data, always do it randomly! Someone might—intentionally or unintentionally—have sorted the dataset for you. If you do not split randomly, your binary classification problem might end up with one class in your training data and the other class in your hold-out validation data . . .*

With the conditions $n_v \geq 1$ and that all data is drawn from $p(\mathbf{x}, y)$, it can be shown that $E_{\text{hold-out}}$ is an unbiased estimate of E_{new} (meaning that if the entire procedure is repeated multiple times, each time with new data, the average value of $E_{\text{hold-out}}$ will be E_{new}). That is reassuring, at least on a theoretical level, but it does not tell us how close $E_{\text{hold-out}}$ will be to E_{new} in a single experiment. However, the variance of $E_{\text{hold-out}}$ decreases when the size of hold-out validation data n_v increases; a small variance of $E_{\text{hold-out}}$ means that we can expect it to be close to E_{new} . Thus, if we take the hold-out validation dataset big enough, $E_{\text{hold-out}}$ will be close to E_{new} . However, setting aside a big validation dataset means that the dataset left for training becomes small. You can presuppose that the more training data, the smaller E_{new} (which we will discuss later in Section 4.3). That is bad news since achieving a small E_{new} is our ultimate goal.

Sometimes there is *a lot* of available data. When we really have a lot of data, we can indeed afford to set aside a few percent data into a reasonably large hold-out validation dataset, without sacrificing the size of the training dataset too much. *In such data-rich situations, the hold-out validation data approach is sufficient.*

If the amount of available data is more limited, this becomes more of a problem. We are in practice faced with the following dilemma: *the better we want to know E_{new}* (more hold-out validation data gives less variance in $E_{\text{hold-out}}$), *the worse we have to make it* (less training data increases E_{new}). That is not very satisfying, and we have to look for an alternative to the hold-out validation data approach.

k -fold cross-validation: $E_{k\text{-fold}} \approx E_{\text{new}}$ without setting aside validation data

To avoid setting aside validation data, but still obtain an estimate of E_{new} , one could suggest a two-step procedure of

- (i) splitting the available data in one training and one hold-out validation set, train the model on the training data and compute $E_{\text{hold-out}}$ using hold-out validation data (as in Figure 4.1), and then
- (ii) training the model again, this time using the entire dataset.

By such a procedure, we get an estimate of E_{new} at the same time as a model trained on the entire dataset. That is not bad, but not perfect either. Why? To achieve small variance in the estimate, we have to put lots of data in the hold-out validation dataset in step (i). Unfortunately that means the model trained in (i) will possibly be very different from step (ii), and the estimate of E_{new} concerns the model from step (i), not the possibly very different model from step (ii). Hence, this will not give us a good estimate of E_{new} . It is, however, rather close to the useful k -fold cross-validation idea.

We would like to use all available data to train a model, and at the same time have a good estimate of E_{new} for that model. By *k -fold cross-validation*, we can achieve (almost) this. The idea of k -fold cross-validation is simply to repeat the hold-out validation dataset approach multiple times with a *different* hold-out dataset each time, in the following way:

- (i) split the dataset in k batches of similar size (see Figure 4.2), and let $\ell = 1$
- (ii) take batch ℓ as the hold-out validation data, and the remaining batches as training data
- (iii) train the model on the training data, and compute $E_{\text{hold-out}}^{(\ell)}$ as the average error on the hold-out validation data, (4.6)
- (iv) if $\ell < k$, set $\ell \leftarrow \ell + 1$ and return to (ii). If $\ell = k$, compute the *k -fold cross-validation error*

$$E_{k\text{-fold}} \triangleq \frac{1}{k} \sum_{\ell=1}^k E_{\text{hold-out}}^{(\ell)} \quad (4.7)$$

- (v) train the model again, this time using the entire dataset

This illustrated in Figure 4.2.

With k -fold cross-validation, we get a model which is trained on all data, as well as an approximation of E_{new} for that model, namely $E_{k\text{-fold}}$. Whereas $E_{\text{hold-out}}$ (Section 4.2) was an unbiased estimate of E_{new} (to the cost of setting aside hold-out validation data), $E_{k\text{-fold}}$ is only approximately so. However, with k large enough, it turns out to often be a sufficiently good approximation. Let us try to understand why k -fold cross-validation works.

First, we have to distinguish between the final model, which is trained on all data in step (v), and the intermediate models which are trained on all except a $1/k$ fraction of the data in step (iii). The key in k -fold cross-validation is that if k is large enough, the intermediate models are quite similar to the final model (since they are trained on almost the same dataset, only a fraction $1/k$ of the data is missing). Furthermore, each intermediate $E_{\text{hold-out}}^{(\ell)}$ is an unbiased but high-variance estimate of E_{new} for the corresponding intermediate model ℓ . Since all intermediate and the final model are similar, $E_{k\text{-fold}}$ (4.7) is approximately the average of k high-variance estimates of E_{new} for the final model. When averaging estimates, the variance decreases and $E_{k\text{-fold}}$ will have a lower variance.

Be aware! For the same reason as with the hold-out validation data approach, it is important to always split the data randomly for cross-validation to work! A simple solution is to first randomly permute the entire dataset, and thereafter split it into batches.



Figure 4.2: Illustration of k -fold cross-validation. The data is split in k batches of similar sizes. When looping over $\ell = 1, 2, \dots, k$, batch ℓ is held out as validation data, and the model is trained on the remaining $k - 1$ data batches. Each time, the trained model is used to compute the average error $E_{k\text{-fold}}^{(\ell)}$ for the validation data. The final model is trained using all available data, and the estimate of E_{new} for that model is $E_{k\text{-fold}}$, the average of all $E_{k\text{-fold}}^{(\ell)}$.

We usually talk about training (or learning) as a procedure that is executed once. However, in k -fold cross-validation the training is repeated k (or even $k + 1$) times. A special case is $k = n$ which is also called *leave-one-out cross-validation*. For methods such as linear regression, the actual training (solving the normal equations) is usually done within milliseconds on modern computers, and doing it an extra n might not be a problem in practice. If the training is computationally intense (as for deep learning), it becomes a rather computationally demanding procedure, and a choice like $k = 10$ might be more practically feasible. If there is much data available, it is also an option to use the computationally less demanding hold-out cross-validation.

Using a test dataset

A very important use of $E_{k\text{-fold}}$ (or $E_{\text{hold-out}}$) in practice is to choose between methods and select different types of hyperparameters such that $E_{k\text{-fold}}$ (or $E_{\text{hold-out}}$) becomes as small as possible. Typical hyperparameters to choose in this way are k in k -NN, tree depths or regularization parameters. However, much like we cannot use the training data error E_{train} to estimate the new data error E_{new} , selecting models and hyperparameters based on $E_{k\text{-fold}}$ (or $E_{\text{hold-out}}$) will invalidate its use as an estimator of E_{new} . If it is important to have a good estimate of the final E_{new} , it is wise to first set aside another hold-out dataset, which we refer to as a *test set*. This test set should be used only once (after selecting models and hyperparameters) to estimate E_{new} for the final model.

In problems where the training data is expensive, it is common to increase the training dataset using more or less artificial techniques. Such techniques can be to duplicate the data and add noise to the duplicated versions, to use simulated data, or to use data from a different but related problem, as we discuss in more depth in Chapter 10. With such techniques (which indeed can be very successful), the training data \mathcal{T} is no longer drawn from $p(\mathbf{x}, y)$. In the worst case (if the artificial training data is very

poor), \mathcal{T} might not provide any information about $p(\mathbf{x}, y)$, and we can not really expect the model to learn anything useful. It can therefore be very useful to have a good estimate of E_{new} if such techniques were used during training, but a reliable estimate of E_{new} can only be achieved from data that we *know* are drawn from $p(\mathbf{x}, y)$ (that is, collected under production-like circumstances). If the training data is extended artificially, it is therefore extra important to set aside a test data set *before* that extension is done.

Remark 4.2 *The error function evaluated on the test data could be called “test error”. To avoid confusion, we do not use that term since it in other places is used ambiguously for both the error on a test dataset as well as E_{new} .*

4.3 The training error–generalization gap decomposition of E_{new}

Designing a method with small E_{new} is the goal in supervised machine learning, and cross-validation helps in estimating E_{new} . However, more can be said to also *understand* E_{new} . To be able to reason about E_{new} , we have to introduce another abstraction level, namely the *training-data averaged* versions of E_{new} and E_{train} ,

$$\bar{E}_{\text{new}} \triangleq \mathbb{E}_{\mathcal{T}}[E_{\text{new}}], \quad (4.8a)$$

$$\bar{E}_{\text{train}} \triangleq \mathbb{E}_{\mathcal{T}}[E_{\text{train}}]. \quad (4.8b)$$

Here, $\mathbb{E}_{\mathcal{T}}$ denotes the expected value when the training dataset $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ (of a fixed size n) is drawn from $p(\mathbf{x}, y)$. Thus \bar{E}_{new} is the average E_{new} if we would train the model multiple times on different training datasets, and similarly for \bar{E}_{train} . The point of introducing these, as it turns out, is that we can say more about the *average* behavior \bar{E}_{new} and \bar{E}_{train} , than we can say about E_{new} and E_{train} when the model is trained on one specific training dataset \mathcal{T} . Even though we most often care about E_{new} in the end (the training data is usually fixed), insights from studying \bar{E}_{new} are still useful. In fact, the k -fold cross-validation is rather an estimate of \bar{E}_{new} instead of E_{new} , since each model in k -fold is trained with a (slightly) different training dataset.

We have already discussed the fact that E_{train} cannot be used in estimating E_{new} . In fact, it usually holds that

$$\bar{E}_{\text{train}} < \bar{E}_{\text{new}}, \quad (4.9)$$

Put in words, this means that on average, a method usually performs worse on new, unseen data, than on training data. A methods ability to perform well on unseen data after being trained on training data can be understood as the method’s ability to *generalize* from training data. We consequently call the difference between \bar{E}_{new} and \bar{E}_{train} the *generalization gap*³, as

$$\text{generalization gap} \triangleq \bar{E}_{\text{new}} - \bar{E}_{\text{train}}. \quad (4.10)$$

The generalization gap is the difference between performance on training data and the performance ‘in production’ on new, previously unseen data. Since we can always compute E_{train} and cross-validation gives an estimate of E_{new} , we can also estimate the generalization gap in practice.

With the decomposition of \bar{E}_{new} into

$$\bar{E}_{\text{new}} = \bar{E}_{\text{train}} + \text{generalization gap}, \quad (4.11)$$

we also have an opening for digging deeper and trying to understand what affects \bar{E}_{new} in practice. We will refer to (4.11) as the *training error–generalization gap decomposition of \bar{E}_{new}* .

³We use a loose terminology and refer also to $E_{\text{new}} - E_{\text{train}}$ as the generalization gap.

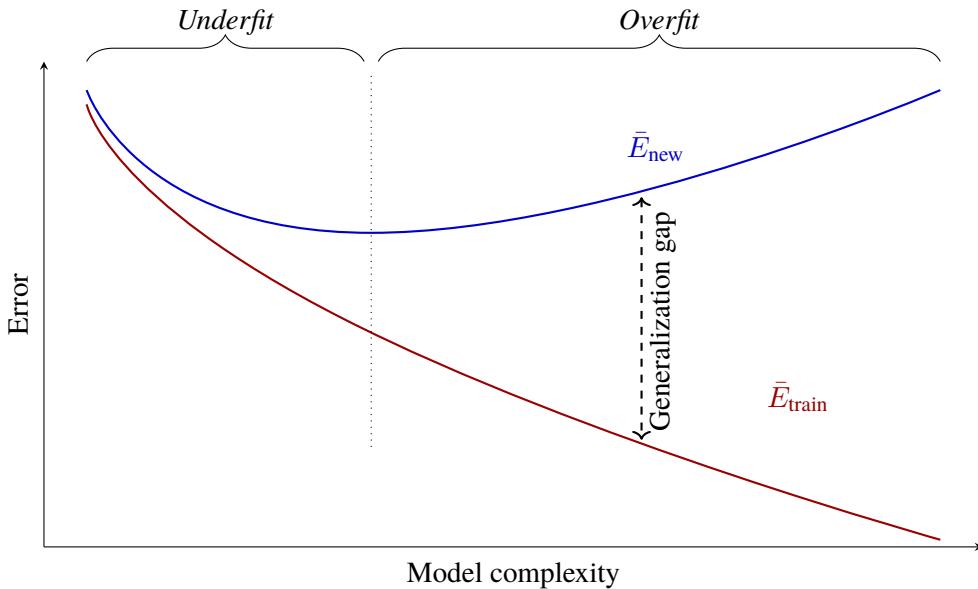


Figure 4.3: Behavior of \bar{E}_{train} and \bar{E}_{new} for many supervised machine learning methods, as a function of model complexity. We have not made a formal definition of complexity, but a rough proxy is the number of parameters that are learned from the data. The difference between the two curves is the generalization gap. The training error \bar{E}_{train} decreases as the model complexity increases, whereas the new data error \bar{E}_{new} typically has a U-shape. If the model is so complex that \bar{E}_{new} is larger than it had been with a less complex model, the term *overfit* is commonly used. Somewhat less commonly is the term *underfit* used for the opposite situation. The level of model complexity which gives the minimum \bar{E}_{new} (at the dotted line) could be called a balanced fit. When we, for example, use cross-validation to select hyperparameters (that is, tuning the model complexity), we are searching for a balanced fit.

What affects the generalization gap

The generalization gap depends on the method and the problem. Concerning the method, one can typically say that *the more a method adapts to training data, the larger the generalization gap*. A theoretical framework for how much method a method adapts to training data is given by the so-called VC dimension. From the VC dimension framework, probabilistic bounds on the generalization gap can be derived, but those bounds are unfortunately rather conservative, and we will not pursue that approach any further. Instead, we only use the vague terms *model complexity* or *model flexibility* (we use them as synonyms), by which we mean the ability of a method to adopt to patterns in the training data. A model with high complexity/flexibility (such as a fully connected deep neural network, deep trees and k -NN with small k) can describe almost arbitrarily complicated relationships, whereas a model with low complexity/flexibility (such as logistic regression) is less flexible in what functions it can describe. For parametric methods, the model complexity is somewhat related to the number of parameters that are trained, but is also affected by regularization techniques. As we will come back to later, the idea of model complexity is an oversimplification and does not capture the full nature of various supervised machine learning methods, but it nevertheless carries some useful intuition.

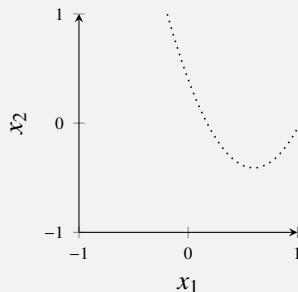
Typically, *higher model complexity implies larger generalization gap*. Furthermore, \bar{E}_{train} decreases as the model complexity increases, whereas \bar{E}_{new} typically attains a minimum for some intermediate model complexity value: too small *and* too high model complexity both raises \bar{E}_{new} . This is illustrated in Figure 4.3. A too high model complexity, meaning that \bar{E}_{new} is higher than it had been with a less complex model, is called *overfit*. The other situation, when the model complexity is too low, is sometimes called *underfit*. In a consistent terminology, the point where \bar{E}_{new} attains its minimum could be referred to as a balanced fit. Since the goal is to minimize \bar{E}_{new} , we are interested in finding this sweet spot. We also illustrate this by Example 4.1.

Remark 4.3 We discuss the usual behavior of \bar{E}_{new} , \bar{E}_{train} and the generalization gap. We use the term ‘usual’ because there are so many supervised machine learning methods and problems that it is

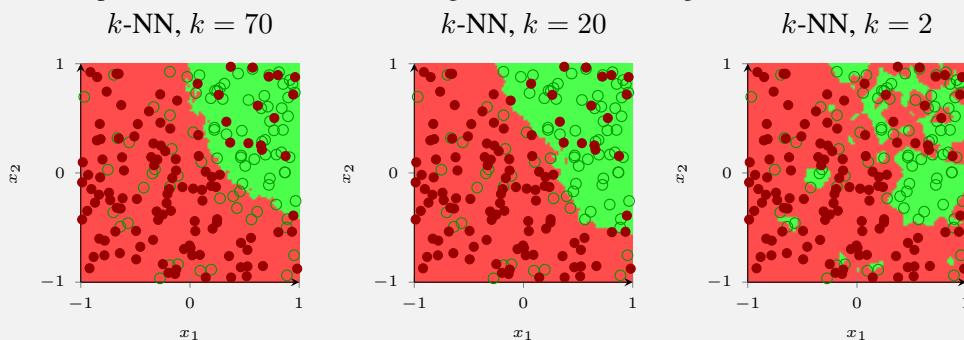
almost impossible to make any claim that is always true for all possible situations, and pathological counter-examples also exist. One should also keep in mind that claims about \bar{E}_{train} and \bar{E}_{new} are about the average behavior, which hopefully is clear in Example 4.1.

Example 4.1: The training error–generalization gap decomposition for k -NN

We consider a simulated binary classification example with a two-dimensional input \mathbf{x} . On the contrary to all real world machine learning problems, in a simulated problem like this we do know $p(\mathbf{x}, y)$ (otherwise we could not make the simulation). In this example, $p(\mathbf{x})$ is a uniform distribution on the square $[-1, 1]^2$, and $p(y | \mathbf{x})$ is defined as follows: all points above the dotted curve in the figure below are green with probability 0.8, and points below the curve are red with probability 0.8. (The optimal classifier, in terms of minimal E_{new} , would have the dotted line as its decision boundary and achieve $E_{\text{new}} = 0.2$.)



We have $n = 200$ in the training data, and learn three classifiers: k -NN with $k = 70$, $k = 20$ and $k = 2$, respectively. In model complexity sense, $k = 70$ gives the least flexible model, and $k = 2$ the most flexible model. We plot their decision boundaries, together with the training data:



Intuitively we see that $k = 2$ (right) adapts too well to the data. With $k = 70$, on the other hand, the model is rigid enough not to adapt to the noise, but appears to possibly be too inflexible to adapt well to the true dotted line above.

We can compute E_{train} by counting the fraction of misclassified training data points in the figures above. From left to right, we get $E_{\text{train}} = 0.27, 0.24, 0.22$. Since this is a simulated example, we can also access E_{new} (or rather estimate it numerically by simulating a lot of test data), and from left to right we get $E_{\text{new}} = 0.26, 0.23, 0.33$. This pattern resembles Figure 4.3, except for the fact that E_{new} is actually smaller than E_{train} for some values of k . This is, however, not unexpected. What we have discussed in the main text is the *average* \bar{E}_{new} and \bar{E}_{train} , *not* the situation with E_{new} and E_{train} for one particular set of training data. To study \bar{E}_{new} and \bar{E}_{train} , we therefore repeat this experiment 100 times, and compute the average over those 100 experiments:

	k -NN with $k = 70$	k -NN with $k = 20$	k -NN with $k = 2$
\bar{E}_{train}	0.24	0.22	0.17
\bar{E}_{new}	0.25	0.23	0.30

This table follows Figure 4.3 well: The generalization gap (difference between \bar{E}_{new} and \bar{E}_{train}) is positive and increases with model complexity (decreasing k in k -NN), whereas \bar{E}_{train} decreases with model complexity. Among these values for k , \bar{E}_{new} has its minimum for $k = 20$. This suggests that k -NN with $k = 2$ suffers from overfitting for this problem, whereas $k = 70$ is a case of underfitting.

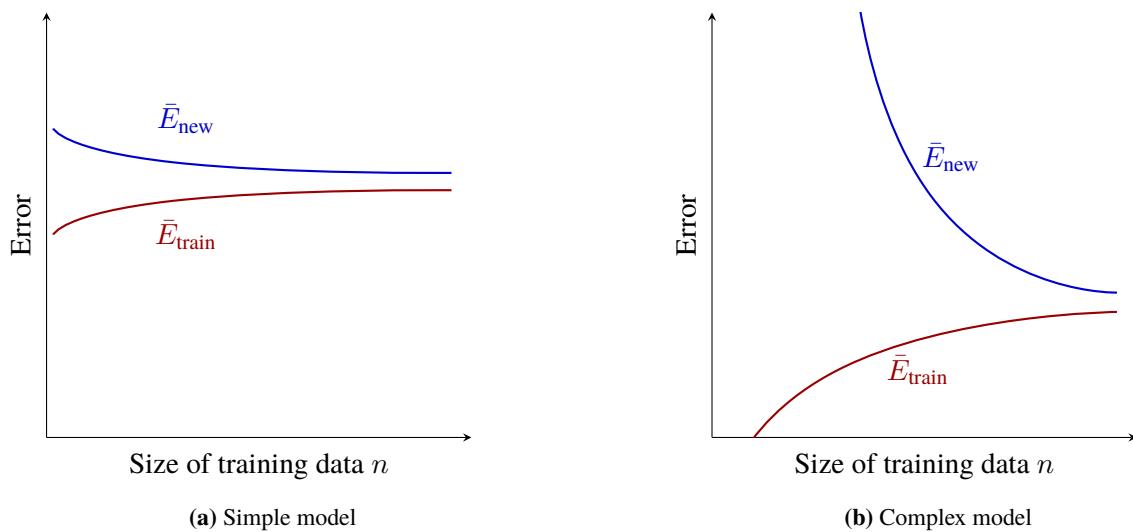


Figure 4.4: Typical relationship between \bar{E}_{new} , \bar{E}_{train} and the number of data points n in the training dataset for a simple model (low model flexibility, left) and a complex model (high model flexibility, right). The generalization gap (difference between \bar{E}_{new} and \bar{E}_{train}) decreases, at the same time as \bar{E}_{train} increases. Typically, a more complex model (right panel) will for large enough n attain a smaller \bar{E}_{new} than a simpler model (left panel) would on the same problem (the axes of the figures are comparable). However, the generalization gap is typically larger for a more complex model, in particular when the training dataset is small.

The have so far been concerned about the relationship between the generalization gap and the model complexity. Another very important aspect is the size of the training dataset, n . We can in general expect that *the more training data, the smaller the generalization gap*. On the other hand, \bar{E}_{train} typically increases as n increases, since most models are not able to fit all training data perfectly if there are too many of them. A typical behavior of \bar{E}_{train} and \bar{E}_{new} is sketched in Figure 4.4.

Reducing E_{new} in practice

Our overall goal is to achieve a small error “in production”, that is, a small E_{new} . To achieve that, according to the decomposition $E_{\text{new}} = E_{\text{train}} + \text{generalization gap}$, we need to have E_{train} as well as the generalization gap small. Let us draw two conclusions from what we have seen so far.

- The new data error E_{new} will (most often) not be smaller than the training error E_{train} . Thus, if E_{train} is much bigger than the E_{new} you need for your application to be successful, you do not even need to waste time on implementing cross-validation for estimating E_{new} . Instead, you should re-think the problem and which method you are using.
 - The generalization gap and E_{new} decreases as good as always as n increases. Thus, increasing the size of the training data may help a lot, and will at least not hurt.

Making the model more flexible decreases E_{train} but increases, often, the generalization gap. Making the model less flexible, on the other hand, typically decreases the generalization gap but increases E_{train} . The optimal tradeoff, in terms of small E_{new} , is often obtained when neither the generalization gap nor the training error E_{train} is zero. Thus, by monitoring E_{train} and estimating E_{new} with cross-validation we also get the following advice:

- If $E_{\text{hold-out}} \approx E_{\text{train}}$ (small generalization gap), it might be beneficial to increase the model flexibility by loosening the regularization, increasing the model order (more parameters to learn), etc.
 - If E_{train} is close to zero, it might be beneficial to decrease the model flexibility by tightening the regularization, decreasing the order (fewer parameters to learn), etc.

Shortcomings of the model complexity scale

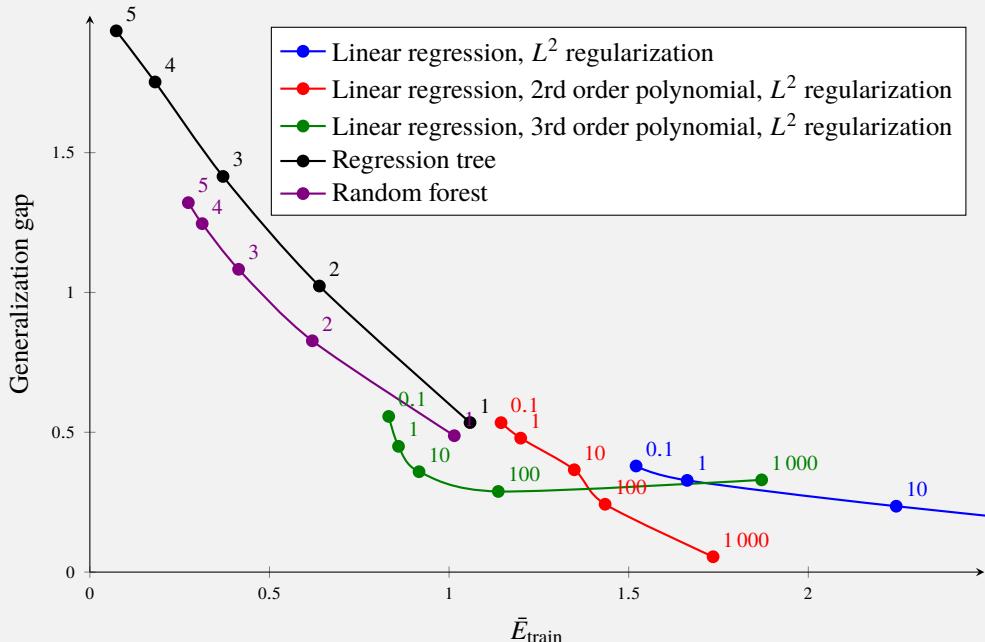
When having only one hyperparameter to choose, the situation sketched in Figure 4.3 is often a relevant picture. However, when having multiple hyperparameters (or even competing methods) to choose, it is important to realize that the one-dimensional model complexity scale in Figure 4.3 does not make justice for the space of all possible choices. For a given problem, one method can have a smaller generalization gap than another method *without* having a larger training error. Some methods are simply better for certain problems. The one-dimensional complexity scale can be particularly misleading for intricate deep learning models, but as we illustrate in Example 4.2 it is not even sufficient for the relatively simple problem of jointly choosing the degree of polynomial regression (higher degree means more flexibility) and regularization parameter (more regularization means less flexibility).

Example 4.2: Training error and generalization gap for a regression problem

To illustrate how the training error and generalization gap can behave, we consider a simulated problem so that we can compute E_{new} . We let $n = 10$ data points be generated as $x \sim \mathcal{U}[-5, 10]$, $y = \min(0.1x^2, 3) + \varepsilon$, and $\varepsilon \sim \mathcal{N}(0, 1)$, and consider the following regression methods:

- Linear regression with L^2 regularization
- Linear regression with a quadratic polynomial and L^2 regularization
- Linear regression with a third order polynomial and L^2 regularization
- Regression tree
- A random forest (Chapter 7) with 10 regression trees

For each of these methods, we try a few different values of the hyperparameters (regularization parameter and tree depth, respectively), compute \bar{E}_{train} and the generalization gap.



For each method the hyperparameter that minimizes \bar{E}_{new} is the value which is closest (in the 1-norm sense) to the origin, since $\bar{E}_{\text{new}} = \bar{E}_{\text{train}} + \text{generalization gap}$. Having decided on a certain model, and only having one hyperparameter left to choose, corresponds well to the situation in Figure 4.3.

However, when we compare the *different* methods, a more complicated situation is revealed than is described by the one-dimensional model complexity scale. Compare, for example, the second (green) to the third order polynomial (red) linear regression: for some values of the regularization parameter, the training error decreases *without* increasing the generalization gap. Similarly is the generalization gap smaller, while the training error remains the same, for the random forest (black) than for the tree (purple) for a maximum tree depth of 2. These type of relationships are quite intricate, problem-dependent and impossible to describe using the simplified picture in Figure 4.3.

In any real problem we cannot make such a plot as in Example 4.2, that is only possible for simulated problems. In practice it is therefore important to use your knowledge about the problem and how different methods are designed, in order to make a good choice. It is good to choose models that are known to work well for a specific type of data and use experience from similar problems. We can also use cross-validation for selecting between different models and choosing hyperparameters. Despite the simplified picture, the intuition about under- and overfit from Figure 4.3 can still be very helpful when deciding on what method or hyperparameter value to explore next with cross-validation.

4.4 The bias-variance decomposition of E_{new}

We will now introduce another decomposition of \bar{E}_{new} into a (squared) *bias* and a *variance* term, as well as an unavoidable component of irreducible noise. This decomposition is somewhat more abstract than the training-generalization gap, but provides some additional insights into E_{new} and how different models behave.

Let us first make a short reminder of the general concepts of bias and variance. Consider an experiment with an unknown constant z_0 , which we would like to estimate. To our help for estimating z_0 we have a random variable z . Think, for example, of z_0 as being the (true) position of an object, and z of being noisy GPS measurements of that position. Since z is a random variable, it has some mean $\mathbb{E}[z]$ which we denote by \bar{z} . We now define

$$\text{Bias: } \bar{z} - z_0 \quad (4.12a)$$

$$\text{Variance: } \mathbb{E}[(z - \bar{z})^2] = \mathbb{E}[z^2] - \bar{z}^2. \quad (4.12b)$$

The *variance* describes how much the experiment varies each time we perform it (the amount of noise in the GPS measurements), whereas the *bias* describes the systematic error in z that remains no matter how many times we repeat the experiment (a possible shift or offset in the GPS measurements). If we consider the expected squared error between z and z_0 as a metric of how good the estimator z is, we can re-write it in terms of the variance and the squared bias,

$$\begin{aligned} \mathbb{E}[(z - z_0)^2] &= \mathbb{E}\left[((z - \bar{z}) + (\bar{z} - z_0))^2\right] = \\ &= \underbrace{\mathbb{E}[(z - \bar{z})^2]}_{\text{Variance}} + 2\underbrace{\mathbb{E}[z] - \bar{z}}_0(\bar{z} - z_0) + \underbrace{(\bar{z} - z_0)^2}_{\text{bias}^2}. \end{aligned} \quad (4.13)$$

In words, the average squared error between z and z_0 is the sum of the squared bias and the variance. The main point here is that to obtain a small expected squared error, we have to consider the bias *and* the variance. Only a small bias *or* little variance in the estimator is not enough, but both aspects are important.

We will now apply the bias and variance concept to our supervised machine learning setting, and in particular to the regression problem. The intuition, however, carries over also to the classification problem⁴. In this setting, z_0 corresponds to the true relationship between inputs and output, and the random variable z corresponds to the model learned from training data. (Since the training data collection includes randomness, the model learned from it will also be random.) Let us spell out the details.

We first make the assumption that the true relationship between input \mathbf{x} and output y can be described as some (possibly very complicated) function $f_0(\mathbf{x})$ plus independent noise ε ,

$$y = f_0(\mathbf{x}) + \varepsilon, \text{ with } \mathbb{E}[\varepsilon] = 0 \text{ and } \text{var}(\varepsilon) = \sigma^2. \quad (4.14)$$

In our notation, $\hat{y}(\mathbf{x}; \mathcal{T})$ represents the model when it is trained on training data \mathcal{T} . This is our random variable, corresponding to z above. We now also introduce the *average trained model*, corresponding to \bar{z} ,

$$\bar{f}(\mathbf{x}) \triangleq \mathbb{E}_{\mathcal{T}} [\hat{y}(\mathbf{x}; \mathcal{T})]. \quad (4.15)$$

⁴For intuition, we may think of classification problems as regression in terms of the decision boundaries.

As before, $\mathbb{E}_{\mathcal{T}}$ denotes the expected value over training data drawn from $p(\mathbf{x}, y)$. Thus, $\bar{f}(\mathbf{x})$ is the (hypothetical) average model we would achieve, if we could re-train the model an infinite number of times on different training datasets and compute the average.

Remember that the definition of \bar{E}_{new} (for regression with squared error) is

$$\bar{E}_{\text{new}} = \mathbb{E}_{\mathcal{T}} \left[\mathbb{E}_{\star} \left[(\hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - y_{\star})^2 \right] \right], \quad (4.16)$$

and assuming these integrals (expressed as expected values) fulfill some technical assumptions, we can change the order of integration and write (4.16) as

$$\bar{E}_{\text{new}} = \mathbb{E}_{\star} \left[\mathbb{E}_{\mathcal{T}} \left[(\hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - f_0(\mathbf{x}_{\star}) - \varepsilon)^2 \right] \right] \quad (4.17)$$

With a slight extension of (4.13) to also include the zero-mean noise term ε (which is independent of $\hat{y}(\mathbf{x}_{\star}; \mathcal{T})$), we can rewrite the expression inside the expected value \mathbb{E}_{\star} in (4.17) as

$$\mathbb{E}_{\mathcal{T}} \left[\underbrace{(\hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - f_0(\mathbf{x}_{\star}))^2}_{\text{"z"}} - \underbrace{\varepsilon^2}_{\text{"z}_0} \right] = (\bar{f}(\mathbf{x}_{\star}) - f_0(\mathbf{x}_{\star}))^2 + \mathbb{E}_{\mathcal{T}} \left[(\hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - \bar{f}(\mathbf{x}_{\star}))^2 \right] + \varepsilon^2. \quad (4.18)$$

This is (4.13) applied to supervised machine learning. In \bar{E}_{new} , which we are interested in decomposing, we also have the expectation over new data points \mathbb{E}_{\star} . By incorporating also that expected value in the expression, we can decompose \bar{E}_{new} as

$$\bar{E}_{\text{new}} = \underbrace{\mathbb{E}_{\star} \left[(\bar{f}(\mathbf{x}_{\star}) - f_0(\mathbf{x}_{\star}))^2 \right]}_{\text{Bias}^2} + \underbrace{\mathbb{E}_{\star} \left[\mathbb{E}_{\mathcal{T}} \left[(\hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - \bar{f}(\mathbf{x}_{\star}))^2 \right] \right]}_{\text{Variance}} + \underbrace{\sigma^2}_{\text{Irreducible error}}. \quad (4.19)$$

The squared bias term $\mathbb{E}_{\star} \left[(\bar{f}(\mathbf{x}_{\star}) - f_0(\mathbf{x}_{\star}))^2 \right]$ now describes how much the average trained model $\bar{f}(\mathbf{x}_{\star})$ differs from the true $f_0(\mathbf{x}_{\star})$, averaged over all possible data points \mathbf{x}_{\star} . In a similar fashion the variance term $\mathbb{E}_{\star} \left[\mathbb{E}_{\mathcal{T}} \left[(\hat{y}(\mathbf{x}_{\star}; \mathcal{T}) - \bar{f}(\mathbf{x}_{\star}))^2 \right] \right]$ describes how much $\hat{y}(\mathbf{x}; \mathcal{T})$ varies each time the model is trained on a different training dataset. The bias term cannot be small unless the model is flexible enough such that $\bar{f}(\mathbf{x})$ can be close to $f_0(\mathbf{x})$. If the variance term is small, the model is not very sensitive to exactly which data points happened to be in the training data, and vice versa. The irreducible error σ^2 is simply an effect of the assumption (4.14) —it is not possible to predict ε since it is truly random. There is not so much more to say about the irreducible error, but we will focus on the bias and variance terms.

What affects the bias and variance

We have not properly defined model complexity, but we can actually use the bias and variance concept to give it a more concrete meaning: A high model complexity means low bias and high variance, and a low model complexity means high bias and low variance, as illustrated by Figure 4.5.

This resonates well with the intuition. The more flexible a model is, the more it will adapt to the training data \mathcal{T} —not only to the interesting patterns, but also to the actual data points and noise that happened to be in \mathcal{T} . That is exactly what is described by the variance term. On the other hand, a model with low flexibility can be too rigid to capture the true relationship $\bar{f}(\mathbf{x})$ between inputs and outputs well. This effect is described by the squared bias term.

Figure 4.5 may very well be compared to Figure 4.3, which builds on the training error–generalization gap decomposition of \bar{E}_{new} instead. From Figure 4.5 we can talk about the challenge of finding the right model complexity level also as the *bias-variance tradeoff*. We give an example of this in Example 4.3.

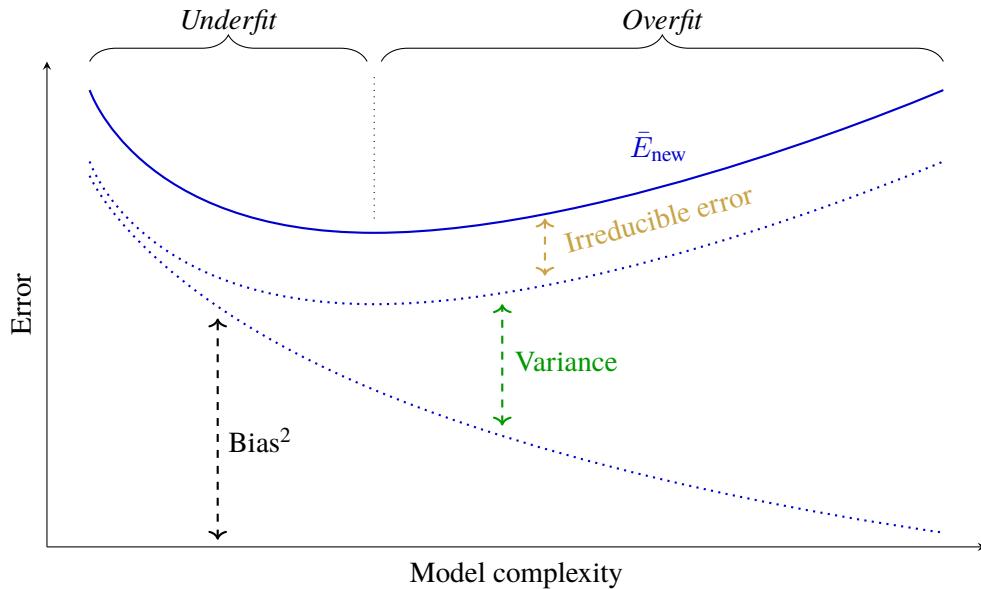


Figure 4.5: The bias-variance decomposition of \bar{E}_{new} (cf. Figure 4.3). Low model complexity means high bias. The more complicated the model is, the more it adapts to training data, and the higher variance. The irreducible error is always constant. The problem of achieving a small \bar{E}_{new} by selecting a good model complexity level is often called the bias-variance tradeoff.

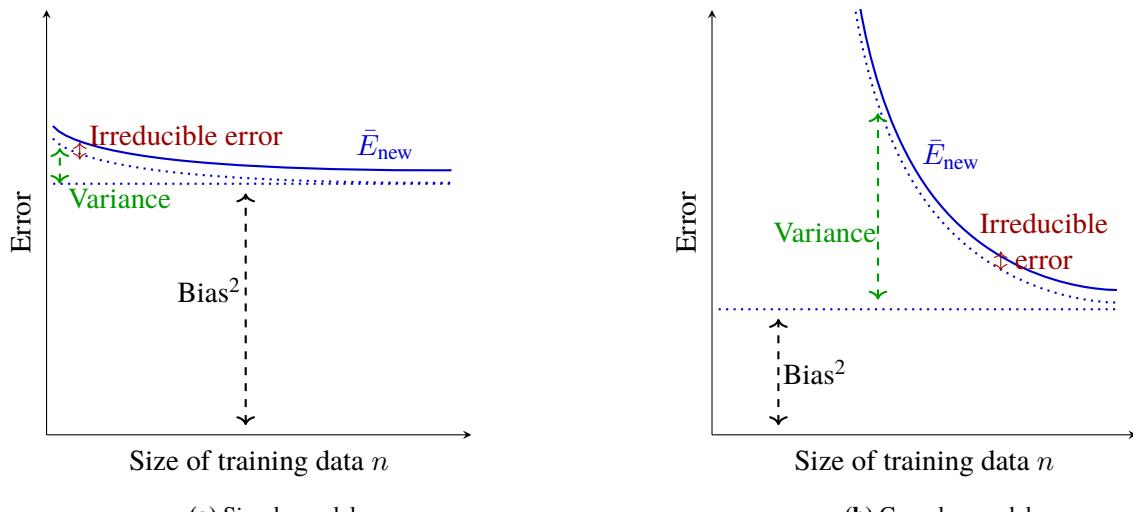


Figure 4.6: The typical relationship between bias, variance and the size n of the training dataset (cf. Figure 4.4). The bias is (approximately) constant, whereas the variance decreases as the size of the training dataset increases.

The squared bias term is mostly a property of the model rather than of the training dataset, and we may think⁵ of the bias term as independent of the number of data points n in the training data. The variance term, on the other hand, varies highly with n . As we know, \bar{E}_{new} typically decreases as n increases, and essentially the entire decline in \bar{E}_{new} is because of the decline in the variance. Intuitively, the more data, the more information about the parameters, meaning less variance. This is summarized by Figure 4.6, which can be compared to Figure 4.4.

⁵This is not exactly true. The average model \bar{f} might indeed be different if all training datasets (which we average over) contain $n = 2$ or $n = 100\,000$ data points, but we neglect that effect here.

Example 4.3: The bias–variance tradeoff for L^2 regularized linear regression

Let us consider a simulated regression example. We let $p(x, y)$ follow from $x \sim \mathcal{U}[0, 1]$ and

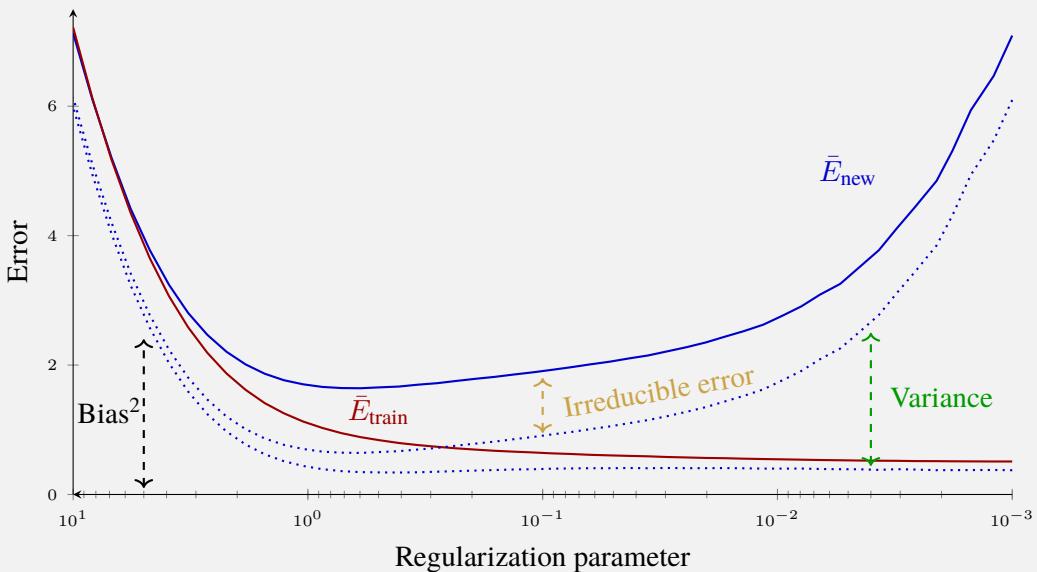
$$y = 5 - 2x + x^3 + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 1). \quad (4.20)$$

We let the training data consist of only $n = 10$ samples. We now try to model the data using linear regression with a 4th order polynomial

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4 + \varepsilon. \quad (4.21)$$

Since (4.20) is a special case of (4.21) and the squared error loss corresponds to Gaussian noise, we do actually have zero bias for this model if we train it using squared error loss. However, learning 5 parameters from only 10 data points leads to very high variance, so we decide to train the model with squared error loss and L^2 regularization, which will decrease the variance (but increase the bias). The more regularization (bigger λ), the more bias and less variance.

Since this is a simulated example, we can repeat the experiment multiple times and estimate the bias and variance terms (since we can simulate as much training and test data as needed). We plot them in the very same style as Figures 4.3 and 4.5 (note the reversed x-axis: a smaller regularization parameter corresponds to a higher model complexity). For this problem the optimal value of λ would have been about 0.7 since \bar{E}_{new} attains its minimum there. Finding this optimal λ is a typical example of the bias–variance tradeoff.



Connections between bias, variance and the generalization gap

The bias and variance are theoretically well defined properties, but often intangible in practice since they are defined in terms of $p(\mathbf{x}, y)$. In practice, we only have an estimate of the generalization gap (for example as $E_{\text{hold-out}} - E_{\text{train}}$), but neither the variance nor the bias. It is therefore interesting to explore what E_{train} and the generalization gap says about the bias and variance.

Considering the regression problem, and assuming that the error function (squared error) also is the loss function and that the global minimum is found during training, we can write

$$\sigma^2 + \text{bias}^2 = \mathbb{E}_\star [(\bar{f}(\mathbf{x}_\star) - y_\star)^2] \approx \frac{1}{n} \sum_{i=1}^n (\bar{f}(\mathbf{x}_i) - y_i)^2 \geq \frac{1}{n} \sum_{i=1}^n (\hat{y}(\mathbf{x}_i; \mathcal{T}) - y_i)^2 = E_{\text{train}}. \quad (4.22)$$

In the approximate equality, we used the “replace the integral with a sum”-trick⁶. In the next step,

⁶Since neither $\bar{f}(\mathbf{x}_\star)$ nor y_\star depends on the training data $\{\mathbf{x}_i, y_i\}_{i=1}^n$, we can use $\{\mathbf{x}_i, y_i\}_{i=1}^n$ for approximating the integral, cf. (4.2)

equality is obtained if $\hat{y} = \bar{f}$ (which we assume is possible), and inequality otherwise. Remembering that $\bar{E}_{\text{new}} = \sigma^2 + \text{bias}^2 + \text{variance}$, and allowing ourselves to write $\bar{E}_{\text{new}} - E_{\text{train}} = \text{generalization gap}$, we have

$$\text{generalization gap} \gtrapprox \text{variance}, \quad (4.23a)$$

$$E_{\text{train}} \lessapprox \text{bias}^2 + \sigma^2. \quad (4.23b)$$

We have made several assumptions in this derivation that are not always met in practice, but it at least gives us some rough idea.

As we discussed previously, the choice of method is crucial for what E_{new} is obtained. Again the one-dimensional scale in Figure 4.5 and the notion of a bias-variance tradeoff is a simplified picture; decreased bias does not *always* lead to increased variance, and vice versa. However, in contrast to the decomposition of E_{new} into training error and generalization gap, the bias and variance decomposition can shed some more light over why E_{new} decreases for different methods: sometimes, the superiority of one method over another can sometimes be attributed to either a lower bias or a lower variance.

A simple (and useless) way to increase the variance without decreasing the bias in linear regression, is to first learn the parameters using the normal equations and thereafter add zero-mean random noise to them. The extra noise does not affect the bias, since the noise has zero mean and hence leaves the average model \bar{f} unchanged, but the variance increases. (This effects also the training error and the generalization gap, but in a less clear way.) This way of training linear regression would be pointless in practice since it increases E_{new} , but it illustrates the fact that increased variance does *not* automatically leads to decreased bias.

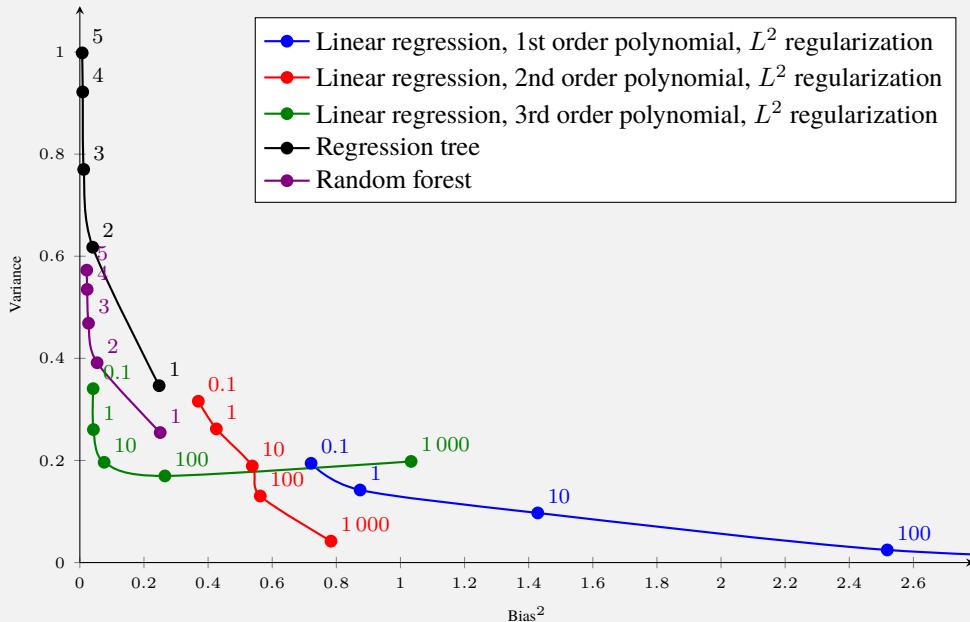
A much more useful way of dealing with bias and variance is the meta-method called bagging, Chapter 7. It makes use of several copies (an ensemble) of a base model, each trained on a slightly different version of the training dataset. Since bagging averages over many base models, it reduces the variance, but the bias remains (essentially) unchanged. Hence, by using bagging instead of the base model, the variance is decreased (almost) without increasing the bias, which consequently decreases E_{new} .

To conclude, the world is more complex than just the one-dimensional model complexity scale used in Figure 4.3 and 4.5, which we illustrate by Example 4.4.

Time to reflect 4.3: Can you modify linear regression such that the bias increases, without decreasing the variance?

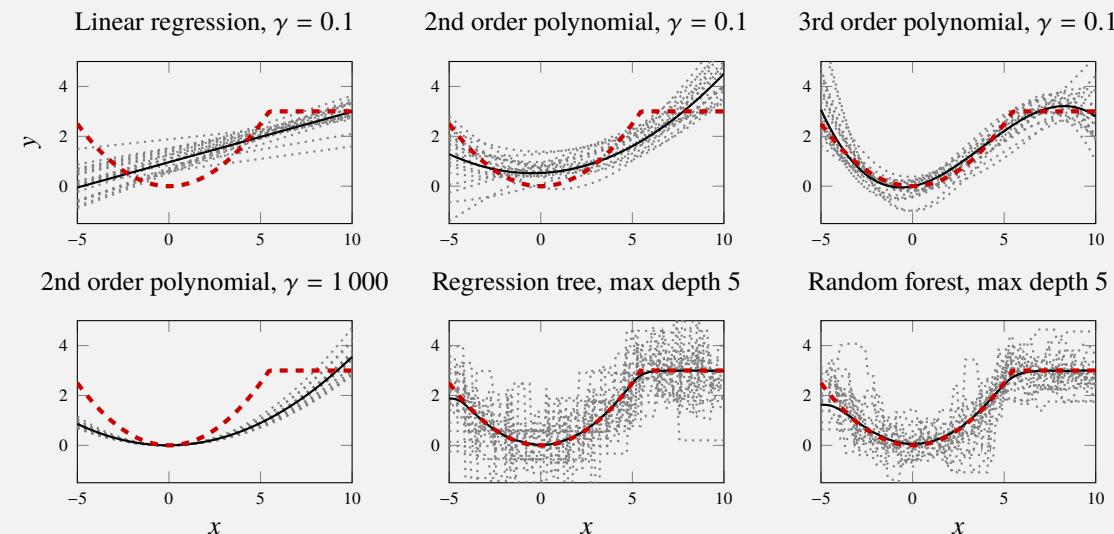
Example 4.4: Bias and variance for a regression problem

We consider the exact same setting as in Example 4.2, but decompose \bar{E}_{new} into bias and variance instead. This gives us the figure below.



There are clear resemblances to Example 4.2, as expected from (4.23). The effect of bagging (random forest) is, however, more clear, namely that it reduces the variance compared to the regression tree with no noteworthy increase in bias.

For another illustration of what bias and variance means, we illustrate some of these cases in more detail. First we plot some of the linear regression models. The dashed red line is the true $f_0(\mathbf{x})$, the dotted gray lines are different models $\hat{y}(\mathbf{x}_*; \mathcal{T})$ learned from different training datasets \mathcal{T} , and the solid black line their mean $\bar{f}(\mathbf{x})$. In these figures, bias is the difference between the dashed red and solid black lines, whereas variance is the spread of the dotted gray lines around the solid black. The variance appears to be roughly the same for all three models, perhaps somewhat smaller for the first order polynomial, whereas the bias is clearly smaller for the higher order polynomials. This can be compared to the figure above.



Comparing the second order polynomial with little ($\gamma = 0.1$) and heavy ($\gamma = 1000$) regularization, it is clear that regularization reduces variance, but also increases bias. Furthermore has random forest a smaller variance than the regression tree, but without any change in the black line $\bar{f}(\mathbf{x})$, and hence no change in bias.

4.5 Evaluation for imbalanced and asymmetric classification problems

We started this chapter by presenting the misclassification as the default error function, but sometimes classification problems are *imbalanced* and/or *asymmetric* and other error functions are more relevant. The typical example is perhaps when binary classification is used to detect the presence of something, such as a disease, an object on the radar, etc. The convention is that $y = 1$ (positive) denotes presence, and $y = -1$ (negative) denotes absence. We say that a problem is

- (i) *imbalanced* if the vast majority of the data points belongs to one class, typically the negative class $y = -1$. This imbalance implies that a (useless) classifier which always predicts $\hat{y}(\mathbf{x}) = -1$ will score very well in terms of misclassification rate (4.1a).
- (ii) *asymmetric* if a missed detection (predicting $\hat{y}(\mathbf{x}) = -1$, when in fact $y = 1$) is considered more severe than a false detection (predicting $\hat{y} = 1$, when in fact $y = -1$), or vice versa. That asymmetry is not taken into account in the misclassification rate (4.1a).

If a classification problem is imbalanced and/or asymmetric, the misclassification rate, and hence E_{new} as we defined it in Section 4.1, is not always very useful for evaluating a classifier. The concepts of generalization gap and bias-variance tradeoff are still applicable, but for those problems other metrics than misclassification rate might be more relevant. We will now briefly introduce some useful tools for asymmetric and/or imbalanced classification problems, which can be used either as an alternative or as a complement to the misclassification error function (4.1a) in E_{new} . For simplicity we consider the binary problem and use a hold-out validation dataset approach, but the ideas can be extended to the multiclass problem as well as to k -fold cross-validation.

The confusion matrix and the F_1 score

If we learn a binary classifier and evaluate it on a hold-out validation dataset, a simple yet useful way to inspect the performance more than just computing $E_{\text{hold-out}}$ is a *confusion matrix*. By separating the validation data in four groups depending on y (the actual output) and $\hat{y}(\mathbf{x})$ (the output predicted by the classifier), we can make the confusion matrix,

	$y = -1$	$y = 1$	<i>total</i>
$\hat{y}(\mathbf{x}) = -1$	True neg (TN)	False neg (FN)	N^*
$\hat{y}(\mathbf{x}) = 1$	False pos (FP)	True pos (TP)	P^*
<i>total</i>	N	P	n

Of course, TN, FN, FP, TP (and also N^* , P^* , N, P and n) should be replaced by the actual numbers, as in Example 4.5. The confusion matrix provides a quick and informative overview of the characteristics of a classifier. For asymmetric problems, it is important to distinguish between false positive (FP, also called *type I error*) and false negative (FN, also called *type II error*). Ideally they both should be 0, but in practice one usually has to decide which ones are considered most important.

There is also a wide body of terminology related to the confusion matrix, which is summarized in Table 4.1. Some particularly common terms are the *recall* (TP/P) and the *precision* (TP/P^*). Recall describes how big proportion among the true positive points that are predicted as positive. A high recall (close to 1) is good, and a low recall (close to 0) indicates a problem with false negatives. Precision describes what the ratio of true positive points are among the ones predicted as positive. A high precision (close to 1) is good, and a low precision (close to 0) indicates a problem with false positives.

For imbalanced (but not asymmetric) problems where the negative class $y = -1$ is the most common class, it makes sense to summarize the precision and recall by their harmonic mean into the F_1 score. Instead of using the misclassification rate for imbalanced problems, the F_1 score can be used instead. Note, however, the scale for the F_1 score: the perfect classifier attains F_1 score equal to 1 whereas 0 is worst (for misclassification rate, 1 is worst).

Ratio	Name
FP/N	False positive rate, Fall-out, Probability of false alarm
TN/N	True negative rate, Specificity, Selectivity
TP/P	True positive rate, Sensitivity, Power, <i>Recall</i> , Probability of detection
FN/P	False negative rate, Miss rate
TP/P*	Positive predictive value, <i>Precision</i>
FP/P*	False discovery rate
TN/N*	Negative predictive value
FN/N*	False omission rate
P/n	Prevalence
(FN+FP)/n	<i>Misclassification rate</i>
(TN+TP)/n	Accuracy
2TP/(P*+P)	F_1 score

Table 4.1: Some common terms related to the quantities (TN, FN, FP, TP) in the confusion matrix. The terms written in italics are discussed in the text.

Example 4.5: The confusion matrix in thyroid disease detection

The thyroid is an endocrine gland in the human body. The hormones it produces influences the metabolic rate and the protein synthesis, and thyroid disorders may have serious implications. We consider the problem of detecting thyroid diseases, using the dataset provided by UCI Machine Learning Repository (Dheeru and Karra Taniskidou 2017). The dataset contains 7200 data points, each with 21 medical indicators as inputs (both qualitative and quantitative). It also contains the qualitative diagnosis {normal, hyperthyroid, hypothyroid}, which we convert into the binary problem with the output classes {normal, not normal}. The dataset is split into a training and hold-out validation part, with 3772 and 3428 samples respectively. The problem is imbalanced since only 7% of the data points are not normal, and possibly asymmetric if false negatives (not indicating the disease) are considered more problematic than false positives (falsely indicating the disease). We train a logistic regression classifier and evaluate it on the validation dataset (using the default decision threshold $r = 0.5$, see (3.36)). We obtain the confusion matrix

		$y = \text{normal}$	$y = \text{not normal}$
		3177	237
$\hat{y}(\mathbf{x}) = \text{normal}$	1	13	
$\hat{y}(\mathbf{x}) = \text{not normal}$			

Most validation data points are correctly predicted as normal, but a large part of the not normal data is also falsely predicted as normal. This might indeed be undesired in the application. The misclassification rate is 0.069 and the F_1 score is 0.106. (The useless predictor of always predicting normal has a very similar misclassification rate of 0.073, but worse F_1 score 0.)

To change the picture, we change the threshold to $r = 0.15$, and obtain new predictions with the following confusion matrix instead:

		$y = \text{normal}$	$y = \text{not normal}$
		3067	165
$\hat{y} = \text{normal}$	111	85	
$\hat{y} = \text{not normal}$			

This change gives more true positives (85 instead of 13 patients are correctly predicted as not normal), but this happens at the expense of more false positives (111, instead of 1, patients are now falsely predicted as not normal). As expected, the misclassification rate is now higher (worse) at 0.081, but the F_1 score is higher (better) at 0.381. Remember, however, that the F_1 score does not take the asymmetry into account, but only the imbalance. We have to decide ourselves whether this classifier is a good tradeoff between the false negative and false positive rates, by considering which type of error has the most severe consequences.

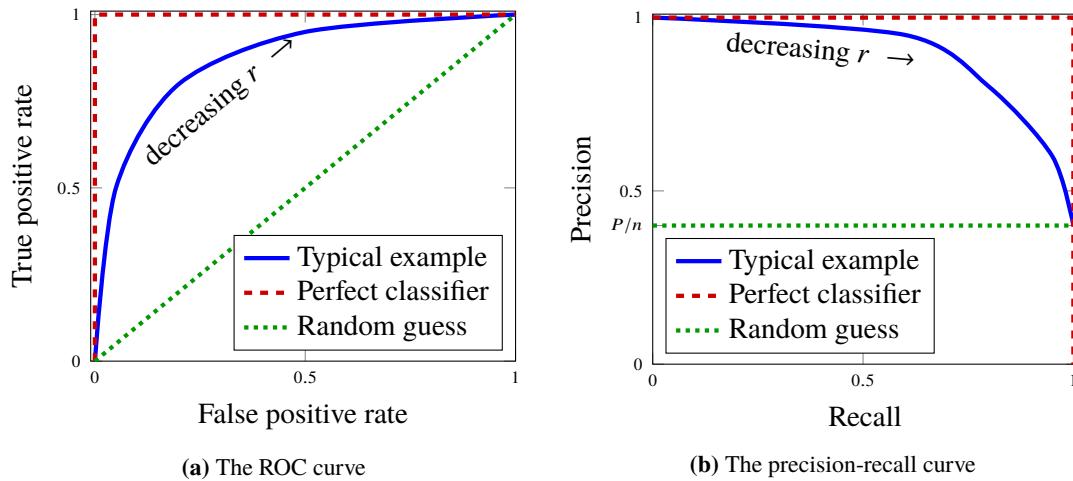


Figure 4.7: The ROC (left) and the precision-recall (right) curve. Both plots summarize the performance of a classifier for *all* decision thresholds r (see (3.36)), but the ROC curve is most relevant for balanced problems, whereas the precision-recall curve is more informative for imbalanced problems.

ROC and precision-recall curve

As suggested by the example above, the tuning of the threshold r in (3.36) can be crucial for the performance in binary classification. If we want to compare different classifiers for a certain problem without specifying a certain decision threshold r , the *ROC curve* can be useful. The abbreviation ROC means "receiver operating characteristics", and is due to its history from communications theory.

To plot an ROC curve, the recall/true positive rate (TP/P , a large value is good) is drawn against the false positive rate (FP/N , a small value is good) for all values of $r \in [0, 1]$. The curve typically looks as shown in Figure 4.7a. An ROC curve for a perfect classifier (always predicting the correct value for all $r \in (0, 1)$) touches the upper left corner, whereas a classifier which only assigns random guesses gives a straight diagonal line.

A compact summary of the ROC curve is the *area under the ROC curve*, *ROC-AUC*. From Figure 4.7a, we conclude that a perfect classifier has ROC-AUC 1, whereas a classifier which only assigns random guesses has ROC-AUC 0.5. The ROC-AUC is thus summarizing the performance of a classifier for *all* possible values of the decision threshold r .

We previously discussed how the misclassification rate might be misleading for imbalanced problems, but we should consider the precision, recall and F_1 score instead. Similarly might the ROC curve be misleading for imbalanced problems, but instead the precision-recall curve can (for imbalanced problems where $y = -1$ is the most common class) be more useful. As the name suggests, the precision-recall curve plots the precision (TP/P^* , a large value is good) against the recall (TP/P , a large value is good) for all values of $r \in [0, 1]$, much like the ROC curve. The precision-recall curve for the perfect classifier touches the upper right corner, and a classifier which only assigns random guesses gives a horizontal line with height P/n , as shown in Figure 4.7b.

Also for the precision-recall curve, we can define *area under the curve*, *precision-recall AUC*. The best possible precision-recall AUC is 1, and the classifier which only makes random guesses has precision-recall AUC equal to P/n .

5 Learning parametric models

This chapter revolves around three different topics, namely loss functions, regularization and optimization. We have touched upon all of them already, mostly in connection to the parametric models in Chapter 3, linear and logistic regression. These topics are however central for many more supervised machine learning methods and deserve a dedicated discussion, which we will give now. The understanding of loss functions is particularly helpful when we introduce boosting (Chapter 7) and support vector machines (Chapter 8) later on. Regularization is crucial in preventing overfit (Chapter 4) and should be in the back of the mind of any machine learning engineer. Finally, it is always helpful to be familiar with the main ideas of optimization, in particular since they are a key enabler for deep learning (Chapter 6).

5.1 Loss functions

Most parametric regression and classification models, including linear and logistic regression (Chapter 3), are trained by minimizing a cost function $J(\theta)$. The cost function is the average of the loss function¹ L evaluated on the training data,

$$\widehat{\theta} = \arg \min_{\theta} \underbrace{\frac{1}{n} \sum_{i=1}^n L(\widehat{y}(\mathbf{x}; \theta), y_i)}_{\text{cost function } J(\theta)}.$$
 (5.1)

The choice of loss function is, in the end, a user choice. Different loss functions give rise to different solutions $\widehat{\theta}$, which in turn results in models with different characteristics. There is never a “right” or “wrong” loss function, but one loss functions may of course perform better than another on a given problem. Certain combinations of models and loss functions have proven particularly fruitful and have historically been branded as different methods carrying different names. For example, the term “linear regression” most often refers to the combination of a linear-in-the-parameter model and the squared error loss, whereas the term “support vector classification” (Chapter 8) refers to a linear-in-the-parameter model trained using the hinge loss. In this section, however, we provide a general discussion about different loss functions and their properties, without connections to a specific method.

One important aspect of a loss function is its *robustness*. Robustness is tightly connected to outliers, meaning spurious data points that do not describe the relationship we are interested in modeling. If outliers in the training data only have a minor impact to the learned model, we say that the loss function is robust. Conversely, a loss function is not robust if the outliers have a major impact to the learned model. Robustness is therefore a very important property in applications where the training data is contaminated with outliers. Some of the common “default” loss-functions, such as the squared error loss, are unfortunately not particularly robust, and it is therefore important to make an informed user choice. Robustness is however not a binary property, instead loss functions can be robust to a greater or smaller degree.

Another important property of loss functions for binary classification is their so-called *asymptotic minimizers*. The asymptotic minimizer refers to the function that minimizes the loss function when $n \rightarrow \infty$, and reveals whether the model will possibly learn the full conditional class probability $p(y = 1 | \mathbf{x})$ or only the decision boundary. We will discuss this in more detail later.

¹As you might already have noticed, the arguments to the loss function (here \widehat{y} and y_*) varies with context. This is unfortunate but unavoidable, since different loss functions are formulated in terms of different quantities, for example the prediction \widehat{y} , the predicted conditional class probability $g(\mathbf{x})$, the classifier margin, etc.

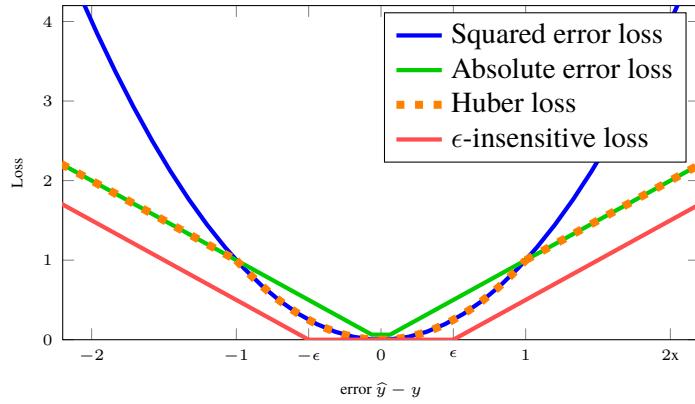


Figure 5.1: The loss functions for regression presented in the text, each as a function of the error $\hat{y} - y$.

Loss functions for regression

In Chapter 3 we introduced the *squared error loss*

$$L(\hat{y}, y) = (\hat{y} - y)^2, \quad (5.2)$$

which is the default choice for linear regression since it simplifies the training to only solving the normal equations. The squared error loss is often used also for other regression models, such as neural networks. Another common, but more robust, loss function is the *absolute error loss*,

$$L(\hat{y}, y) = |\hat{y} - y|. \quad (5.3)$$

In Chapter 3 we introduced the maximum likelihood motivation of the squared error loss by assuming that the output y is measured with additive uncorrelated noise ε from a Gaussian distribution. We can similarly motivate the absolute error loss by instead assuming ε to have a Laplace distribution, $\varepsilon \sim \mathcal{L}(0, b_\varepsilon)$. This statistical perspective is one way to understand the fact that the absolute error loss is more robust (less sensitive to outliers) than the squared error loss, since the Laplace distribution has a thicker tail compared to the Gaussian distribution, see Figure 5.1. The Laplace distribution therefore encodes sporadic large noise values (i.e., outliers) as more probable, compared to the Gaussian distribution.

If we have statistical knowledge about how the noise ε in a regression problem behaves (for example a skewed normal distribution), it is perfectly possible to use that information and make a maximum likelihood derivation like (3.17)-(3.21) to derive a loss function for regression tailored to that information. There are, however, a few additional off-the-shelf loss functions commonly used for regression which are less natural to derive from a maximum likelihood perspective.

It is sometimes argued that the squared error loss is a good choice because of its quadratic shape which penalizes small errors ($\varepsilon < 1$) less than linearly. After all, the Gaussian distribution appears (at least approximately) quite often in nature. However, the quadratic shape for large errors ($\varepsilon > 1$) is the reason for its non-robustness, and the *Huber loss* has therefore been suggested as a hybrid between the absolute loss and squared error loss:

$$L(\hat{y}, y) = \begin{cases} (\hat{y} - y)^2 & \text{if } |\hat{y} - y| < 1, \\ |\hat{y} - y| & \text{otherwise.} \end{cases} \quad (5.4)$$

Another extension to the absolute error loss is the *ϵ -insensitive loss*,

$$L(\hat{y}, y) = \begin{cases} 0 & \text{if } |\hat{y} - y| < \epsilon, \\ |\hat{y} - y| - \epsilon & \text{otherwise,} \end{cases} \quad (5.5)$$

where ϵ is a user-chosen design parameter. This loss places a “tolerance” of width 2ϵ around the true y and behaves like the absolute error loss outside this region. Furthermore the robustness properties of the

ϵ -insensitive loss are very similar to those of the absolute error loss. The ϵ -insensitive loss turns out to be a useful loss function for support vector regression in Chapter 8. We illustrate all these loss functions for regression in Figure 5.1.

Loss functions for binary classification

An intuitive loss function for binary classification is provided by the *misclassification loss*,

$$L(\hat{y}, y) = \begin{cases} 1 & \text{if } y_\star = y, \\ 0 & \text{otherwise.} \end{cases} \quad (5.6)$$

However, even though a small misclassification loss often is the ultimate goal in practice, this loss function is rarely used when training models. There are (at least) two reasons for this. The most apparent reason is that the resulting cost function is a piecewise constant function, which is a poor optimization objective since the gradient is zero everywhere, except where it is undefined. Furthermore, a more subtle reason is that for many models the final prediction \hat{y} does not reveal all aspects of the classifier. Thinking intuitively in terms of decision boundaries, we may prefer to not have the decision boundary close to the training data points, but instead push the boundary further away to have some “safety leeway” if possible. The misclassification loss (5.6) cannot achieve this, since it only encourages training data points to be on the right side of the decision boundary.

For binary classifiers that predicts conditional class probabilities $p(y = 1 | \mathbf{x})$ in terms of a function $g(\mathbf{x})$, the *cross entropy loss* is, as introduced in Chapter 3,

$$L(g(\mathbf{x}), y) = \begin{cases} \ln g(\mathbf{x}) & \text{if } y = 1, \\ \ln(1 - g(\mathbf{x})) & \text{if } y = -1. \end{cases} \quad (5.7)$$

This loss was derived from a maximum likelihood perspective, but unlike regression (where we had to specify a distribution for ε) there are no user choices left in the cross entropy loss, other than what model to use for $g(\mathbf{x})$.

The misclassification and cross entropy loss are however not the only loss functions for binary classification. To introduce an entire family of loss functions for logistic regression, let us first introduce the concept of *margins* in binary classification. Many binary classifiers $\hat{y}(\mathbf{x})$ can be constructed by thresholding some real-valued function $c(\mathbf{x})$ at 0. That is, we can write

$$\hat{y}(\mathbf{x}) = \text{sign}\{c(\mathbf{x})\}. \quad (5.8)$$

Logistic regression, for example, can be brought into this form by simply using $c(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x}$. The decision boundary for any classifier on the form (5.8) is given by the values of \mathbf{x} for which $c(\mathbf{x}) = 0$. To simplify our discussion we will assume that none of the data points fall exactly on the decision boundary (which always gives rise to an ambiguity). This will imply that we can assume that $\hat{y}(\mathbf{x})$ as defined above is always either -1 or $+1$.

Based on the function $c(\mathbf{x})$, we say that

$\text{the margin of a classifier is } y \cdot c(\mathbf{x}).$

(5.9)

It follows that if y and $c(\mathbf{x})$ have the same sign, meaning the classification is correct, then the margin is positive. Analogously, for an incorrect classification y and $c(\mathbf{x})$ will have different signs and the margin is negative. The margin can be viewed as a measure of certainty in a prediction, where data points with small margins in some sense (not necessarily Euclidian) are close to the decision boundary. The margin plays a similar role for binary classification as the error $\hat{y} - y$ does for regression.

We can now define loss functions for binary classification in terms of the margin, by assigning a small loss to positive margins (correct classifications) and a large loss to negative margins (misclassifications). We can, for instance, re-formulate the logistic loss (3.34) as

$$L(y \cdot c(\mathbf{x})) = \ln(1 + \exp(-y \cdot c(\mathbf{x}))), \quad (5.10)$$

with $c(\mathbf{x}) = \theta^\top \mathbf{x}$. That is, learning θ by minimizing the loss function (5.10) and making predictions according to (5.8) is equivalent to logistic regression as described in Chapter 3, except for the fact that we have lost the notion of a conditional class probability estimate $g(\mathbf{x})$ and only have a “hard” prediction $\hat{y}(\mathbf{x}_\star)$. We will, however, recover the class probability estimate later when we discuss the asymptotic minimizer of (5.10).

We can also formulate the misclassification loss in terms of the margin,

$$L(y \cdot c(\mathbf{x})) = \begin{cases} 1 & \text{if } y \cdot c(\mathbf{x}) < 0, \\ 0 & \text{otherwise.} \end{cases} \quad (5.11)$$

The benefit of formulating a loss function in terms of the margin is that it is easy to come up with new loss functions. Let us now introduce the *exponential loss* defined as

$$L(y \cdot c(\mathbf{x})) = \exp(-y \cdot c(\mathbf{x})), \quad (5.12)$$

which turns out to be a useful loss function when we later will derive the AdaBoost in Chapter 7 due to its properties. The downside of the exponential loss is that it is not particularly robust against outliers, compared to the logistic loss. We also have the *hinge loss*, which we will use later for support vector classification in Chapter 8,

$$L(y \cdot c(\mathbf{x})) = \begin{cases} 1 - y \cdot c(\mathbf{x}) & \text{for } y \cdot c(\mathbf{x}) \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (5.13)$$

As we will see in Chapter 8, the hinge loss has an attractive so-called support-vector property. However, as we soon will see its asymptotic minimizer is somewhat unfortunate. As a remedy to that, one may instead consider the *squared hinge loss*

$$L(y \cdot c(\mathbf{x})) = \begin{cases} (1 - y \cdot c(\mathbf{x}))^2 & \text{for } y \cdot c(\mathbf{x}) \leq 1, \\ 0 & \text{otherwise,} \end{cases} \quad (5.14)$$

which on the other hand is much less robust than the hinge loss. A more elaborate alternative might therefore be the *Huberized squared hinge loss*

$$L(y \cdot c(\mathbf{x})) = \begin{cases} -4y \cdot c(\mathbf{x}) & \text{for } y \cdot c(\mathbf{x}) \leq -1, \\ (1 - y \cdot c(\mathbf{x}))^2 & \text{for } -1 \leq y \cdot c(\mathbf{x}) \leq 1, \\ 0 & \text{otherwise,} \end{cases} \quad (\text{squared hinge loss}) \quad (5.15)$$

whose name refers to its similarities to the Huber loss for regression, namely that the quadratic function is replaced with a linear function for margins < -1 . The three loss functions presented above are all particularly interesting for support vector classification, due to the fact that they are all exactly 0 for margins > 1 .

We summarize this cascade of loss functions for binary classification in Figure 5.2, which illustrates all these losses as a function of the margin.

We have already made some claims about robustness. Let us motivate them using Figure 5.2. One characterization of an outlier is as a data point on the “wrong” side of and far away from the decision boundary. In a margin perspective, that is equivalent to a large negative margin. The robustness of a loss function is therefore tightly connected to the shape of the loss function for large negative margins. The steeper slope and heavier penalization of large negative margins, the more sensitive it is to outliers. We can therefore tell from Figure 5.2 that the exponential loss is most sensitive to outliers, due to its exponential left asymptote, whereas the squared hinge loss is somewhat more robust with a quadratic asymptote instead. However, even more robust are the Huberized squared hinge loss, the hinge loss and the logistic loss which all have an asymptotic behavior which is linear. Most robust is the misclassification loss, but as already discussed that loss has other disadvantages.

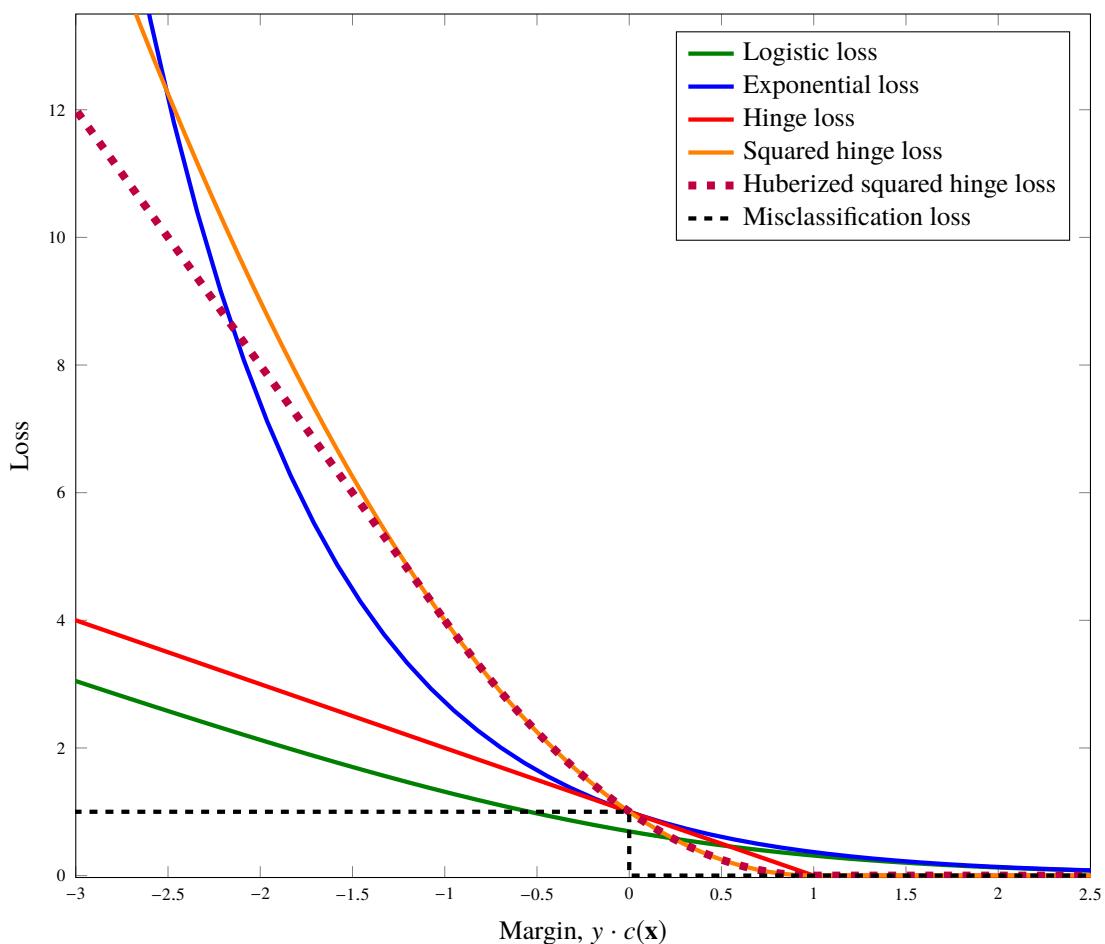


Figure 5.2: Comparison of some common loss functions for classification, plotted as a function of the margin.

The asymptotic minimizer in binary classification

As mentioned earlier the *asymptotic minimizer* is a useful theoretical concept un understanding loss functions for classification. The asymptotic minimizer is the function ($g(\mathbf{x})$ or $c(\mathbf{x})$, depending on the parametrization of the loss function) which minimizes the cost function when the number of training data $n \rightarrow \infty$, hence the name asymptotic. It contains valuable information whether the model will, eventually as $n \rightarrow \infty$, learn the conditional class probability $p(y = 1 | \mathbf{x})$ correctly or not.

In the following discussion we will make heavy use of the notion of a ground truth conditional class probability $p(y = 1 | \mathbf{x})$, which we assume the training data is drawn from. Learning $p(y = 1 | \mathbf{x})$, or some aspect of it, is the ultimate goal with binary classification in supervised machine learning. We assume that the model is flexible enough such that it can describe $p(y = 1 | \mathbf{x})$ well. This assumption is not always met in practice, and we will discuss the implications of this assumption at the end.

There is a notion of a *strictly proper* loss function. A strictly proper loss function is a loss function for which the asymptotic minimizer (i) is unique and (ii) is an invertible transformation of $p(y = 1 | \mathbf{x})$. That is, a strictly proper loss function forces the model to learn the full² distribution $p(y = 1 | \mathbf{x})$ if given enough training data. All models can be used for predicting classes $\hat{y}(\mathbf{x}_*)$, but only a model trained with a strictly proper loss function can be used for predicting conditional class probabilities (like $g(\mathbf{x})$ in logistic regression). Deriving the asymptotic minimizer of a loss function is most often a straightforward calculation, but for brevity we do not include the derivations here.

Starting with the binary maximum likelihood loss (5.7), its asymptotic minimizer can be shown to be $g(\mathbf{x}) = p(y = 1 | \mathbf{x})$. In other words, when $n \rightarrow \infty$ the loss function (5.7) trains the model such that $g(\mathbf{x})$ becomes the conditional class probability. Since we derived (5.7) with this intention in Chapter 3, this was indeed expected.

The asymptotic minimizer of the logistic loss (5.10) is $c(\mathbf{x}) = \ln \frac{p(y=1 | \mathbf{x})}{1-p(y=1 | \mathbf{x})}$. This is an invertible expression and hence the logistic loss is a strictly proper loss function. By inverting $c(\mathbf{x})$ we obtain $p(y = 1 | \mathbf{x}) = \frac{\exp c(\mathbf{x})}{1+\exp c(\mathbf{x})}$ which shows how conditional class probability predictions can be obtained from $c(\mathbf{x})$. (With the “margin formulation” of logistic regression, we seemingly lost the class probability predictions $g(\mathbf{x})$. We have now recovered it.)

For the misclassification loss (5.11) the asymptotic minimizer is not unique but only fulfills

$$\text{sign}\{c(\mathbf{x})\} = \begin{cases} 1 & \text{if } p(y = 1 | \mathbf{x}) > 0.5, \\ -1 & \text{if } p(y = 1 | \mathbf{x}) < 0.5. \end{cases}$$

From this expression we see that a model trained with the misclassification loss has the correct decision boundary (when $n \rightarrow \infty$). However, the misclassification loss does not train the model such that $p(y = 1 | \mathbf{x})$ can be inferred from $c(\mathbf{x})$, since the asymptotic minimizer is not an invertible transformation of it. The misclassification loss is therefore not a strictly proper loss function, and $c(\mathbf{x})$ can not be used for predicting conditional class probabilities. This is yet another downside of the misclassification loss.

Turning to the exponential loss (5.12), its asymptotic minimizer is $c(\mathbf{x}) = \frac{1}{2} \ln \frac{p(y=1 | \mathbf{x})}{1-p(y=1 | \mathbf{x})}$, almost like the logistic loss. The exponential loss is therefore also strictly proper, and $c(\mathbf{x})$ can be transformed and used for predicting conditional class probabilities.

The asymptotic minimizer of the hinge loss is

$$c(\mathbf{x}) = \begin{cases} 1 & \text{if } p(y = 1 | \mathbf{x}) > 0.5, \\ -1 & \text{if } p(y = 1 | \mathbf{x}) < 0.5. \end{cases}$$

This is a non-invertible transformation of $p(y = 1 | \mathbf{x})$, which means that $p(y = 1 | \mathbf{x})$ is not possible to infer from the asymptotic minimizer $c(\mathbf{x})$. This implies that a classifier learned using hinge loss (such as support vector classification, Section 8.5) can *not* predict conditional class probabilities.

The squared hinge loss, on the other hand, is a strictly proper loss function, since its asymptotic minimizer is $c(\mathbf{x}) = 2p(y | \mathbf{x}) - 1$. This also holds for the Huberized square hinge loss. Recalling our

²Loss functions that are proper, but not strictly proper, may for example learn only the correct decision boundary, that is $p(y = 1 | \mathbf{x}) = 0.5$, which is not a full characterization of $p(y = 1 | \mathbf{x})$.

robustness discussion, we see that by squaring the hinge loss we make it strictly proper but at the same time we impact its robustness. However, the “huberization” (replacing the quadratic curve with a linear one for margins < -11) improves the robustness while keeping the property of being strictly proper.

We have now seen that some (but not all) of the loss functions are strictly proper, meaning they could potentially predict conditional class probabilities correctly. However, this is only under the assumption that the model is sufficiently flexible such that $g(\mathbf{x})$ or $c(\mathbf{x})$ actually can take the shape of the asymptotic minimizer. This is possibly problematic; remember that $c(\mathbf{x})$ is a linear function in logistic regression, whereas $p(y = 1 | \mathbf{x})$ can be almost arbitrarily complicated in real world applications. It is therefore not sufficient to use a strictly proper loss function in order to predict conditional class probabilities, but our model also has to be flexible enough. This discussion is also only valid in the limit as $n \rightarrow \infty$. However, in practice n is always bounded, and we may ask how large n has to be for a flexible enough model to at least approximately learn the asymptotic minimizer? Unfortunately we cannot give any general numbers, but following the same principles as the overfit discussion in Chapter 4, the more flexible the model the larger n is required. If n is not large enough, the predicted conditional class probabilities tend to be “overfitted” to the training data. We therefore unfortunately have to conclude that the concept of the asymptotic minimizer only give limited guidance on how to interpret the learned $c(\mathbf{x})$ in practice.

Multiclass classification

So far we have only discussed the binary classification problem with $M = 2$. The maximum likelihood loss is indeed possible to generalize to the multiclass problem, that is $M > 2$, as we did in Chapter 3 for logistic regression. The generalization of the other loss functions is, however, a more complicated problem and would require a generalization of the margin to the multiclass problem to start with. That is possible but we do not discuss it in this book. A more pragmatic approach is instead to reformulate the problem as several binary problems. This re-formulation can be done using either a one-versus-rest or one-versus-one scheme.

The one-versus-rest (or one-versus-all or binary relevance) idea is to train M binary classifiers. Each classifier in this scheme is trained for predicting one class against all the other classes. To make a prediction for a test data point, all M classifiers are used, and the class which, for example, is predicted with the largest margin is taken as the predicted class. This approach is a pragmatic solution, which may turn out to work well for some problems.

The one-versus-one idea is instead to train one classifiers for each pair of classes. If there are M classes in total, there are $\frac{1}{2}M(M - 1)$ such pairs. To make a prediction each classifier predicts either of its two classes, and the class which overall obtains most “votes” is chosen as the final prediction. The predicted margins can be used to break a tie if that happens. Compared to one-versus-rest, the one-versus-one approach has the disadvantage of involving $\frac{1}{2}M(M - 1)$ classifiers, instead of only M . On the other hand each of these classifiers is trained on much smaller datasets (only the data points that belong to either of the two classes) compared to one-versus-rest which uses the entire original training dataset for all M classifiers.

5.2 Regularization

We will now take a closer look at regularization, which was briefly introduced in Section 3.3 as a useful tool for avoiding overfit if the model was too flexible, such as a polynomial of high degree. We have also discussed thoroughly in Chapter 4 the need for tuning the model flexibility, which is the purpose of regularization. Finding the right level of flexibility, and thereby avoiding overfit, is very important in practice.

Regularization is applicable to any parametric model, and the idea is to ‘keep the parameters $\hat{\theta}$ small unless the data really convinces us otherwise’, or alternatively ‘if a model with small values of the parameters $\hat{\theta}$ fits the data almost as well as a model with larger parameter values, the one with small parameter values should be preferred’. We will now have a closer look at the so-called L^2 and L^1 regularization.

L^2 regularization

The L^2 regularization (also known as *Tikhonov regularization*, *ridge regression* and *weight decay*) amounts to adding an extra penalty term $\|\theta\|_2^2$ to the cost function. Linear regression with squared error loss and L^2 regularization, as an example, amounts to solving

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \|\mathbf{X}\theta - \mathbf{y}\|_2^2 + \lambda \|\theta\|_2^2. \quad (5.16)$$

By choosing the regularization parameter $\lambda \geq 0$, a trade-off between the original cost function (fitting the training data as well as possible) and the regularization term (keeping the parameters $\hat{\theta}$ close to zero) is made. In the setting $\lambda = 0$ we recover the original least squares problem (3.12), whereas $\lambda \rightarrow \infty$ will force all parameters $\hat{\theta}$ to 0. A good choice of λ is usually somewhere in between and depends on the actual problem, and can be selected using cross-validation.

It is actually possible to derive a version of the normal equations for (5.16), namely

$$(\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}_{p+1})\hat{\theta} = \mathbf{X}^\top \mathbf{y}, \quad (5.17)$$

where \mathbf{I}_{p+1} is the identity matrix of size $(p+1) \times (p+1)$. For $\lambda > 0$, the matrix $\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}_{p+1}$ is always invertible, and we have the closed form solution

$$\hat{\theta} = (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}_{p+1})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (5.18)$$

This also reveals another reason for using regularization in linear regression, namely if $\mathbf{X}^\top \mathbf{X}$ is not invertible. When $\mathbf{X}^\top \mathbf{X}$ is not invertible, the ordinary normal equations (3.13) have no unique solution $\hat{\theta}$, whereas the L^2 -regularized version always has the unique solution (5.18) if $\lambda > 0$.

L^1 regularization

With L^1 regularization (also called *LASSO*, an abbreviation for Least Absolute Shrinkage and Selection Operator), the penalty term $\|\theta\|_1$ is added to the cost function. Here $\|\theta\|_1$ is the 1-norm or ‘taxicab norm’ $\|\theta\|_1 = |\theta_0| + |\theta_1| + \dots + |\theta_p|$. The L^1 regularized cost function for linear regression (with squared error loss) (3.12) then becomes

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \|\mathbf{X}\theta - \mathbf{y}\|_2^2 + \lambda \|\theta\|_1. \quad (5.19)$$

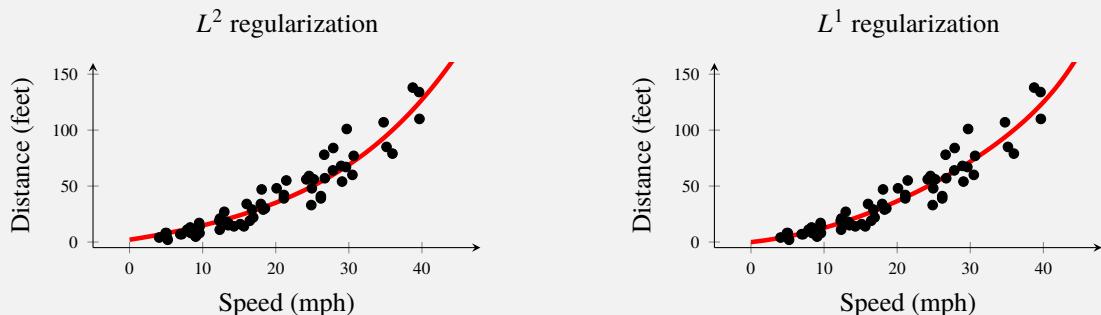
Contrary to linear regression with L^2 regularization (3.48), there is no closed-form solution available for (5.19). However, as we will see in Section 5.3, it is possible to design an efficient numerical optimization algorithm for solving (5.19).

As for L^2 regularization, the regularization parameter λ has to be chosen by the user, and has a similar meaning: $\lambda = 0$ gives the ordinary least squares solution and $\lambda \rightarrow \infty$ gives $\hat{\theta} = 0$. Between these extremes,

however, L^1 and L^2 tend to give different solutions. Whereas L^2 regularization pushes all parameters towards small values (but not necessarily exactly zero), L^1 tends to favor so-called *sparse* solutions where only a few of the parameters are non-zero, and the rest are exactly zero. Thus, L^1 regularization can effectively ‘switch off’ some inputs (by setting the corresponding parameter θ_k to zero) and it can therefore be used as an input (or feature) selection method.

Example 5.1: Regularization for car stopping distance

Consider again Example 2.2 with the car stopping distance regression problem. We use the 10th order polynomial that was considered meaningless in Example 3.5 and apply L^2 and L^1 regularization to it, respectively. With manually chosen λ , we obtain the following models



Both models suffer less from overfit than the non-regularized 10th order polynomial in Example 3.5. The two models here are, however, not identical. Whereas all parameters are relatively small but non-zero in the L^2 -regularized model (left panel), only 4 (out of 11) parameters are non-zero in the L^1 -regularized model (right panel). It is typical for L^1 regularization to give sparse models, where some parameters are set to exactly zero.

General cost function regularization

The L^1 and L^2 regularization are two commonly used regularization methods, and they are both formulated as additions to the cost function. Together they suggest a more general regularization scheme

$$\widehat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \underbrace{J(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y})}_{(i)} + \underbrace{\lambda}_{(iii)} \underbrace{R(\boldsymbol{\theta})}_{(ii)}. \quad (5.20)$$

This expression contains three important elements:

- (i) the cost function, which encourages a good fit to training data,
- (ii) the regularization term, which encourages small parameter values, and
- (iii) the regularization parameter λ , which determines the trade-off between (i) and (ii).

The regularization term can be designed in many ways. As a combination of the L^1 and L^2 terms, one option is $R(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1 + \|\boldsymbol{\theta}\|_2^2$, which often is referred to as elastic net regularization. Regardless of the exact expression of the regularization term, its purpose is to encourage small parameter values and thereby decrease the flexibility of the model. As we discussed in depth in Chapter 4, a too flexible (or complex) model may overfit to training data and not be able to generalize well to previously unseen test data. By using regularization we get a systematic and practically useful way of controlling the model flexibility and thereby counteracting overfit.

Implicit regularization

Any supervised machine learning method that is trained by minimizing a cost function can be regularized as (5.20). There are, however, alternative ways to achieve similar effects. One such example of *implicit*

5 Learning parametric models

regularization is early stopping. Early stopping is applicable to any method that is trained using numerical optimization, and amounts to aborting the optimization before it has reached the minimum of the cost function. Even though it might appear counter-intuitive to pre-maturely abort an optimization procedure, early stopping has shown to be a of practical importance to avoid overfit for some models, most notably deep learning Chapter 6.

5.3 Parameter optimization

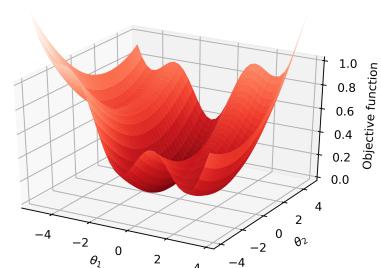
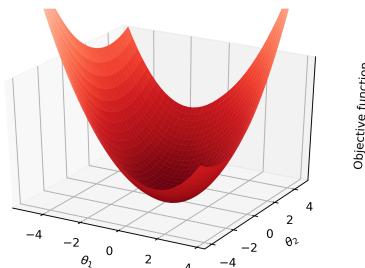
Many supervised machine learning methods, including linear and logistic regression, involves one (or more) optimization problems, such as (3.12), (3.35) or (5.19). It is therefore good to be familiar with some of the main strategies for how to solve optimization problems fast. Starting in the optimization problems from linear and logistic regression, we will introduce the ideas behind some common methods.

Optimization is about finding the minimum (or maximum) of an *objective function*. Since the maximization problem can be formulated as minimization of the negative objective function, it is sufficient to describe only minimization. When we use optimization for training models in machine learning, that objective function is the cost function J . Optimization is, however, often used also in combination with cross-validation (Chapter 4) for tuning hyperparameters, such as selecting the regularization parameter λ by minimizing the prediction error on a held-out validation dataset.

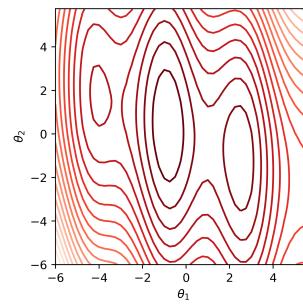
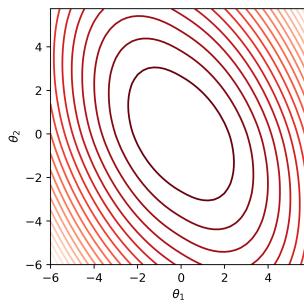
An important class of objective functions are *convex* functions. Optimization is often easier to carry out for convex objective functions, and it is a general advise to spend some extra effort to consider whether a non-convex optimization problem can be re-formulated into a convex problem (which sometimes, but not always, is possible). The most important property of a convex function, for this discussion, is that a convex function has a unique minimum³, and no other local minima. Examples of convex functions are the cost functions for logistic regression, linear regression and L^1 -regularized linear regression. Examples of non-convex functions will come later in the book, including the cost function for a deep neural network. Two examples of convex and non-convex functions are given in Example 5.2 below.

Example 5.2: Example of objective functions

These figures are examples of what an objective function can look like.



Both examples are functions of a two-dimensional parameter vector $\theta = [\theta_1 \theta_2]^\top$. The left is convex and has a finite unique global minimum, whereas the right is non-convex and has three local minima (of which only one is the global minimum). We will in the following examples illustrate these objective functions using contour plots instead, as shown below.



³The minimum does, however, not have to be finite. Consider for example the exponential function, which is convex but attains its minimum in $-\infty$. Convexity is, however, a relatively strong property, and a function may have only one minimum even if it is not convex.

Time to reflect 5.1: After reading the rest of this book, return here and try to fill out this table, summarizing how optimization is used by the different methods.

Method	What is optimization used for?			What type of optimization?			
	Training	Hyper-parameters	Nothing	Closed-form*	Grid search	Gradient-based	Stochastic gradient
Linear regression							
Linear regression with L2-regularization							
Linear regression with L1-regularization							
Logistic regression							
k-NN							
Trees							
Random forests							
AdaBoost							
Gradient boosting							
Deep learning							
Gaussian processes							
*including coordinate descent							

Optimization using closed-form expressions

For linear regression with squared error loss, training the model amounts to solving the optimization problem (3.12)

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2.$$

As we have discussed, and also prove in Appendix 3.A, the solution (3.14) to this problem can (under the assumption that $\mathbf{X}^\top \mathbf{X}$ is invertible) be derived analytically. If we only once spend some time to efficiently implement (3.14), for example using the Cholesky or QR factorization, we can use that every time we want to train a linear regression model with squared error loss. Each time we use it we know that we have found the optimal solution in a computationally efficient way.

If we instead want to learn the L^1 -regularized version, we instead have to solve (5.19)

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1.$$

This problem can, unfortunately, not be solved analytically. Instead we have to use computer power to solve it, by constructing an iterative procedure for seeking the solution. With a certain choice of such an optimization algorithm, we can make use of some analytical expressions along the way, which turns out to yield an efficient way to solve (5.19). Remember that $\boldsymbol{\theta}$ is a vector containing $p + 1$ parameters we want to learn from the training data. As it turns out, if we seek the minimum for only one of these parameters, say θ_j , while keeping the other parameters fix, we can find the optimum as

$$\arg \min_{\theta_j} \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1 = \text{sign}(t)(|t| - \lambda), \text{ where } t = \sum_{i=1}^n x_{ij}(y_i - \sum_{k \neq j} x_{ik}\theta_k). \quad (5.21)$$

It turns out that making repeated “sweeps” through the vector $\boldsymbol{\theta}$ and updating one-parameter-at-a-time according to (5.21) is a good way to solve (5.19). This type of algorithm, where we update one-parameter-at-a-time, is referred to as *coordinate descent*, and we illustrate it in Example 5.3

It can be shown that the cost function in (5.19) is convex. Only convexity is not sufficient to guarantee that coordinate descent will find its (global) minimum, but for the L^1 -regularized cost function (5.19) it can be shown that coordinate descent actually finds the (global) minimum. In practice we know that we have found the global minimum when no parameters have changed during a full “sweep” of the parameter vector.

It turns out that coordinate descent is a very efficient method for L^1 -regularized linear regression (5.19). The keys are (i) that (5.21) exists and is cheap to compute, and (ii) many updates will simply set $\theta_j = 0$ due to the sparsity of the optimal $\hat{\boldsymbol{\theta}}$. This makes the algorithm fast. For most machine learning optimization problems it can, however, *not* be said that coordinate descent is the preferred method. We will now have a look at some more general families of optimization methods that are widely used in machine learning.

Example 5.3: Coordinate descent

We apply coordinate descent to the objective functions from Example 5.2. For coordinate descent to be an efficient alternative in practice, closed-form solutions for updating one-parameter-at-a-time, similar to (5.21), have to be available.



The figures show how the parameters are updated in the coordinate descent algorithm, for two different initial parameter vectors (blue and green trajectory, respectively). It is clear from the figures that only one parameter is updated each time, which gives the trajectory a characteristic shape. The obtained minimum is marked with a yellow dot. Note how the different initializations lead to different (local) minima in the non-convex (right) case.

Grid search

Most often, we do not have closed-form solutions to the optimization problems we need to solve. Sometimes, we can only evaluate the objective function in certain points, and we have no knowledge about its derivatives. This typically happens when selecting hyperparameters, such as the regularization parameter λ in (5.19). The objective is then to minimize the prediction error for a held-out validation data set (or some other version of cross-validation). That is a very complicated objective function to write down, and it is even harder to find its derivative (in fact, this objective function includes an optimization problem itself—learning $\hat{\theta}$ for a given value of λ), but we can nevertheless evaluate it for any given λ (that is, run the entire learning procedure and see how good the predictions become for the validation data set).

The perhaps simplest way to solve an optimization problem, which also is the idea of grid search, is to “try a few different parameter values, and pick the one which works best”. That is the idea of *grid search*. The term “grid” here refers to some (more or less arbitrary chosen) set of different parameter values to try out, and we illustrated it in Example 5.4.

Although simple to implement, grid search can be computationally quite inefficient, in particular if the parameter vector has a high dimension. As an example, having a grid with a resolution of 10 grid points per dimension (which is a very coarse-grained grid) for a 5-dimensional parameter vector requires $10^5 = 100\,000$ evaluations of the objective function. If possible, one should avoid using grid search for this reason. However, with low-dimensional hyperparameters (in L^1 and L^2 regularization, λ is 1-dimensional, for example), grid search can be feasible. We summarize grid search in Algorithm 7, where we use it for determining a regularization parameter λ .

Algorithm 7: Grid search for regularization parameter λ

Input: Training data $\{\mathbf{x}_i, y_i\}_{i=1}^{n_t}$, validation data $\{\mathbf{x}_j, y_j\}_{j=1}^{n_v}$
Result: $\hat{\lambda}$

- 1 **for** $\lambda = 10^{-3}, 10^{-2}, \dots, 10^3$ **do**
- 2 Learn $\hat{\theta}$ with regularization parameter λ from training data
- 3 Compute error on validation data $E_{\text{val}}(\lambda) \leftarrow \sum_{j=1}^{n_v} (\hat{y}(\mathbf{x}_j; \hat{\theta}) - y_j)^2$
- 4 **end**
- 5 **return** $\hat{\lambda}$ as $\arg \min_{\lambda} E_{\text{val}}(\lambda)$

Example 5.4: Grid search

We apply grid search to the objective functions from Example 5.2, with an arbitrary chosen grid indicated below by blue marks. The found minimum, which is the grid point with the smallest value of the objective functions, is marked with a yellow dot.



Due to the unfortunate selection of the grid, the global minimum is not found in the non-convex problem (right). That problem could be handled by increasing the resolution of the grid, which however requires more computations (more evaluations of the objective function).

Some hyperparameters (for example k in k -NN, Chapter 2) are integers, and sometimes it is feasible to simply try all reasonable integer values in grid search. However, most of the time the major challenge in grid search is to select a good grid. The grid used in Algorithm 7 is logarithmic between 0.001 and 1 000, but that is of course only an example. Of course, one could do some manual work by first selecting a coarse grid to get an initial guess, and then refine the grid only around the promising candidates, etc. In practice, if the problem has more than one dimension, it can also be useful to select the grid points randomly instead of using an equally spaced linear or logarithmic grid.

The manual procedure of choosing a grid might, however, become quite tedious, and one could wish for an automated method. That is, in fact, possible by treating the grid point selection problem as a machine learning problem itself. If we consider the points in which the objective function already has been evaluated as a training data set, we can use a regression method to learn a model for the objective function. That model can, in turn, be used to answer questions on where to evaluate the objective function next, and thereby automatically selecting the next grid point. A concrete method built from this idea is the Gaussian process optimization method, which uses Gaussian processes (Chapter 9) for learning a model of the objective function.

Gradient methods

In many situations, we do have access to more than just the value of the objective function in certain points. For example when training parameters it is most often possible to compute the derivative (or rather gradient) of the objective function, and sometimes even the second derivative (or rather Hessian). In those situations, it is often a good idea to use a *gradient descent* method, or even *Newton's method*. In certain situations Newton's method offers a great speedup, but may also break if the circumstances are less fortunate. In practice we therefore have to modify Newton's method, for example by using *trust regions*. We will now introduce the fundamentals of these methods now.

Gradient descent

Gradient descent is typically used for learning parameters, where the parameter vector θ often has a high dimension and where the objective function $J(\theta)$ is simple enough such that its gradient is possible to compute. Let us therefore consider the parameter learning problem

$$\widehat{\theta} = \arg \min_{\theta} J(\theta) \quad (5.22)$$

(even though gradient descent possibly can be used for hyperparameters as well). We will assume⁴ that the gradient of the cost function $\nabla_{\theta} J(\theta)$ exists for all θ . As an example, the gradient of the cost function for logistic regression (3.34) is

$$\nabla_{\theta} J(\theta) = -\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{y_i \theta^T \mathbf{x}_i}} \right) y_i \mathbf{x}_i. \quad (5.23)$$

Note that $\nabla_{\theta} J(\theta)$ is a vector of the same dimension as θ , which describes the direction in which $J(\theta)$ increases. Consequently, and more useful for us, $-\nabla_{\theta} J(\theta)$ describes the direction in which $J(\theta)$ decreases. That is,

$$J(\theta - \gamma \nabla_{\theta} J(\theta)) \leq J(\theta) \quad (5.24)$$

for some (possibly very small) $\gamma > 0$. If $J(\theta)$ is convex, the inequality in (5.24) is strict except at the minimum (where $\nabla_{\theta} J(\theta)$ is zero). This suggests that if we have $\theta^{(t)}$ and want to select $\theta^{(t+1)}$ such that $J(\theta^{(t+1)}) \leq J(\theta^{(t)})$, we should

$$\text{update } \theta^{(t+1)} \text{ as } \theta^{(t)} - \gamma \nabla_{\theta} J(\theta^{(t)}).$$

(5.25)

Repeating (5.25) gives the gradient descent Algorithm 8.

Algorithm 8: Gradient descent

Input: Objective function $J(\theta)$, initial $\theta^{(0)}$, learning rate γ
Result: $\widehat{\theta}$

- 1 Set $t \leftarrow 0$
- 2 **while** $\|\theta^{(t)} - \theta^{(t-1)}\|$ not small enough **do**
- 3 Update $\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma \nabla_{\theta} J(\theta^{(t)})$
- 4 Update $t \leftarrow t + 1$
- 5 **end**
- 6 **return** $\widehat{\theta} \leftarrow \theta^{(t-1)}$

⁴This assumption is primarily made for the theoretical discussion. In practice, there are successful examples of gradient descent being applied to objective functions not differentiable everywhere, such as neural networks with ReLu activation functions (Chapter 6).

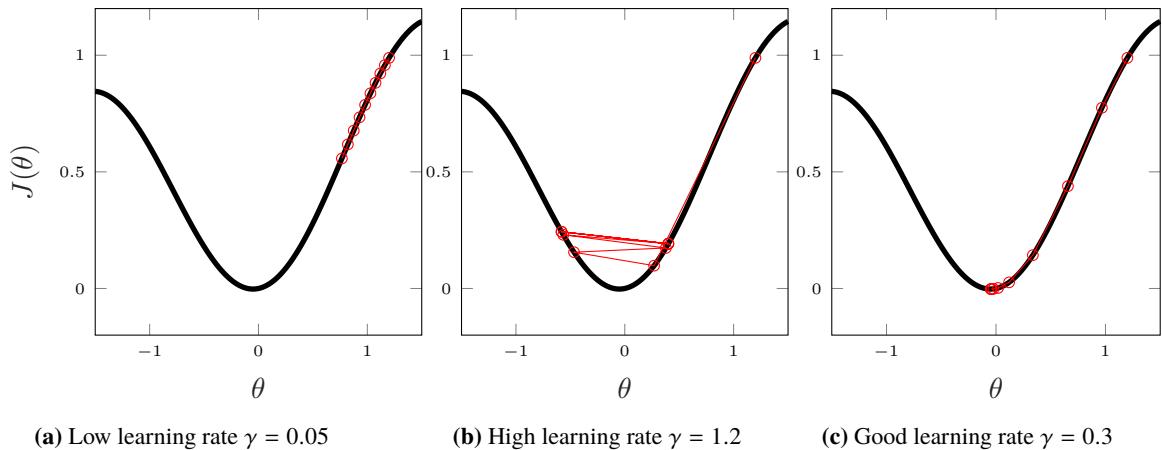


Figure 5.3: Optimization using gradient descent of a cost function $J(\theta)$ where θ is a scalar parameter. In the different subfigures we use a too low learning rate (a), a too high learning rate (b), and a good learning rate (c). Remember that a good value of γ is very much related to the shape of the cost function; $\gamma = 0.3$ might be too small (or high) for a different $J(\theta)$.

In practice, we do not know γ , which determines how big the θ -step is at each iteration. It is possible to formulate the selection of γ as an optimization problem itself to be solved at each iteration, a so-called line-search problem. Here we will consider the simpler solution where we leave the choice of γ to the user. When left as a user choice, γ is often referred to as the *learning rate* or step-size. Note that the gradient $\nabla_{\theta} J(\theta)$ will typically decrease, and eventually attain 0 at a stationary point (possibly, but not necessarily, a minimum), so Algorithm 8 may converge if γ is kept constant. This is in contrast to what you later will learn about the *stochastic* gradient algorithm.

The choice of learning rate γ is important. Some typical situations with too small, too high and a good choice of learning rate are shown in Figure 5.3. With the intuition from these figures, we advise to monitor $J(\theta^{(t)})$ during your optimization, and

- decrease the learning rate γ if the cost function values $J(\theta^{(t)})$ are getting worse or oscillates widely (cf. Figure 5.3b),
- increase the learning rate γ if the cost function values $J(\theta^{(t)})$ are fairly constant and only slowly decreasing (cf. Figure 5.3a).

No general convergence guarantees can be given for gradient descent, basically because a bad learning rate γ may break the method. However, with the “right” choice of γ , the value of $J(\theta)$ will decrease for each iteration (cf. (5.24)) until a point with zero gradient is found, a so-called stationary point. That is, however, not necessarily a minimum, but can also be a maximum or a saddle-point of the objective function. In practice, one typically monitor the value of $J(\theta)$ and terminate the algorithm when it seems not to decrease anymore, and hope it has arrived at a minimum.

In non-convex problem with multiple local minima, we can not expect gradient descent to always find the global one either. The initialization is usually critical for determining which minimum (or stationary point) is found, as illustrated by Example 5.5. It can therefore be a good practice (if time and computational resources permits) to run the optimization multiple times with different initializations. For computationally heavy non-convex problems such as training a deep neural network (Chapter 6) when we cannot afford to re-run the training, we usually employ method-specific heuristics and tricks to find a good initialization point.

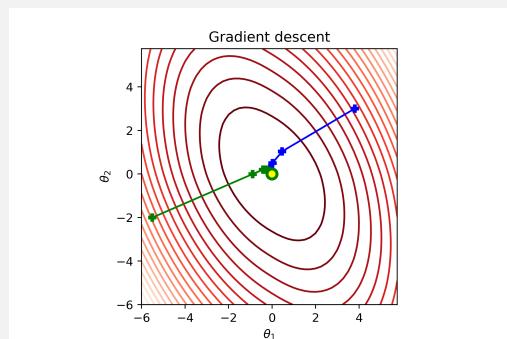
For convex problems there is only one stationary point, which also is the global minimum. Hence, the initialization for convex problem can be done arbitrarily. However, by *warm-starting* the optimization with a good initial guess, we may still save valuable computational time. Sometimes, such as when doing k -fold cross validation (Chapter 4), we have to train k models on similar (but not identical) datasets. In

situations like that we can typically make use of that by initializing Algorithm 8 with the parameters learned for the previous model.

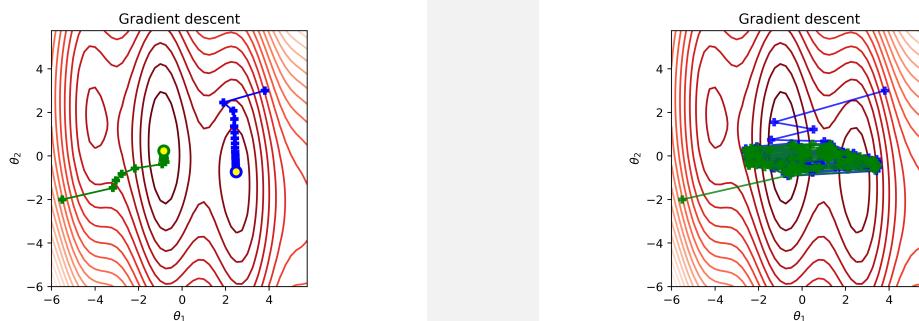
For learning logistic regression (3.35), gradient descent can be used. If we remember that its cost function is convex, we know that once gradient descent has converged to a minimum, it has reached the global minimum and we are done. However, for logistic regression, there are more advanced alternatives that usually performs even better.

Example 5.5: Gradient descent

We first consider the convex objective function from Example 5.2, and apply gradient descent to it with a seemingly reasonable learning rate. Note that each step is perpendicular to the level curves at the point where it starts, which is a property of the gradient. As expected, we find the (global) minimum with both of the two different initializations.



For the non-convex objective function from Example 5.2 we apply gradient descent with two different learning rates. In the left plot, the learning rate seems well chosen and the optimization converges nicely, albeit to different minima depending on the initialization. Note that it *could* have converged also to one of the saddle points between the different minima. In the right plot the learning rate is too big, and the procedure does not seem to converge.



Second order methods

We can think of gradient descent as approximating $J(\boldsymbol{\theta})$ with a first order Taylor expansion around $\boldsymbol{\theta}^{(t)}$, that is, a (hyper-)plane. The next parameter $\boldsymbol{\theta}^{(t+1)}$ is selected by taking a step in the steepest direction of the (hyper-)plane. Let us now see what happens if we instead use a second order Taylor expansion,

$$J(\boldsymbol{\theta} + \mathbf{v}) \approx J(\boldsymbol{\theta}) + \underbrace{\mathbf{v}^\top [\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})] + \frac{1}{2} \mathbf{v}^\top [\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta})] \mathbf{v}}_{\triangleq s(\boldsymbol{\theta}, \mathbf{v})}, \quad (5.26)$$

where \mathbf{v} is a vector of same dimension as $\boldsymbol{\theta}$. This expression does not only contain the gradient of the cost function $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$, but also the Hessian matrix of the cost function $\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta})$. Remember that we are searching for the minimum of $J(\boldsymbol{\theta})$. If the Hessian $\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta})$ is positive definite, that minimum is obtained where the derivative of $s(\boldsymbol{\theta}, \mathbf{v})$ is zero,

$$\frac{\partial}{\partial \mathbf{v}} s(\boldsymbol{\theta}, \mathbf{v}) = \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) + [\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta})] \mathbf{v} = 0 \Leftrightarrow \mathbf{v} = -[\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta})]^{-1} [\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})]. \quad (5.27)$$

This suggests to update

$$\boldsymbol{\theta}^{(t+1)} \text{ as } \boldsymbol{\theta}^{(t)} - [\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^{(t)})]^{-1} [\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(t)})], \quad (5.28)$$

which is *Newton's method* for minimization. Unfortunately, no general convergence guarantees can be given for Newton's method either. For certain cases, Newton's method can be much faster than gradient descent. In fact, if the cost function $J(\boldsymbol{\theta})$ is a quadratic function in $\boldsymbol{\theta}$ then (5.26) is exact and Newton's method (5.28) will find the optimum in only one iteration! Quadratic objective functions are, however, rare in machine learning⁵. It is not even guaranteed that the Hessian $\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta})$ always is positive definite in practice, which may result in rather strange parameter updates in (5.28). To still make use of the potentially valuable second order information, but at the same time also have a robust and practically useful algorithm, we have to introduce some modification of Newton's method. There are multiple options, and we will look at so-called *trust regions*.

We derived Newton's method using the second order Taylor expansion (5.26) as a model for how $J(\boldsymbol{\theta})$ behave around $\boldsymbol{\theta}^{(t)}$. We should perhaps not trust the Taylor expansion to be a good model for all values of $\boldsymbol{\theta}$, but only in the vicinity of $\boldsymbol{\theta}^{(t)}$. One natural restriction is to trust the second order Taylor expansion (5.26) only within a ball of radius D around $\boldsymbol{\theta}^{(t)}$, which we refer to as our trust region. This suggests that we could make a Newton update (5.28) of the parameters, unless the step is longer than D , in which case we downscale the step to never leave our trust region. In the next iteration, the trust region is moved to be centered around the updated $\boldsymbol{\theta}^{(t+1)}$, and another step is taken from there. We can express this as

$\text{update } \boldsymbol{\theta}^{(t+1)} \text{ as } \boldsymbol{\theta}^{(t)} - \eta [\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^{(t)})]^{-1} [\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(t)})],$

(5.29)

where $\eta \leq 1$ is chosen as large as possible such that $\|\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)}\| \leq D$. The radius of the trust region D can be updated and adapted as the optimization proceeds, but for simplicity we consider here D to be a user choice (much like the step size for gradient descent). We summarize this by Algorithm 9. The trust-region Newton method, with a certain set of rules on how to update D is one of the methods commonly used for training logistic regression in practice.

It can be computationally expensive or even impossible to compute the inverse of the Hessian matrix $[\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta}^{(t)})]^{-1}$. To this end, there is an entire class of methods called *quasi-Newton methods* that all use different ways to *approximate* the inverse of the Hessian matrix $[\nabla_{\boldsymbol{\theta}}^2 J(\boldsymbol{\theta})]^{-1}$ in (5.28). This class includes, among other, the Broyden method and the BFGS method (an abbreviation of Broyden, Fletcher, Goldfarb and Shanno). A further approximation of the latter, called limited-memory BFGS or L-BFGS, has proven to be another good choice for the logistic regression problem.

⁵For regression, we often use the squared error loss $L(\hat{y}, y) = (\hat{y} - y)^2$, which indeed is a quadratic function in \hat{y} . That does not imply that $J(\boldsymbol{\theta})$ (the objective function) is a quadratic function in $\boldsymbol{\theta}$.

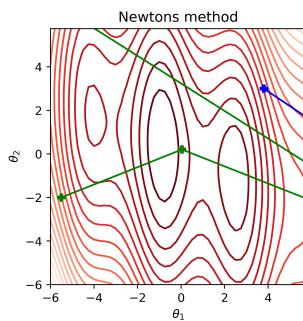
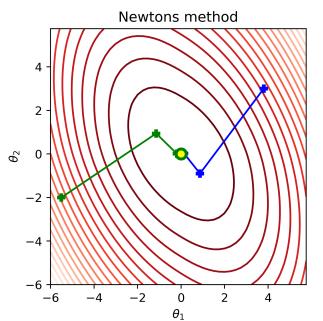
Algorithm 9: Truncated Newton's method

Input: Objective function $J(\theta)$, initial $\theta^{(0)}$, trust region radius D
Result: $\hat{\theta}$

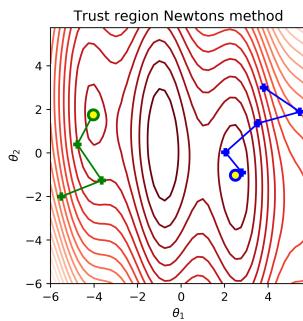
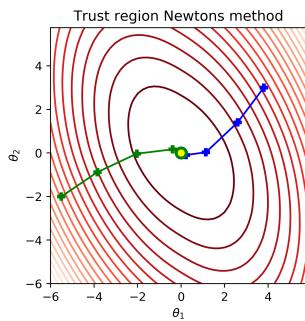
- 1 Set $t \leftarrow 0$
- 2 **while** $\|\theta^{(t)} - \theta^{(t-1)}\|$ not small enough **do**
- 3 Compute $v \leftarrow [\nabla_{\theta}^2 J(\theta^{(t)})]^{-1} [\nabla_{\theta} J(\theta^{(t)})]$
- 4 Compute $\eta \leftarrow \frac{D}{\max(\|v\|, D)}$
- 5 Update $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta v$
- 6 Update $t \leftarrow t + 1$
- 7 **end**
- 8 **return** $\hat{\theta} \leftarrow \theta^{(t-1)}$

Example 5.6: Newton's method

We first apply Newton's method to the cost functions from Example 5.2. Since the convex cost (left) function also happens to be close to a quadratic function, the Newton's method works well and finds, for both initializations, the minimum in only two iterations. For the non-convex problem (right), Newton's method diverges for both initializations, since the second order Taylor expansion (5.26) is a poor approximation of this function and leads the method wrong.



We apply also the trust-region Newton's method to both problems. Note that the first step direction is identical to the non-truncated version above, but the steps are now limited to stay within the trust region (here a circle of radius 2). This prevents the severe divergence problems for the non-convex case, and all cases converge nicely. Indeed, the convex case (left) requires more iterations than for the non-truncated version above, but that is a price we have to pay in order to have a method which also works for the non-convex case shown to the right.



5.4 Optimization with large datasets

In machine learning the training data may have $n = \text{millions}$ (or more) of data samples, the so-called *big data* problem. Computing, for example, the gradient of the cost function

$$\nabla_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L(\mathbf{x}_i, y_i, \theta) \quad (5.30)$$

thus involves summing a million of elements. Besides taking lot of time to sum, it can also be an issue to keep all data samples in the computer memory at the same time. However, with so many data samples it is usually the case that many of them are similar, and in practice we do not always need to consider them all at once, but it might be sufficient to have a look at only a subset of them. This is a general idea called *subsampling*, and we will have a closer look at how subsampling can be combined with gradient descent into a very useful optimization method called *stochastic gradient*. It is, however, possible to combine the subsampling idea also with other methods.

With n very big, we can expect the gradient computed only for the first half of the dataset $\nabla_{\theta} J(\theta) \approx \sum_{i=1}^{n/2} \nabla_{\theta} L(\mathbf{x}_i, y_i, \theta)$ to be almost identical to the gradient based on the second half of the dataset $\nabla_{\theta} J(\theta) \approx \sum_{i=n/2+1}^n \nabla_{\theta} L(\mathbf{x}_i, y_i, \theta)$. Consequently, it might be a waste of time to compute the gradient based on the whole training data set for each iteration of gradient descent. Instead, we could compute the gradient based on the first half of the training data set, update the parameters according to the gradient descent method Algorithm 8, and then compute the gradient for the new parameters based on the second half of the training data,

$$\theta^{(t+1)} = \theta^{(t)} - \gamma \frac{1}{n/2} \sum_{i=1}^{\frac{n}{2}} \nabla_{\theta} L(\mathbf{x}_i, y_i, \theta^{(t)}), \quad (5.31a)$$

$$\theta^{(t+2)} = \theta^{(t+1)} - \gamma \frac{1}{n/2} \sum_{i=\frac{n}{2}+1}^n \nabla_{\theta} L(\mathbf{x}_i, y_i, \theta^{(t+1)}). \quad (5.31b)$$

In other words, we use only a *subsample* of the training when we compute the gradient. These two steps would only require roughly half the computational time compared to using the entire training data set for each gradient computation, and this computational saving illustrates well the benefit of the subsampling idea.

The extreme version of this strategy would be to use only one single data sample each time when computing the gradient. However, most commonly we do something in between, using more than one but not all data samples when computing the gradient. We call this small subsample of data a *mini-batch*, which typically can contain $n_b = 10$, $n_b = 100$ or $n_b = 1000$ data samples. One complete pass through the training data is called an *epoch*, and consists consequently of n/n_b iterations.

When using mini-batches it is important to ensure that the different mini-batches are balanced and representative for the whole data set. For example, if we have a big training data set with a few different output classes and the dataset is sorted with respect to the output, the mini-batch with the first n_b samples would only include one class and hence not give a good approximation of the gradient for the full data set. For this reason, the mini-batches should be formed randomly. One implementation of this is to first randomly shuffle the training data, and thereafter dividing it into mini-batches in an ordered manner. When we have completed one epoch, we do another random reshuffling of the training data and do another pass through the data set. We summarize gradient descent with, *stochastic gradient*, by Algorithm 10.

Stochastic gradient is widely used in machine learning, and there are many extensions tailored for different methods. For training deep neural networks (Chapter 6), some commonly used methods include automatic adaption of the learning rate and an idea called momentum to counteract the subsampling noise. The AdaGrad (short for adaptive gradient), RMSProp (short for root mean square propagation) and Adam (short for adaptive moments) methods are such examples. For logistic regression in the big data setting, the stochastic average gradient (SAG) method, which averages over all previous gradient estimates, has proven useful, to only mention a few.

Algorithm 10: Stochastic gradient

Input: Objective function $J(\theta) = \sum_{i=1}^n L(\mathbf{x}_i, y_i; \theta)$, initial $\theta^{(0)}$, learning rate $\gamma^{(t)}$

Result: $\hat{\theta}$

- 1 Set $t \leftarrow 0$
- 2 **while** Convergence criteria not met **do**
- 3 **for** $i = 1, 2, \dots, E$ **do**
- 4 Randomly shuffle the training data $\{\mathbf{x}_i, y_i\}_{i=1}^n$
- 5 **for** $j = 1, 2, \dots, \frac{n}{n_b}$ **do**
- 6 Approximate the gradient using the mini-batch $\{(\mathbf{x}_i, y_i)\}_{i=(j-1)n_b+1}^{jn_b}$,
- 7 $\hat{\mathbf{d}}^{(t)} = \frac{1}{n_b} \sum_{i=(j-1)n_b+1}^{jn_b} \nabla_{\theta} L(\mathbf{x}_i, y_i, \theta^{(t)})$.
- 8 Update $\theta^{(t+1)} \leftarrow \theta^{(t)} - \gamma^{(t)} \hat{\mathbf{d}}^{(t)}$
- 9 Update $t \leftarrow t + 1$
- 10 **end**
- 11 **end**
- 12 **return** $\hat{\theta} \leftarrow \theta^{(t-1)}$

Learning rate and convergence for stochastic gradient

Standard gradient descent converges if the learning rate is constant, since the gradient itself is zero at the minimum (or any other stationary point). For stochastic gradient, on the other hand, we do *not* obtain convergence if we use a constant learning rate. The reason is that we only use an *estimate* of the true gradient, and this estimate will not necessarily be zero at the minimum of the objective function, but there might still be a considerable amount of “noise” in the gradient estimate due to the subsampling. As a consequence, the stochastic gradient algorithm with a fixed learning rate will not converge towards a point, but walk around somewhat randomly. However, if we do not use a constant learning rate, but gradually decrease it towards zero, the parameter updates will be smaller and smaller, and eventually converge. We hence start at $t = 0$ with a fairly high learning rate $\gamma^{(t)}$ (meaning that we take big steps) and then decay $\gamma^{(t)}$ as t increases. Under certain regularity conditions of the cost function and with a learning rate fulfilling the Robbins-Monro conditions $\sum_{t=0}^{\infty} \gamma^{(t)} = \infty$ and $\sum_{t=0}^{\infty} (\gamma^{(t)})^2 < \infty$, the stochastic gradient algorithm can be shown to converge almost surely to a local minimum. The Robbins-Monro conditions are, for example, fulfilled if using $\gamma^{(t)} = \frac{1}{t}$.

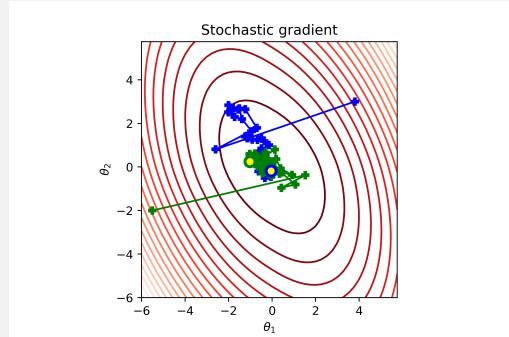
For many machine learning problems, however, it has been found that better performance in practice often is obtained if not letting $\gamma^{(t)} \rightarrow 0$, but to cap it at some small value $\gamma_{\min} > 0$. This will cause stochastic gradient not to exactly converge, but the algorithm will in fact walk around indefinitely (theoretically, in practice until the algorithm is aborted by the user). For practical purposes this seemingly undesired property does usually not cause any major issue if γ_{\min} is only small enough, and one possible rule for setting the learning rate in practice is

$$\gamma^{(t)} = \gamma_{\min} + (\gamma_{\max} - \gamma_{\min}) e^{-\frac{t}{\tau}}. \quad (5.32)$$

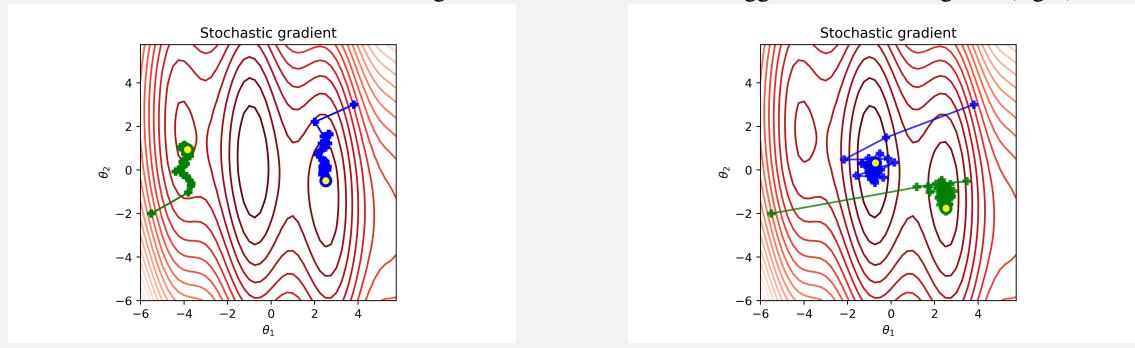
Now the learning rate $\gamma^{(t)}$ starts at γ_{\max} and goes to γ_{\min} as $t \rightarrow \infty$. How to pick the parameters γ_{\min} , γ_{\max} and τ is more of an art than science. As a rule of thumb γ_{\min} can be chosen approximately as 1% of γ_{\max} . The parameter τ depends on the size of the dataset and the complexity of the problem, but should be chosen such that multiple epochs have passed before we reach γ_{\min} . The strategy to pick γ_{\max} can be done by monitoring the cost function as for standard gradient descent in Figure 5.3.

Example 5.7: Optimization

We apply the stochastic gradient method to the objective functions from Example 5.2. For the convex function below, the choice of learning rate is not very crucial. Note, however, that the algorithm does not converge as nicely as, for example, gradient descent, due to the “noise” in the gradient estimate caused by the subsampling. This is the price we have to pay for the substantial computational savings offered by the subsampling.



For the objective function with multiple local minima, we apply stochastic gradient with two decaying learning rates, but with different initial $\gamma^{(0)}$. With a smaller learning rate, left, stochastic gradient converges to the closest minima, whereas it converges to other minima with a bigger initial learning rate (right).



6 Neural networks and deep learning

Neural networks can be used for both regression and classification, and they can be seen as an extension of linear regression and logistic regression, respectively. Traditionally neural networks with *one* so-called hidden layer have been used and analysed, and several success stories came in the 1980s and early 1990s. In the 2000s it was, however, realized that *deep* neural networks with *several* hidden layers, or simply *deep learning*, are even more powerful. With the combination of a lot of training data, new software, hardware and parallel algorithms for training, deep learning has made a major contribution to machine learning and several other fields. Deep learning has excelled in many applications, including image classification, speech recognition and language translation. New applications, analysis, and algorithmic developments to deep learning are published literally every day.

We will start in Section 6.1 by generalizing linear regression to a two-layer neural network (i.e., a neural network with one hidden layer), and then generalize it further to a deep neural network. We thereafter leave regression and look at the classification setting in Section 6.1. In Section 6.2 we present a special neural network tailored for images and finally we look into some details on how to train neural networks in Section 6.3.

6.1 Neural networks

A neural network is a nonlinear function that describes a prediction of the output \hat{y} as a nonlinear function of its input variables

$$\hat{y} = f_{\theta}(x_1, \dots, x_p), \quad (6.1)$$

where the function f is parametrized by θ . Such a nonlinear function can be parametrized in many ways. In a neural network, the strategy is to use several *layers* of linear regression models and nonlinear *activation functions*. We will explain this carefully in turn below.

Generalized linear regression

We start the description with a graphical illustration of the linear regression model

$$\hat{y} = W_1 x_1 + W_2 x_2 + \dots + W_p x_p + b, \quad (6.2)$$

which is shown in Figure 6.1a. Each input variable x_j is represented with a node and each parameter W_j with a link. Furthermore, the output z is described as the sum of all terms $W_j x_j$. Note that we use 1 as the input variable corresponding to the offset term b .

To describe *nonlinear* relationships between $\mathbf{x} = [1 \ x_1 \ x_2 \ \dots \ x_p]^T$ and \hat{y} we introduce a nonlinear scalar function called the *activation function* $h : \mathbb{R} \rightarrow \mathbb{R}$. The linear regression model (6.2) is now modified into a *generalized* linear regression model where the linear combination of the inputs is squashed through the (scalar) activation function

$$\hat{y} = h(W_1 x_1 + W_2 x_2 + \dots + W_p x_p + b). \quad (6.3)$$

This extension to the generalized linear regression model is visualized in Figure 6.1b.

Common choices for activation function are the *logistic function* and the *rectified linear unit* (ReLU). These are illustrated in Figure 6.2a and Figure 6.2b, respectively. The logistic (or sigmoid) function has already been used in the context of logistic regression (Section 3.2). The logistic function is linear close to



Figure 6.1: Graphical illustration of a linear regression model (Figure 6.1a), and a generalized linear regression model (Figure 6.1b). In Figure 6.1a, the output z is described as the sum of all terms b and $\{W_j x_j\}_{j=1}^p$, as in (6.2). In Figure 6.1b, the circle denotes addition and also transformation through the activation function h , as in (6.3).

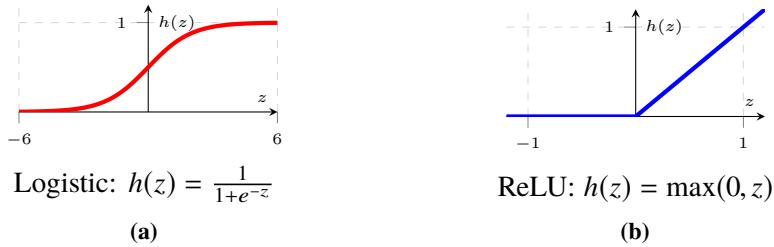


Figure 6.2: Two common activation functions used in neural networks. The logistic (or sigmoid) function (Figure 6.2a), and the rectified linear unit (Figure 6.2b).

$z = 0$ and saturates at 0 and 1 as z decreases or increases. The ReLU is even simpler. The function is the identity function for positive inputs and equal to zero for negative inputs.

The logistic function used to be the standard choice of activation function in neural networks for many years, whereas the ReLU is now the standard choice (despite its simplicity!) in most neural network models.

The generalized linear regression model (6.3) is very simple and is itself not capable of describing very complicated relationships between the input \mathbf{x} and the output \hat{y} . Therefore, we make two further extensions to increase the generality of the model: We will first make use of *several* generalized linear regression models to build a layer (which will lead us to the *two-layer* neural network) and then stack these layers in a *sequential* construction (which will result in a *deep* neural network, or simply *deep learning*).

Two-layer neural network

In (6.3), the output is constructed by one scalar regression model. To increase its flexibility and turn it into a two-layer neural network, we instead let the output be a sum of U such generalized linear regression models, each of which has its own parameters. The parameter for the i th regression model are b_i, \dots, W_{ip} and we denote its output by q_i ,

$$q_i = h(W_{i1}x_1 + W_{i2}x_2 + \dots + W_{ip}x_p + b_i), \quad i = 1, \dots, U. \quad (6.4)$$

These intermediate outputs q_i are so-called *hidden units*, since they are not the output of the whole model. The U different hidden units $\{q_i\}_{i=1}^U$ instead act as input variables to an additional linear regression model

$$\hat{y} = W_1 q_1 + W_2 q_2 + \dots + W_U q_U + b. \quad (6.5)$$

To distinguish the parameters in (6.4) and (6.5) we add the superscripts (1) and (2), respectively. The equations describing this two-layer neural network (or equivalently, a neural network with one layer of

Input variables	Hidden units	Output
-----------------	--------------	--------

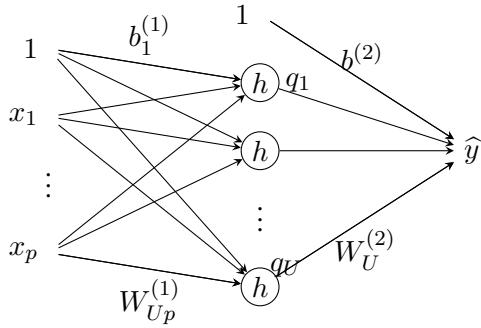


Figure 6.3: A two-layer neural network, or equivalently, a neural network with one intermediate layer of hidden units.

hidden units) are thus

$$\begin{aligned} q_1 &= h\left(W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + \dots + W_{1p}^{(1)}x_p + b_1^{(1)}\right), \\ q_2 &= h\left(W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + \dots + W_{2p}^{(1)}x_p + b_2^{(1)}\right), \\ &\vdots \\ q_U &= h\left(W_{U1}^{(1)}x_1 + W_{U2}^{(1)}x_2 + \dots + W_{Up}^{(1)}x_p + b_U^{(1)}\right), \end{aligned} \quad (6.6a)$$

$$\hat{y} = W_1^{(2)}q_1 + W_2^{(2)}q_2 + \dots + W_U^{(2)}q_U + b^{(2)}. \quad (6.6b)$$

Extending the graphical illustration from Figure 6.1, this model can be depicted as a graph with two layers of links (illustrated using arrows), see Figure 6.3. As before, each link has a parameter associated with it. Note that we include an offset term not only in the input layer, but also in the hidden layer.

Matrix notation

The two-layer neural network model in (6.6) can also be written more compactly using matrix notation, where the parameters in each layer are stacked in a *weight matrix* \mathbf{W} and an *offset vector*¹ \mathbf{b} as

$$\mathbf{W}^{(1)} = \begin{bmatrix} W_{11}^{(1)} & \dots & W_{1p}^{(1)} \\ \vdots & \ddots & \vdots \\ W_{U1}^{(1)} & \dots & W_{Up}^{(1)} \end{bmatrix}, \quad \mathbf{b}^{(1)} = \begin{bmatrix} b_1^{(1)} \\ \vdots \\ b_U^{(1)} \end{bmatrix}, \quad \mathbf{W}^{(2)} = \begin{bmatrix} W_1^{(2)} & \dots & W_U^{(2)} \end{bmatrix}, \quad \mathbf{b}^{(2)} = [b^{(2)}]. \quad (6.7)$$

The full model can then be written as

$$\mathbf{q} = h(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}), \quad (6.8a)$$

$$\hat{y} = \mathbf{W}^{(2)}\mathbf{q} + \mathbf{b}^{(2)}, \quad (6.8b)$$

where we have also stacked the components in \mathbf{x} and \mathbf{q} as $\mathbf{x} = [x_1, \dots, x_p]^\top$ and $\mathbf{q} = [q_1, \dots, q_U]^\top$. The activation function h acts element-wise. The two weight matrices and the two offset vectors will be the parameters in the model, which can be written as

$$\boldsymbol{\theta} = [\text{vec}(\mathbf{W}^{(1)})^\top \quad \text{vec}(\mathbf{W}^{(2)})^\top \quad \mathbf{b}^{(1)\top} \quad \mathbf{b}^{(2)\top}]^\top. \quad (6.9)$$

¹The word ‘‘bias’’ is often used for the offset vector in the neural network literature, but this is really just a model parameter and not a bias in the statistical sense. To avoid confusion we refer to it as an offset instead.

By this we have described a nonlinear regression model on the form $\hat{y} = f(\mathbf{x}; \boldsymbol{\theta})$ according to above. Note that the predicted output \hat{y} in (6.8b) depends on all the parameters in $\boldsymbol{\theta}$ even though it is not explicitly stated in the notation.

Deep neural network

The two-layer neural network is a useful model on its own, and a lot of research and analysis has been done for it. However, the real descriptive power of a neural network is realized when we stack multiple such layers of generalized linear regression models, and thereby achieve a *deep* neural network. Deep neural networks can model complicated relationships (such as the one between an image and its class), and is one of the state-of-the-art methods in machine learning as of today.

We enumerate the layers with index l . Each *layer* is parametrized with a weight matrix $\mathbf{W}^{(l)}$ and an offset vector $\mathbf{b}^{(l)}$, as for the two-layer case. For example, $\mathbf{W}^{(1)}$ and $\mathbf{b}^{(1)}$ belong to layer $l = 1$, $\mathbf{W}^{(2)}$ and $\mathbf{b}^{(2)}$ belong to layer $l = 2$ and so forth. We also have multiple *layers of hidden units* denoted by $\mathbf{q}^{(l-1)}$. Each such layer consists of U_l hidden units $\mathbf{q}^{(l)} = [q_1^{(l)}, \dots, q_{U_l}^{(l)}]^T$, where the dimensions U_1, U_2, \dots can be different across the various layers.

Each layer maps a hidden layer $\mathbf{q}^{(l-1)}$ to the next hidden layer $\mathbf{q}^{(l)}$ according to

$$\mathbf{q}^{(l)} = h(\mathbf{W}^{(l)}\mathbf{q}^{(l-1)} + \mathbf{b}^{(l)}). \quad (6.10)$$

This means that the layers are stacked such that the output of the first layer $\mathbf{q}^{(1)}$ (the first layer of hidden units) is the input to the second layer, the output of the second layer $\mathbf{q}^{(2)}$ (the second layer of hidden units) is the input to the third layer, etc. By stacking multiple layers we have constructed a *deep* neural network. A deep neural network of L layers can mathematically be described as (cf. (6.8))

$$\begin{aligned} \mathbf{q}^{(1)} &= h(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}), \\ \mathbf{q}^{(2)} &= h(\mathbf{W}^{(2)}\mathbf{q}^{(1)} + \mathbf{b}^{(2)}), \\ &\vdots \\ \mathbf{q}^{(L-1)} &= h(\mathbf{W}^{(L-1)}\mathbf{q}^{(L-2)} + \mathbf{b}^{(L-1)}), \\ \hat{\mathbf{y}} &= \mathbf{W}^{(L)}\mathbf{q}^{(L-1)} + \mathbf{b}^{(L)}. \end{aligned} \quad (6.11)$$

A graphical representation of this model is provided in Figure 6.4.

The weight matrix $\mathbf{W}^{(1)}$ for the first layer $l = 1$ has the dimension $U_1 \times p$ and the corresponding offset vector $\mathbf{b}^{(1)}$ has the dimension U_1 . In deep learning it is common to consider applications where also the output is multi-dimensional $\hat{\mathbf{y}} = [\hat{y}_1, \dots, \hat{y}_M]^T$. This means that for the last layer the weight matrix $\mathbf{W}^{(L)}$ has the dimension $M \times U_{L-1}$ and the offset vector $\mathbf{b}^{(L)}$ has the dimension M . For all intermediate layers $l = 2, \dots, L-1$, $\mathbf{W}^{(l)}$ has the dimension $U_l \times U_{l-1}$ and the corresponding offset vector U_l .

The number of inputs p and the number of outputs M are given by the problem, but the number of layers L and the dimensions U_1, U_2, \dots are user design choices that will determine the flexibility of the model.

Learning the network from data

Analogously to the parametric models presented earlier (e.g. linear regression and logistic regression) we need to learn all the parameters in order to use the model. For deep neural networks the parameters are

$$\boldsymbol{\theta} = [\text{vec}(\mathbf{W}^{(1)})^T \quad \text{vec}(\mathbf{W}^{(2)})^T \quad \dots \text{vec}(\mathbf{W}^{(L)})^T \quad \mathbf{b}^{(1)\top} \quad \mathbf{b}^{(2)\top} \quad \dots \quad \mathbf{b}^{(L)\top}]^T \quad (6.12)$$

The wider and deeper the network is, the more parameters there are. Practical deep neural networks can easily have in the order of millions of parameters and these models are therefore also extremely flexible. Hence, some mechanism to avoid overfitting is needed. Regularization such as ridge regression is common (cf. ??), but there are also other techniques specific to deep learning; see further Section 6.3. Furthermore,

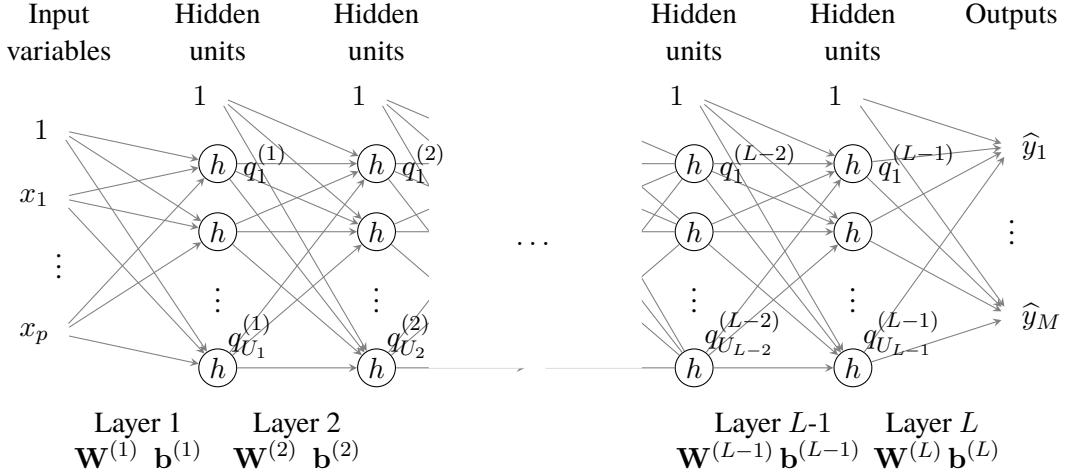


Figure 6.4: A deep neural network with L layers. Each layer l is parameterized by $\mathbf{W}^{(l)}$ and $\mathbf{b}^{(l)}$.

the more parameters there are, the more computational power is needed to train the model. As before, the training data $\mathcal{T} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ consists of n samples of the input \mathbf{x} and the output \mathbf{y} .

For a regression problem we typically start with maximum likelihood and assume Gaussian noise $\varepsilon \sim \mathcal{N}(0, h_\varepsilon^2)$, and thereby obtain the square error loss function as in Section 1,

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \mathbf{y}_i, \boldsymbol{\theta}), \quad \text{where} \quad L(\mathbf{x}_i, \mathbf{y}_i, \boldsymbol{\theta}) = \|\mathbf{y}_i - f(\mathbf{x}_i; \boldsymbol{\theta})\|^2 = \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|^2. \quad (6.13)$$

This problem can be solved with numerical optimization, and more precisely stochastic gradient. This is described in more detail in Section 6.3.

From the parameters $\boldsymbol{\theta}$ and inputs $\{\mathbf{x}_i\}_{i=1}^n$ we can compute the predicted outputs $\{\hat{\mathbf{y}}_i\}_{i=1}^n$ using the model $\hat{\mathbf{y}}_i = f(\mathbf{x}_i; \boldsymbol{\theta})$. For example, for the two-layer neural network presented in Section 6.1 we have

$$\mathbf{q}_i^\top = h(\mathbf{x}_i^\top \mathbf{W}^{(1)\top} + \mathbf{b}^{(1)\top}), \quad (6.14a)$$

$$\hat{\mathbf{y}}_i^\top = \mathbf{q}_i^\top \mathbf{W}^{(2)\top} + \mathbf{b}^{(2)\top}, \quad i = 1, \dots, n \quad (6.14b)$$

In (6.14) the equations are transposed in comparison to the model in (6.8). This is a small trick such that we easily can extend (6.14) to include multiple data points i . Similar to (3.5) we stack all data points in matrices, where each data point represents one row

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}_1^\top \\ \vdots \\ \mathbf{y}_n^\top \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix}, \quad \hat{\mathbf{Y}} = \begin{bmatrix} \hat{\mathbf{y}}_1^\top \\ \vdots \\ \hat{\mathbf{y}}_n^\top \end{bmatrix}, \quad \text{and} \quad \mathbf{Q} = \begin{bmatrix} \mathbf{q}_1^\top \\ \vdots \\ \mathbf{q}_n^\top \end{bmatrix}. \quad (6.15)$$

We can then conveniently write (6.14) as

$$\mathbf{Q} = h(\mathbf{X} \mathbf{W}^{(1)\top} + \mathbf{b}^{(1)\top}), \quad (6.16a)$$

$$\hat{\mathbf{Y}} = \mathbf{Q} \mathbf{W}^{(2)\top} + \mathbf{b}^{(2)\top}, \quad (6.16b)$$

where we have also stacked the predicted output and the hidden units in matrices. Note that the transposed offset vectors \mathbf{b}_1^\top and \mathbf{b}_2^\top are added and broadcasted to each row in this notation.

The vectorized equations in (6.16) is also how the model would typically be implemented in languages that support array programming. For the implementation you might want to consider using the transposed version of \mathbf{W} and \mathbf{b} as your weight matrix and offset vector to avoid taking transpose in each layer.

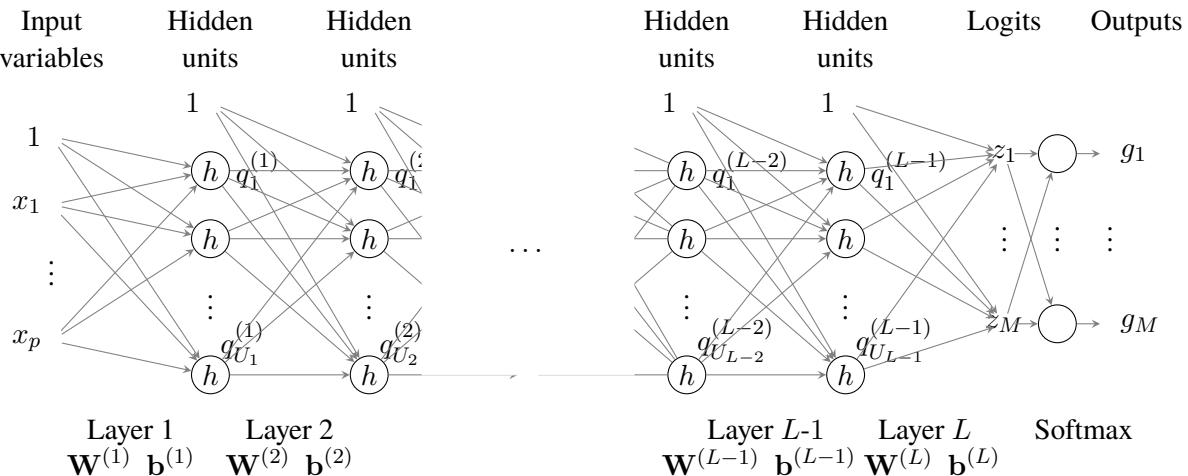


Figure 6.5: A deep neural network with L layers for classification. The only difference to regression (Figure 6.4) is the softmax transformation after layer L .

Neural networks for classification

Neural networks can also be used for classification where we have qualitative outputs $y \in \{1, \dots, M\}$ instead of quantitative. In Section 3.2 we extended linear regression to logistic regression by simply adding the logistic function to the output. In the same manner we can extend the neural network presented in the previous section to a neural network for classification. In doing this extension, we will use the multi-class version of logistic regression presented in Section 3, and more specifically the softmax parametrization given in (3.41), repeated here for convenience

$$\text{softmax}(\mathbf{z}) \triangleq \frac{1}{\sum_{j=1}^M e^{z_j}} \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ \vdots \\ e^{z_M} \end{bmatrix}. \quad (6.17)$$

The softmax function now becomes an additional activation function acting on the final layer of the neural network. In addition to the regression network in (6.11) we add the softmax function at the end of the network as

$$\mathbf{q}^{(1)} = h(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}), \quad (6.18a)$$

⋮

$$\mathbf{q}^{(L-1)} = h(\mathbf{W}^{(L-1)} \mathbf{q}^{(L-2)} + \mathbf{b}^{(L-1)}), \quad (6.18b)$$

$$\mathbf{z} = \mathbf{W}^{(L)} \mathbf{q}^{(L-1)} + \mathbf{b}^{(L)}, \quad (6.18c)$$

$$\mathbf{g} = \text{softmax}(\mathbf{z}). \quad (6.18d)$$

The softmax function maps the output of the last layer $\mathbf{z} = [z_1, \dots, z_M]^\top$ to $\mathbf{g} = [g_1, \dots, g_M]^\top$ where g_m is a model of the class probability $p(y = m | \mathbf{x}_i)$, see also Figure 6.5 for a graphical illustration. The inputs to the softmax function, i.e. the variables z_1, \dots, z_M , are referred to as *logits*.

Note that the softmax function does not come as a layer with additional parameters, it merely acts as a transformation from the output into the modeled class probabilities. By construction, the outputs of the softmax function will always be in the interval $g_m \in [0, 1]$ and sum to $\sum_{m=1}^M g_m = 1$, otherwise they could not be interpreted as probabilities.

	$m = 1$	$m = 2$	$m = 3$
y_{im}	0	1	0
$g_{im}(\theta_A)$	0.1	0.8	0.1

Cross-entropy:

$$L(\mathbf{x}_i, \mathbf{y}_i, \theta_A) = -1 \cdot \ln 0.8 = 0.22$$

	$m = 1$	$m = 2$	$m = 3$
y_{im}	0	1	0
$g_{im}(\theta_B)$	0.8	0.1	0.1

Cross-entropy:

$$L(\mathbf{x}_i, \mathbf{y}_i, \theta_B) = -1 \cdot \ln 0.1 = 2.30$$

Figure 6.6: Illustration of the cross-entropy between a data point \mathbf{y}_i and two different prediction outputs $\mathbf{g}_i(\theta_A) = [g_{i1}(\theta_A), g_{i2}(\theta_A), g_{i3}(\theta_A)]$ and $\mathbf{g}_i(\theta_B) = [g_{i1}(\theta_B), g_{i2}(\theta_B), g_{i3}(\theta_B)]$.

Learning classification networks from data

As before, the training data consists of n samples of inputs and outputs $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$. For the classification problem we use one-hot encoding scheme for the output \mathbf{y}_i . This means that for a problem with M different classes, \mathbf{y}_i consists of M elements $\mathbf{y}_i = [y_{i1} \dots y_{iM}]^\top$. If a data point i belongs to class m then $y_{im} = 1$ and $y_{ij} = 0$ for all $j \neq m$. See more about the one-hot encoding in Section 3.

For a neural network which has the softmax activation function in the final layer we typically use the negative log-likelihood, which is also commonly referred to as the *cross-entropy* loss function, to train the model (cf. (3.44))

$$\widehat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \mathbf{y}_i, \boldsymbol{\theta}), \quad \text{where} \quad L(\mathbf{x}_i, \mathbf{y}_i, \boldsymbol{\theta}) = - \sum_{m=1}^M y_{im} \ln p(y = m | \mathbf{x}_i; \boldsymbol{\theta}). \quad (6.19)$$

To motivate the use of this cost function, we note that cross-entropy is close to its minimum if the predicted probability $p(m | \mathbf{x}_i; \boldsymbol{\theta})$ is close to 1 for the class m for which $y_{im} = 1$. For example, if the i th data point belongs to class $m = 2$ out of a total of $M = 3$ classes we have $\mathbf{y}_i = [0 \ 1 \ 0]^\top$. Assume that we have a set of parameters for the network that we denote by $\boldsymbol{\theta}_A$, and with these parameters we predict $p(y = 1 | \mathbf{x}_i; \boldsymbol{\theta}_A) = 0.1$, $p(y = 2 | \mathbf{x}_i; \boldsymbol{\theta}_A) = 0.8$ and $p(y = 3 | \mathbf{x}_i; \boldsymbol{\theta}_A) = 0.1$ indicating that we are quite sure that data point i actually belongs to class $m = 2$. This would generate a low cross-entropy $L(\mathbf{x}_i, \mathbf{y}_i, \boldsymbol{\theta}_A) = -(0 \cdot \ln 0.1 + 1 \cdot \ln 0.8 + 0 \cdot \ln 0.1) \approx 0.22$. If we instead for another set of parameters $\boldsymbol{\theta}_B$ predict $p(y = 1 | \mathbf{x}_i; \boldsymbol{\theta}_B) = 0.8$, $p(y = 2 | \mathbf{x}_i; \boldsymbol{\theta}_B) = 0.1$ and $p(y = 3 | \mathbf{x}_i; \boldsymbol{\theta}_B) = 0.1$, the cross-entropy would be much higher $L(\mathbf{x}_i, \mathbf{y}_i, \boldsymbol{\theta}_B) = -(0 \cdot \ln 0.8 + 1 \cdot \ln 0.1 + 0 \cdot \ln 0.1) \approx 2.30$. For this case, we would indeed prefer the parameters $\boldsymbol{\theta}_A$ over $\boldsymbol{\theta}_B$. The above reasoning is summarized in Figure 6.6.

Computing the loss function explicitly via the logarithm could lead to numerical problems when $p(y = m | \mathbf{x}_i; \boldsymbol{\theta})$ is close to zero since $\ln(x) \rightarrow -\infty$ as $x \rightarrow 0$. This can be avoided since the logarithm in the cross-entropy loss function (6.19) can “undo” the exponential in the softmax function (6.17),

$$\begin{aligned} L(\mathbf{x}_i, \mathbf{y}_i, \boldsymbol{\theta}) &= - \sum_{m=1}^M y_{im} \ln p(y = m | \mathbf{x}_i; \boldsymbol{\theta}) = - \sum_{m=1}^M y_{im} \ln g_{im} \\ &= - \sum_{m=1}^M y_{im} \left(z_{im} - \ln \left\{ \sum_{j=1}^M e^{z_{ij}} \right\} \right), \end{aligned} \quad (6.20)$$

$$= - \sum_{m=1}^M y_{im} \left(z_{im} - \max_j z_{ij} - \ln \left\{ \sum_{j=1}^M e^{z_{ij} - \max_j z_{ij}} \right\} \right), \quad (6.21)$$

where z_{im} denote the logits.

Image	Data representation	Input variables																																																																								
	<table border="1"> <tbody> <tr><td>0.0</td><td>0.0</td><td>0.8</td><td>0.9</td><td>0.6</td><td>0.0</td></tr> <tr><td>0.0</td><td>0.9</td><td>0.6</td><td>0.0</td><td>0.8</td><td>0.0</td></tr> <tr><td>0.0</td><td>0.0</td><td>0.0</td><td>0.0</td><td>0.9</td><td>0.0</td></tr> <tr><td>0.0</td><td>0.0</td><td>0.0</td><td>0.9</td><td>0.6</td><td>0.0</td></tr> <tr><td>0.0</td><td>0.0</td><td>0.9</td><td>0.0</td><td>0.0</td><td>0.0</td></tr> <tr><td>0.0</td><td>0.8</td><td>0.9</td><td>0.9</td><td>0.9</td><td>0.9</td></tr> </tbody> </table>	0.0	0.0	0.8	0.9	0.6	0.0	0.0	0.9	0.6	0.0	0.8	0.0	0.0	0.0	0.0	0.0	0.9	0.0	0.0	0.0	0.0	0.9	0.6	0.0	0.0	0.0	0.9	0.0	0.0	0.0	0.0	0.8	0.9	0.9	0.9	0.9	<table border="1"> <tbody> <tr><td>$x_{1,1}$</td><td>$x_{1,2}$</td><td>$x_{1,3}$</td><td>$x_{1,4}$</td><td>$x_{1,5}$</td><td>$x_{1,6}$</td></tr> <tr><td>$x_{2,1}$</td><td>$x_{2,2}$</td><td>$x_{2,3}$</td><td>$x_{2,4}$</td><td>$x_{2,5}$</td><td>$x_{2,6}$</td></tr> <tr><td>$x_{3,1}$</td><td>$x_{3,2}$</td><td>$x_{3,3}$</td><td>$x_{3,4}$</td><td>$x_{3,5}$</td><td>$x_{3,6}$</td></tr> <tr><td>$x_{4,1}$</td><td>$x_{4,2}$</td><td>$x_{4,3}$</td><td>$x_{4,4}$</td><td>$x_{4,5}$</td><td>$x_{4,6}$</td></tr> <tr><td>$x_{5,1}$</td><td>$x_{5,2}$</td><td>$x_{5,3}$</td><td>$x_{5,4}$</td><td>$x_{5,5}$</td><td>$x_{5,6}$</td></tr> <tr><td>$x_{6,1}$</td><td>$x_{6,2}$</td><td>$x_{6,3}$</td><td>$x_{6,4}$</td><td>$x_{6,5}$</td><td>$x_{6,6}$</td></tr> </tbody> </table>	$x_{1,1}$	$x_{1,2}$	$x_{1,3}$	$x_{1,4}$	$x_{1,5}$	$x_{1,6}$	$x_{2,1}$	$x_{2,2}$	$x_{2,3}$	$x_{2,4}$	$x_{2,5}$	$x_{2,6}$	$x_{3,1}$	$x_{3,2}$	$x_{3,3}$	$x_{3,4}$	$x_{3,5}$	$x_{3,6}$	$x_{4,1}$	$x_{4,2}$	$x_{4,3}$	$x_{4,4}$	$x_{4,5}$	$x_{4,6}$	$x_{5,1}$	$x_{5,2}$	$x_{5,3}$	$x_{5,4}$	$x_{5,5}$	$x_{5,6}$	$x_{6,1}$	$x_{6,2}$	$x_{6,3}$	$x_{6,4}$	$x_{6,5}$	$x_{6,6}$
0.0	0.0	0.8	0.9	0.6	0.0																																																																					
0.0	0.9	0.6	0.0	0.8	0.0																																																																					
0.0	0.0	0.0	0.0	0.9	0.0																																																																					
0.0	0.0	0.0	0.9	0.6	0.0																																																																					
0.0	0.0	0.9	0.0	0.0	0.0																																																																					
0.0	0.8	0.9	0.9	0.9	0.9																																																																					
$x_{1,1}$	$x_{1,2}$	$x_{1,3}$	$x_{1,4}$	$x_{1,5}$	$x_{1,6}$																																																																					
$x_{2,1}$	$x_{2,2}$	$x_{2,3}$	$x_{2,4}$	$x_{2,5}$	$x_{2,6}$																																																																					
$x_{3,1}$	$x_{3,2}$	$x_{3,3}$	$x_{3,4}$	$x_{3,5}$	$x_{3,6}$																																																																					
$x_{4,1}$	$x_{4,2}$	$x_{4,3}$	$x_{4,4}$	$x_{4,5}$	$x_{4,6}$																																																																					
$x_{5,1}$	$x_{5,2}$	$x_{5,3}$	$x_{5,4}$	$x_{5,5}$	$x_{5,6}$																																																																					
$x_{6,1}$	$x_{6,2}$	$x_{6,3}$	$x_{6,4}$	$x_{6,5}$	$x_{6,6}$																																																																					

Figure 6.7: Data representation of a grayscale image with 6×6 pixels. Each pixel is represented with a number encoding the grayscale color. We denote the whole image as \mathbf{X} (a matrix), and each pixel value is an input variable $x_{j,k}$ (element in the matrix \mathbf{X}).

6.2 Convolutional neural networks

Convolutional neural networks (CNN) are a special kind of neural networks originally tailored for problems where the input data has a grid-like topology. In this text we will focus on images, which have a 2D-topology of pixels. Images are also the most common type of input data in applications where CNNs are applied. However, CNNs can be used for any input data on a grid, also in 1D (e.g. audio waveform data) and 3D (volumetric data e.g. CT scans or video data). We will focus on grayscale images, but the approach can easily be extended to color images as well.

Data representation of an image

Digital grayscale images consist of pixels ordered in a matrix. Each pixel can be represented as a range from 0 (total absence, black) to 1 (total presence, white) and values between 0 and 1 represent different shades of gray. In Figure 6.7 this is illustrated for an image with 6×6 pixels. In an image classification problem, an image is the input \mathbf{x} and the pixels in the image are the input variables $x_{1,1}, x_{1,2}, \dots, x_{6,6}$. The two indices j and k determine the position of the pixel in the image, as illustrated in Figure 6.7.

If we put all input variables representing the image pixels in a long vector, we can use the network architecture presented in Section 6.1 and 6.1. However, by doing that, a lot of the structure present in the image data will be lost. For example, we know that two pixels close to each other typically have more in common than two pixels further apart. This information would be destroyed by such a vectorization. In contrast, CNNs preserve this information by representing the input variables as well as the hidden layers as matrices. The core component in a CNN is the convolutional layer, which will be explained next.

The convolutional layer

Following the input layer, we use a hidden layer with as many hidden units as there are input variables. For the image with 6×6 pixels we consequently have $6 \times 6 = 36$ hidden units. We choose to order the hidden units in a 6×6 matrix, i.e. in the same manner as we did for the input variables, see Figure 6.8a.

The network layers presented in earlier sections (like the one in Figure 6.3) have been *dense layers*. This means that each input variable is connected to all hidden units in the subsequent layer, and each such connection has a unique parameter W_{jk} associated to it. These layers have empirically been found to provide too much flexibility for images and we might not be able to capture the patterns of real importance, and hence not generalize and perform well on unseen data. Instead, a convolutional layer appears to exploit the structure present in images to find a more efficiently parameterized model. In contrast to a dense layer, a convolutional layer leverages two important concepts – *sparse interactions* and *parameter sharing* – to achieve such a parametrization.

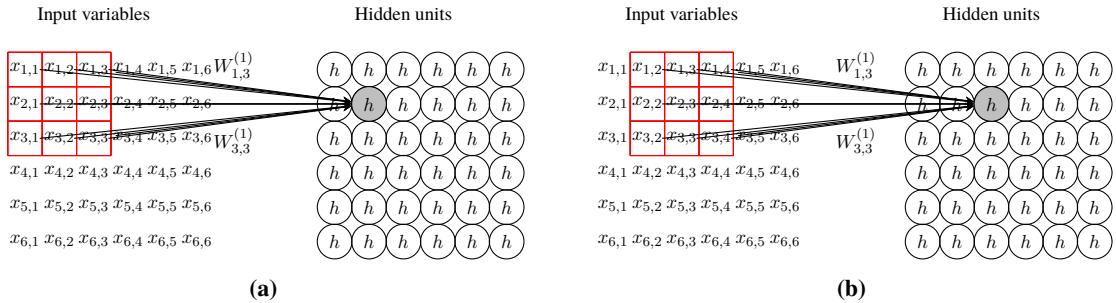


Figure 6.8: An illustration of the interactions in a convolutional layer: Each hidden unit (circle) is only dependent on the pixels in a small region of the image (red boxes), here of size 3×3 pixels. The location of the hidden unit corresponds to the location of the region in the image: if we move to a hidden unit one step to the right, the corresponding region in the image also moves one step to the right, compare Figure 6.8a and Figure 6.8b. Furthermore, the nine parameters $W_{1,1}^{(1)}, W_{1,2}^{(1)}, \dots, W_{3,3}^{(1)}$ are the *same* for all hidden units in the layer.

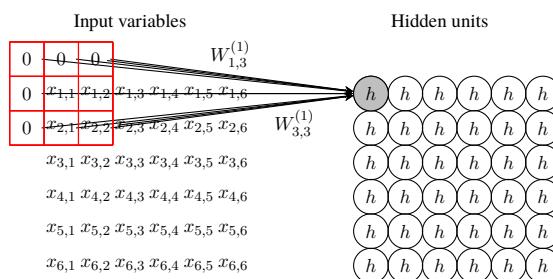


Figure 6.9: An illustration of zero-padding used when the region is partly outside the image. With zero-padding, the size of the image can be preserved in the following layer.

Sparse interactions

With sparse interactions we mean that most of the parameters in a corresponding dense layer are forced to be equal to zero. More specifically, a hidden unit in a convolutional layer only depends on the pixels in a small region of the image and not on all pixels. In Figure 6.8 this region is of size 3×3 . The position of the region is related to the position of the hidden unit in its matrix topology. If we move to the hidden unit one step to the right, the corresponding region in the image also moves one step to the right, as displayed by comparing Figure 6.8a and Figure 6.8b. For the hidden units on the border, the corresponding region is partly located outside the image. For these border cases, we typically use zero-padding where the missing pixels are simply replaced with zeros. Zero-padding is illustrated in Figure 6.9.

Parameter sharing

In a dense layer each link between an input variable and a hidden unit has its own unique parameter. With parameter sharing we instead let the same parameter be present in multiple places in the network. In a convolutional layer the set of parameters for the different hidden units are all the *same*. For example, in Figure 6.8a we use the same set of parameters to map the 3×3 region of pixels to the hidden unit as we do in Figure 6.8b. Instead of learning separate sets of parameters for every position we only learn one set of a few parameters, and use it for all links between the input layer and the hidden units. We call this set of parameters a *filter*. The mapping between the input variables and the hidden units can be interpreted as a convolution between the input variables and the filter, hence the name convolutional neural network.

The sparse interactions and parameter sharing in a convolutional layer makes the CNN fairly invariant to translations of objects in the image. If the parameters in the filter are sensitive to a certain detail (such as a corner, an edge, etc.) a hidden unit will react to this detail (or not) *regardless of where in the image that detail is present!* Furthermore, a convolutional layer uses significantly fewer parameters compared to the corresponding dense layer. In Figure 6.8 only $3 \cdot 3 + 1 = 10$ parameters are required (including the

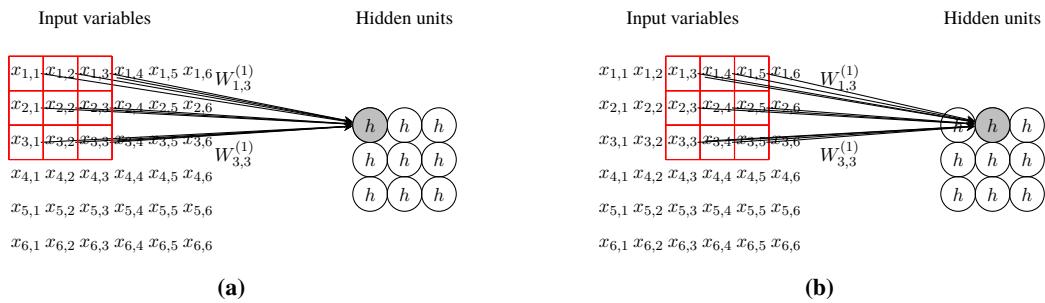


Figure 6.10: A convolutional layer with stride [2,2] and a filter of size 3×3 .

offset parameter). If we instead had used a dense layer $(36 + 1) \cdot 36 = 1332$ parameters would have been needed! Another way of interpreting this is: with the same amount of parameters, a convolutional layer can encode more properties of an image than a dense layer.

Condensing information with strides

In the convolutional layer presented above we have equally many hidden units as we have pixels in the image. As we add more layers to the CNN we usually want to condense the information by reducing the number of hidden units in each layer. One way of doing this is by not applying the filter to every pixel but to say every two pixels. If we apply the filter to every two pixels both row-wise and column-wise, the hidden units will only have half as many rows and half as many columns. For a 6×6 image we get 3×3 hidden units. This concept is illustrated in Figure 6.10.

The *stride* controls how many pixels the filter shifts over the image at each step. In Figure 6.8 the stride is [1,1] since the filter moves by one pixel both row- and column-wise. In Figure 6.10 the stride is [2,2] since it moves by two pixels row- and column-wise. Note that the convolutional layer in Figure 6.10 still requires 10 parameters, as the convolutional layer in Figure 6.8 does. Another way of condensing the information after a convolutional layer is by subsampling the data, so-called *pooling*. The interested can read further about pooling in Goodfellow, Bengio, and Courville 2016.

Multiple channels

The networks presented in Figure 6.8 and 6.10 only have 10 parameters each. Even though this parameterization comes with several important advantages, one filter is probably not sufficient to encode all interesting properties of the images in our dataset. To extend the network, we add multiple filters, each with their own set of parameters. Each filter produces its own set of hidden units—a so-called *channel*—using the same convolution operation as explained in Section 6.2. Hence, each layer of hidden units in a CNN is organized into a tensor with the dimensions (rows \times columns \times channels). In Figure 6.11, the first layer of hidden units has four channels and that hidden layer consequently has dimension $6 \times 6 \times 4$.

When we continue to stack convolutional layers, each filter depends not only on one channel, but on all the channels in the previous layer. This is displayed in the second convolutional layer in Figure 6.11. As a consequence, each filter is a tensor of dimension (filter rows \times filter columns \times input channels). For example, each filter in the second convolutional layer in Figure 6.11 is of size $3 \times 3 \times 4$. If we collect all filter parameters in one weight tensor \mathbf{W} , that tensor will be of dimension (filter rows \times filter columns \times input channels \times output channels). In the second convolutional layer in Figure 6.11, the corresponding weight matrix $\mathbf{W}^{(2)}$ is a tensor of dimension $3 \times 3 \times 4 \times 6$. With multiple filters in each convolutional layer, each of them can be sensitive to different features in the image, such as certain edges, lines or circles enabling a rich representation of the images in our training data.

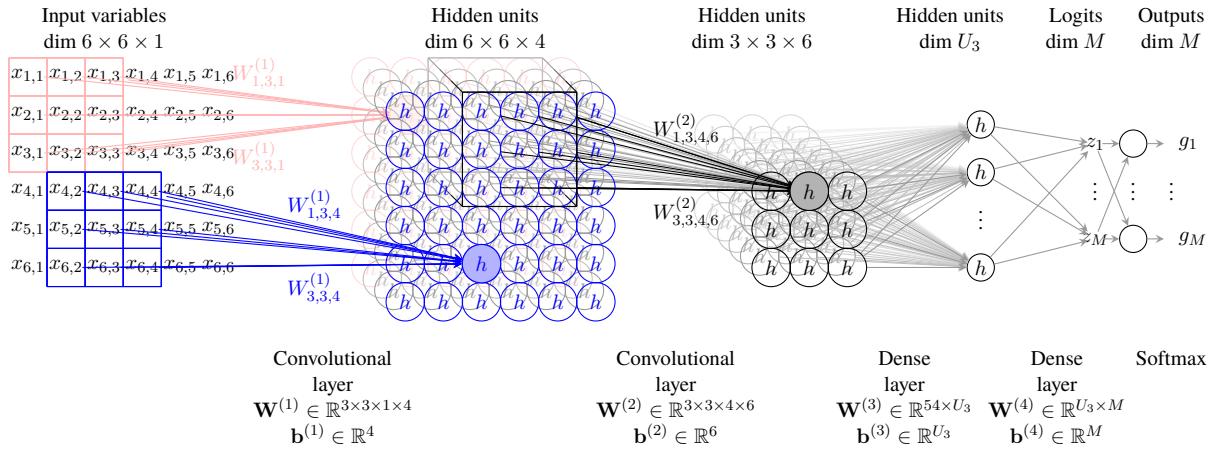


Figure 6.11: A full CNN architecture for classification of grayscale 6×6 images. In the first convolutional layer four filters each of size 3×3 produce a hidden layer with four channels. The first channel (in the back) is visualized in red and the forth channel (in the front) is visualized in blue. We use stride [1,1] which maintains the number of rows and columns. In the second convolutional layer, six filters of size $3 \times 3 \times 4$ and the stride [2,2] are used. They produce a hidden layer with 3 rows, 3 columns and 6 channels. After the two convolutional layers follows a dense layer where all $3 \cdot 3 \cdot 6 = 54$ hidden units in the second hidden layer are densely connected to the third layer of hidden units where all links have their unique parameters. We add an additional dense layer mapping down to the M logits. The network ends with a softmax function to provide predicted class probabilities as output.

Full CNN architecture

A full CNN architecture consists of multiple convolutional layers. Typically, we decrease the number of rows and columns in the hidden layers as we proceed through the network, but instead increase the number of channels to enable the network to encode more high level features. After a few convolutional layers we usually end the network with one or more dense layers. If we consider an image classification task, we place a softmax layer at the very end to get outputs in the range [0,1]. The loss function when training a CNN will be the same as in the regression and classification networks explained earlier, depending on which type of problem we have at hand. In Figure 6.11 a small example of a full CNN architecture is displayed.

6.3 Training a neural network

To use a neural network for prediction we need to find suitable values for its parameters θ . To do that we solve an optimization problem on the form

$$\widehat{\theta} = \arg \min_{\theta} J(\theta) \quad \text{where } J(\theta) = \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \mathbf{y}_i, \theta). \quad (6.22)$$

We denote $J(\theta)$ as the *cost function* and $L(\mathbf{x}_i, \mathbf{y}_i, \theta)$ as the loss function. The functional form of the loss function depends on the characteristics of the problem at hand, see e.g. (6.13) for regression and (6.19) for classification.

These optimization problems can not be solved in closed form, so numerical optimization has to be used. In all numerical optimization algorithms the parameters are updated in an iterative manner. In deep learning we typically use various versions of gradient based search:

1. Pick an initialization θ_0 .
2. Update the parameters as $\theta_{t+1} = \theta_t - \gamma \nabla_{\theta} J(\theta_t)$ for $t = 1, 2, \dots$
3. Terminate when some criterion is fulfilled, and take the last θ_t as $\widehat{\theta}$.

This problem has two main computational challenges. The first computational challenge is the big data problem. For deep learning application the number of data points n is typically very big making the computation of the cost function and its gradient very costly since it requires a sum over all data points. As a consequence, we cannot afford to compute the exact gradient $\nabla_{\theta}J(\theta_t)$ at each iteration. Instead, we compute an approximation of this gradient by considering a random subset for the training data at each iteration. These so-called stochastic gradient algorithms are further explained in Section 5.4.

The second computational challenge is that the number of parameters $\dim(\theta)$, which is also very big for deep learning problems. To efficiently compute the gradient $\nabla_{\theta}J(\theta_t)$ we apply the chain rule of calculus and reuse partial derivatives needed to compute this gradient. This is called the back-propagation algorithm and is not further explained here. The interested reader can for example consult Goodfellow, Bengio, and Courville 2016.

Initialization

Most of the previous optimization problems (such as L^1 regularization and logistic regression) that we have so far encountered have all been convex. This means that we can guarantee global convergence regardless of what initialization θ_0 we use. In contrast, the cost functions for training neural networks is usually non-convex. This means that the training is sensitive to the value of the initial parameters. Typically, we initialize all the parameters to small random numbers to enable the different hidden units to encode different aspects of the data. If the ReLU activation functions are used, the offset elements b_0 are typically initialized to a small positive value such that it operates in the non-negative range of the ReLU.

Dropout

Like all models presented in this course, neural network models can suffer from overfitting if we have a too flexible model in relation to the complexity of the data. Bagging (Section 7.1) is one way to reduce the variance and by that also reducing the risk of overfitting. In bagging we train an entire *ensemble* of models. We train all models (ensemble members) on a different dataset each, which has been bootstrapped (sampled with replacement) from the original training dataset. To make a prediction, we first make one prediction with each model (ensemble member), and then average over all models to obtain the final prediction.

Bagging is also applicable to neural networks. However, it comes with some practical problems; a large neural network model usually takes quite some time to train and it also has quite some parameters to store. To train not just one, but an entire ensemble of many large neural networks would thus be very costly, both in terms of runtime and memory. *Dropout* is a bagging-like technique that allows us to combine many neural networks without the need to train them separately. The trick is to let the different models share parameters with each other, which reduces the computational cost and memory requirement.

Ensemble of sub-networks

Consider a neural network like the one in Figure 6.12a. In dropout we construct the equivalent to an ensemble member by randomly removing some of the hidden units. We say that we drop the units, hence the name dropout. Via this process we obtain a sub-network of our original network. Two such sub-networks are displayed in Figure 6.12b. We randomly sample with a pre-defined probability which units to drop, and the collection of dropped units in one sub-network is independent from the collection of dropped units in another sub-network. When a unit is removed, we also remove all of its incoming and outgoing connections. Not only hidden units can be dropped, but also input variables.

Since all sub-networks stem from the very same original network, the different sub-networks share some parameters with each other. For example, in Figure 6.12b the parameter $W_{55}^{(1)}$ is present in both sub-networks. The fact that they share parameters with each other allow us to train the ensemble of sub-networks in an efficient manner.

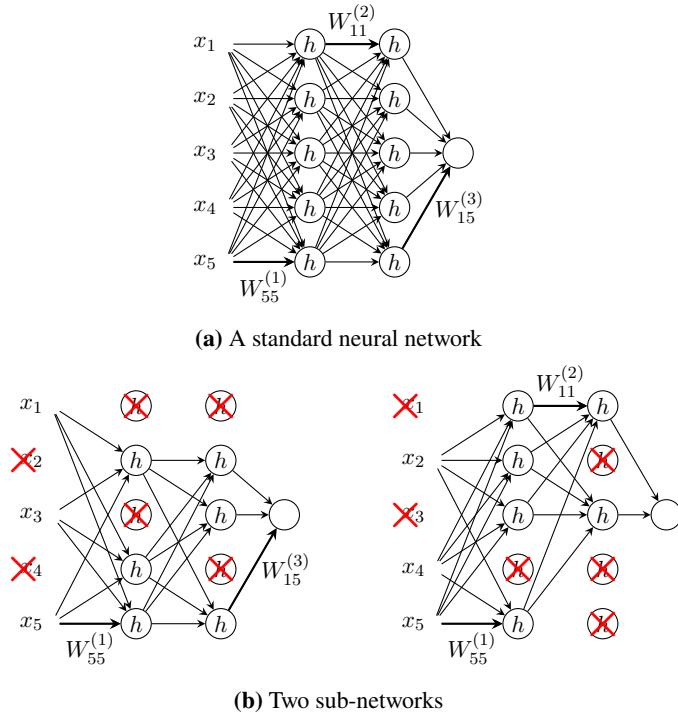


Figure 6.12: A neural network with two hidden layers (a), and two sub-networks with dropped units (b). The collection of units that have been dropped are independent between the two sub-networks.

Training with dropout

To train with dropout we use the stochastic gradient algorithm described in Algorithm 10. In each gradient step a mini-batch of data is used to compute an approximation of the gradient, as before. However, instead of computing the gradient for the full network, we generate a random sub-network by randomly dropping units as described above. We compute the gradient for that sub-network as if the dropped units were not present and then do a gradient step. This gradient step only updates the parameters present in the sub-network. The parameters that are not present are left untouched. In the next gradient step we grab another mini-batch of data, remove another randomly selected collection of units and update the parameters present in that sub-network. We proceed in this manner until some terminal condition is fulfilled.

Dropout vs bagging

The dropout procedure to generate an ensemble of models differs from bagging in a few ways:

- In bagging all models are independent in the sense that they have their own parameters. In dropout the different models (the sub-networks) share parameters.
- In bagging each model is trained until convergence. In dropout each sub-network is only trained for a single gradient step. However, since they share parameters all models will be updated also when the other networks are trained.
- Similar to bagging, in dropout we train each model on a dataset that has been randomly selected from our training data. However, in bagging we usually do it on a bootstrapped version of the whole dataset whereas in dropout each model is trained on a randomly selected mini-batch of data.

Even though dropout differs from bagging in some aspects it has empirically been shown to enjoy similar properties as bagging in terms of avoiding overfitting and reducing the variance of the model.

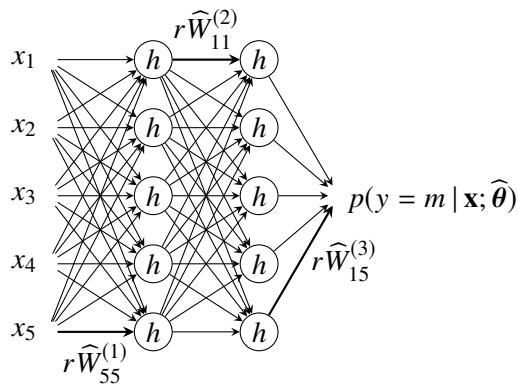


Figure 6.13: The network used for prediction after being trained with dropout. All units and links are present (no dropout) but the weights going out from a certain unit are multiplied with the probability of that unit being included during training. This is to compensate for the fact that some of them were dropped during training. Here all units have been kept with probability r during training (and consequently dropped with probability $1 - r$).

Prediction at test time

After we have trained the sub-networks, we want to make a prediction based on an unseen input data point \mathbf{x}_* . In bagging we evaluate all the different models in the ensemble and combine their results. This would be infeasible in dropout due to the very large (combinatorial) number of possible sub-networks. However, there is a simple trick to approximately achieve the same result. Instead of evaluating all possible sub-networks we simply evaluate the full network containing all the parameters. To compensate for the fact that the model was trained with dropout, we multiply each estimated parameter going out from a unit with the probability of that unit being included during training. This ensures that the expected value of the input to a unit is the same during training and testing, as during training only a fraction of the incoming links were active. For instance, assume that we during training kept a unit with probability p in all layers, then during testing we multiply all estimated parameters with p before we do a prediction based on the network. This is illustrated in Figure 6.13. This procedure of approximating the average over all ensemble members has been shown to work surprisingly well in practice even though there is not yet any solid theoretical argument for the accuracy of this approximation.

Dropout as a regularization method

As a way to reduce the variance and avoid overfitting, dropout can be seen as a regularization method. There are plenty of other regularization methods for neural networks including parameter penalties (analog to ridge regression and LASSO in Section ??), early stopping (the training is stopped before the parameters have converged, and thereby the risk of overfitting is reduced) and various sparse representations (for example CNNs can be seen as a regularization method where most parameters are forced to be zero), just to mention a few. Since its invention, dropout has become one of the most popular regularization techniques due to its simplicity, the fact that it is computationally cheap and its good performance. In fact, a good practice of designing a neural network is often to extend the network until it overfits, then extend it a bit more and finally add a regularization like dropout to avoid that overfitting.

6.4 Perspective and further reading

Although the first conceptual ideas of neural networks date back to the 1940s (McCulloch and Pitts 1943), they had their first main success stories in the late 1980s and early 1990s with the use of the so-called back-propagation algorithm. At that stage, neural networks could, for example, be used to classify handwritten digits from low-resolution images (LeCun, Boser, et al. 1990). However, in the late 1990s neural networks were largely forsaken because it was widely believed that they could not be used to

solve any challenging problems in computer vision and speech recognition. In these areas, neural networks could not compete with hand-crafted solutions based on domain specific prior knowledge.

This situation has changed dramatically since the late 2000s, with multiple layers under the name deep learning. Progress in software, hardware and algorithm parallelization made it possible to address more complicated problems, which were unthinkable only a couple of decades ago. For example, in image recognition, these deep models are now the dominant methods of use and they reach human or even super-human performance on some specific tasks (LeCun, Bengio, and Hinton 2015). Recent advances based on deep neural networks have generated algorithms that can learn how to play computer games based on pixel information only (Mnih et al. 2015), how to beat the world champion in the board game of Go (Silver et al. 2016) and automatically understand the situation in images for automatic caption generation (Xu et al. 2015).

An accessible introduction and overview of deep learning is provided by LeCun, Bengio, and Hinton (2015), and via the textbook by Goodfellow, Bengio, and Courville (2016).

7 Ensemble methods: Bagging and boosting

In Chapter 3 and 2 we introduced four fundamental models for supervised machine learning. In this chapter we will introduce ensemble methods, a type of *meta-algorithm*, which makes use of multiple copies of some fundamental model. We refer to a set of multiple copies of a fundamental model as an *ensemble of base models*, and the key idea is to train each such base model in a slightly different way. To obtain a prediction from an ensemble, we let each base model make its own prediction and then use a (possibly weighted) average or majority vote to obtain the final prediction. With a carefully constructed ensemble, the prediction obtained in this way is better than the prediction of a single base model.

We start in Section 7.1 by introducing a general technique referred to as bootstrap aggregating, or *bagging* for short. The bagging idea is to first create multiple slightly different “versions” of the training data by, essentially, randomly sample overlapping subsets of the training data (the so-called bootstrap). Thereafter, one base model is trained from each such “version” of the training data. In this way, an ensemble of similar, but not identical, base models is obtained. With this procedure it is possible to *reduce the variance* (without any notable increase in bias) compared to using only a single base model learned from the entire training dataset. In practice this means that by using bagging the risk of overfit decreases, compared to using the base model itself. In Section 7.2 we introduce an extension to bagging only applicable when the base model is a classification or regression tree, which results in a powerful off-the-shelf method called random forests. In random forests, each tree is randomly perturbed in order to obtain additional variance reduction, beyond what is already obtained by the bagging procedure itself.

In Section 7.3-7.4 we introduce another ensemble method known as *boosting*. Boosting is different from bagging and random forests, since its base models are trained sequentially, one after the other, where each model tries to “correct” for the “mistakes” made by the previous ones. On the contrary to bagging, the main effect of boosting is *bias reduction* compared to the base model. Thus, boosting is able to turn an ensemble of “weak” base models (e.g., linear classifiers) into one “strong” ensemble model (e.g., a heavily non-linear classifier), and has also shown very useful in practice.

7.1 Bagging

As discussed already in Chapter 4, a central concept in machine learning is the bias–variance trade–off. Roughly speaking, the more flexible a model is, the lower its bias will be. That is, a flexible model is capable of representing complicated input–output relationships. Examples of simple yet flexible models are k -NN with a small value of k and a classification tree that is grown deep. Such highly flexible models are sometimes needed for solving real-world machine learning problems, where relationships are far from linear. The downside, however, is the risk of overfitting¹, or equivalently, high model variance. Despite their high variance, those models are not useless. By using them as base models in bootstrap aggregating, or *bagging*, we can

reduce the variance of the base model, without increasing its bias.

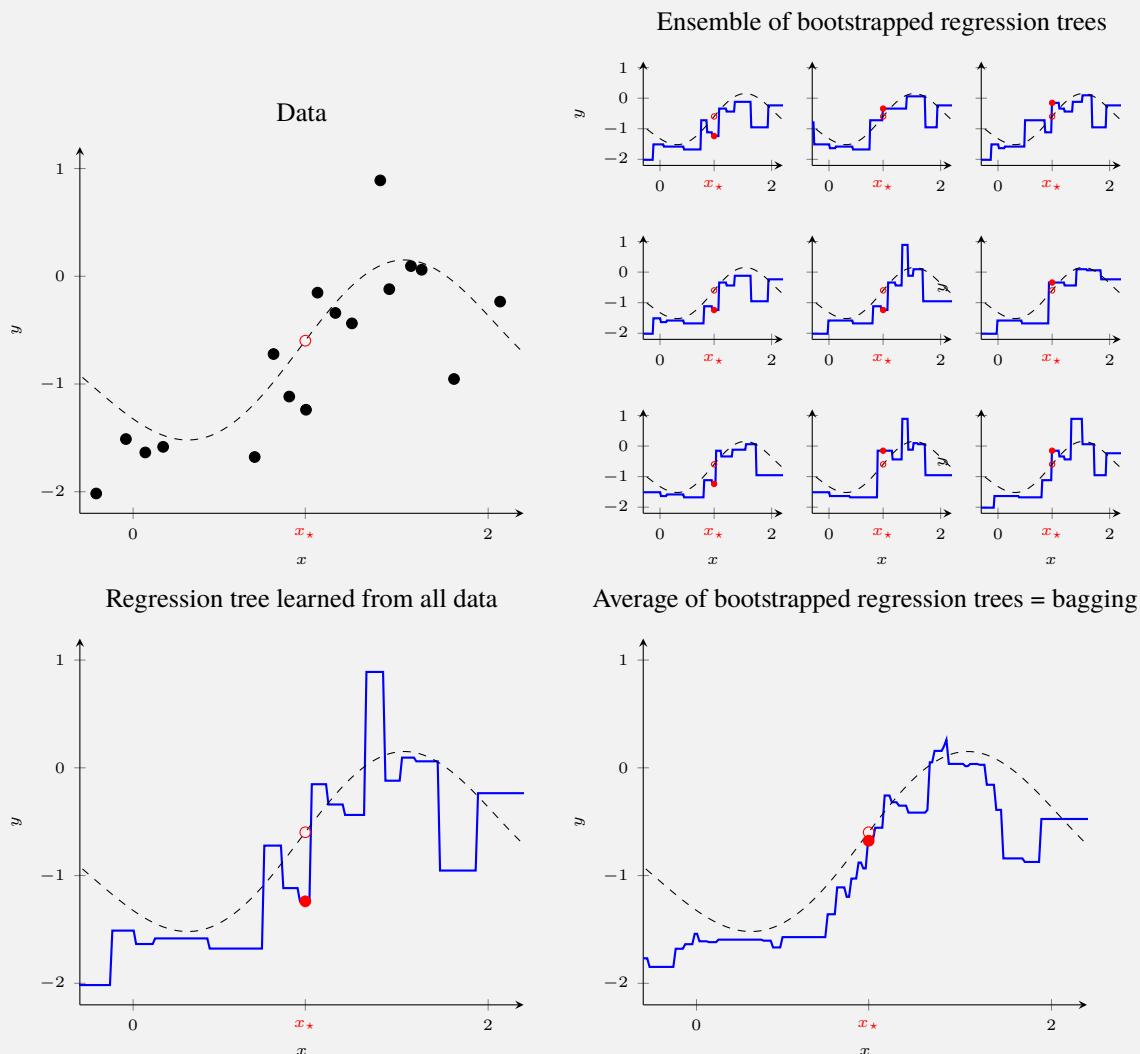
We outline the main idea of bagging with the example below.

¹Both a k -NN model with $k = 1$ and a classification tree with a single data point per leaf node will result in zero training error, typical cases of severe over-fitting.

Example 7.1: Using bagging for a regression problem

Consider the data (black dots) below that are drawn from a function (dashed line) plus noise. As always in supervised machine learning, we want to learn a model from the data which is able to predict new data points well. Being able to predict new data points well means, among other things, that the model should predict the dotted line at x_* (the red circle) well.

For solving this problem, we could use any regression method. Here, we use a regression tree which is grown until each leaf node only contains one data point, whose prediction is shown to the lower left (blue line and red dot). This is a typical low-bias high-variance model, and the overfit to the training data is apparent from the figure. We could decrease its variance, and hence the overfit, by using a more shallow (less deep) tree, but that would on the other hand increase the bias. Instead, we lower the variance (without increasing the bias much) by using bagging with the regression tree as base model.



The rationale behind bagging goes as follows: Because of the noise in the training data, we may think of the prediction $\hat{y}(x_*)$ (the red dot) as a random variable. In bagging, we learn an ensemble of base models (upper right panel), where each base model is trained on a different “version” of the training data obtained using the bootstrap. We may therefore think of each base model to be a different realization of the random variable $\hat{y}(x_*)$. It is well-known that the average of multiple realizations of a random variable has a lower variance than the random variable itself, which means that by taking the average (lower right) of all base models we obtain a prediction with less variance than the base model itself. That is, the bagged regression tree (lower right) has lower variance than a single prediction tree (lower left). Since the base model itself also has low bias, the averaged prediction will have low bias *and* low variance. We can visually confirm that the prediction is better (red dot and circle are closer to each other) for bagging than for the single regression tree.

The bootstrap

As outlined in Example 7.1, the idea of bagging is to average over multiple base models, each learned from a different training dataset. First we therefore have to construct different training datasets. In the best of worlds we would just collect multiple datasets, but most often we cannot do that and instead we have to make the most out of the limited data available. For this purpose, the bootstrap is useful.

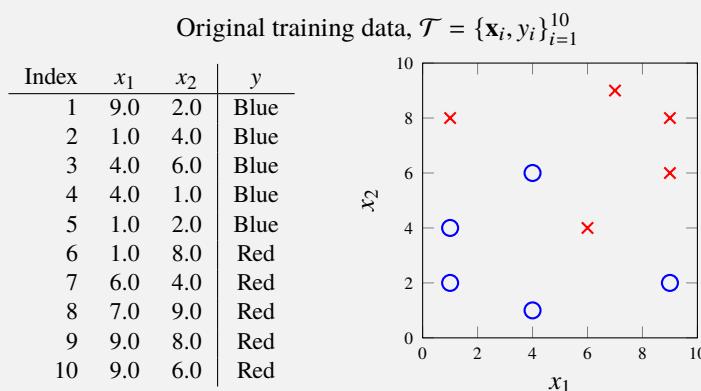
The bootstrap is a method for artificially creating multiple datasets (of size n) out of one dataset (also of size n). The traditional usage of the bootstrap is to quantify uncertainties in statistical estimators (such as confidence intervals), but it turns out to be useful also for machine learning. We denote the original dataset $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$, and assume that \mathcal{T} provides a good representation of the real-world data generating process, in the sense that if we were to collect more training data, these data points would likely be similar to the training data points already contained in \mathcal{T} . We can thus argue that randomly picking data points from \mathcal{T} is a reasonable way to simulate a “new” training dataset. In statistical terms, instead of sampling from the population (collecting more data), we sample from the available training data which is assumed to provide a good representation of the population.

The bootstrap is stated in Algorithm 11 and illustrated in Example 7.2 below. Note that the sampling is done with replacement, meaning that the resulting bootstrapped dataset may contain multiple copies of some of the original training data points, whereas other data points are not included at all.

Time to reflect 7.1: What would happen if the sampling was done without replacement in the bootstrap?

Example 7.2: The bootstrap

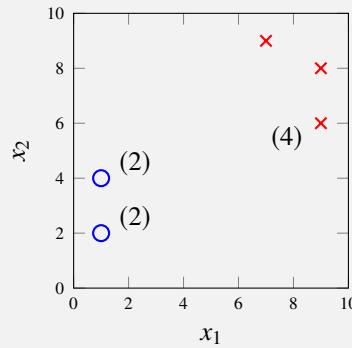
We have a small training dataset with $n = 10$ data points with a two-dimensional input $\mathbf{x} = [x_1 \ x_2]$ and a binary output $y \in \{\text{Blue, Red}\}$.



To generate a bootstrapped dataset $\tilde{\mathcal{T}} = \{\tilde{\mathbf{x}}_i, \tilde{y}_i\}_{i=1}^{10}$, we simulate 10 times with replacement from the index set $\{1, \dots, 10\}$, resulting in the indices $\{2, 10, 10, 5, 9, 2, 5, 10, 8, 10\}$. Thus, $(\tilde{\mathbf{x}}_1, \tilde{y}_1) = (\mathbf{x}_2, y_2)$, $(\tilde{\mathbf{x}}_2, \tilde{y}_2) = (\mathbf{x}_{10}, y_{10})$, etc. We end up with the following dataset, where the numbers in parentheses in the right panel indicate that there are multiple copies of some of the original data points in the bootstrapped data.

Bootstrapped data, $\tilde{\mathcal{T}} = \{\tilde{\mathbf{x}}_i, \tilde{y}_i\}_{i=1}^{10}$

Index	\tilde{x}_1	\tilde{x}_2	\tilde{y}
2	1.0	4.0	Blue
10	9.0	6.0	Red
10	9.0	6.0	Red
5	1.0	2.0	Blue
9	9.0	8.0	Red
2	1.0	4.0	Blue
5	1.0	2.0	Blue
10	9.0	6.0	Red
8	7.0	9.0	Red
10	9.0	6.0	Red



Algorithm 11: The bootstrap.

Data: Training dataset $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$

Result: Bootstrapped data $\tilde{\mathcal{T}} = \{\tilde{\mathbf{x}}_i, \tilde{y}_i\}_{i=1}^n$

```

1 for  $i = 1, \dots, n$  do
2   | Sample  $\ell$  uniformly on the set of integers  $\{1, \dots, n\}$ 
3   | Set  $\tilde{\mathbf{x}}_i = \mathbf{x}_\ell$  and  $\tilde{y}_i = y_\ell$ 
4 end

```

Algorithm 12: Bagging.

Data: Training dataset $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$

Result: A prediction $\hat{y}_{\text{bag}}(x_\star)$ or $\mathbf{g}_{\text{bag}}(x_\star)$

```

1 for  $b = 1, \dots, B$  do
2   | Run Algorithm 11 to obtain a bootstrapped training dataset  $\tilde{\mathcal{T}}^{(b)}$ 
3   | Learn a base model from  $\tilde{\mathcal{T}}^{(b)}$ 
4   | Use the base model to predict  $\tilde{y}^b(x_\star)$ 
5 end
6 Obtain  $\hat{y}_{\text{bag}}(x_\star)$  or  $\mathbf{g}_{\text{bag}}(x_\star)$  by averaging (7.1).

```

Variance reduction by averaging

By running the bootstrap (Algorithm 11) repeatedly B times we obtain B identically distributed bootstrapped datasets $\tilde{\mathcal{T}}^1, \dots, \tilde{\mathcal{T}}^B$. We can then use those bootstrapped datasets to train an ensemble of B base models. We thereafter average their predictions

$$\hat{y}_{\text{bag}}(x_\star) = \frac{1}{B} \sum_{b=1}^B \tilde{y}^b(x_\star) \quad \text{or} \quad \mathbf{g}_{\text{bag}}(x_\star) = \frac{1}{B} \sum_{b=1}^B \tilde{\mathbf{g}}^b(x_\star), \quad (7.1)$$

depending on whether we are concerned with regression (predicting an output value $\hat{y}_{\text{bag}}(x_\star)$) or classification (predicting class probabilities $\mathbf{g}_{\text{bag}}(x_\star)$). In (7.1), $\tilde{y}^1(x_\star), \dots, \tilde{y}^B(x_\star)$ and $\tilde{\mathbf{g}}^1(x_\star), \dots, \tilde{\mathbf{g}}^B(x_\star)$ denote the predictions from the individual ensemble members. The averaged prediction, denoted $\hat{y}_{\text{bag}}(x_\star)$ or $\mathbf{g}_{\text{bag}}(x_\star)$, is the final prediction obtained from bagging. We summarize this by Algorithm 12. (For classification, the prediction could alternatively be decided by majority vote among the ensemble members, but that typically degrades the performance slightly compared to averaging the predicted class probabilities.)

We will now give some more details on the variance reduction that happens in (7.1), which is the entire point of bagging. We focus on regression, but the intuition works also for classification.

Let us point out a basic property of random variables, namely that averaging reduces variance. To formalize this, let z_1, \dots, z_B be a collection of identically distributed (but possibly dependent) random variables with mean value $\mathbb{E}[z_b] = \mu$ and variance $\text{Var}[z_b] = \sigma^2$ for $b = 1, \dots, B$. Furthermore, assume that the average correlation² between any pair of variables is ρ . Then, computing the mean and the variance of the average $\sum_{b=1}^B z_b$ of these variables we get

$$\mathbb{E}\left[\frac{1}{B} \sum_{b=1}^B z_b\right] = \mu, \quad (7.2a)$$

$$\text{Var}\left[\frac{1}{B} \sum_{b=1}^B z_b\right] = \frac{1-\rho}{B} \sigma^2 + \rho \sigma^2. \quad (7.2b)$$

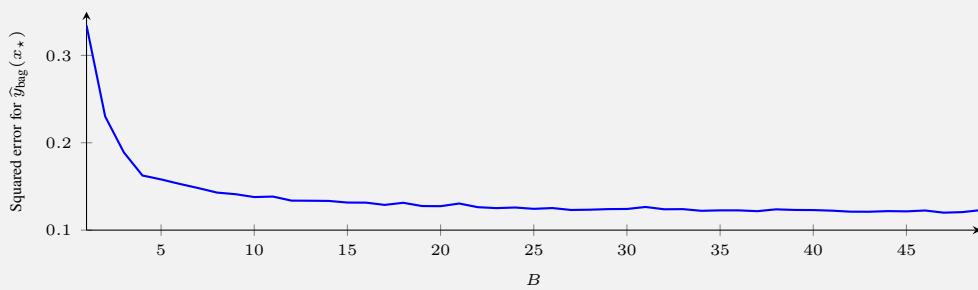
The first equation (7.2a) tells us that the mean is unaltered by averaging a number of identically distributed random variables. Furthermore, the second equation (7.2b) tells us that the variance is reduced by averaging if the correlation $\rho < 1$. The first term in the variance expression (7.2b) can be made arbitrarily small by increasing B , whereas the second term is only determined by the correlation ρ .

To make a connection between bagging and (7.2), consider the predictions $\tilde{y}^b(x_\star)$ from the base models as random variables (in other words, imagine z_b being the red dots in Example 7.1). All base models, and hence their predictions, originate from the same data \mathcal{T} (via the bootstrap), and $\tilde{y}^b(x_\star)$ are therefore identically distributed but correlated. By averaging the predictions we decrease the variance, according to (7.2b). If we only chose B large enough, the achieved variance reduction will be limited by the correlation ρ . Experience has shown that ρ is often small enough such that the computational complexity of bagging (compared to only using the base model itself) pays off well in terms of decreased variance. To summarize, by averaging the identically distributed predictions from several base models as in (7.1), each with a low bias, the *bias remains low*³ (cf. (7.2a)) and *the variance is reduced* (cf. (7.2b)).

At first glance, one might *think* that a bagging model (7.1) becomes more “complex” as the number of ensemble members B increase, and that we therefore run a risk of overfitting if we use many ensemble members B . However, there is nothing in (7.2) which indicate any such problem (bias remains low, variance decreases), and we confirm this by Example 7.3.

Example 7.3: Bagging for regression (cont.)

We consider the problem from Example 7.1 again, and explore how the number of base models B affects the result. We measure the squared error between the “true” function value at x_\star and the predicted $\hat{y}^{\text{bag}}(x_\star)$ when using different B . (Because of the bootstrap, there is a certain amount of randomness in the bagging algorithm itself. To avoid that “noise”, we average the result over multiple runs of the bagging algorithm.)



What we see here is that the squared error eventually reaches a plateau as $B \rightarrow \infty$. Had there been an overfit issue with $B \rightarrow \infty$, the squared error would have started do increase again for some large value of B .

²That is $\frac{1}{B(B-1)} \sum_{b \neq c} \mathbb{E}[(z_b - \mu)(z_c - \mu)] = \rho \sigma^2$.

³Strictly speaking, (7.2a) implies that the bias is identical for a single ensemble member and the ensemble average. The use of the bootstrap might however affect the bias, in that a base model trained on the original data might have a smaller bias than a base model trained on a bootstrapped version of the training data. Most often, this is not an issue in practice.

Despite the fact that the number of parameters in the model increases as B increases, the lack of overfit as $B \rightarrow \infty$ according to Example 7.3 is the expected (and intended) behavior. It is important to understand that by the construction of bagging, *more ensemble members does not make the resulting model more flexible*, but only reduces the variance. With this in mind, in practice the choice of B is mainly guided by computational constraints. The larger B the better, but increasing B when there is no further reduction in test error is computationally wasteful.

Remark 7.1 *Bagging can still overfit, we cannot prevent from that. The only claim we have is that the problem never gets worse when $B \rightarrow \infty$.*

Out-of-bag error estimation

When using bagging (or random forests), it turns out that there is a way to estimate the expected new data error E_{new} *without* using cross-validation. The first observation we have to make is that not all data points from the original dataset \mathcal{T} will have been used for training all ensemble members. It is actually possible to show that with the bootstrap, on average only 63% of the original training data points in $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ will be present in a bootstrapped training dataset $\tilde{\mathcal{T}} = \{\tilde{\mathbf{x}}_i, \tilde{y}_i\}_{i=1}^n$. Roughly speaking, this means that for any given $\{\mathbf{x}_i, y_i\}$ in \mathcal{T} , one third of the ensemble members will not have seen that data point yet. We refer to these roughly $B/3$ ensemble members as being out-of-bag for sample i , and we let them form their own ensemble, the out-of-bag-ensemble i . Note that the out-of-bag-ensemble is different for each data point $\{\mathbf{x}_i, y_i\}$.

The next key insight is that for the out-of-bag-ensemble i , the data point $\{\mathbf{x}_i, y_i\}$ can actually act as a test data point since it has not yet been seen by any of its ensemble members. By computing the squared or misclassification error when the out-of-bag-ensemble i predicts $\{\mathbf{x}_i, y_i\}$, we thus get an estimate of E_{new} for this out-of-bag-ensemble, which we denote $E_{\text{OOB}}^{(i)}$. Since $E_{\text{OOB}}^{(i)}$ is based on only one data point, it will be a fairly poor estimate of E_{new} . If we however repeat this for all data points $\{\mathbf{x}_i, y_i\}$ in the training data \mathcal{T} and average $E_{\text{OOB}} = \frac{1}{n} \sum_{i=1}^n E_{\text{OOB}}^{(i)}$, we get a better estimate of E_{new} . Indeed, E_{OOB} will be an estimate of E_{new} for an ensemble with only $B/3$ (and not B) members, but as we have seen (Example 7.3), the performance of bagging plateaus after a certain number of ensemble members. Hence, if B is large enough so that ensembles with B and $B/3$ members perform roughly equally, E_{OOB} provides an estimate of E_{new} which can be at least as good as the estimate $E_{k\text{-fold}}$ from k -fold cross-validation. Most importantly, however, E_{OOB} comes almost for free in bagging, whereas $E_{k\text{-fold}}$ requires much more computations when re-training k times.

7.2 Random forests

In bagging we reduce the variance by averaging over an ensemble of models. That reduction, however, is limited by the correlation between the individual ensemble members (cf. the role of ρ in (7.2b)). Using a simple trick, it turns out to be possible to reduce that correlation further, beyond what is achieved by using the bootstrap. This is known as random forests.

While bagging is a general technique that in principle can be used to reduce the variance of any base model, random forests assumes that these base models are classification or regression trees. The idea is to inject additional randomness when constructing each tree, in order to further reduce the correlation among the base models. At first this might seem like a silly idea: randomly perturbing the training of a model should intuitively degrade its performance. There is a rationale for this perturbation, however, which we will discuss below, but first we present the details of the algorithm.

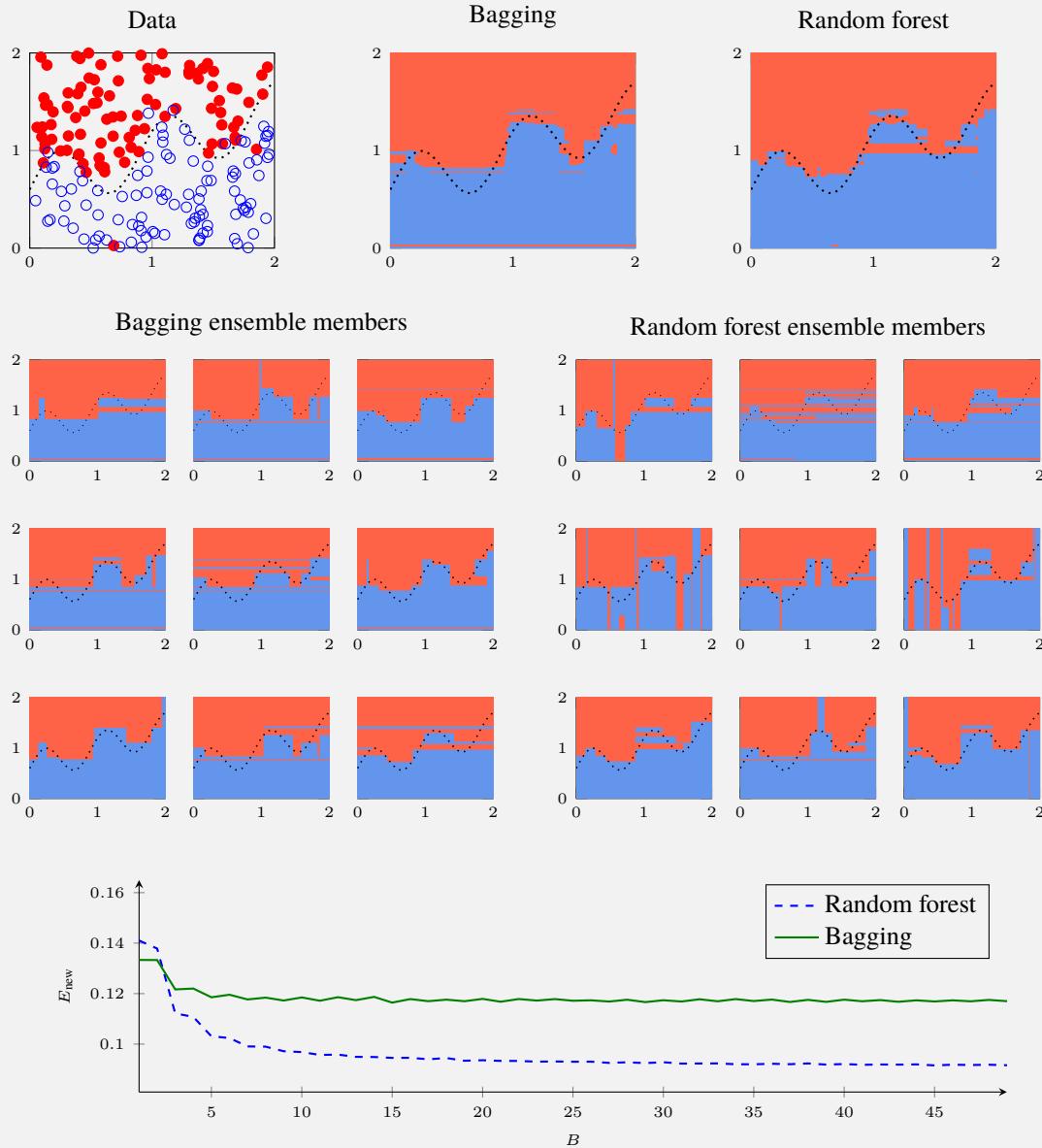
Let $\tilde{\mathcal{T}}^{(b)}$ be one of the B bootstrapped datasets in bagging. To train a classification or regression tree on this data we proceed as usual (see Section 2.3), but with one difference. Throughout the training, whenever we are about to split a node we do not consider all possible input variables x_1, \dots, x_p as splitting variables. Instead, we pick a random subset consisting of $q \leq p$ inputs, and only consider these q variables as possible splitting variables. At the next splitting point we draw a new random subset of q inputs to use as possible splitting variables, and so on. Naturally, this random subset selection is done independently for each of the B ensemble members, so that we (with high probability) end up using different subsets for the different trees. This additional random constraint when training is what turns bagging into *random forests*. This will cause the B trees to be less correlated and averaging their predictions can therefore result in larger variance reduction compared to bagging. It should be noted, however, that this random perturbation of the training procedure will increase the variance⁴ of each *individual tree*. In the notation of Equation (7.2b), random forests decreases ρ (good) but increases σ^2 (bad) compared to bagging. Experience has however shown that the reduction in correlation is the dominant effect, so that the averaged prediction variance is often reduced. We illustrate this in Example 7.4 below.

⁴And possibly also the bias, in a similar manner as the bootstrap might increase the bias, see Footnote 3, page 105.

Example 7.4: Random forests and bagging for a binary classification problem

Consider the binary classification with $p = 2$ using the data given below. The different classes are the blue circles and the red dots. The input values were randomly sampled from $[0, 2] \times [0, 2]$, and labeled red with probability 0.98 if above the dotted line, and vice versa. We use two different classifiers: bagging with classification trees (which is equivalent to a random forest with $q = p = 2$) and a random forest with $q = 1$, each with $B = 9$ ensemble members. Below we plot the decision boundary for each ensemble member as well as the majority-voted final decision boundary.

The most apparent difference is that the variation among the ensemble members of the random forest than those of bagging. Roughly half of the random forest ensemble members have been forced to make the first split along the horizontal axis, which has lead to increased variance and decreased correlation compared to bagging where all ensemble members made the first split along the vertical axis.



While it is hard to visually compare the final decision boundaries for bagging and random forest (top right), we also compute E_{new} for different numbers of ensemble members B . Since the learning itself has a certain amount of randomness, we average over multiple learned models to not be confused by that random effect. Indeed we see that the random forest performs better than bagging, except for very small B , and we conclude that the positive effect of the reduced correlation between the ensemble members outweighs the negative effect of additional variance. The poor performance of random forest with only one ensemble member is expected, since this lonely model has higher variance and no averaging is taking place when $B = 1$.

To understand why it can be a good idea to only consider a subset of inputs as splitting variables, recall that tree-building is based on recursive binary splitting which is a greedy algorithm. This means that the algorithm can make choices early on that appear to be good, but which nevertheless turn out to be suboptimal further down the splitting procedure. For instance, consider the case when there is one dominant input variable. If we construct an ensemble of trees using plain bagging, it is then very likely that all of the ensemble members will pick this dominant variable up as the first splitting variable, making all trees identical (i.e., perfectly correlated) after the first split. If we instead apply random forests, some of the ensemble members will not even have access to this dominant variable at the first split, since it most likely will not be present in the random subset of q inputs selected at the first split for some of the ensemble members. This will force those members to split according to some other variable. While there is no reason for why this would improve the performance of the individual tree, it *could* prove to be useful further down the splitting process, and since we average over many ensemble members the overall performance can be improved.

User aspects

Since random forest is a bagging method, the tools and properties from Section 7.1 applies also to random forests. One such example is the out-of-bag error estimation, which is applicable also to random forests. Also for random forests, $B \rightarrow \infty$ does not lead to overfit. Hence, there is no reason to choose B small, other than the computational load. Compared to using a single tree, a random forest requires approximately B times as much computations. Since all trees are identically distributed, it is however possible to parallelize the random forest learning over multiple nodes, where each node learns a few ensemble members.

The choice of q is a tuning parameter, where for $q = p$ we recover the basic bagging method described previously. As a rule-of-thumb we can set $q = \sqrt{p}$ for classification problems and $q = p/3$ for regression problems (values rounded down to closest integer). A more systematic way of selecting q is to use out-of-bag error estimation or cross-validation and select q such that E_{OOB} or $E_{k\text{-fold}}$ is minimized.

7.3 Boosting and AdaBoost

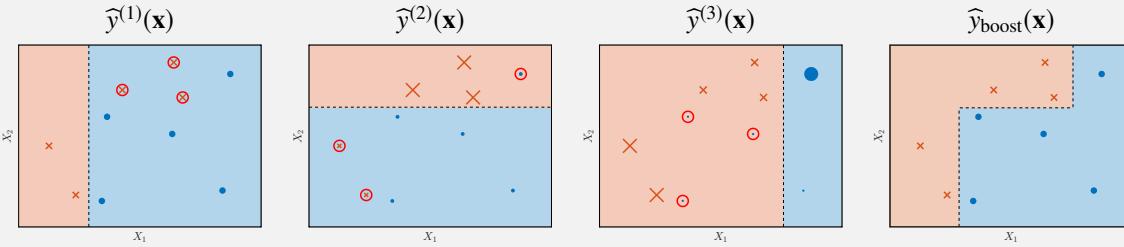
Whereas bagging primarily is an ensemble method for reducing variance in high-variance base models, boosting is rather an ensemble method for reducing bias in high-bias base models. A typical example of a simple (or, in other words, weak) high-bias model is a classification tree of depth one (sometimes called a classification stump). Boosting is built on the idea that even a weak high-bias model often can capture *some* of the relationship between the inputs and the output. Thus, by training multiple weak models, each describing part of the input-output relationship, it might be possible to combine the predictions of these models into an overall better prediction. Hence, the intention is to *reduce the bias* by turning an ensemble of weak models into one strong model.

Boosting shares some similarities with bagging. Both are ensemble methods, in the sense that they are based on combining the predictions from multiple models (an ensemble). Both bagging and boosting can also be viewed as meta-algorithms, in the sense that they can be used to combine essentially any regression or classification algorithm—they are algorithms built on top of other algorithms. However, there are also important differences between boosting and bagging which we will discuss below.

The main difference is how the base models are learned. In bagging we learn B identically distributed models in parallel. Boosting, on the other hand, uses a *sequential* construction of the ensemble members. Informally, this is done in such a way that each model tries to correct the mistakes made by the previous one. This is accomplished by modifying the training dataset at each iteration in order to put more emphasis on the data points for which the model (so far) has performed poorly. The final prediction is obtained from a weighted average or a weighted majority vote among the models. We look at a simple example to illustrate this idea.

Example 7.5: Boosting illustration

We consider a binary classification problem with a two-dimensional input $\mathbf{x} = [x_1 \ x_2]$. The training data consists of $n = 10$ data points, 5 from each of the two classes. We use a decision stump, a classification tree of depth one, as a simple (weak) base classifier. A decision stump means that we select one of the input variables, x_1 or x_2 , and split the input space into two half spaces, in order to minimize the training error. This results in a decision boundary that is perpendicular to one of the axes. The left panel below shows the training data, illustrated by red crosses and blue dots for the two classes, respectively. The colored regions shows the decision boundary for a decision stump $\hat{y}^{(1)}(\mathbf{x})$ trained on this data.



The model $\hat{y}^{(1)}(\mathbf{x})$ misclassifies three data points (red crosses falling in the blue region), which are encircled in the figure. To improve the performance of the classifier we want to find a model that can distinguish these three points from the blue class. To this end, we train another decision stump, $\hat{y}^{(2)}(\mathbf{x})$, on the same data. To put emphasis on the three misclassified points, however, we assign *weights* $\{w_i^{(2)}\}_{i=1}^n$ to the data. Points correctly classified by $\hat{y}^{(1)}(\mathbf{x})$ are down-weighted, whereas the three points misclassified by $\hat{y}^{(1)}(\mathbf{x})$ are up-weighted. This is illustrated in the second panel of the figure above, where the marker sizes have been scaled according to the weights. The classifier $\hat{y}^{(2)}(\mathbf{x})$ is then found by minimizing the *weighted* misclassification error, $\frac{1}{n} \sum_{i=1}^n w_i^{(2)} \mathbb{I}\{y_i \neq \hat{y}^{(2)}(\mathbf{x}_i)\}$, resulting in the decision boundary shown in the second panel. This procedure is repeated for a third and final iteration: we update the weights based on the hits and misses of $\hat{y}^{(2)}(\mathbf{x})$ and train a third decision stump $\hat{y}^{(3)}(\mathbf{x})$ shown in the third panel. The final classifier, $\hat{y}^{\text{boost}}(\mathbf{x})$ is then obtained as a majority vote of the three decision stumps.

The decision boundary of the boosted classifier is shown in the right panel. Note that this decision boundary is nonlinear, whereas the decision boundary for each ensemble member is linear. This illustrates the concept of turning an ensemble of three weak (high-bias) base models into a stronger (low-bias) model.

This example illustrates the idea of boosting, but there are several important details left to specify in order to have a useful off-the-shelf algorithm. We will soon have a look at a specific boosting algorithm called AdaBoost, and thereafter consider the more general framework of gradient boosting. We will restrict our attention to binary classification ($K = 2$, with $y \in \{+1, -1\}$), but boosting is possible to apply also for the multiclass problem and for regression.

AdaBoost

What we have discussed so far is a general idea, but there are still a few technical design choices left. Let us now derive an actual boosting method, the AdaBoost (Adaptive Boosting) algorithm for binary classification. AdaBoost was the first successful practical implementation of the boosting idea and lead the way for its popularity.

As we outlined in Example 7.5 boosting attempts to construct a sequence of B (weak) binary classifiers $\hat{y}^{(1)}(\mathbf{x}), \hat{y}^{(2)}(\mathbf{x}), \dots, \hat{y}^{(B)}(\mathbf{x})$. We will in this procedure only consider the final ‘hard’ prediction $\hat{y}(\mathbf{x})$ from the base models, and not their predicted class probabilities $g(\mathbf{x})$. Any classification model can in principle be used as base classifier—shallow classification trees are common in practice. The individual predictions of the B ensemble members are then combined into a final prediction. Unlike bagging, all ensemble members are not treated equally. Instead, we assign some positive coefficients $\{\alpha^{(b)}\}_{b=1}^{(B)}$ and construct the

boosted classifier using a *weighted* majority vote

$$\hat{y}_{\text{boost}}^{(B)}(\mathbf{x}) = \text{sign} \left\{ \sum_{b=1}^B \alpha^{(b)} \hat{y}^{(b)}(\mathbf{x}) \right\}. \quad (7.3)$$

Each ensemble member votes either -1 or $+1$, and the output from the boosted classifier is $+1$ if the weighted sum of the individual votes is positive and -1 if it is negative.

In AdaBoost, the ensemble members and their coefficients $\alpha^{(b)}$ are trained greedily by minimizing the exponential loss of the boosted classifier at each iteration. That is, one ensemble member is added iteratively at a time. When member b is added, it is trained such that the exponential loss (5.12) of the entire ensemble of b members (that is, the boosted classifier we have so far) is minimized. The reason for choosing the exponential loss is that the problem becomes easier to solve (much like the squared error loss in linear regression), which we will now see when we derive a mathematical expression for this procedure.

Let us write the boosted classifier after b iterations as $\hat{y}_{\text{boost}}^{(b)}(\mathbf{x}) = \text{sign}\{c^{(b)}(\mathbf{x})\}$ where $c^{(b)}(\mathbf{x}) = \sum_{j=1}^b \alpha^{(j)} \hat{y}^{(j)}(\mathbf{x})$. Since $c^{(b)}(\mathbf{x})$ is a sum (cf. (7.3)), we can express it iteratively as

$$c^{(b)}(\mathbf{x}) = c^{(b-1)}(\mathbf{x}) + \alpha^{(b)} \hat{y}^{(b)}(\mathbf{x}), \quad (7.4)$$

initialized with $c^0(\mathbf{x}) = 0$. The ensemble members are constructed sequentially, meaning that at iteration b of the procedure the function $c^{(b-1)}(\mathbf{x})$ is known and fixed. This is what makes this construction “greedy”. What remains to be learned at iteration b is the ensemble member $\hat{y}^{(b)}(\mathbf{x})$ and its coefficient $\alpha^{(b)}$. We do this by minimizing the exponential loss of the training data,

$$(\alpha^{(b)}, \hat{y}^{(b)}) = \arg \min_{(\alpha, \hat{y})} \sum_{i=1}^n L(y_i \cdot c^{(b)}(\mathbf{x}_i)) \quad (7.5a)$$

$$= \arg \min_{(\alpha, \hat{y})} \sum_{i=1}^n \exp \left(-y_i \left(c^{(b-1)}(\mathbf{x}_i) + \alpha \hat{y}(\mathbf{x}_i) \right) \right) \quad (7.5b)$$

$$= \arg \min_{(\alpha, \hat{y})} \sum_{i=1}^n \underbrace{\exp \left(-y_i c^{(b-1)}(\mathbf{x}_i) \right)}_{=w_i^{(b)}} \exp (-y_i \alpha \hat{y}(\mathbf{x}_i)), \quad (7.5c)$$

where for the first equality we have used the definition of the exponential loss function (5.12) and the sequential structure of the boosted classifier (7.4). The last equality is where the convenience of the exponential loss appears, namely the fact that $\exp(a + b) = \exp(a) \exp(b)$. This allows us to define the quantities

$$w_i^{(b)} \stackrel{\text{def}}{=} \exp \left(-y_i c^{(b-1)}(\mathbf{x}_i) \right), \quad (7.6)$$

which can be interpreted as *weights* for the individual data points in the training dataset. Note that the weights $w_i^{(b)}$ are independent of α and \hat{y} . That is, when learning $\hat{y}^{(b)}(\mathbf{x})$ and its coefficient $\alpha^{(b)}$ by solving (7.5c) we can regard $\{w_i^{(b)}\}_{i=1}^n$ as constants.

To solve (7.5) we start by rewriting the objective function as

$$\sum_{i=1}^n w_i^{(b)} \exp (-y_i \alpha \hat{y}(\mathbf{x}_i)) = e^{-\alpha} \underbrace{\sum_{i=1}^n w_i^{(b)} \mathbb{I}\{y_i = \hat{y}(\mathbf{x}_i)\}}_{=W_c} + e^{\alpha} \underbrace{\sum_{i=1}^n w_i^{(b)} \mathbb{I}\{y_i \neq \hat{y}(\mathbf{x}_i)\}}_{=W_e}, \quad (7.7)$$

where we have used the indicator function to split the sum into two sums: the first ranging over all training data points correctly classified by \hat{y} and the second ranging over all points erroneously classified by \hat{y} . (Remember that \hat{y} is the ensemble member we are to learn at this step.) Furthermore, for notational

simplicity we define W_c and W_e for the sum of weights of correctly classified and erroneously classified data points, respectively. Furthermore, let $W = W_c + W_e$ be the total weight sum, $W = \sum_{i=1}^n w_i^{(b)}$.

Minimizing (7.7) is done in two stages, first w.r.t. \hat{y} and then w.r.t. α . This is possible since the minimizing argument in \hat{y} turns out to be independent of the actual value of $\alpha > 0$, another convenient effect of using the exponential loss function. To see this, note that we can write the objective function (7.7) as

$$e^{-\alpha}W + (e^\alpha - e^{-\alpha})W_e. \quad (7.8)$$

Since the total weight sum W is independent of \hat{y} and since $e^\alpha - e^{-\alpha} > 0$ for any $\alpha > 0$, minimizing this expression w.r.t. \hat{y} is equivalent to minimizing W_e w.r.t. \hat{y} . That is,

$$\hat{y}^{(b)} = \arg \min_{\hat{y}} \sum_{i=1}^n w_i^{(b)} \mathbb{I}\{y_i \neq \hat{y}(\mathbf{x}_i)\}. \quad (7.9)$$

In words, the b^{th} ensemble member should be trained by minimizing the *weighted misclassification loss*, where each data point (\mathbf{x}_i, y_i) is assigned a weight $w_i^{(b)}$. The intuition for these weights is that, at iteration b , we should focus our attention on the data points previously misclassified in order to “correct the mistakes” made by the ensemble of the first $b - 1$ classifiers.

Time to reflect 7.2: In AdaBoost, we use the exponential loss for training the boosting ensemble. How come that we end up training the individual ensemble members using a weighted misclassification loss (and not the unweighted exponential loss) then?

How the problem (7.9) is solved in practice depends on the choice of base classifier that we use, i.e. on the specific restrictions that we put on the function \hat{y} (for example a shallow classification tree). However, solving (7.9) is almost our standard classification problem, except for the weights $w_i^{(b)}$. Training the ensemble member b on a *weighted* classification problem is, for most base classifiers, straightforward. Since most classifiers are trained by minimizing some cost function, this simply boils down to weighting the individual terms of the cost function and solve that slightly modified problem instead.

When the b^{th} ensemble member, $\hat{y}^{(b)}(\mathbf{x})$, has been trained for solving the weighted classification problem (7.9) it remains to learn its coefficient $\alpha^{(b)}$. This is done by solving (7.5), which amounts to minimizing (7.8) once \hat{y} has been trained. By differentiating (7.8) w.r.t. α and setting the derivative to zero we get the equation

$$-\alpha e^{-\alpha}W + \alpha(e^\alpha + e^{-\alpha})W_e = 0 \Leftrightarrow W = \left(e^{2\alpha} + 1\right)W_e \Leftrightarrow \alpha = \frac{1}{2} \ln\left(\frac{W}{W_e} - 1\right).$$

Thus, by defining

$$E_{\text{train}}^{(b)} \stackrel{\text{def}}{=} \frac{W_e}{W} = \sum_{i=1}^n \frac{w_i^{(b)}}{\sum_{j=1}^n w_j^{(b)}} \mathbb{I}\{y_i \neq \hat{y}^{(b)}(\mathbf{x}_i)\} \quad (7.10)$$

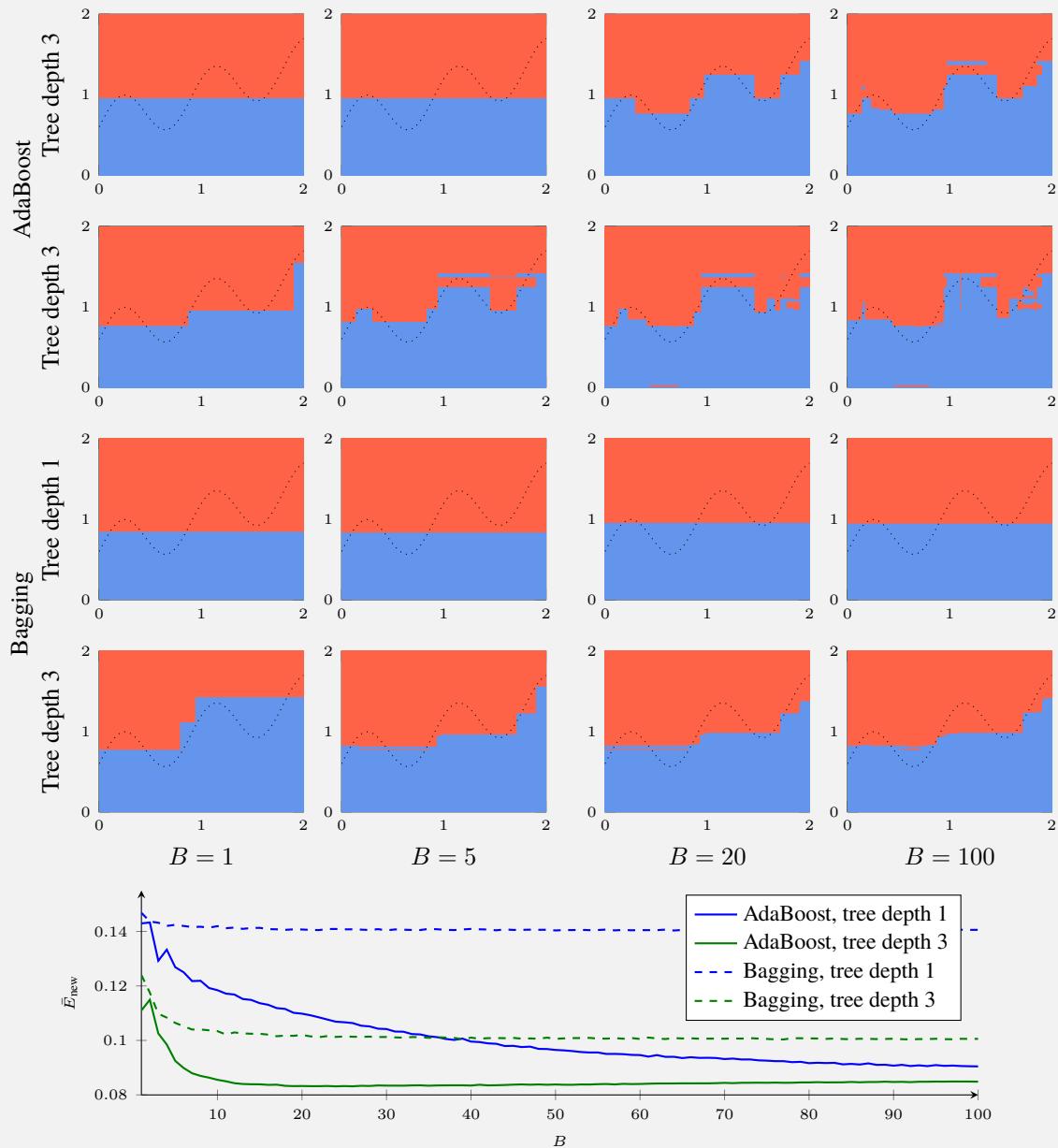
to be the weighted misclassification error for the b^{th} classifier, we can express the optimal value for its coefficient as

$$\alpha^{(b)} = \frac{1}{2} \ln\left(\frac{1 - E_{\text{train}}^{(b)}}{E_{\text{train}}^{(b)}}\right). \quad (7.11)$$

This completes the derivation of the AdaBoost algorithm, which is summarized in Algorithm 13. In the algorithm we exploit the fact that the weights (7.6) can be computed recursively by using the expression (7.4) in line 5. Furthermore, we have added an explicit weight normalization (line 6) which is convenient in practice and which does not affect the derivation of the method above.

Example 7.6: AdaBoost and bagging for a binary classification example

Consider the same binary classification problem as in Example 7.4. We now compare how AdaBoost and bagging performs on this problem, when using trees of depth one (decision stumps) and three, respectively. The decision boundaries for each method with $B = 1, 5, 20$ and 100 ensemble members are shown below. Despite using quite weak ensemble members (a shallow tree has high bias), AdaBoost adapts quite well to the data (Example 7.4). This is in contrast to bagging, where the decision boundary does not become much more flexible despite using many ensemble members. In other words, AdaBoost reduces the bias of the base model, whereas bagging only has minor effect on the bias.



We also numerically compute \bar{E}_{new} for this problem, as a function of B , which is shown above. Remember that \bar{E}_{new} depends on both the bias and the variance. As discussed, the main effect of bagging is variance reduction, but that does not help much since the base model already is quite low-variance (but high-bias). Boosting, on the other hand, reduces bias, which has a much bigger effect in this case. Furthermore, bagging does not overfit as $B \rightarrow \infty$, but that is *not* the case for boosting! We can indeed see that for trees of depth 3, the smallest \bar{E}_{new} is obtained for $B \approx 25$, and there is actually a slight increase in \bar{E}_{new} for larger values of B . Hence, AdaBoost with depth-3 trees suffers from a (minor) overfit as $B \gtrsim 25$ in this problem.

Algorithm 13: AdaBoost

Data: Training dataset $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$
Result: A prediction $\hat{y}^{\text{boost}}(\mathbf{x}_*)$

- 1 Assign weights $w_i^{(1)} = 1/n$ to all data points.
- 2 **for** $b = 1, \dots, B$ **do**
- 3 Train a weak classifier $\hat{y}^{(b)}(\mathbf{x})$ on the weighted training data $\{(\mathbf{x}_i, y_i, w_i^{(b)})\}_{i=1}^n$.
- 4 Compute $E_{\text{train}}^{(b)} = \sum_{i=1}^n w_i^{(b)} \mathbb{I}\{y_i \neq \hat{y}^{(b)}(\mathbf{x}_i)\}$
- 5 Compute $\alpha^{(b)} = 0.5 \ln((1 - E_{\text{train}}^{(b)})/E_{\text{train}}^{(b)})$
- 6 Compute $w_i^{(b+1)} = w_i^{(b)} \exp(-\alpha^{(b)} y_i \hat{y}^{(b)}(\mathbf{x}_i)), i = 1, \dots, n$
- 7 Set $w_i^{(b+1)} \leftarrow w_i^{(b+1)} / \sum_{j=1}^n w_j^{(b+1)}$, for $i = 1, \dots, n$
- 8 **end**
- 9 Output $\hat{y}_{\text{boost}}^{(B)}(\mathbf{x}) = \text{sign}\left\{\sum_{b=1}^B \alpha^{(b)} \hat{y}^{(b)}(\mathbf{x})\right\}$

Remark 7.2 The derivation of AdaBoost assumes that all coefficients $\{\alpha^{(b)}\}_{b=1}^{(B)}$ are positive. To see that this is indeed the case when the coefficients are computed according to (7.11), note that the function $\ln((1 - x)/x)$ is positive for any $0 < x < 0.5$. Thus, $\alpha^{(b)}$ will be positive as long as the weighted training error for the b^{th} classifier, $E_{\text{train}}^{(b)}$, is less than 0.5. That is, the classifier just has to be slightly better than a coin flip, which is always the case in practice (note that $E_{\text{train}}^{(b)}$ is the training error). (Indeed, if $E_{\text{train}}^{(b)} > 0.5$, then we could simply flip the sign of all predictions made by $\hat{y}^{(b)}(\mathbf{x})$ to reduce the error below 0.5.)

User aspects of AdaBoost

AdaBoost, and in fact any boosting algorithm, has two important design choices, (i) which base classifier to use, an (ii) how many iterations B to run the boosting algorithm for. As previously pointed out, we can use essentially any classification method as base classifier. However, the most common choice in practice is to use a shallow classification tree, or even a decision stump (a tree of depth one; see Example 7.5). This choice is guided by the fact that boosting reduces bias efficiently, and can thereby learn good models despite using a very weak (high-bias) base model. Since shallow trees can be trained quickly, they are a good default choice. Practical experience suggests that trees with roughly 6 terminal nodes often work well as base models, but trees of depth one (only $M = 2$ terminal nodes in binary classification) are perhaps even more commonly used. In fact, using deep classification trees (high-variance models) as base classifiers typically deteriorates performance.

The base models are learned sequentially in boosting, each iteration introduces a new base model aiming at reducing the errors made by the current model. As an effect, the boosting model becomes more and more flexible as the number of iterations B increases and using too many base models can result in overfitting (in contrast to bagging, where increased B cannot lead to overfit). It has been observed in practice, however, that this overfitting often occurs slowly and the performance tends to be rather insensitive to the choice of B . Nevertheless, it is a good practice to select B in some systematic way, for instance using early stopping during training, similarly to how it is used for neural networks (recall Section 6.3). Another unfortunate aspect of the sequential nature of boosting is that it is not possible to parallelize the learning.

In the method discussed above we have assumed that each base classifier outputs a class prediction, $\hat{y}^{(b)}(\mathbf{x}) \in \{-1, 1\}$. However, many classification models output $g(\mathbf{x})$, which is an estimate of the class probability $p(y = 1 | \mathbf{x})$. In AdaBoost it is possible to use the predicted probabilities $g(\mathbf{x})$ (instead of the binary prediction $\hat{y}(\mathbf{x})$) when constructing the prediction, however at the cost of a more complicated expression than (7.3). This extension of Algorithm 13 is referred to as Real AdaBoost.

7.4 Gradient boosting

It has been seen in practice that AdaBoost performs well if there is little noise (few mislabeled data points) in the training data. If the training data has more noise and hence contains more outliers, the performance of AdaBoost typically deteriorates. That is not an artifact of the boosting idea, but of the exponential loss function that we discussed in Section 5.1. It is possible to construct more robust boosting algorithms by choosing another loss function, but the training becomes more computationally involved, as we will see.

A first gradient boosting algorithm

If the squared error loss in linear regression is replaced with, say, absolute error loss, the simple analytical solution (the normal equations) is lost. We can, however, still learn the model by using numerical optimization. The situation is quite similar for boosting. The simplicity of AdaBoost is that the boosting problem is turned into a sequence of weighted classification problems, which we can solve (at least approximately) by almost any off-the-shelf classification method. If the exponential loss in (7.5a) is replaced, the problem does not separate into weighted classification problems anymore, and the situation becomes more intricate.

It turns out to be possible to approximately solve (7.5a) for rather general loss functions using a method reminiscent of gradient descent ((5.25) in Chapter 5). The resulting method is referred to as *gradient boosting*.

Like AdaBoost (7.4), we construct the classifier sequentially as

$$c^{(b)}(\mathbf{x}) = c^{(b-1)}(\mathbf{x}) + \alpha^{(b)}\tilde{y}^{(b)}(\mathbf{x}), \quad (7.12)$$

for $b = 1, \dots, B$, and the goal is to sequentially select $\{\alpha^{(b)}, \tilde{y}^{(b)}(\mathbf{x})\}$ such that the final $c^{(B)}(\mathbf{x})$ minimizes

$$J(c) = \frac{1}{n} \sum_{i=1}^n L(y_i \cdot c(\mathbf{x}_i)) \quad (7.13)$$

for any arbitrary differentiable loss function L , such as the logistic loss.

The idea of gradient boosting is to think of

$$\underbrace{\begin{bmatrix} c^{(B)}(\mathbf{x}_1) \\ \vdots \\ c^{(B)}(\mathbf{x}_n) \end{bmatrix}}_{c^{(B)}(\mathbf{X})}$$

as an n -dimensional optimization variable, which by construction is build up in the additive fashion (7.12). We then observe that the boosted model (7.12) actually looks like a gradient descent update (recall (5.25) in Chapter 5), if (7.13) would be the objective function $J(c(\mathbf{X}))$ and $\tilde{y}^{(b)}(\mathbf{x})$ its gradient $\nabla_c J(c^{(b-1)}(\mathbf{X}))$.

With this inspiration from gradient descent, we would like to choose $\tilde{y}^{(b)}(\mathbf{x})$ in (7.12) such that

$$\begin{bmatrix} \tilde{y}^{(b)}(\mathbf{x}_1) \\ \vdots \\ \tilde{y}^{(b)}(\mathbf{x}_n) \end{bmatrix} \text{ is close to } \underbrace{\begin{bmatrix} \frac{\partial L(y_1 \cdot c^{(b-1)}(\mathbf{x}_1))}{\partial c} \\ \vdots \\ \frac{\partial L(y_n \cdot c^{(b-1)}(\mathbf{x}_n))}{\partial c} \end{bmatrix}}_{\nabla_c J(c^{(b-1)}(\mathbf{X}))}. \quad (7.14)$$

By the arguments of gradient descent this would as $b \rightarrow \infty$ lead to a $c^{(B)}(\mathbf{x})$ which is, at least approximately, a local minimizer of (7.13).

We have now outlined the idea of gradient boosting, but it is not yet a concrete algorithm. Specifically we have to be more precise about $\tilde{y}^{(b)}$ in (7.14). In gradient boosting, we simply handle (7.14) as a

regression problem where $\hat{y}^{(b)}$ is learned with input \mathbf{x}_i and output $\frac{\partial L(y_i \cdot c^{(b-1)}(\mathbf{x}_i))}{\partial c}$ ($i = 1, \dots, n$). (We have assumed that the loss function L is differentiable, meaning that we can always compute the necessary derivatives.) That is, gradient boosting builds a classifier by using regression. The type of base regression model \hat{y} is in principle arbitrary, but regression trees are often used in practice.

We have not yet discussed how to choose $\alpha^{(b)}$, which corresponds to γ (the step size or learning rate) in gradient descent. In the simplest version of gradient descent it is considered a tuning choice left to the user, but it can also be formulated as a line-search optimization problem itself to choose the optimal value at each iteration. For gradient boosting, it is most often handled in the latter way. When using trees as base models optimizing $\alpha^{(b)}$ can be done jointly with learning \hat{y} , resulting in a more efficient implementation. If multiplying the optimal $\alpha^{(b)}$ with a parameter < 1 , a regularizing effect is obtained which has proven useful in practice.

We summarize gradient boosting by Algorithm 14.

Algorithm 14: A gradient boosting algorithm

1. Initialize (as a constant), $c^0(\mathbf{x}) \equiv \arg \min_c \sum_{i=1}^n L(y_i, c)$.

2. For $b = 1$ to B

- (a) Compute the negative gradient of the loss function,

$$g_i^{(b)} = - \left[\frac{\partial L(y_i, c)}{\partial c} \right]_{c=c^{(b-1)}(\mathbf{x}_i)}, \quad i = 1, \dots, n.$$

- (b) Train a base *regression* model $\hat{f}^{(b)}(\mathbf{x})$ to fit the gradient values,

$$\hat{f}^{(b)} = \arg \min_f \sum_{i=1}^n \left(f(\mathbf{x}_i) - g_i^{(b)} \right)^2.$$

- (c) Update the boosted model,

$$c^{(b)}(\mathbf{x}) = c^{(b-1)}(\mathbf{x}) + \gamma \hat{f}^{(b)}(\mathbf{x})$$

3. Output $\hat{y}_{\text{boost}}^{(B)}(\mathbf{x}) = \text{sign}\{c^{(B)}(\mathbf{x})\}$.

While presented for classification in Algorithm 14, gradient boosting can also be used for regression with minor modifications. As mentioned above, gradient boosting requires a certain amount of smoothness in the loss function. A minimal requirement is that it is almost everywhere differentiable, so that it is possible to compute the gradient of the loss function. However, some implementations of gradient boosting require stronger conditions, such as second order differentiability. The logistic loss (see Figure 5.2) is in this respect a “safe choice” which is infinitely differentiable and strongly convex, while still enjoying good statistical properties. As a consequence, the logistic loss is one of the most commonly used loss functions in practice.

8 Nonlinear input transformations and kernels

In this chapter we will continue to develop the idea from Chapter 3 of creating new input features by using nonlinear transformations $\phi(\mathbf{x})$. It turns out that by the so-called *kernel trick*, we can have *infinitely* many such nonlinear transformations and we can extend our basic methods, such as linear regression and k -NN, into more versatile and flexible ones. When we also change the loss function of linear regression, we obtain support vector machines, another powerful off-the-shelf machine learning method. The concept of kernels is important also to the next chapter (Chapter 9), where a Bayesian perspective of linear regression and kernels leads us to the Gaussian process.

8.1 Creating features by nonlinear input transformations

The reason for the word ‘linear’ in the name ‘linear regression’ is that the output is modelled as a *linear* combination of the inputs. We have, however, not made a clear definition of what an input is: if the speed is an input in Example 2.2, then why could not also the kinetic energy—the square of the speed—be considered as another input? The answer is yes, it can. We can in fact make use of arbitrary nonlinear transformations of the “original” input variables in any model, including linear regression. If we, for example, only have a one-dimensional input x , the vanilla linear regression model is

$$y = \theta_0 + \theta_1 x + \varepsilon. \quad (8.1)$$

Starting from this we can extend the model with x^2, x^3, \dots, x^{d-1} as inputs (d is a user-choice), and thus obtain a linear regression model which is a polynomial in x ,

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_{d-1} x^{d-1} + \varepsilon = \boldsymbol{\theta}^\top \phi(x) + \varepsilon. \quad (8.2)$$

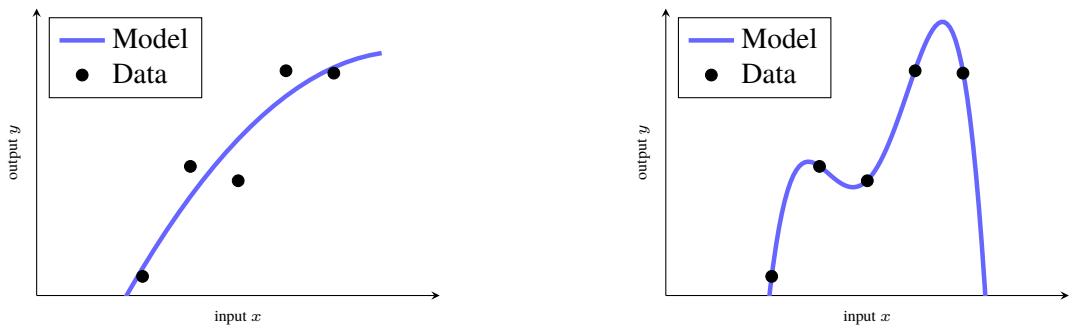
Since x is known, we can directly compute x^2, \dots, x^{d-1} . Note that this is still a linear regression model since the parameters $\boldsymbol{\theta}$ appear in a linear fashion with $\phi(x) = [1 \ x \ x^2 \ \dots \ x^{d-1}]^\top$ as a new input vector. We refer to a transformation of \mathbf{x} as a *feature*¹ and the vector of transformed inputs $\phi(\mathbf{x})$, a vector of dimension $d \times 1$, as a *feature vector*. The parameters $\hat{\boldsymbol{\theta}}$ are still learned the same way, but we

replace the original $\mathbf{X} = \underbrace{\begin{bmatrix} -\mathbf{x}_1^\top - \\ -\mathbf{x}_2^\top - \\ \vdots \\ -\mathbf{x}_n^\top - \end{bmatrix}}_{n \times p+1}$	with the transformed $\Phi(\mathbf{X}) = \underbrace{\begin{bmatrix} -\phi(\mathbf{x}_1)^\top - \\ -\phi(\mathbf{x}_2)^\top - \\ \vdots \\ -\phi(\mathbf{x}_n)^\top - \end{bmatrix}}_{n \times d}$
--	--

in the cost function. For linear regression, this means that we can learn the parameters by doing the substitution (8.3) directly in the normal equations (3.13).

The idea of nonlinear input transformations is not unique to linear regression, but any choice of nonlinear transformation $\phi(\cdot)$ can be used with any supervised machine learning methods: the nonlinear transformation is first applied to the input, like a pre-processing step, and the transformed input is thereafter used when training, evaluating and using the model. We illustrated this for regression already by Example 3.5 in Chapter 3, and for classification by Example 8.1.

¹The original input \mathbf{x} is sometimes also referred to as a feature.



(a) A linear regression model with a 2nd order polynomial, trained with squared error loss. The line is no longer straight (cf. Figure 3.1), but this is merely an artifact of the plot: in a three-dimensional plot with each feature (here, x and x^2) on a separate axis, it would still be an affine model.

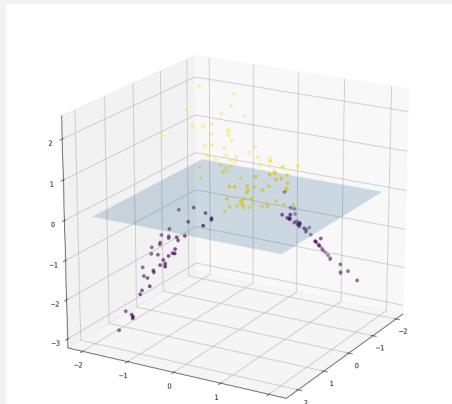
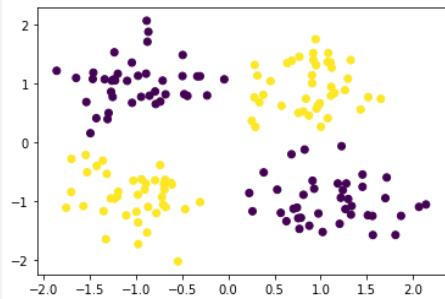
(b) A linear regression model with a 4th order polynomial, trained with squared error loss. Note that a 4th order polynomial implies 5 unknown parameters, which roughly means that we can expect the learned model to fit 5 data points exactly, a typical case of overfit.

Figure 8.1: A linear regression model with 2nd and 4th order polynomials in the input x , as shown in (8.2).

Time to reflect 8.1: Figure 8.1 shows an example of two linear regression models with transformed (polynomial) inputs. When studying the figure one may ask how a linear regression model can result in a curved line? Are linear regression models not restricted to linear (or affine) straight lines?

Example 8.1: Nonlinear feature transformations for classification

This example is not yet complete.



Polynomials are only one out of (infinitely) many possible choices of features $\phi(\mathbf{x})$. One should take care when using polynomials higher than second order in practice, because of their behavior outside the range where the data is observed (recall Figure 8.1b). Instead, there are several alternatives that are often more useful in practice, such as Fourier series, corresponding to (for scalar x) something like $\phi(x) = [1 \sin(x) \cos(x) \sin(2x) \cos(2x) \dots]^T$, step functions, regression splines, etc. The use of nonlinear input transformations $\phi(\mathbf{x})$ arguably makes simple models more flexible and applicable to real-world problems with nonlinear characteristics. In order to obtain a good performance, it is important to choose $\phi(\mathbf{x})$ such that enough flexibility is obtained, but overfit avoided. With a very careful choice of $\phi(\mathbf{x})$ good performance can be obtained for certain problems, but that choice is problem-specific and perhaps more of an art form than a science. Let us instead explore the conceptual idea of letting the number of features $d \rightarrow \infty$, and combine it with regularization. In a sense this will automate the choice of features, and leads us to a family of powerful off-the-shelf machine learning tools called kernel methods.

8.2 Kernel ridge regression

A carefully engineered transformation $\phi(x)$ in linear regression, or any other method for that matter, may indeed perform well for a specific machine learning problem. However, we would like $\phi(x)$ to contain “all” transformations that could possibly be of interest for most problems, in order to obtain a general off-the-shelf method. We will therefore explore the idea of choosing d really big, much bigger than the number of data points n , and eventually even let $d \rightarrow \infty$. The derivation and reasoning will be done using L^2 regularized linear regression, but we will later see that the idea is applicable also to other models.

Re-formulating linear regression

First of all we have to use some kind of regularization if we are going to increase d in linear regression, in order to avoid overfit when $d \gg n$. For reasons that we will discuss later we chose to use L^2 regularization. Let us repeat the equation for L^2 -regularized linear regression,

$$\widehat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \left(\underbrace{\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}_i)}_{\widehat{v}(\mathbf{x}_i)} - y_i \right)^2 + \lambda \|\boldsymbol{\theta}\|_2^2 = (\boldsymbol{\Phi}(\mathbf{X})^\top \boldsymbol{\Phi}(\mathbf{X}) + n\lambda \mathbf{I})^{-1} \boldsymbol{\Phi}(\mathbf{X})^\top \mathbf{y}, \quad (8.4a)$$

$$\hat{y}(\mathbf{x}_\star) = \hat{\boldsymbol{\theta}}^\top \phi(\mathbf{x}_\star) \quad (8.4b)$$

We have not fixed the nonlinear transformations $\phi(\mathbf{x})$ to anything specific yet, but we are preparing for choosing $d \gg n$ of these transformations. The downside of choosing d , the dimension of $\phi(\mathbf{x})$, large is that we also have to learn d parameters when training. In linear regression, we usually first learn and store the d -dimensional vector $\hat{\theta}$, and thereafter we use it for computing a prediction. To be able to choose d really large, perhaps even $d \rightarrow \infty$, we have to re-formulate the model such that there are no computations or storage demands that scales with d . The first step is to realize that the prediction $\hat{y}(\mathbf{x}_*)$ can be rewritten as

(where the underbraces provides the size of the corresponding vectors and matrices). This expression for $\widehat{y}(\mathbf{x}_*)$ suggests that instead of computing and storing the d -dimensional $\widehat{\boldsymbol{\theta}}$ once (independently of \mathbf{x}_*) we could, for each test input \mathbf{x}_* , compute the n -dimensional vector $\Phi(\mathbf{X})(\Phi(\mathbf{X})^\top \Phi(\mathbf{X}) + n\lambda \mathbf{I})^{-1} \boldsymbol{\phi}(\mathbf{x}_*)$. By doing so, we avoid storing a d -dimensional vector, but this would still require the inversion of a $d \times d$ matrix. Hence we have some more work to do before we have a practically useful method where we can select d really large.

The push-through matrix identity says that $\mathbf{A}(\mathbf{A}^\top \mathbf{A} + \mathbf{I})^{-1} = (\mathbf{A}\mathbf{A}^\top + \mathbf{I})^{-1}\mathbf{A}$ holds for any matrix \mathbf{A} . By using it in (8.5), we can write $\hat{\mathbf{y}}(\mathbf{x}_*)$ as

$$\hat{y}(\mathbf{x}_\star) = \underbrace{\mathbf{y}^\top}_{1 \times n} \underbrace{(\Phi(\mathbf{X})\Phi(\mathbf{X})^\top + n\lambda\mathbf{I})^{-1}}_{n \times n} \underbrace{\Phi(\mathbf{X})\phi(\mathbf{x}_\star)}_{n \times 1}. \quad (8.6)$$

It now seems as if we can compute $\widehat{y}(\mathbf{x}_\star)$ without having to deal with any d -dimensional vectors or matrices. That requires, however, that the matrix multiplication $\Phi(\mathbf{X})\Phi(\mathbf{X})^\top$ and $\Phi(\mathbf{X})\phi(\mathbf{x}_\star)$ in (8.6) can

be computed. Let us therefore have a closer look at these,

$$\Phi(\mathbf{X})\Phi(\mathbf{X})^\top = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_n) \\ \phi(\mathbf{x}_2)^\top \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)^\top \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_2)^\top \phi(\mathbf{x}_n) \\ \vdots & & \ddots & \vdots \\ \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_1) & \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_n) \end{bmatrix} \quad (8.7)$$

$$\text{and } \Phi(\mathbf{X})\phi(\mathbf{x}_*) = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_*) \\ \phi(\mathbf{x}_2)^\top \phi(\mathbf{x}_*) \\ \vdots \\ \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_*) \end{bmatrix}. \quad (8.8)$$

Remember that $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$ is an inner product between the two d -dimensional vectors $\phi(\mathbf{x})$ and $\phi(\mathbf{x}')$. The key insight here is to note that the transformed inputs $\phi(\mathbf{x})$ only enters into (8.6) as inner products $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$. That is, if we are able to compute the inner product $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$ directly, without first explicitly computing $\phi(\mathbf{x})$, we are safe.

As a concrete illustration, let us use polynomials since they are familiar and have no complicated expressions. With $p = 1$, meaning \mathbf{x} is a scalar x , and $\phi(x)$ is a third order polynomial ($d = 4$) with the second and third term scaled² by $\sqrt{3}$, we have

$$\phi(x)^\top \phi(x') = [1 \quad \sqrt{3}x \quad \sqrt{3}x^2 \quad x^3] \begin{bmatrix} 1 \\ \sqrt{3}x' \\ \sqrt{3}x'^2 \\ x'^3 \end{bmatrix} = 1 + 3xx' + 3x^2x'^2 + x^3x'^3 = (1 + xx')^3. \quad (8.9)$$

It can generally be shown that if $\phi(x)$ is a (scaled) polynomial of order $d - 1$, then $\phi(x)^\top \phi(x') = (1 + xx')^{d-1}$. The point we want to make is that instead of first computing the two d -dimensional vectors $\phi(x)$ and $\phi(x')$ and thereafter computing their inner product, we could just evaluate the expression $(1 + xx')^{d-1}$ directly instead. With a second or third order polynomial this might not make much of a difference, but imagine a situation where it is of interest to use d in the hundreds or thousands.

The key insight we are getting at is that *if we only make the choice of $\phi(\mathbf{x})$ such that the inner product $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$ can be computed without first computing $\phi(\mathbf{x})$, we can let d be arbitrary big*. Since it is possible to define inner products also between infinite-dimensional vectors, there is nothing preventing us from letting $d \rightarrow \infty$.

We have now derived a version of L^2 -regularized linear regression that we can use in practice also with an unbounded number of features d in $\phi(\mathbf{x})$, if we only restrict ourselves to $\phi(\mathbf{x})$ such that its inner product $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$ has a closed-form expression (or can, at least, be computed in such a way that it does not scale with d). This might appear to be of rather limited interest for a machine learning engineer, since one still has to come up with a nonlinear transformation $\phi(\mathbf{x})$, choose d (possibly ∞) and thereafter make a pen-and-paper derivation (like (8.9)) of $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$. To alleviate this, let us introduce the concept of a *kernel*.

Introducing the kernel idea

A kernel $\kappa(\mathbf{x}, \mathbf{x}')$ is (in this book) any function that takes two arguments \mathbf{x} and \mathbf{x}' from the same space, and returns a scalar. Throughout this book, we will limit ourselves to kernels that are real-valued and symmetric, that is, $\kappa(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x}', \mathbf{x}) \in \mathbb{R}$ for all \mathbf{x} and \mathbf{x}' . The inner product of two nonlinear input transformations is an example of a kernel

$$\kappa(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}'). \quad (8.10)$$

²The scaling $\sqrt{3}$ can be compensated with an inverse scaling of the second and third element in θ .

Equation (8.9), for example, is such a kernel. The key insight at this stage is that since $\phi(\mathbf{x})$ only appears in the linear regression model (8.6) via inner products, we do not have to engineer a d -dimensional vector $\phi(\mathbf{x})$ and derive its inner product. Instead, we can just choose a kernel $\kappa(\mathbf{x}, \mathbf{x}')$ directly. This is known as the *kernel trick*,

If \mathbf{x} only enters the model as $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$, we can choose a kernel $\kappa(\mathbf{x}, \mathbf{x}')$ instead of choosing $\phi(\mathbf{x})$.

To be clear on what this means in practice, we rewrite (8.6) using the kernel (8.10),

$$\hat{y}(\mathbf{x}_*) = \underbrace{\mathbf{y}^\top}_{1 \times n} \underbrace{(\mathbf{K}(\mathbf{X}, \mathbf{X}) + n\lambda\mathbf{I})^{-1}}_{n \times n} \underbrace{\mathbf{K}(\mathbf{X}, \mathbf{x}_*)}_{n \times 1}, \quad (8.11a)$$

$$\mathbf{K}(\mathbf{X}, \mathbf{X}) = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \kappa(\mathbf{x}_1, \mathbf{x}_2) & \dots & \kappa(\mathbf{x}_1, \mathbf{x}_n) \\ \kappa(\mathbf{x}_2, \mathbf{x}_1) & \kappa(\mathbf{x}_2, \mathbf{x}_2) & \dots & \kappa(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \ddots & \ddots & \vdots \\ \kappa(\mathbf{x}_n, \mathbf{x}_1) & \kappa(\mathbf{x}_n, \mathbf{x}_2) & \dots & \kappa(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}, \quad (8.11b)$$

$$\mathbf{K}(\mathbf{X}, \mathbf{x}_*) = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_*) \\ \kappa(\mathbf{x}_2, \mathbf{x}_*) \\ \vdots \\ \kappa(\mathbf{x}_n, \mathbf{x}_*) \end{bmatrix}. \quad (8.11c)$$

These equations describes linear regression with L^2 regularization using a kernel $\kappa(\mathbf{x}, \mathbf{x}')$. Since L^2 regularization also is called ridge regression, we refer to (8.11) as kernel ridge regression. We initially argued that linear regression with a possibly infinite-dimensional nonlinear transformation vector $\phi(\mathbf{x})$ could be an interesting model, and (8.11) is (for certain choices of $\phi(\mathbf{x})$ and $\kappa(\mathbf{x}, \mathbf{x}')$) equivalent to that. The design choice for the user is now to select a kernel $\kappa(\mathbf{x}, \mathbf{x}')$ instead of $\phi(\mathbf{x})$. In practice, choosing $\kappa(\mathbf{x}, \mathbf{x}')$ is a much less tedious problem than choosing $\phi(\mathbf{x})$.

As users we may in principle choose the kernel $\kappa(\mathbf{x}, \mathbf{x}')$ arbitrarily, as long as we can compute (8.11a). This requires the inverse of $(\mathbf{K}(\mathbf{X}, \mathbf{X}) + n\lambda\mathbf{I})$ to exist. We are therefore on the safe side if we restrict ourselves to kernels for which the matrix $\mathbf{K}(\mathbf{X}, \mathbf{X})$ (8.11b) is always positive semidefinite. Such kernels are called³ positive semidefinite kernels.

The user of kernel ridge regression chooses a positive semidefinite kernel $\kappa(\mathbf{x}, \mathbf{x}')$, and does not have to select (and even less compute) $\phi(\mathbf{x})$. However, a corresponding $\phi(\mathbf{x})$ always exists for a positive semidefinite kernel as we will discuss in Section 8.4.

There is a number of positive semidefinite kernels that a user can choose between. A commonly used positive semidefinite kernel is the squared exponential kernel

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right), \quad (8.12)$$

where the hyperparameter $\ell > 0$ is a design choice left to the user, for example to be chosen using cross validation. Another example of a positive semidefinite kernel mentioned earlier is the polynomial kernel $\kappa(\mathbf{x}, \mathbf{x}') = (c + \mathbf{x}^\top \mathbf{x}')^{d-1}$. A special case thereof is the linear kernel $\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$. We will give more examples later.

From the formulation (8.11), it may seem as if we have to compute the inverse of $(\mathbf{K}(\mathbf{X}, \mathbf{X}) + n\lambda\mathbf{I})$ every time we want to make a prediction. That is, however, not necessary since it does not depend on the test input \mathbf{x}_* . It is therefore wise to introduce the n -dimensional vector

$$\boldsymbol{\alpha} = \begin{bmatrix} \hat{\alpha}_1 \\ \hat{\alpha}_2 \\ \vdots \\ \hat{\alpha}_n \end{bmatrix} \quad (8.13a)$$

³Confusingly enough, such kernels are called positive definite in some texts.

which allows us to rewrite kernel ridge regression (8.11) as

$$\hat{y}(\mathbf{x}_*) = \hat{\alpha}^\top \mathbf{K}(\mathbf{X}, \mathbf{x}_*), \quad (8.13b)$$

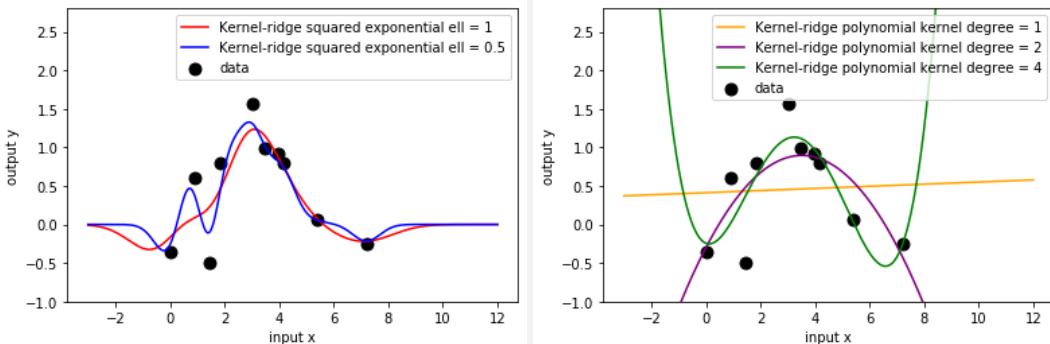
$$\hat{\alpha} = \mathbf{y}^\top (\mathbf{K}(\mathbf{X}, \mathbf{X}) + n\lambda \mathbf{I})^{-1}. \quad (8.13c)$$

That is, instead of computing and storing a d -dimensional vector $\hat{\theta}$ as we intended initially, we now instead compute and store an n -dimensional vector $\hat{\alpha}$. However, we also need to store \mathbf{X} , since we have to compute $\mathbf{K}(\mathbf{X}, \mathbf{x}_*)$ for every prediction.

We summarize kernel ridge regression in Algorithm 15. Kernel ridge regression is in itself a practically useful method, but we will next take a step back and discuss what we have derived, in order to prepare for more kernel methods. We will also come back to kernel ridge regression in Chapter 9, where it is used as a stepping stone in deriving the Gaussian process.

Example 8.2: Linear regression with kernels

This example is not yet complete.



Time to reflect 8.2: Verify that you retrieve L^2 -regularized linear regression, without any nonlinear transformations, by using the linear kernel $\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$ in (8.11).

Algorithm 15: Kernel ridge regression

1 To be completed

8.3 Support vector regression

Kernel ridge regression, as we just derived, is our first kernel method for regression, which is a useful method on its own. We will now develop kernel ridge regression into support vector regression by replacing the loss function. First, however, we take a step back and make an interesting observation that suggests the so-called representer theorem, which will be useful later in this chapter.

Preparing for more kernel methods: the representer theorem

The formulation (8.13) is not only practical for implementation, it is also important for theoretical understanding. We can interpret (8.13) as a *dual formulation* of linear regression, where we have the *dual parameters* α instead of the primal formulation (8.4) with primal parameters θ . Remember that d , the (possibly infinite) number of primal parameters in θ , is a user design choice, whereas n , the number of dual parameters in α , is the number of data points.

By comparing (8.5) and (8.4b), we have that

$$\hat{y}(\mathbf{x}_*) = \hat{\theta}^\top \phi(\mathbf{x}_*) = \hat{\alpha}^\top \underbrace{\Phi(\mathbf{X})\phi(\mathbf{x}_*)}_{K(\mathbf{X}, \mathbf{x}_*)} \quad (8.14)$$

for all \mathbf{x}_* , which suggests that

$$\hat{\theta} = \Phi(\mathbf{X})^\top \hat{\alpha}. \quad (8.15)$$

This relationship between the primal parameters θ and the dual parameters α is not specific for kernel ridge regression, but (8.15) a general result called *the representer theorem*.

In essence the representer theorem states that if $\hat{y}(\mathbf{x}) = \theta^\top \phi(\mathbf{x})$, the equation (8.15) holds when θ is learned using (almost) any loss functions and L^2 regularization. A full treatment is beyond the scope of this chapter, but we give a complete statement of the theorem in Section 8.A. An implication of the representer theorem is that L^2 regularization is crucial in order to obtain kernel ridge regression (8.13), and we could not have achieved it using, say, L^1 regularization instead. The representer theorem is a cornerstone for most kernel methods, since it tells us that *we can express some models in terms of dual parameters α (of finite length n) and a kernel $\kappa(\mathbf{x}, \mathbf{x}')$ instead of the primal parameters θ (possibly of infinite length d) and a nonlinear feature transformation $\phi(\mathbf{x})$* , just like we did with linear regression in (8.13).

Support vector regression

We will now look at support vector regression, another off-the-shelf kernel method for regression. From a model perspective the only difference to kernel ridge regression is a change of loss function. This new loss function has an interesting effect in that the dual parameter vector $\hat{\alpha}$ in support vector regression becomes *sparse*, meaning that several elements of $\hat{\alpha}$ are exactly zero. The training data points corresponding to the non-zero elements of $\hat{\alpha}$ are referred to as *support vectors*, since $\hat{y}(\mathbf{x}_*)$ will depend only on those (in contrast to kernel ridge regression (8.11), where all training data points are needed to compute $\hat{y}(\mathbf{x}_*)$). This makes support vector regression an example of a so-called *support vector machine* (SVM), a family of methods with sparse dual parameter vectors.

The loss function we will use for support vector regression is the ϵ -insensitive loss

$$L(\hat{y}(\mathbf{x}_*), y_*) = \begin{cases} 0 & \text{if } |\hat{y}(\mathbf{x}_*) - y_*| < \epsilon, \\ |\hat{y}(\mathbf{x}_*) - y_*| - \epsilon & \text{otherwise,} \end{cases} \quad (8.16)$$

which was introduced in Section 5.1. The parameter ϵ is a user design choice. In its primal formulation, support vector regression also makes use of the linear regression model

$$\hat{y}(\mathbf{x}_*) = \theta^\top \phi(\mathbf{x}_*), \quad (8.17a)$$

but instead of the least square cost function in (8.4), we have

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, \underbrace{|\theta^\top \phi(\mathbf{x}_i) - y_i| - \epsilon}_{\hat{y}(\mathbf{x}_i)} \right\} + \lambda \|\theta\|_2^2. \quad (8.17b)$$

As with kernel ridge regression, we reformulate the primal formulation (8.17) into a dual formulation with α instead of θ and use the kernel-trick. For the dual formulation, we cannot repeat the convenient closed-form derivation along the lines of (8.4)-(8.13) since there is no closed-form solution for $\hat{\theta}$. Instead we have to use optimization theory and introduce slack variables and construct the Lagrangian of (8.17b)). We do not give the full derivation here, but as it turns out the dual formulation becomes

$$\hat{y}(\mathbf{x}_*) = \hat{\alpha}^\top K(\mathbf{X}, \mathbf{x}_*), \quad (8.18a)$$

where $\hat{\alpha}$ is the solution to the optimization problem

$$\hat{\alpha} = \arg \min_{\alpha} \frac{1}{2} \alpha^T K(\mathbf{X}, \mathbf{X}) \alpha - \alpha^T \mathbf{y} + \epsilon \|\alpha\|_1 \quad (8.18b)$$

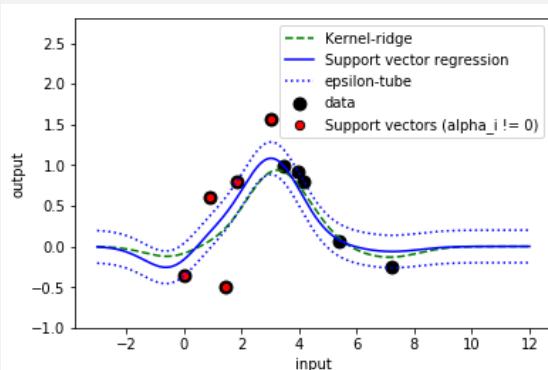
$$\text{subject to } |\alpha_i| \leq \frac{2}{n\lambda} \quad (8.18c)$$

Note that (8.18a) is identical to the corresponding expression for kernel ridge regression (8.13b). This is a consequence of the representer theorem. The only difference to kernel ridge regression is how the dual parameters α are learned; by numerically solving the optimization problem (8.18b) instead of using the closed-form solution in (8.13c).

The ϵ -insensitive loss function could be used for any regression model but it is particularly interesting in combination in this kernel context since the dual parameter vector α becomes sparse (meaning that only some elements are non-zero). Remember that α has one entry per training data point. This implies that the prediction (8.18a) depends only on *some* of the training data points, namely those whose corresponding α_i is non-zero. All training data is indeed used at training time (solving (8.18b)), but when making predictions (using (8.18a)) only the data points, the support vectors, with non-zero α_i contributes, which significantly reduces the computational burden. The larger ϵ chosen by the user, the fewer support vectors and the fewer computations needed in making predictions. It can also be argued that ϵ has a regularizing effect, in the sense that the more/fewer support vectors that are used the more/less complicated model. We illustrate support vector regression with Example 8.3, and summarize it in Algorithm 16.

Example 8.3: support vector regression and kernel ridge regression

This example is not yet complete.



Algorithm 16: Support vector regression

1 To be completed

The ϵ -insensitive loss indeed makes the dual parameter vector α sparse. Note however that it does *not* mean that the corresponding primal parameter vector θ is sparse (their relationship is given by (8.15)). Also note that (8.18b) is a *constrained* optimization problem (there is a constraint given by (8.18c)), and more theory than what we presented in Section 5.3 is needed to derive a good solver.

Remark 8.1 *The nonlinear feature vector $\Phi(\mathbf{x})$ corresponding to some kernels, such as the squared exponential kernel (8.12), does not have a constant/offset term. Therefore an additional θ_0 is sometimes included in (8.18a) for support vector regression which adds the constraint $\sum_i \alpha_i = 0$ to the optimization problem in (8.18b). The same addition could be made also to kernel ridge regression (8.13b), but that would break the closed-form calculation of α (8.13b).*

Conclusions on using kernels for regression

With kernel ridge regression and support vector regression we have actually been dealing with the interplay between three different concepts, each of them interesting on its own. To clarify this, we repeat them in an ordered list:

1. We have considered the **primal and dual formulation** of a model. The primal formulation expresses the model in terms of θ (fixed size d), whereas the dual formulation uses α (one α_i per data sample i , hence the size is n no matter what d happens to be). Both formulations are mathematically equivalent, but more or less useful in practice depending on whether $d > n$ or $n > d$.
2. We have introduced **kernels** κ , which allows us to let $d \rightarrow \infty$. This idea is practically useful only when using the dual formulation with α , since d is the dimension of θ .
3. We have used **different loss functions**. Kernel ridge regression makes use of squared error loss, whereas support vector regression uses the ϵ -insensitive loss. The ϵ -insensitive loss is particularly interesting in the dual formulation, since it gives *sparse* α . (We will later also use the hinge loss for support vector classification in Section 8.5, which has a similar effect.)

We will now spend some additional effort on understanding the kernel concept in Section 8.4, and thereafter introduce support vector classification (SVC) in Section 8.5.

8.4 Kernel theory

We have defined a kernel as being any function taking two arguments from the same space, and returning a scalar. We have also suggested that we should (at least sometimes) restrict ourselves to positive semidefinite kernels, and presented two practically useful algorithms—kernel ridge regression and support vector regression. Before we introduce also support vector classification, we will discuss what a kernel means in practice, and also give a flavor of the available theory behind it. To make the discussion more concrete, let us start by introducing the kernel version of k -NN.

Introducing kernel k -NN

As you know from Chapter 2, k -NN constructs the prediction for \mathbf{x}_* by taking the average or a majority vote among the k nearest neighbors to \mathbf{x}_* . In its standard formulation, “nearest” is defined by the Euclidian distance. Since the Euclidian distance always is positive, we can equivalently consider the squared Euclidian distance instead, which can be written in terms of the linear kernel $\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$ as

$$\|\mathbf{x} - \mathbf{x}'\|_2^2 = (\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}') = \mathbf{x}^\top \mathbf{x} + \mathbf{x}'^\top \mathbf{x}' - 2\mathbf{x}^\top \mathbf{x}' = \kappa(\mathbf{x}, \mathbf{x}) + \kappa(\mathbf{x}', \mathbf{x}') - 2\kappa(\mathbf{x}, \mathbf{x}'). \quad (8.19)$$

To generalize the k -NN algorithm to use kernels, we allow the linear kernel to be replaced by any, say, positive semidefinite kernel $\kappa(\mathbf{x}, \mathbf{x}')$ in (8.19). *Kernel k*-NN thus works as standard k -NN, but determines proximity between the data points using the right hand side of (8.19) (with a user-chosen kernel $\kappa(\mathbf{x}, \mathbf{x}')$) instead of the left hand side of (8.19).

For many (but not all) kernels, it holds that $\kappa(\mathbf{x}, \mathbf{x}) = \kappa(\mathbf{x}', \mathbf{x}')$ for all \mathbf{x} and \mathbf{x}' , suggesting that the most interesting part of the right hand side of (8.19) is the term $-2\kappa(\mathbf{x}, \mathbf{x}')$. Thus, if $\kappa(\mathbf{x}, \mathbf{x}')$ takes a large value, the two data points \mathbf{x} and \mathbf{x}' are considered to be close, and vice versa if $\kappa(\mathbf{x}, \mathbf{x}')$ takes a small value. That is, *the kernel determines how close any two data points are*.

Furthermore, kernel k -NN allows us to use k -NN also for data where the Euclidian distance has no natural meaning. As long as we have a kernel which acts on the input space, we can apply kernel k -NN even if the Euclidian distance is not defined for that input type. We can thereby apply kernel k -NN to input data that is neither numerical nor categorical, such as text snippets, as illustrated by Example 8.4.

Example 8.4: Kernel k -NN for interpreting words

This example illustrates how kernel k -NN can be applied to text data, where the Euclidian distance has no meaning and standard k -NN therefore can not be applied. In this example, the input is single words (or more technically: character strings) and we use the so-called Levenshtein distance to construct a kernel. The Levenshtein distance is the number of single-character edits needed to transform one word (string) into another. It takes two strings and returns a non-negative integer, which is zero only if the two strings are equivalent. We can thereby use it to construct a kernel for example as $\kappa(x, x') = \frac{1}{1 + \text{LD}(x, x')}$ (where LD is the Levenshtein distance). This defines a positive semidefinite kernel since it will always yield real-valued symmetric and diagonal-dominant (and thereby positive semidefinite) matrices $\mathbf{K}(\mathbf{X}, \mathbf{X})$.

In this very small example, we consider a training data set of 10 adjectives shown below (x_i), each labeled (y_i) Positive or Negative, according to their meaning. We will now use kernel k -NN (with the kernel defined above) to predict whether the word “horrendous” (x_\star) is a positive or negative word. In the third column below we have therefore computed the Levenshtein distance (LD) between each known word (x_i) and “horrendous” (x_\star). The rightmost column shows the value of the right hand side of (8.19), which is the value that kernel k -NN uses to determine how close two data points are.

Word, x_i	Meaning, y_i	Levenshtein dist. from x_i to x_\star	$\kappa(x_i, x_i) + \kappa(x_\star, x_\star) - 2\kappa(x_i, x_\star)$
Awesome	Positive	8	1.79
Excellent	Positive	10	1.81
Spotless	Positive	9	1.80
Terrific	Positive	8	1.79
Tremendous	Positive	4	1.60
Awful	Negative	9	1.80
Dreadful	Negative	6	1.71
Horrific	Negative	6	1.71
Outrageous	Negative	6	1.71
Terrible	Negative	8	1.79

According to the rightmost column, the closest word to horrendous is the positive word tremendous. Thus, if we use $k = 1$ the conclusion would be that horrendous is a positive word. However, the second, third and fourth closest words are all negative (dreadful, horrific, outrageous), and with $k = 3$ or $k = 4$ the conclusion thereby becomes that horrendous is a negative word (which happens to be the correct conclusion, for this particular example).

The purpose of this example is to illustrate how a kernel allows a basic method such as k -NN to be used for a problem where the input has a more intricate structure than just being numerical. For the particular application of predicting word semantics, more elaborate machine learning methods exist.

The meaning of a kernel

From kernel k -NN we got (at least) two lessons about kernels that are generally applicable to all supervised machine learning methods that use kernels:

- The kernel defines how close/similar any two data points are. If, say, $\kappa(\mathbf{x}_i, \mathbf{x}_\star) > \kappa(\mathbf{x}_j, \mathbf{x}_\star)$, then \mathbf{x}_\star is considered to be more similar to \mathbf{x}_i than \mathbf{x}_j . Intuitively speaking, for most methods the prediction $\hat{y}(\mathbf{x}_\star)$ is most influenced by the training data points that are closest/most similar to \mathbf{x}_\star . The kernel thereby plays an important role in determining the individual influence of each training data point when making a prediction.
- Even though we started by introducing kernels from the inner product $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$, no inner product has to be defined for the space in which \mathbf{x} itself lives! As we saw in Example 8.4, we can apply a positive semidefinite kernel method also to text string without worrying about inner products, as long as we have a kernel for that type of data.

In addition to this, the kernel also plays a somewhat more subtle role in methods that builds on the representer theorem (such as kernel ridge regression, support vector regression and support vector classification, but not kernel k -NN). Remember that the primal formulation of those methods, by virtue

of the representer theorem, contains the L^2 regularization term $\lambda\|\theta\|_2^2$. Even though we do not solve the primal formulation explicitly when using kernels (we solve the dual instead), it is nevertheless an equivalent representation, and we may ask what impact the regularization $\lambda\|\theta\|_2^2$ has on the solution?

The L^2 regularization means that primal parameter values θ close to zero are favored. Besides the regularization term, θ only appear in the expression $\theta^\top \phi(\mathbf{x})$. The solution $\hat{\theta}$ to the primal problem is therefore an interplay between the feature vector $\phi(\mathbf{x})$ and the L^2 regularization term. Consider two different choices of feature vectors $\phi_1(\mathbf{x})$ and $\phi_2(\mathbf{x})$. If they both span the same space of functions, there exists θ_1 and θ_2 such that $\theta_1^\top \phi_1(\mathbf{x}) = \theta_2^\top \phi_2(\mathbf{x})$ for all \mathbf{x} , and it might appear irrelevant which feature vector is being used. However, the L^2 regularization complicates the situation because it acts directly on θ , and it therefore matters whether $\phi_1(\mathbf{x})$ or $\phi_2(\mathbf{x})$ is being used. In the dual formulation we choose the kernel $\kappa(\mathbf{x}, \mathbf{x}')$ instead of feature vector $\phi(\mathbf{x})$, but since that choice implicitly corresponds to a feature vector the effect is still present, and we may add one more bullet point about the meaning of a kernel:

- The choice of kernel corresponds to a choice of a regularization functional. That is, the kernel implies a preference for certain functions in the space of all functions that are spanned by the feature vector. For example, the squared exponential kernel enforce smooth solutions.

Using a kernel makes a method quite flexible, and one could perhaps expect it to suffer heavily from overfit. However, the regularizing role of the kernel explains why that rarely is the case in practice.

All three bullet points above are central to understand the usefulness and versatility of kernel methods. They also highlight the importance for the machine learning engineer to choose the kernel wisely, and not only resort to “default” choices.

Valid choices of kernels

We introduced kernels as a way to compactly work with nonlinear feature transformations as (8.10). A direct consequence of this is that it is now sufficient $\kappa(\mathbf{x}, \mathbf{x}')$, and not $\phi(\mathbf{x})$. A natural question to ask is now if an arbitrary kernel $\kappa(\mathbf{x}, \mathbf{x}')$ always corresponds to a feature transformation $\phi(\mathbf{x})$, such that it can be written as the inner product

$$\kappa(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')? \quad (8.20)$$

Before answering the question, we have to be aware that this question is merely of theoretical nature. As long as we can use $\kappa(\mathbf{x}, \mathbf{x}')$ when computing predictions it serves its purpose, no matter if it admits the factorization (8.20) or not. The specific requirements on $\kappa(\mathbf{x}, \mathbf{x}')$ are different for different methods, for example the inverse $(\mathbf{K}(\mathbf{X}, \mathbf{X}) + n\lambda\mathbf{I})^{-1}$ is needed for kernel ridge regression but not for support vector regression. Furthermore, whether a kernel admits the factorization (8.20) or not has no direct correspondence to how well it performs in terms of, say, prediction accuracy. For any practical machine learning problem the performance still has to be evaluated using cross validation or similarly.

That being said, we will now have a closer look at the important family of positive semidefinite kernels. A kernel is said to be positive semidefinite if the matrix $\mathbf{K}(\mathbf{X}, \mathbf{X})$ (8.11b) is positive semidefinite (has no negative eigenvalues) for any choice of \mathbf{X} .

First, it holds that any kernel $\kappa(\mathbf{x}, \mathbf{x}')$ that is defined as an inner product between feature vectors $\phi(\mathbf{x})$, as in (8.20), always is positive semidefinite. It can, for example, be shown by using Cauchy-Schwartz inequality which implies that $\mathbf{K}(\mathbf{X}, \mathbf{X})$ is diagonally dominant. Since $\mathbf{K}(\mathbf{X}, \mathbf{X})$ also is real-valued, symmetric and has no negative diagonal entries, it follows that $\mathbf{K}(\mathbf{X}, \mathbf{X})$ is always positive semidefinite.

Less trivially the other direction also holds, that is, for any positive semidefinite kernel $\kappa(\mathbf{x}, \mathbf{x}')$ there always exists a feature vector $\phi(\mathbf{x})$ such that $\kappa(\mathbf{x}, \mathbf{x}')$ can be written as an inner product (8.20). Technically it can be shown that for any positive semidefinite kernel $\kappa(\mathbf{x}, \mathbf{x}')$ it is possible to construct a function space, more specifically a Hilbert space, that is spanned by a feature vector $\phi(\mathbf{x})$ for which (8.20) holds. The dimensionality of the Hilbert space, and thereby also the dimension of $\phi(\mathbf{x})$, can however be infinite.

Given a kernel $\kappa(\mathbf{x}, \mathbf{x}')$ there are multiple ways to construct a Hilbert space spanned by $\phi(\mathbf{x})$, and we will only mention some directions here. One alternative is to consider the so-called reproducing kernel map.

The reproducing kernel map is obtained by consider one argument, say the latter, to $\kappa(\mathbf{x}, \mathbf{x}')$ fix and let $\kappa(\cdot, \mathbf{x}')$ span the Hilbert space and define the inner product $\langle \cdot, \cdot \rangle$ such that $\langle \kappa(\cdot, \mathbf{x}), \kappa(\cdot, \mathbf{x}') \rangle = \kappa(\mathbf{x}, \mathbf{x}')$. This inner product has a so-called reproducing property, and it is the main building block for the so-called reproducing kernel Hilbert space. Another alternative is to use the so-called Mercer kernel map, which constructs the Hilbert space using eigenfunctions to an integral operator which is related to the kernel.

A given Hilbert space uniquely defines a kernel, but for a given kernel there exists multiple Hilbert spaces which corresponds to it. In practice this means that given a kernel $\kappa(\mathbf{x}, \mathbf{x}')$ the corresponding feature vector $\phi(\mathbf{x})$ is not unique, in fact not even its dimensionality. As a simple example consider the linear kernel $\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$, which can either be expressed as an inner product between $\phi(\mathbf{x}) = \mathbf{x}$ (one dimensional $\phi(\mathbf{x})$) or as an inner product between $\phi(\mathbf{x}) = \begin{bmatrix} \frac{1}{\sqrt{2}}\mathbf{x} & \frac{1}{\sqrt{2}}\mathbf{x} \end{bmatrix}^T$ (two dimensional $\phi(\mathbf{x})$).

Examples of kernels

We will now give a list of some of the commonly used kernels, of which we already have introduced some. These examples are only for the case when \mathbf{x} is a numeric variable. For other types of input variables (such as Example 8.4), we have to resort to the more application specific literature. We start with some positive semidefinite kernels, where the *linear kernel* might be the simplest one

$$\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'. \quad (8.21)$$

A generalization thereof, still positive semidefinite, is the *polynomial kernel*

$$\kappa(\mathbf{x}, \mathbf{x}') = (c + \mathbf{x}^T \mathbf{x}')^{d-1} \quad (8.22)$$

with hyperparameter $c \geq 0$ and polynomial order $d - 1$ (integer). The polynomial kernel corresponds to a finite-dimensional feature vector $\phi(\mathbf{x})$ of monomials up to order $d - 1$. The polynomial kernel does therefore not conceptually enable anything that could not be achieved by instead implementing the primal formulation and the finite-dimensional $\phi(\mathbf{x})$ explicitly. The other positive semidefinite kernels below, on the other hand, all corresponds to infinite-dimensional feature vectors $\phi(\mathbf{x})$.

We have previously also mentioned the *squared exponential kernel*⁴,

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\ell^2}\right), \quad (8.23)$$

with hyperparameter $\ell > 0$ (usually called lengthscale). As we saw in Example 8.2 and 8.3, this kernel has more of a “local” nature compared to the polynomial since $\kappa(\mathbf{x}, \mathbf{x}') \rightarrow 0$ as $\|\mathbf{x} - \mathbf{x}'\| \rightarrow \infty$. This property makes sense in many problems, and is perhaps the reason why this might be the most commonly used kernel.

Somewhat related to the squared exponential we have the family of *Matérn kernels*

$$\kappa(\mathbf{x}, \mathbf{x}') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{\ell} \right)^\nu \zeta_\nu \left(\frac{\sqrt{2\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{\ell} \right), \quad (8.24)$$

with hyperparameters $\ell > 0$ and $\nu > 0$. Here Γ is the Gamma function and ζ_ν is a modified Bessel function. All Matérn kernels are positive semidefinite. Of particular interest are the cases $\nu = \frac{1}{2}, \frac{3}{2}$ and $\frac{5}{2}$, when (8.24) simplifies to

$$\nu = \frac{1}{2} \Rightarrow \kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2}{\ell}\right), \quad (8.25)$$

$$\nu = \frac{3}{2} \Rightarrow \kappa(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{\sqrt{3} \|\mathbf{x} - \mathbf{x}'\|_2}{\ell}\right) \exp\left(-\frac{\sqrt{3} \|\mathbf{x} - \mathbf{x}'\|_2}{\ell}\right), \quad (8.26)$$

$$\nu = \frac{5}{2} \Rightarrow \kappa(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{\sqrt{5} \|\mathbf{x} - \mathbf{x}'\|_2}{\ell} + \frac{5 \|\mathbf{x} - \mathbf{x}'\|_2^2}{3\ell^2}\right) \exp\left(-\frac{\sqrt{5} \|\mathbf{x} - \mathbf{x}'\|_2}{\ell}\right). \quad (8.27)$$

⁴Also known as the Gaussian radial basis function kernel.

The Matérn kernel with $\nu = \frac{1}{2}$ is also called the *exponential kernel*. It can furthermore be shown that the Matérn kernel (8.24) equals the squared exponential (8.23) when $\nu \rightarrow \infty$.

As a final example of a positive semidefinite kernel, we mention the *rational quadratic kernel*

$$\kappa(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2a\ell^2} \right)^{-a} \quad (8.28)$$

with hyperparameters $\ell > 0$ and $a > 0$. The squared exponential, Matérn and rational quadratic kernel are all examples of *isotropic* kernels, since they are a function only of $\|\mathbf{x} - \mathbf{x}'\|_2$.

Going back to the discussion in connection to (8.20), positive semidefinite kernels are a subset of all kernels, for which we know that certain theoretical properties holds. In practice, however, a kernel is potentially useful as long as we can compute a prediction using it, no matter its theoretical properties. One (at least historically) popular kernel for SVMs which is not positive semidefinite is the *Sigmoid kernel*

$$\kappa(\mathbf{x}, \mathbf{x}') = \tanh(a\mathbf{x}^\top \mathbf{x}' + b), \quad (8.29)$$

where $a > 0$ and $b < 0$ are hyperparameters. The fact that it is not positive semidefinite can, for example, be seen by computing the eigenvalues of $\mathbf{K}(\mathbf{X}, \mathbf{X})$ with $a = 1$, $b = -1$ and $\mathbf{X} = [1 \ 2]^\top$. Since this kernel is not positive semidefinite the inverse $(\mathbf{K}(\mathbf{X}, \mathbf{X}) + n\lambda\mathbf{I})^{-1}$ does not always exists, and it is therefore not suitable for kernel ridge regression. It can, however, be used in support vector regression and classification, where that inverse is not needed. For certain values of b it can be shown to be a so-called conditional positive semidefinite kernel (a weaker property than positive semidefinite).

It is possible to construct “new” kernels by modifying or combining existing ones. There are in particular a set of operations that preserves the positive semidefinite property: If $\kappa(\mathbf{x}, \mathbf{x}')$ is a positive semidefinite kernel, then so is $a\kappa(\mathbf{x}, \mathbf{x}')$ if $a > 0$ (scaling). Furthermore if both $\kappa_1(\mathbf{x}, \mathbf{x}')$ and $\kappa_2(\mathbf{x}, \mathbf{x}')$ are positive semidefinite kernels, then so are $\kappa_1(\mathbf{x}, \mathbf{x}') + \kappa_2(\mathbf{x}, \mathbf{x}')$ (addition) as well as $\kappa_1(\mathbf{x}, \mathbf{x}')\kappa_2(\mathbf{x}, \mathbf{x}')$ (multiplication).

Most kernels contains a few hyperparameters that are left to the user to chose. Much like cross validation can give valuable help in choosing between different kernels, cross validation can also help choosing hyperparameters (and regularization parameter λ) with grid search.

8.5 Support vector classification

We have so far spent most time deriving two kernel versions of linear regression: kernel ridge regression and support vector regression. We will now focus on classification. Unfortunately the derivations are now more technically intricate than for regression, and we have therefore placed the details in the chapter appendix. However, the intuition carries over from regression, as well as the main ideas of the dual formulation, the kernel trick and the change of loss function.

It is possible to derive a kernel version of logistic regression with L^2 regularization. The derivation can be made by first replacing \mathbf{x} with $\phi(\mathbf{x})$, and then use the representer theorem to derive its dual formulation and apply the kernel trick. Kernel logistic regression is a classification counterpart to kernel ridge regression. Due to the fact that kernel logistic regression os very scarcely used in practice, we will instead go straight to support vector classification. As the name suggests, support vector classification is the classification counterpart of support vector regression. Both support vector regression and classification are referred to as support vector machines (SVM) since they both have sparse dual parameter vectors.

We consider the binary classification problem $y \in \{-1, 1\}$, and start with the margin formulation (see (5.8) in Chapter 5) of the logistic regression classifier

$$\hat{y}(\mathbf{x}_*) = \text{sign}\{\boldsymbol{\theta}^\top \phi(\mathbf{x})\}. \quad (8.30)$$

If we now were to learn $\boldsymbol{\theta}$ using the logistic loss (5.10), we would obtain logistic regression with a nonlinear feature transformation $\phi(\mathbf{x})$, from which kernel logistic regression eventually would follow. Inspired by

support vector regression, we will instead make use of the hinge loss function (5.13),

$$L(\mathbf{x}, y, \boldsymbol{\theta}) = \begin{cases} 1 - y\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}) & \text{if } y\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}) < 1 \\ 0, & \text{otherwise.} \end{cases} \quad (8.31)$$

From Figure 5.2, it is not immediately clear what advantages the hinge loss has over the logistic loss. Analogously to the ϵ -insensitive loss, the main advantage of the hinge loss comes when we look at the dual formulation using α instead of the primal formulation with $\boldsymbol{\theta}$. But before introducing a dual formulation we first have to spell out the primal one. Since the representer theorem is lurking behind all this, we have to use L^2 regularization, meaning which together with (8.31) gives the primal formulation

$$\widehat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - y_i \boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}_i)\} + \lambda \|\boldsymbol{\theta}\|_2^2. \quad (8.32)$$

The primal formulation does not allow for the kernel trick, since the feature vector does not show up as $\boldsymbol{\phi}(\mathbf{x})^\top \boldsymbol{\phi}(\mathbf{x})$. By using optimization theory and constructing the Lagrangian (the details are found in Appendix 8.B), we can arrive at the dual formulation⁵ of (8.32),

$$\widehat{\boldsymbol{\alpha}} = \arg \min_{\boldsymbol{\alpha}} \frac{1}{2} \boldsymbol{\alpha}^\top \mathbf{K}(\mathbf{X}, \mathbf{X}) \boldsymbol{\alpha} - \boldsymbol{\alpha}^\top \mathbf{y} \quad (8.33a)$$

$$\text{subject to } |\alpha_i| \leq \frac{2}{n\lambda} \text{ and } 0 \leq \alpha_i y_i. \quad (8.33b)$$

with

$$\widehat{y}(\mathbf{x}_*) = \text{sign}(\widehat{\boldsymbol{\alpha}} \mathbf{K}(\mathbf{X}, \mathbf{x}_*)) \quad (8.33c)$$

instead of (8.30). Note that we here also have made use of the kernel trick by replacing $\boldsymbol{\phi}(\mathbf{x})\boldsymbol{\phi}(\mathbf{x}')$ with $\kappa(\mathbf{x}, \mathbf{x}')$, which gives $\mathbf{K}(\mathbf{X}, \mathbf{X})$ in (8.33a) and $\mathbf{K}(\mathbf{X}, \mathbf{x}_*)$ in (8.33c).

Because of the representer theorem, the formulation (8.33c) should come as no surprise to us, since it simply corresponds to inserting (8.15) into (8.30). However, the representer theorem only tells us that this dual formulation exists, but (8.33a)-(8.33b) does not follow automatically from the representer theorem but requires its own derivation.

The perhaps most interesting property of the constrained optimization problem (8.33) is that its solution $\widehat{\boldsymbol{\alpha}}$ turns out to be sparse. This is exactly the same phenomena as with support vector regression, and explains why (8.33) is an SVM. More specifically, (8.33) is called support vector classification. The strength of this method, like support vector regression, is that the model has the full flexibility of being a kernel method, and yet the prediction (8.33c) only explicitly depends on some of the training data points (the support vectors). It is however important to realize that all training data points are needed when solving (8.33a).

The support vector property is due to the fact that the loss function is exactly equal to zero when the margin (see Section 5.1) is > 1 . In the dual formulation, the parameter α_i becomes nonzero only if the margin for datapoint i is < 1 , which makes $\widehat{\boldsymbol{\alpha}}$ sparse.

As a consequence of using the hinge loss, as we discussed in Section 5.1, support vector classification does not provide probability estimates $g(\mathbf{x})$ but only a “hard” classification $\widehat{y}(\mathbf{x}_*)$. The margin, in this case $y\widehat{\boldsymbol{\alpha}} \mathbf{K}(\mathbf{X}, \mathbf{x}_*)$, is not possible to interpret as a class probability estimate, because of the asymptotic minimizer of the hinge loss. As an alternative it is possible to instead use the squared hinge loss or the Huberized squared hinge loss which allows for a probability interpretation of the margin. Since all these loss functions are exactly zero for margins > 1 , they retain the support vector property with a sparse $\widehat{\boldsymbol{\alpha}}$. However, by using a different loss function than the hinge loss we will obtain a different optimization problem since the derivation would no longer start in (8.32).

⁵In other texts, it is common to let the dual variables instead be $\alpha' = (\alpha_i y_i)$ (which happens to be the Lagrange multipliers). It is mathematically equivalent, but this formulation better highlights the similarities to the other methods and the importance of the representer theorem.

It is not *necessary* to make use of kernels in support vector classification. It is perfectly possible to use the linear kernel $\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$ or any other kernel corresponding to a finite dimensional $\phi(\mathbf{x})$ in (8.33). That would indeed limit the flexibility of the classifier, using the linear kernel would limit the classifier to linear decision boundaries, for instance. The possible benefit, however, of not making full use of kernels is that it suffices to implement and solve the primal (and not the dual) formulation (8.32), since θ is of finite dimension. The support vector property would still be present, but much less visible since α not are explicitly computed. If only using the primal formulation and no kernels, the representer theorem is not needed, and one could therefore also replace the L^2 regularization with any other regularization method.

Algorithm 17: Support vector classification

1 To be completed

Example 8.5:

ex:svc This example is not yet complete.

The support vector classifier is most often formulated as a solution to the binary classification problem. The generalization to the multiclass problem is unfortunately not straightforward, since it requires a multiclass generalization of the loss function, as discussed in Section 5.1. In practice it is common to construct a multiclass classifier from multiple binary ones using either the one-versus-rest or the one-versus-one strategy (see page 69).

8.A The representer theorem

8.B Derivation of support vector classification

9 The Bayesian approach and Gaussian processes

So far, learning a parametric models has in all previous chapters amounted to somehow finding a parameter value $\hat{\theta}$. With the *Bayesian approach* (also called the *probabilistic approach*) learning amounts to instead finding the *distribution* of parameter values $p(\theta | \mathbf{y})$ conditioned on the observed training data $\mathcal{T} = \{\mathbf{X}, \mathbf{y}\}$ (to ease the notation we do not write out \mathbf{X} in the conditioning). Also the prediction is a distribution $p(y_\star | \mathbf{y})$, instead of a single value. Before we become more concrete and use equations, let us just say that on a theoretical (or even philosophical) level the Bayesian approach is rather different to what we previously have seen in this book. As it however opens up for a family of new, versatile and practically useful methods, the extra effort required to understand this somewhat different approach pays off well and provides another interesting perspective on supervised machine learning. As the Bayesian approach makes repeated use of probability distribution, it is also natural that this chapter will be more heavy on the probability theory side compared to the rest of the book.

We will start this chapter by first giving a general introduction to the Bayesian idea. We thereafter go back to basics and applies the Bayesian approach to linear regression, which we thereafter extend to the non-parametric Gaussian process model.

9.1 The Bayesian idea

In this chapter we will use the term *model* with a very specific meaning. A model, in this chapter, refers to a set of equations that contain parameters θ and defines a probability distribution over output \mathbf{y} given some input \mathbf{X} and parameters θ , that is, $p(\mathbf{y} | \theta)$ (we omit \mathbf{X} from the probability notation since it is not a random variable). Remember that $p(\mathbf{y} | \theta)$ was used also for the maximum likelihood perspective; one example of a model $p(\mathbf{y} | \theta)$ is linear regression (3.17)-(3.18).

In the Bayesian approach the parameters θ of any model are consistently treated as being random variables. This implies that θ is always subject to a probability distribution. Consequently, learning (in the Bayesian setting) amounts to computing the distribution of θ *conditional on* training data $\mathcal{T} = \{\mathbf{X}, \mathbf{y}\}$. To ease the notation we consistently omit \mathbf{X} in the conditioning (mathematically we motivate this by the fact that we only consider y , and not x , to be a random variable), and we therefore denote the distribution of θ conditioned on training data as $p(\theta | \mathbf{y})$. The computation of $p(\theta | \mathbf{y})$ is done using the laws of probabilities, and notably Bayes' theorem,

$$p(\theta | \mathbf{y}) = \frac{p(\mathbf{y} | \theta)p(\theta)}{p(\mathbf{y})}, \quad (9.1)$$

which is the reason why it is called the Bayesian approach. The left hand side of (9.1) is the sought after distribution $p(\theta | \mathbf{y})$. The right hand side of (9.1) contains some important elements: $p(\mathbf{y} | \theta)$ is the model, and $p(\theta)$ is the distribution of θ before any observations are made (that is, not conditional on training data). By nature $p(\theta)$ cannot be computed, but has to be postulated by the user (much like the model is chosen by the user as well). Finally $p(\mathbf{y})$ can, by the laws of probabilities, be rewritten as

$$p(\mathbf{y}) = \int p(\mathbf{y}, \theta) d\theta = \int p(\mathbf{y} | \theta)p(\theta) d\theta \quad (9.2)$$

which is an integral that can be computed. In other words, learning a parametric model (in the Bayesian fashion) amounts to conditioning $p(\theta)$ on \mathbf{y} , that is, compute $p(\theta | \mathbf{y})$. After learning, the next step with a parametric model is to compute predictions. Again this is a matter of computing a distribution $p(y_\star | \mathbf{y})$

(rather than a point prediction \hat{y}) for a test input \mathbf{x}_\star , which connects to $p(\theta | \mathbf{y})$ via the integral

$$p(y_\star | \mathbf{y}) = \int p(y_\star | \cdot, \theta) p(\theta | \mathbf{y}) d\theta. \quad (9.3)$$

Here $p(y_\star | \theta)$ encodes the model for a test data input \mathbf{x}_\star (again omitted in the notation).

The model is denoted by $p(\mathbf{y} | \theta)$ (with the training data) as well as $p(y_\star | \cdot, \theta)$ (with the test data). Sometimes $p(\mathbf{y} | \theta)$ is referred to, somewhat sloppily, as the ‘likelihood’. The other elements involved in the Bayesian approach are traditionally given the names

- $p(\theta)$ prior,
- $p(\theta | \mathbf{y})$ posterior,
- $p(y_\star | \mathbf{y})$ posterior predictive,

These names are useful for communication, but it is important to remember that they are nothing but probability distributions. In addition $p(\mathbf{y})$ is often called the *marginal likelihood* or *evidence*.

A representation of beliefs

The main feature of the Bayesian approach is its repeated use of probability distributions. It is possible to interpret those distributions as representing *beliefs* in the following sense: The prior $p(\theta)$ represents what beliefs there are about θ *a priori*, that is, before any data has been observed. The model, encoded as $p(\mathbf{y} | \theta)$, defines how data \mathbf{y} relates to the parameter θ . Using Bayes’ theorem (9.1) we *update* the belief about θ to the posterior $p(\theta | \mathbf{y})$ which takes also the observed data \mathbf{y} into account. In an everyday language, these distributions could be said to represent uncertainty about θ before and after observing the data \mathbf{y} respectively.

An interesting and practically relevant consequence is that the Bayesian approach is less prone to overfit, compared to using a maximum-likelihood based method. This is due to the fact that the prior (which has a similar effect as regularization) is taken into account, as well as the fact that an entire distribution $p(\theta | \mathbf{y})$ is obtained instead of a single value $\hat{\theta}$. In particular in the regime of small data sets \mathbf{y} , the “uncertainty” seen in the posterior $p(\theta | \mathbf{y})$ reflects how much (or little) that can be said about θ , under the assumed conditions.

The posterior $p(\theta | \mathbf{y})$ is therefore a combination of the prior belief $p(\theta)$ and the information about θ carried by \mathbf{y} through the model. Without a meaningful prior $p(\theta)$, the posterior $p(\theta | \mathbf{y})$ is no meaningful either. In some applications it can be hard to make a choice of $p(\theta)$ that is not influenced by the personal experience of the machine learning engineer, which sometimes is emphasized by saying that $p(\theta)$, and thereby also $p(\theta | \mathbf{y})$, represents a *subjective* belief. This notion is meant to reflect that it is not independent of the human engineer that designed the solution. However, also the model (no matter if the Bayesian approach is used or not) is often chosen based on the personal experience of the machine learning engineer, meaning that most machine learning results are in that sense subjective.

An interesting situation for the Bayesian approach is when more data arrives after each other. Say that we have two sets of data, \mathbf{y}_1 and \mathbf{y}_2 . Starting with a prior $p(\theta)$ we can condition on \mathbf{y}_1 by computing $p(\theta | \mathbf{y}_1)$ using Bayes’ theorem (9.1). However, if we thereafter want to condition on all data, \mathbf{y}_1 and \mathbf{y}_2 , that is compute $p(\theta | \mathbf{y}_1, \mathbf{y}_2)$, we do not have to start over again. We can instead replace the prior $p(\theta)$ with $p(\theta | \mathbf{y}_1)$ in Bayes’ theorem to compute $p(\theta | \mathbf{y}_1, \mathbf{y}_2)$. In a sense, the “old posterior” can become the “new prior” when data arrives sequentially.

The marginal likelihood as a model selection tool

When using the Bayesian approach there are often some hyperparameters in the model or the prior, say η , that needs to be chosen. It is indeed possible to assume a prior $p(\eta)$ and compute the posterior also for the hyperparameters, $p(\eta | \mathbf{y})$. That would be the fully Bayesian solution, but sometimes that is too computationally challenging.

A more pragmatic solution is to select a value of $\hat{\eta}$ using cross-validation instead. It is perfectly possible to use cross-validation, but the Bayesian approach also comes with an alternative for selecting hyperparameters η by maximizing the marginal likelihood (9.2) as

$$\hat{\eta} = \arg \max_{\eta} p_{\eta}(\mathbf{y}), \quad (9.4)$$

where we have added an index η to emphasize the fact that the value of $p(\mathbf{y})$ depends on η . This approach is sometimes referred to as *empirical Bayes*. Choosing hyperparameter η is, in a sense, a selection of a model (and/or prior), and we can therefore understand the marginal likelihood as a tool for selecting a model. Maximizing the marginal likelihood is, however, not equivalent to using cross-validation (the obtained hyperparameter value might differ), and unlike cross-validation does the marginal likelihood not give an estimate of E_{new} . In many situations, however, it is relatively easy to compute (and maximize) the marginal likelihood, compared to employing a full cross-validation procedure.

In the previous section we argued that the Bayesian approach was less prone to overfit compared to using the maximum likelihood approach. However, maximizing the marginal likelihood is, in a sense, a kind of a maximum likelihood approach, and one may ask if there is a risk of overfitting when maximizing the marginal likelihood. To some extent that might be the case, but the key point is that handling one (or, at most, a few) hyperparameters η with maximum (marginal) likelihood does typically not cause any severe overfit, much like there rarely are overfit issues when learning a straight line with plain linear regression. In other words, we can usually “afford” to learn one or a few (hyper)parameters by maximizing the (marginal) likelihood; maximum likelihood typically becomes problematic only when learning a large number of (hyper)parameters.

9.2 Bayesian linear regression

As a first example of the Bayesian approach, we will apply it to linear regression. In itself Bayesian linear regression is perhaps not the most useful and versatile method, but just like ordinary linear regression it is a good starting point and illustrates the main concepts well. Just like ordinary linear regression it is possible to extend in various directions, and it also opens up for the arguably more interesting and useful Gaussian process. Before we work out the details of the Bayesian approach applied to linear regression, we will however repeat some facts about the multivariate Gaussian distribution that will be useful later.

The multivariate Gaussian distribution

A central mathematical object for Bayesian linear regression (and later also the Gaussian process) is the multivariate Gaussian distribution. We assume the reader already has some familiarity with multivariate random variables, or equivalently random vectors, and repeat only the most important properties of the multivariate Gaussian distribution here.

Let \mathbf{z} denote a q -dimensional multivariate Gaussian random vector $\mathbf{z} = [z_1 \ z_2 \ \dots \ z_q]^T$. The multivariate Gaussian distribution is parametrized by a q -dimensional mean vector $\boldsymbol{\mu}$ and a $q \times q$ covariance matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1q} \\ \sigma_{21} & \sigma_2^2 & & \sigma_{2q} \\ \vdots & & & \vdots \\ \sigma_{q1} & \sigma_{q2} & \dots & \sigma_q^2 \end{bmatrix}.$$

The covariance matrix is a real-valued positive semidefinite matrix, that is, a symmetric matrix with nonnegative eigenvalues. As a shorthand notation we write $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ or $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$. Note that we use the same symbol \mathcal{N} to denote the univariate as well as the multivariate Gaussian distribution. The reason is that the former is only a special case of the latter.

The expected value of \mathbf{z} is $\mathbb{E}[\mathbf{z}] = \boldsymbol{\mu}$ and the variance of z_1 is $\text{var}(z_1) = \mathbb{E}[(z_1 - \mathbb{E}[z_1])^2] = \sigma_1^2$, and similarly for z_2, \dots, z_q . Moreover the covariance between z_1 and z_2 is $\text{cov}(z_1, z_2) = \mathbb{E}[(z_1 - \mathbb{E}[z_1])(z_2 - \mathbb{E}[z_2])] = \sigma_{12}$.

$\sigma_{12} = \sigma_{21}$, and similarly for z_3, \dots, z_q . All these properties can be derived from the probability density function of the multivariate Gaussian distribution, which is

$$\mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{q}{2}} \det(\boldsymbol{\Sigma})^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu})\right) \quad (9.5)$$

If all off-diagonal elements of $\boldsymbol{\Sigma}$ are 0, the elements of \mathbf{z} are just independent univariate Gaussian random variables. However, if some off-diagonal element, say σ_{ij} ($i \neq j$), is nonzero there is a correlation between z_i and z_j . Intuitively the correlation means that z_i carries information also about z_j , and vice versa. Some important results on how the multivariate Gaussian distribution can be manipulated are summarized in Appendix 9.A.

Linear regression with the Bayesian approach

We will now apply the Bayesian approach to the linear regression model. We will first spend some effort on mapping the elements of linear regression from Chapter 3 to the Bayesian terminology, and thereafter derive the solution.

From Chapter 3 we have that the linear regression model is

$$y = \boldsymbol{\theta}^\top \mathbf{x} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2) \quad (9.6)$$

which equivalently can be written as

$$p(y | \boldsymbol{\theta}) = \mathcal{N}\left(y; \boldsymbol{\theta}^\top \mathbf{x}, \sigma^2\right). \quad (9.7)$$

This is an expression for one output data point y , and for the entire vector of all training data outputs \mathbf{y} we can write

$$p(\mathbf{y} | \boldsymbol{\theta}) = \prod_{i=1}^n p(y_i | \boldsymbol{\theta}) = \prod_{i=1}^n \mathcal{N}\left(y_i; \boldsymbol{\theta}^\top \mathbf{x}_i, \sigma^2\right) = \mathcal{N}(\mathbf{y}; \mathbf{X}\boldsymbol{\theta}, \mathbf{I}\sigma^2) \quad (9.8)$$

where we in the last step used the fact that a n -dimensional Gaussian random vector with a diagonal covariance matrix is equivalent to n scalar Gaussian random variables.

In the Bayesian approach there is also a need for a prior $p(\boldsymbol{\theta})$ for the unknown parameters $\boldsymbol{\theta}$. In Bayesian linear regression the prior distribution is chosen as a Gaussian with mean $\boldsymbol{\mu}_0$ and covariance $\boldsymbol{\Sigma}_0$,

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0). \quad (9.9)$$

with $\boldsymbol{\Sigma}_0 = \mathbf{I}\sigma_0^2$. The reason for this choice is frankly that it simplifies the calculations (much like the squared error loss simplifies the computations for plain linear regression).

The next step is now to compute the posterior. It is perfectly possible to derive it using Bayes' theorem, but since $p(\mathbf{y} \text{ mod } \boldsymbol{\theta})$ as well as $p(\boldsymbol{\theta})$ are multivariate Gaussian distributions, Corollary 9.1 in Section 9.A directly gives us that

$$p(\boldsymbol{\theta} | \mathbf{y}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n) \quad (9.10a)$$

$$\boldsymbol{\mu}_n = \boldsymbol{\Sigma}_n \left(\frac{1}{\sigma_0^2} \boldsymbol{\mu}_0 + \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{y} \right) \quad (9.10b)$$

$$\boldsymbol{\Sigma}_n = \left(\frac{1}{\sigma_0^2} \mathbf{I} + \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} \right)^{-1} \quad (9.10c)$$

From (9.10) we can also derive the posterior predictive by Corollary 9.2 in Section 9.A,

$$p(y_\star | \mathbf{y}) = \mathcal{N}(y_\star; m_\star, s_\star) \quad (9.11a)$$

$$m_\star = \mathbf{x}_\star^\top \boldsymbol{\mu}_n \quad (9.11b)$$

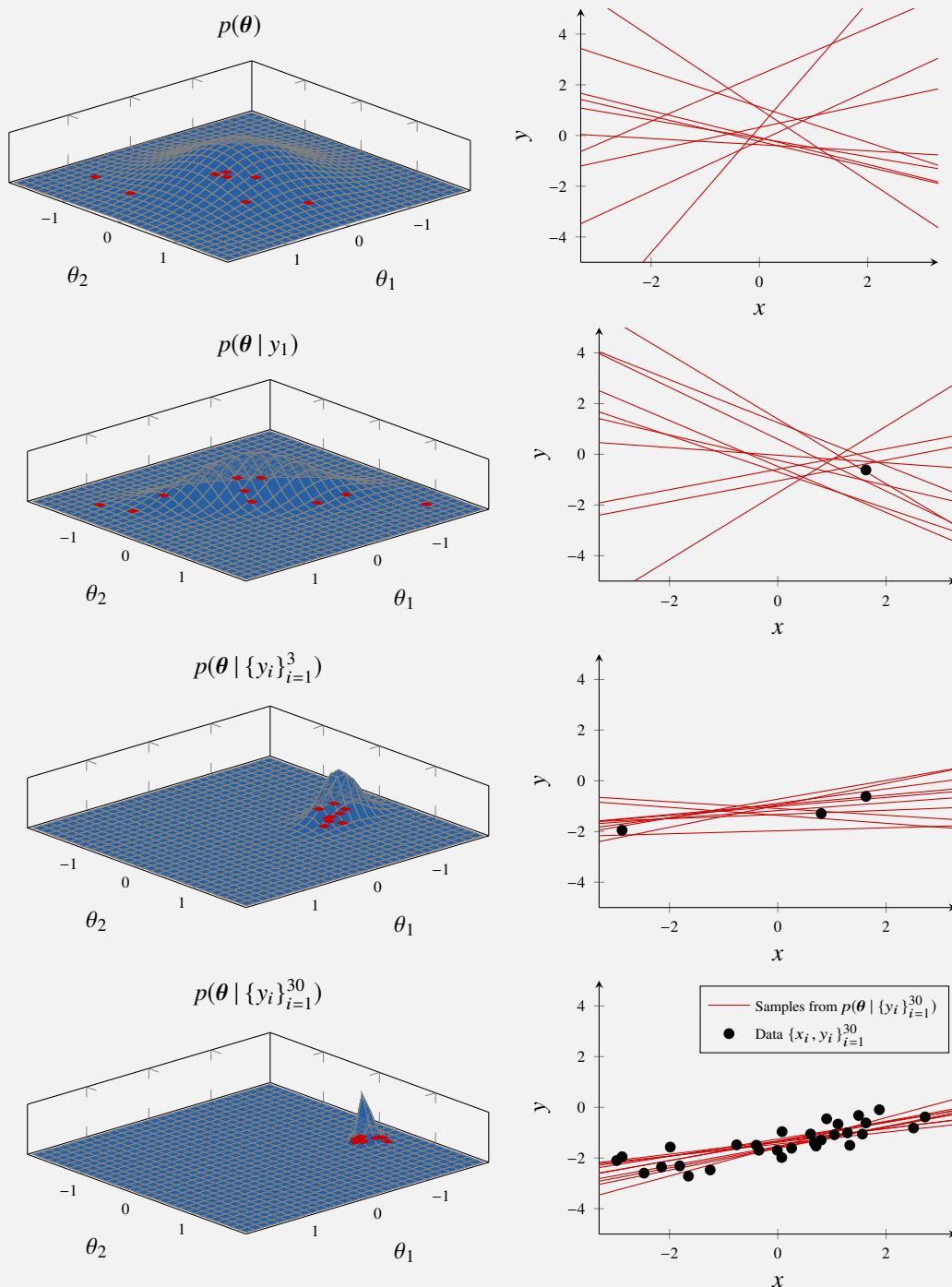
$$s_\star = \sigma^2 + \mathbf{x}_\star^\top \boldsymbol{\Sigma}_n \mathbf{x}_\star. \quad (9.11c)$$

We now have all pieces of Bayesian linear regression in place. The main difference to plain linear regression is that we compute a posterior $p(\boldsymbol{\theta} | \mathbf{y})$ (instead of a value $\hat{\boldsymbol{\theta}}$) and a posterior predictive $p(y_\star | \mathbf{y})$ instead of a prediction \hat{y} . We illustrate it with Example 9.1.

We have so far assumed that the noise variance σ^2 is fix. Most often σ^2 is also a parameter the user has to decide, which can be done by maximizing the marginal likelihood. Corollary 9.2 gives us the marginal likelihood

$$p(\mathbf{y}) = \mathcal{N}\left(\mathbf{y}; \mathbf{X}\boldsymbol{\mu}_0, \sigma^2 \mathbf{I} + \mathbf{X}\boldsymbol{\Sigma}_0\mathbf{X}^\top\right). \quad (9.12)$$

It is also possible to chose the prior variance σ_0^2 in $\boldsymbol{\Sigma}_0 = \mathbf{I}\sigma_0^2$ by maximizing (9.12).

Example 9.1: Bayesian linear regression


TODO: Describe

It is perfectly possible to use nonlinear input transformations, such as polynomials, as we introduced for plain linear regression. We give an example of that in Example 9.2, where we return to the running regression example of car stopping distances. We can, however, take one step further and also use the kernel trick from Chapter 8. That will lead us to the Gaussian process, which is the topic of Section 9.3.

Example 9.2: Car stopping distances

To be completed.

Connection to regularized linear regression

The main feature of the Bayesian approach is that it provides a full distribution $p(\boldsymbol{\theta} | \mathbf{y})$ over the parameters $\boldsymbol{\theta}$, instead of only a point estimate $\hat{\boldsymbol{\theta}}$. There is, however, also an interesting connection between the Bayesian approach and regularization. We will make this concrete by considering the posterior $p(\boldsymbol{\theta} | \mathbf{y})$ from Bayesian linear regression and the point estimate $\hat{\boldsymbol{\theta}}_{L^2}$ obtained from L^2 -regularized linear regression with squared error loss. Let us extract the so called *maximum a posteriori (MAP)* point estimate $\hat{\boldsymbol{\theta}}_{\text{MAP}}$ from the posterior $p(\boldsymbol{\theta} | \mathbf{y})$. The MAP estimate is the value of $\boldsymbol{\theta}$ for which the posterior reaches its maximum,

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \mathbf{y}) = \arg \max_{\boldsymbol{\theta}} p(\mathbf{y} | \boldsymbol{\theta})p(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} [\ln p(\mathbf{y} | \boldsymbol{\theta}) + \ln p(\boldsymbol{\theta})] \quad (9.13)$$

where the second equality follows from the fact that $p(\boldsymbol{\theta} | \mathbf{y}) = \frac{p(\mathbf{y} | \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})}$ and that $p(\mathbf{y})$ does not depend on $\boldsymbol{\theta}$. Remember that L^2 -regularized linear regression can be understood as using the cost function ((3.20) and (3.48))

$$\hat{\boldsymbol{\theta}}_{L^2} = \arg \max_{\boldsymbol{\theta}} [\ln p(\mathbf{y} | \boldsymbol{\theta}) + \tilde{\lambda} \|\boldsymbol{\theta}\|_2^2] \quad (9.14)$$

with some regularization parameter $\tilde{\lambda}$. When comparing (9.14) to (9.13), we realize that if only $\ln p(\boldsymbol{\theta}) \propto \|\boldsymbol{\theta}\|_2^2$, the MAP estimate and the L^2 regularized estimate of $\boldsymbol{\theta}$ are identical for some value of $\tilde{\lambda}$. By (9.9), that is indeed the case. The MAP estimate is indeed identical to $\hat{\boldsymbol{\theta}}_{L^2}$.

This connection between MAP estimates and regularized maximum likelihood estimates holds as long as the regularization is proportional to the logarithm of the prior. For example, if we instead would chose a Laplace prior for $\boldsymbol{\theta}$, the MAP estimate would be identical to L^1 regularization. In general there are many regularization methods that can be interpreted as implicitly choosing a certain prior. This connection to regularization gives another perspective on why the Bayesian approach is less prone to overfit.

Remark 9.1 *The connection between the Bayesian approach and the use of regularization does not imply that the two approaches are equivalent. The main point with the Bayesian approach is still that a posterior distribution is computed for $\boldsymbol{\theta}$ (instead of just a point estimate $\hat{\boldsymbol{\theta}}$).*

9.3 The Gaussian process

We introduced the Bayesian approach as the idea of considering unknown parameters θ as random variables, and consequently learning a posterior distribution $p(\theta | \mathbf{y})$ instead of a single value $\hat{\theta}$. The Bayesian idea does, however, not only apply to models with parameters, but also to non-parametric models. We will now introduce the Gaussian process, where instead of considering only parameters θ as being random variables we effectively consider an entire function f function to be a stochastic process, and compute the posterior $p(f | \mathbf{y})$. The Gaussian process is a very interesting and useful Bayesian non-parametric model. In a nutshell it is the Bayesian approach applied to kernel ridge regression (Chapter 8). We will present the Gaussian process as a method for handling regression problems, but it is possible to also extend it to deal with classification problems as well (somewhat similar to how linear regression is modified into logistic regression).

In this section we will discuss the fundamentals of the Gaussian process, and thereafter in Section 9.4 discuss more on how it can be used as a supervised machine learning method. We will first discuss what a Gaussian process is, at all, and thereafter see how we can construct a Gaussian process that connects closely to kernel ridge regression from Section 8.2.

Remember that for regression in general we have

$$y_\star = f(\mathbf{x}_\star) + \varepsilon. \quad (9.15)$$

As for Bayesian linear regression, we assume $\varepsilon \sim \mathcal{N}(0, \sigma^2)$. Throughout the Gaussian process discussion we will however consider the posterior predictive for the (theoretical) “noise-free” $f(\mathbf{x}_\star)$ instead of the noise-corrupted test output y_\star . In mathematical terms we will compute $p(\mathbf{f}(\mathbf{x}_\star) | \mathbf{y})$ instead of $p(y_\star | \mathbf{y})$. In the previous chapters with only a point prediction $\hat{y}(\mathbf{x}_\star)$ there was no difference in predicting $f(\mathbf{x}_\star)$ or y_\star , since ε by assumption has mean 0. Now, when considering the posterior predictive, there is a difference since we are less certain (there is more variance) about y_\star than $f(\mathbf{x}_\star)$ because of the random noise ε added to it. With the Gaussian process, $p(\mathbf{f}(\mathbf{x}_\star) | \mathbf{y})$ will be a Gaussian distribution, and we will later discuss how it can be turned into $p(y_\star | \mathbf{y})$ (another Gaussian distribution). Note, however, that we still assume that the training data \mathbf{y} contains noise ε .

What is a Gaussian process?

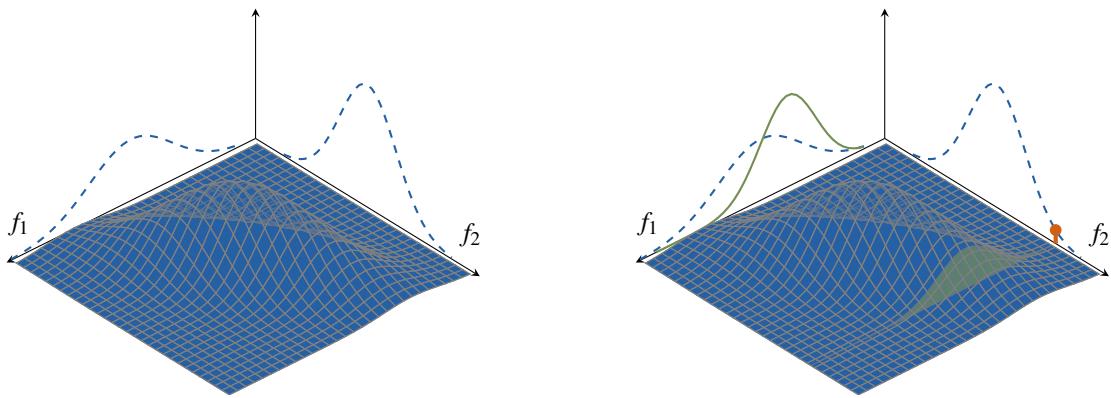
A Gaussian process is a stochastic process. A stochastic process, in turn, is a generalization of a random variable. Whereas a random variable has a finite dimension (even though it may take infinitely many different values), a stochastic process is not limited to that. We will now introduce the Gaussian process by extending a multivariate Gaussian random variable into a Gaussian process.

Consider a q -dimensional random vector $\mathbf{f} = [f_1 \ \cdots \ f_q]^\top$ with a multivariate Gaussian distribution

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (9.16)$$

with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. The reason why we denote the random variable by \mathbf{f} (instead of, say, \mathbf{z}) is that it eventually will represent function values (like $f(\mathbf{x}_1)$ and $f(\mathbf{x}_2)$) in the regression problem. We partition \mathbf{f} into two vectors \mathbf{f}_1 and \mathbf{f}_2 such that $\mathbf{f} = [\mathbf{f}_1^\top \ \mathbf{f}_2^\top]^\top$, and $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ similarly, and then write

$$p\left(\begin{bmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right). \quad (9.17)$$



(a) A two-dimensional Gaussian distribution for the random variables f_1 and f_2 , with a blue surface plot for the density, and the marginal distribution for each component sketched using dashed blue lines along each axis. Note that the marginal distributions do *not* contain all information about the distribution of f_1 and f_2 , since the covariance information is lacking in that representation.

(b) The conditional distribution of f_1 (green line), when f_2 is observed (orange dot). The conditional distribution of f_1 is given by (9.18), which (apart from a normalizing constant) in this graphical representation also is the green ‘slice’ of the joint distribution (blue surface). The marginals of the joint distribution from Figure 9.1a are kept for reference (blue dashed lines).

Figure 9.1: A two-dimensional multivariate Gaussian distribution for f_1 and f_2 in (a), and the conditional distribution for f_1 , when a particular value of f_2 is observed, in (b).



(a) The marginal distributions for f_1 and f_2 from Figure 9.1a.

(b) The distribution for f_1 (green line) when f_2 is observed (orange dot), as in Figure 9.1b.

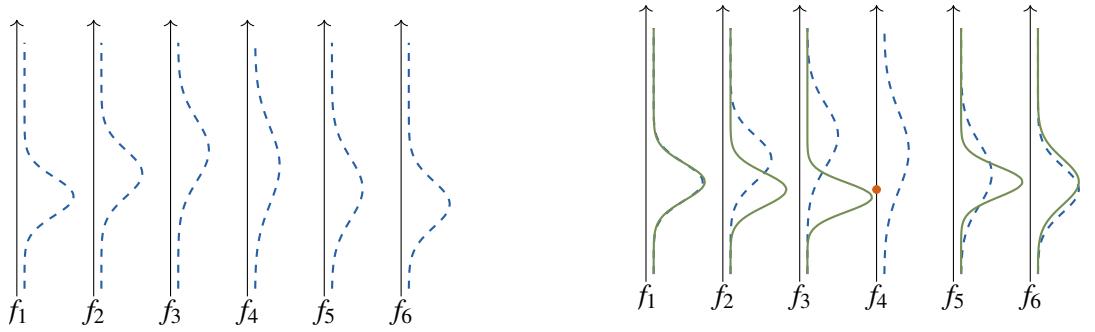
Figure 9.2: The marginals of the distributions in Figure 9.1, here plotted slightly differently. Note that this more compact plot comes with the cost of missing the information about the covariance between f_1 and f_2 .

If some elements of \mathbf{f} , let us say the ones in \mathbf{f}_2 , are observed, the conditional distribution for \mathbf{f}_1 given the observation of \mathbf{f}_2 is, by Theorem 9.2,

$$p(\mathbf{f}_1 | \mathbf{f}_2) = \mathcal{N}\left(\mathbf{f}_1; \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{f}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}\right). \quad (9.18)$$

The conditional distribution is nothing but another Gaussian distribution with closed-form expressions for the mean and covariance. This is particularly useful.

Figure 9.1 shows a 2-dimensional example (\mathbf{f}_1 is a scalar f_1 , and \mathbf{f}_2 is a scalar f_2) where a multivariate Gaussian distribution is conditioned on an observation of f_2 . In Figure 9.2 we have plotted the marginal distributions from Figure 9.1. We have made this plot to prepare for the generalization to the Gaussian process. In a similar fashion to Figure 9.2 we can plot a 6-dimensional multivariate Gaussian distribution by its marginal distributions in Figure 9.3. Bear in mind that to fully illustrate the joint distribution for f_1, \dots, f_6 , a 6-dimensional surface plot would be needed, whereas Figure 9.3a only contains the marginal distributions for each component. We may also condition the 6-dimensional distribution underlying



(a) A 6-dimensional Gaussian distribution, plotted in the same way as Figure 9.2a, i.e., only its marginals are illustrated.

(b) The conditional distribution \$f_1, f_2, f_3, f_5\$ and \$f_6\$ when \$f_4\$ is observed (orange dot), illustrated by its marginals (green lines), cf Figure 9.2b.

Figure 9.3: A 6-dimensional Gaussian distribution, illustrated in the same fashion as Figure 9.2.

Figure 9.3a on an observation of, say, \$f_4\$. Once again, the conditional distribution is another Gaussian distribution, and the marginal distributions of the 5-dimensional random variable \$[f_1, f_2, f_3, f_5, f_6]^\top\$ are plotted in Figure 9.3b.

In Figure 9.2 and 9.3 we illustrated the marginal distributions for a finite-dimensional multivariate Gaussian random variable. However, our aim is the Gaussian process, which is a stochastic process on a continuous space.

The extension of the Gaussian distribution (defined on a finite set) to the Gaussian process (defined on a continuous space) is achieved by replacing the index set \$\{1, 2, 3, 4, 5, 6\}\$ in Figure 9.3 by a variable \$\mathbf{x}\$ taking values on a continuous space, for example the real line. We then also have to replace the random variables \$f_1, f_2, \dots, f_6\$ with a random function (a stochastic process) \$f\$ which can be evaluated at any \$\mathbf{x}\$ as \$f(\mathbf{x})\$. Furthermore, in the Gaussian multivariate distribution \$\boldsymbol{\mu}\$ is a vector with \$q\$ components, and \$\boldsymbol{\Sigma}\$ is a \$q \times q\$ matrix. In the Gaussian process we replace \$\boldsymbol{\mu}\$ by a mean *function* \$\mu(\mathbf{x})\$ into which we can insert any \$\mathbf{x}\$, and the covariance matrix \$\boldsymbol{\Sigma}\$ by a covariance *function* \$\kappa(\mathbf{x}, \mathbf{x}')\$ into which we can insert any pair \$\mathbf{x}\$ and \$\mathbf{x}'\$.

Let us from this define the Gaussian process. If, for any finite set of points \$\{\mathbf{x}_1, \dots, \mathbf{x}_n\}\$, it holds that

$$p\left(\begin{bmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_n) \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_n) \end{bmatrix}; \begin{bmatrix} \mu(\mathbf{x}_1) \\ \vdots \\ \mu(\mathbf{x}_n) \end{bmatrix}, \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ \kappa(\mathbf{x}_n, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}\right), \quad (9.19)$$

then \$f\$ is a Gaussian process.

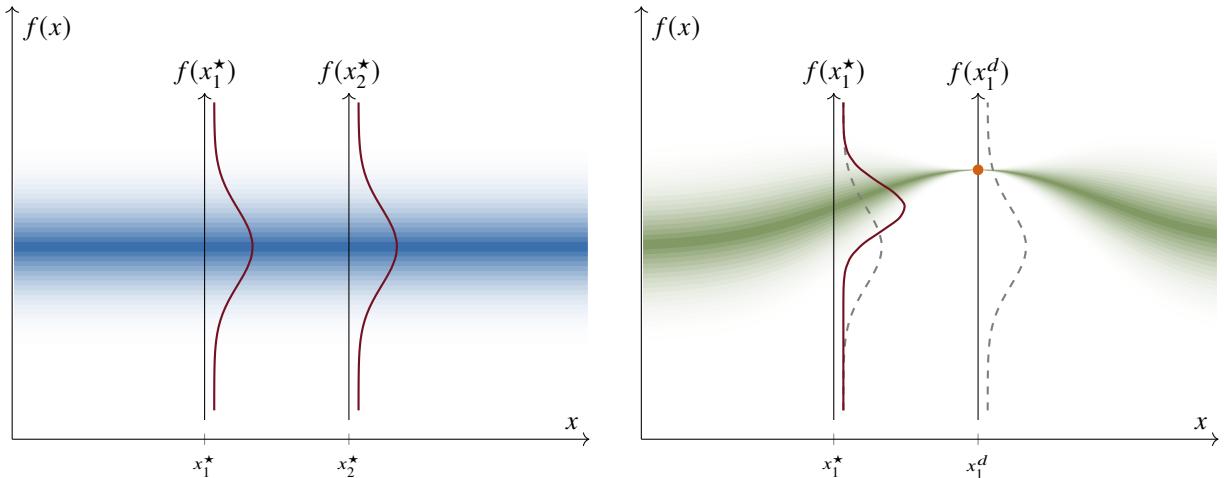
That is, with a Gaussian process \$f\$ and any choice of \$\{\mathbf{x}_1, \dots, \mathbf{x}_n\}\$, the vector of functions values \$[f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)]\$ have a multivariate Gaussian distribution, just like the one in Figure 9.3. Since \$\{\mathbf{x}_1, \dots, \mathbf{x}_n\}\$ can be chosen arbitrarily from the continuous space on which it lives, the Gaussian process defines a distribution for *all* points in that space. Indeed, for this definition to make sense, \$\kappa(\mathbf{x}, \mathbf{x}')\$ has to be such that a positive semidefinite covariance matrix is obtained for any choice of \$\{\mathbf{x}_1, \dots, \mathbf{x}_n\}\$.

We will use the notation

$$f \sim \mathcal{GP}(\boldsymbol{\mu}, \kappa) \quad (9.20)$$

to express that the function \$f\$ is distributed according to a Gaussian process with mean function \$\mu(\mathbf{x})\$ and covariance function \$\kappa(\mathbf{x}, \mathbf{x}')\$. If we want to illustrate a Gaussian process, which we do in Figure 9.4, we can choose \$\{\mathbf{x}_1, \dots, \mathbf{x}_n\}\$ to correspond to the pixels on the screen or the printer dots on the paper, and print the marginal distributions for each \$\{\mathbf{x}_1, \dots, \mathbf{x}_n\}\$, so that it appears as a continuous line to the eye (despite that we actually can access the distribution only in a finite, however arbitrary, set of points).

It is no coincidence that we use the same symbol \$\kappa\$ for covariance functions as we used for kernels in Chapter 8. As we soon will discuss, applying the Bayesian approach to kernel ridge regression will give us a Gaussian process in which the covariance function is the kernel.



(a) A Gaussian process defined on the real line parameterized by x , not conditioned on any observations. The intensity of the blue color is proportional to the (marginal) density, and the marginal distributions for some x_1^* and x_2^* are pictured in red. Akin to Figure 9.3, we only plot the marginal distribution for each \mathbf{x}_\star , but the Gaussian process defines a full joint distribution for all points on the x -axis, even though it is hard to illustrate.

(b) The conditional Gaussian process distribution given the observation of $f(x_1^d)$ in the point x_1^d corresponding to x_2^* in (a). The prior distribution from Figure (a) is dashed gray. Note how the conditional distribution adjusts to the observation, both in terms of mean (closer to the observation) and (marginal) variance (smaller in the proximity of the observation, but it remains unchanged in areas distant from it).

Figure 9.4: A Gaussian process. Figure (a) shows the prior distribution (shaded blue), whereas (b) shows the posterior distribution (shaded green) after conditioning on one observation (orange dot).

We can also condition the Gaussian process on some observations $\{f(\mathbf{x}_i), \mathbf{x}_i\}_{i=1}^n$, the Gaussian process counterpart to Figure 9.1b, 9.2b and 9.3b. As usually we stack the observed inputs in \mathbf{X} , and let $f(\mathbf{X})$ denote the observed outputs (we assume for now that the observations are made without any noise). We use the notation with $\mathbf{K}(\mathbf{X}, \mathbf{X})$ and $\mathbf{K}(\mathbf{X}, \mathbf{x}_\star)$ as defined by (8.13b) to write the joint distribution between the observed values $f(\mathbf{X})$ and the test value $f(\mathbf{x}_\star)$ as

$$p\left(\begin{bmatrix} f(\mathbf{x}_\star) \\ f(\mathbf{X}) \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} f(\mathbf{x}_\star) \\ f(\mathbf{X}) \end{bmatrix}; \begin{bmatrix} \mu(\mathbf{x}_\star) \\ \mu(\mathbf{X}) \end{bmatrix}, \begin{bmatrix} \kappa(\mathbf{x}_\star, \mathbf{x}_\star) & \mathbf{K}(\mathbf{X}, \mathbf{x}_\star)^T \\ \mathbf{K}(\mathbf{X}, \mathbf{x}_\star) & \mathbf{K}(\mathbf{X}, \mathbf{X}) \end{bmatrix}\right). \quad (9.21)$$

Now, as we have observed $f(\mathbf{X})$, we use the by now well-known expressions for the Gaussian distribution to write the distribution for $f(\mathbf{x}_\star)$ conditional on the observations of $f(\mathbf{X})$ as

$$\begin{aligned} p(f(\mathbf{x}_\star) | f(\mathbf{X})) = \\ \mathcal{N}\left(f(\mathbf{x}_\star); \mu(\mathbf{x}_\star) + \mathbf{K}(\mathbf{X}, \mathbf{x}_\star)^T \mathbf{K}(\mathbf{X}, \mathbf{X})^{-1} (f(\mathbf{X}) - \mu(\mathbf{X})), \kappa(\mathbf{x}_\star, \mathbf{x}_\star) - \mathbf{K}(\mathbf{X}, \mathbf{x}_\star)^T \mathbf{K}(\mathbf{X}, \mathbf{X})^{-1} \mathbf{K}(\mathbf{X}, \mathbf{x}_\star)\right), \end{aligned} \quad (9.22)$$

which is another Gaussian distribution for any test input \mathbf{x}_\star . We illustrate this by Figure 9.4.

We have now introduced the somewhat abstract concept of a Gaussian process. In some subjects, such as signal processing, are so-called white Gaussian processes common. A white Gaussian process has a white covariance function

$$\kappa(\mathbf{x}, \mathbf{x}') = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}', \\ 0 & \text{otherwise,} \end{cases} \quad (9.23)$$

which implies that $f(\mathbf{x})$ is uncorrelated to $f(\mathbf{x}')$ unless $\mathbf{x} = \mathbf{x}'$. White Gaussian processes are of less use in supervised machine learning, but we will instead have a look at how kernel ridge regression can be turned into a Gaussian process, where the mean function becomes zero and the covariance function becomes the kernel from Chapter 8.

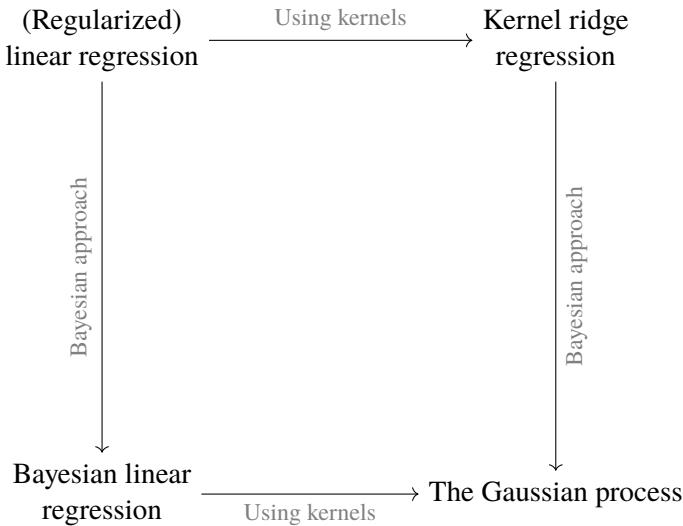


Figure 9.5: A graphical summary on the connections between linear regression, Bayesian linear regression, kernel ridge regression and the Gaussian process.

Extending kernel ridge regression into a Gaussian process

In order to obtain a Gaussian process construction that we can use as a practical tool for supervised machine learning, we apply the kernel trick from Section 8.2 to Bayesian linear regression (9.11). The connection between linear regression, Bayesian linear regression, kernel ridge regression and the Gaussian process is summarized by Figure 9.5. This will lead us essentially to (9.22) again, with the kernel being the covariance function $\kappa(\mathbf{x}, \mathbf{x}')$ and the mean function being $\mu(\mathbf{x}) = 0$.

Let us now repeat the posterior predictive for Bayesian linear regression (9.11), however with three changes. The first change is that we assume the prior mean and covariance for θ is $\mu_0 = \mathbf{0}$ and $\sigma_0 = 1$, respectively. This assumption is not strictly needed for our purposes, but simplifies the expressions. The second change is that we, as in the previous section, consider the prediction of $f(\mathbf{x}_\star)$ rather than y_\star (see (9.15)), which simply amounts to not adding σ^2 in (9.11c). Lastly, as in Section 8.1, we introduce nonlinear feature transformations $\phi(\mathbf{x})$ of the input variable \mathbf{x} in the linear regression model. We therefore replace \mathbf{X} with $\Phi(\mathbf{X})$ in the notation. Altogether (9.11) becomes

$$p(f(\mathbf{x}_\star) | \mathbf{y}) = \mathcal{N}(f(\mathbf{x}_\star); m_\star, s_\star) \quad (9.24a)$$

$$m_\star = \phi(\mathbf{x}_\star)^\top \left(\sigma^2 \mathbf{I} + \Phi(\mathbf{X})^\top \Phi(\mathbf{X}) \right)^{-1} \Phi(\mathbf{X})^\top \mathbf{y} \quad (9.24b)$$

$$s_\star = \phi(\mathbf{x}_\star)^\top \left(\mathbf{I} + \frac{1}{\sigma^2} \Phi(\mathbf{X})^\top \Phi(\mathbf{X}) \right)^{-1} \phi(\mathbf{x}_\star) \quad (9.24c)$$

In a similar fashion to the derivation of kernel ridge regression, we use the push-through matrix identity $\mathbf{A}(\mathbf{A}^\top \mathbf{A} + \mathbf{I})^{-1} = (\mathbf{A}\mathbf{A}^\top + \mathbf{I})^{-1}\mathbf{A}$ to re-write m_\star with the aim to have $\phi(\mathbf{x})$ only entering through inner products,

$$m_\star = \phi(\mathbf{x}_\star)^\top \Phi(\mathbf{X})^\top \left(\sigma^2 \mathbf{I} + \Phi(\mathbf{X}) \Phi(\mathbf{X})^\top \right)^{-1} \mathbf{y}. \quad (9.25a)$$

To re-write s_\star in a similar fashion, we have to use the matrix inversion lemma $(\mathbf{I} - \mathbf{U}\mathbf{V})^{-1} = \mathbf{I} - \mathbf{U}(\mathbf{I} + \mathbf{V}\mathbf{U})^{-1}\mathbf{V}$ (which holds for any matrices \mathbf{U}, \mathbf{V} of compatible dimensions),

$$s_\star = \phi(\mathbf{x}_\star)^\top \phi(\mathbf{x}_\star) - \phi(\mathbf{x}_\star)^\top \Phi(\mathbf{X})^\top \left(\sigma^2 \mathbf{I} + \Phi(\mathbf{X}) \Phi(\mathbf{X})^\top \right)^{-1} \Phi(\mathbf{X}) \phi(\mathbf{x}_\star). \quad (9.25b)$$

In complete analogy to the derivation of kernel ridge regression as in (8.11), we are now ready to apply the kernel trick and replace all instances of $\phi(\mathbf{x})^\top \phi(\mathbf{x}')$ with a kernel $\kappa(\mathbf{x}, \mathbf{x}')$. With the same notation as in

(8.11), we get

$$m_\star = \mathbf{K}(\mathbf{X}, \mathbf{x}_\star)^\top \left(\sigma^2 \mathbf{I} + \mathbf{K}(\mathbf{X}, \mathbf{X}) \right)^{-1} \mathbf{y} \quad (9.26a)$$

$$s_\star = \kappa(\mathbf{x}_\star, \mathbf{x}_\star) - \mathbf{K}(\mathbf{X}, \mathbf{x}_\star)^\top \left(\sigma^2 \mathbf{I} + \mathbf{K}(\mathbf{X}, \mathbf{X}) \right)^{-1} \mathbf{K}(\mathbf{X}, \mathbf{x}_\star). \quad (9.26b)$$

The posterior predictive that is defined by (9.24a) and (9.26) is the Gaussian process model again, identical to (9.22) if $\mu(\mathbf{x}_\star) = 0$ and $\sigma^2 = 0$. The reason for $\mu(\mathbf{x}_\star) = 0$ is that we made this derivation starting with $\mu_0 = 0$. When we derived (9.22) we assumed that we observed $f(\mathbf{x}_\star)$ (rather than $y_\star = f(\mathbf{x}_\star) + \varepsilon$), which is the reason why $\sigma^2 = 0$ in (9.22). The Gaussian process is hence a kernel version of Bayesian linear regression, much like kernel ridge regression is a kernel version of (regularized) linear regression, as pictured by Figure 9.5.

The fact that the kernel plays the role of a covariance function in the Gaussian process gives us another interpretation of the kernel in addition to the ones in Section 8.4, namely that the kernel $\kappa(\mathbf{x}, \mathbf{x}')$ determines how much $f(\mathbf{x})$ and $f(\mathbf{x}')$ are assumed to correlate.

We have so far only derived the posterior predictive for $f(\mathbf{x}_\star)$. Since we have $y_\star = f(\mathbf{x}_\star) + \varepsilon$, where ε is assumed to be independent with variance σ^2 , we get

$$p(y_\star | \mathbf{y}) = \mathcal{N}(y_\star; m_\star, s_\star + \sigma^2) \quad (9.27)$$

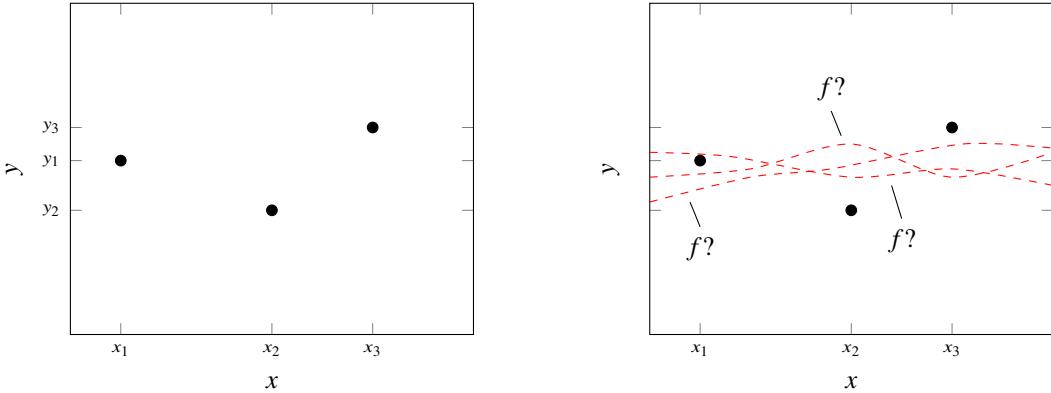
if we are interested in the posterior predictive for y_\star .

Time to reflect 9.1: Verify that you retrieve Bayesian linear regression when using the linear kernel (8.21) in the Gaussian process. Why is that?

A nonparametric distribution over functions

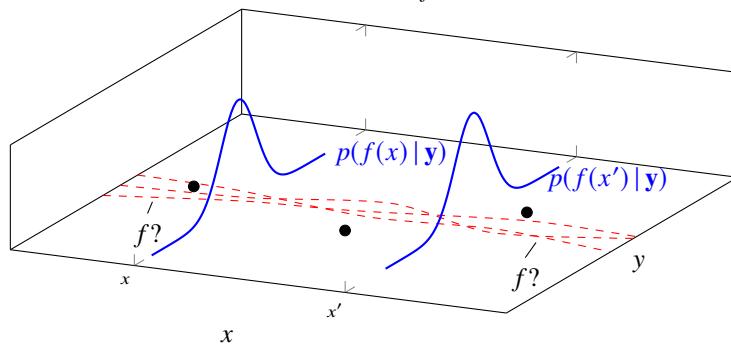
The Gaussian process is a somewhat abstract object, but as a supervised machine learning tool we use it for the very concrete purpose of predictions, that is, computing the posterior predictive $p(f(\mathbf{x}_\star) | \mathbf{y})$ (or $p(y_\star | \mathbf{y})$). However, unlike most other methods which only delivers a point prediction $\hat{y}(\mathbf{x}_\star)$, the posterior predictive is a distribution. Since we can compute the posterior predictive for any \mathbf{x}_\star , the Gaussian process actually defines a *distribution over functions*, as we illustrate by Figure 9.6.

Much like we could derive a connection between Gaussian linear regression and L^2 regularized linear regression, there is a connection between the Gaussian process and kernel ridge regression as well. If we only consider the mean m_\star of the posterior predictive, we retain the kernel ridge regression. To take full advantage of the Bayesian perspective, we have to consider also the posterior predictive variance s_\star . With most kernels, the predictive variance is smaller if there is a training data point nearby, and larger if the closest training data point is distant. The predictive variance hence give a quantification of the “uncertainty” in the prediction . Altogether the Gaussian process is an interesting tool for regression problems, as we illustrate by Figure 9.7.



(a) The training data $\{x_i, y_i\}_{i=1}^3$ of a regression problem is given to us.

(b) The underlying assumption when we do regression is that there exists *some* function f , which describes the data as $y_i = f(x_i) + \varepsilon$. It is unknown to us, but the purpose of regression (no matter which method is used) is, in a sense, to nail f .



(c) The Gaussian process defines a distribution over f . We can condition that distribution on data (that is, learn) and access it for any input, say x and x' . That is, we make a prediction for x and x' . The Gaussian process gives us a Gaussian distribution for $f(x)$ and $f(x')$, illustrated by solid blue lines. That distribution is heavily influenced by the choice of kernel, which is a design choice in the Gaussian process.

Figure 9.6: The Gaussian process defines a distribution over functions, which we can condition on training data and access in arbitrary points (such as x and x') in order to compute predictions.

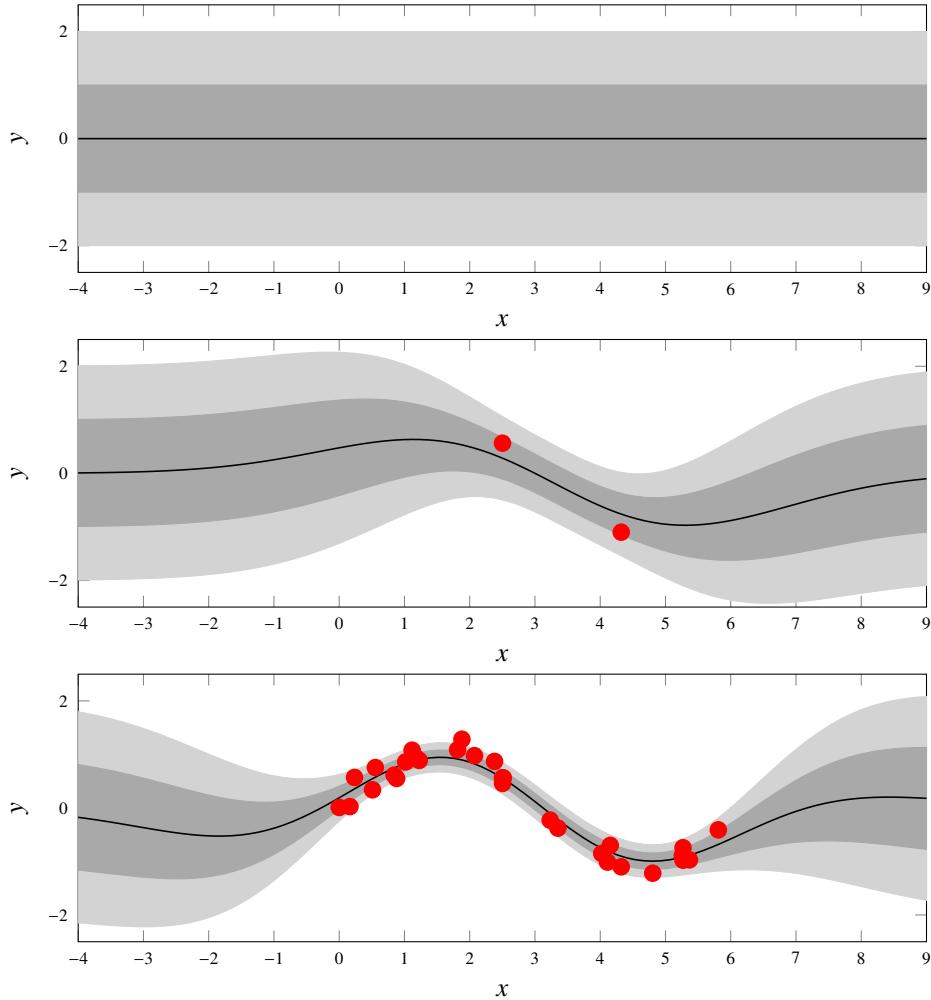


Figure 9.7: The Gaussian process as a supervised machine learning method: we can learn (that is, compute (9.24a) and (9.26)) the posterior predictive for $f(x)$ (shaded gray; dark gray for one standard deviation, light gray for 2 standard deviations, and solid black line for the mean) learned from 0 (upper), 2 (middle), 30 (bottom) observations (blue stars). We see how the model adapts to training data, and note in particular that the variance shrinks in the regions where observations are made, but remains larger in regions where no observations are made.

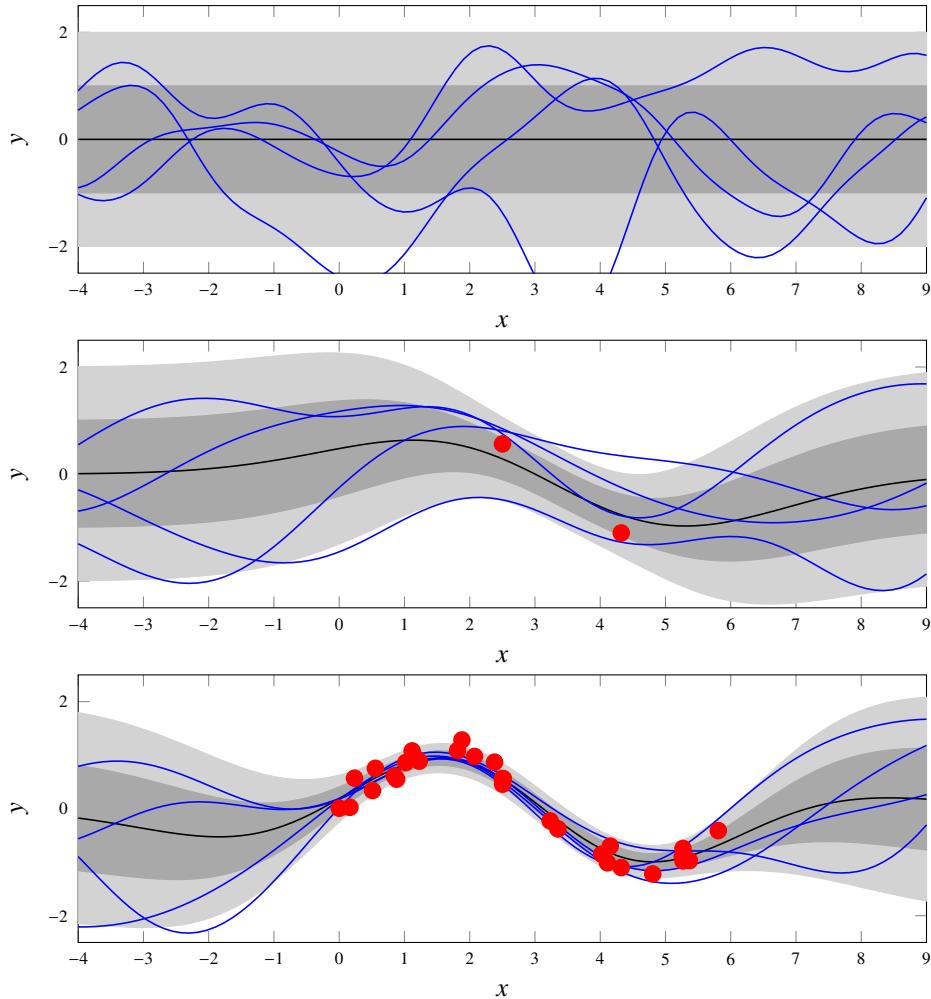


Figure 9.8: Figure 9.7 again, this time appended also with samples from $p(f(\mathbf{x}_*) | \mathbf{y})$.

Drawing samples from a Gaussian process

When making a prediction of $f(\mathbf{x}_*)$ for a single input point \mathbf{x}_* , the posterior predictive $p(f(\mathbf{x}_*) | \mathbf{y})$ captures all information the Gaussian process has about $f(\mathbf{x}_*)$. However, if we are interested not only in predicting $f(\mathbf{x}_*)$, but also $f(\mathbf{x}_* + \delta)$, the Gaussian process contains more information than is contained in the two posterior predictive distributions $p(f(\mathbf{x}_*) | \mathbf{y})$ and $p(f(\mathbf{x}_* + \delta) | \mathbf{y})$. This is because the Gaussian process also contains information on the correlation between $f(\mathbf{x}_*)$ and $f(\mathbf{x}_* + \delta)$, and the pitfall is that $p(f(\mathbf{x}_*) | \mathbf{y})$ and $p(f(\mathbf{x}_* + \delta) | \mathbf{y})$ are only the marginal distributions of the joint distribution $p(f(\mathbf{x}_*), f(\mathbf{x}_* + \delta) | \mathbf{y})$. As we saw in connection to Figure 9.1 and 9.2, the joint distribution contain more information than only its marginals.

If we are interested in computing predictions for a larger set of input values it can be rather cumbersome to grasp and visualize the resulting high-dimensional posterior predictive. It can therefore be an interesting alternative to visualize the Gaussian process posterior by samples from it. Technically this simply amounts to drawing a sample from the posterior predictive distribution, which we illustrate in Figure 9.8.

9.4 Practical usage of the Gaussian process

When using the Gaussian process as a method for supervised machine learning, there are a few important design choices left to the user. Like the methods presented in Chapter 8, the Gaussian process is a kernel method, and the choice of kernel is very important. Most kernels also contain a few hyperparameters, which have to be chosen as well. Since the Gaussian process belongs to the Bayesian framework, that choice can be done by maximizing the marginal likelihood, which we will discuss now.

Kernel choice

Since the Gaussian process can be understood as the Bayesian version of kernel ridge regression, also the Gaussian process requires a positive semidefinite kernel. Any of the positive semidefinite kernels presented in Section 8.4 can therefore be used also for Gaussian processes.

Among all kernel methods presented in this book, the Gaussian process could be the method where the exact choice of kernel has the biggest importance since the Gaussian posterior predictive $p(f(\mathbf{x}_\star | \mathbf{y}))$ has a mean and a variance, which both are heavily affected by the kernel. It is therefore important to make a good choice, and besides the discussion in Section 8.4 it can also be instructive to visually compare different kernel choices as in Figure 9.9. In the end it is a design choice left to the machine learning engineer. In the Bayesian perspective the kernel is an important part of the prior that makes crucial assumptions about the function f . For example, as can be seen from Figure 9.9 the squared exponential and the Matérn kernel with $\nu = \frac{1}{2}$ corresponds to drastically different assumptions about the smoothness of $f(\mathbf{x})$.

Hyperparameter tuning

Most kernels $\kappa(\mathbf{x}, \mathbf{x}')$ contain a few hyperparameters η (such as ℓ and α , as well as the noise variance σ^2), that needs to be determined. Cross-validation can indeed be used, but the Bayesian approach also comes with the option to maximize the marginal likelihood as a method for selecting hyperparameters, (9.4). From the construction of the Gaussian process, we have that $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})$, and consequently

$$\ln p(\mathbf{y}) = -\frac{1}{2}\mathbf{y}^\top (\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y} - \frac{1}{2} \ln \det (\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}) - \frac{n}{2} \log 2\pi. \quad (9.28)$$

Remember that the hyperparameters enter through $\mathbf{K}(\mathbf{X}, \mathbf{X})$. In other words, the hyperparameters of the Gaussian process kernel can be chosen by solving the optimization problem of maximizing (9.28). If using this approach, solving the optimization problem can be seen as a part of learning the Gaussian process. We illustrate this in Figure 9.10.

When selecting hyperparameters η of a kernel, it is important to be aware that (9.28) (as a function of the hyperparameters) may have multiple local maxima, as we illustrate by Figure 9.11. It is therefore important to carefully choose the initialization of the optimization procedure. The challenge with local maxima is not unique to using the marginal likelihood approach, but can arise also if using cross-validation for choosing hyperparameters.

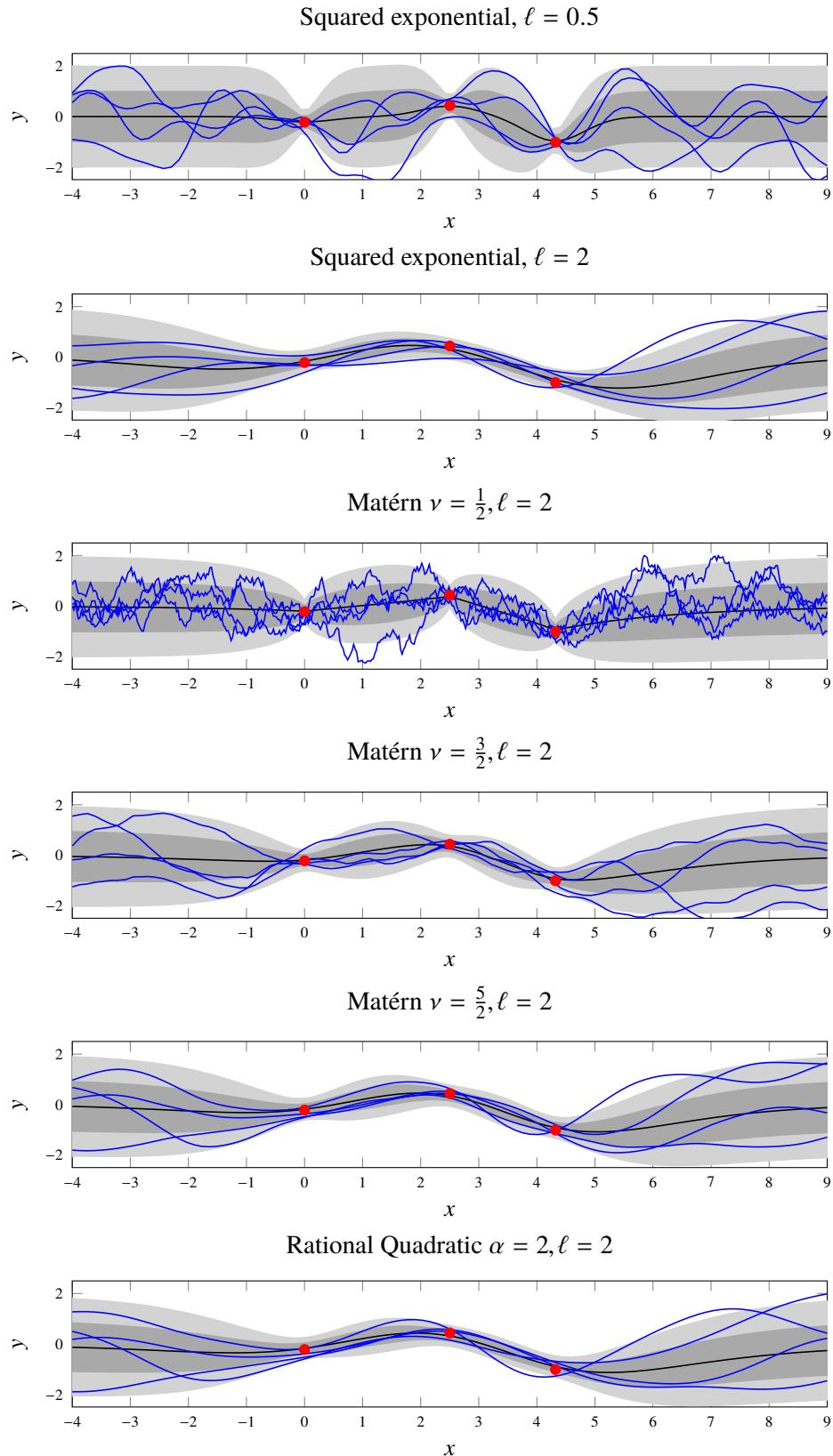


Figure 9.9: The Gaussian process applied with different kernels and hyperparameters to the same data, in order to illustrate what assumptions are made by the different kernels. The data is marked with red dots, and the Gaussian process is illustrated by its mean (black line), variance (gray areas; dark gray for one standard deviation and light gray for two standard deviations) and samples (blue lines).

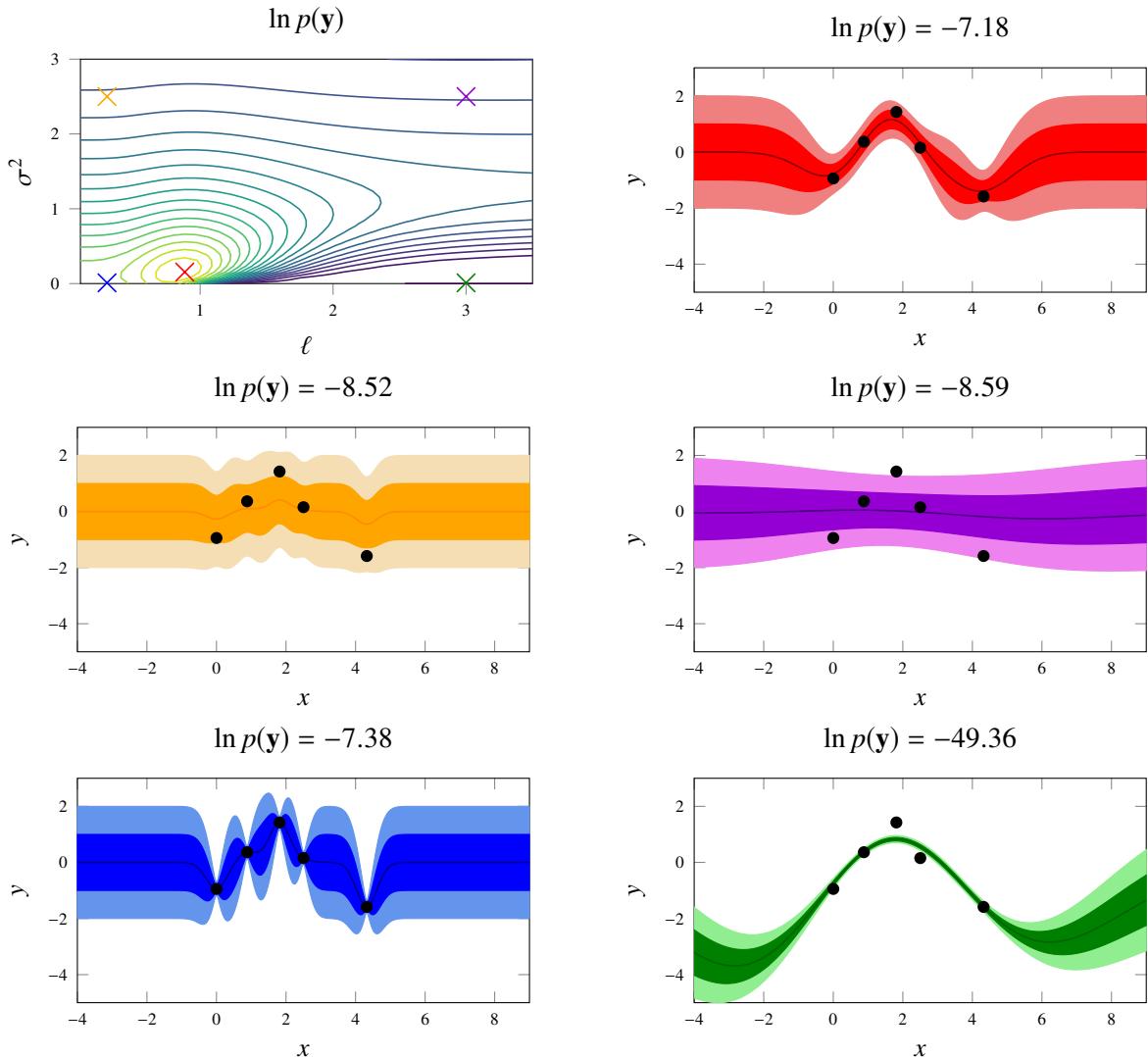


Figure 9.10: To choose hyperparameters for the Gaussian process kernel, the marginal likelihood $p(\mathbf{y})$ can be maximized. For a given dataset of five points (black dots) and the squared exponential kernel, the landscape of the logarithm of the marginal likelihood (as a function of the hyperparameters ℓ , lengthscales, and σ^2 , noise variance) is shown as a contour plot in the upper left panel. Each point in that plot corresponds to a certain selection of the hyperparameters ℓ, σ^2 . For five such points (the yellow, purple, blue, green and red crosses) the corresponding Gaussian process is shown. In particular is the red cross at the maximum of $p(\mathbf{y})$. The red upper right plot therefore corresponds to a Gaussian process where the hyperparameters have been maximized using marginal likelihood. It is clear that optimizing $p(\mathbf{y})$ does *not* mean selecting hyperparameters such that the Gaussian process follows the data as close as possible (as the blue one does).

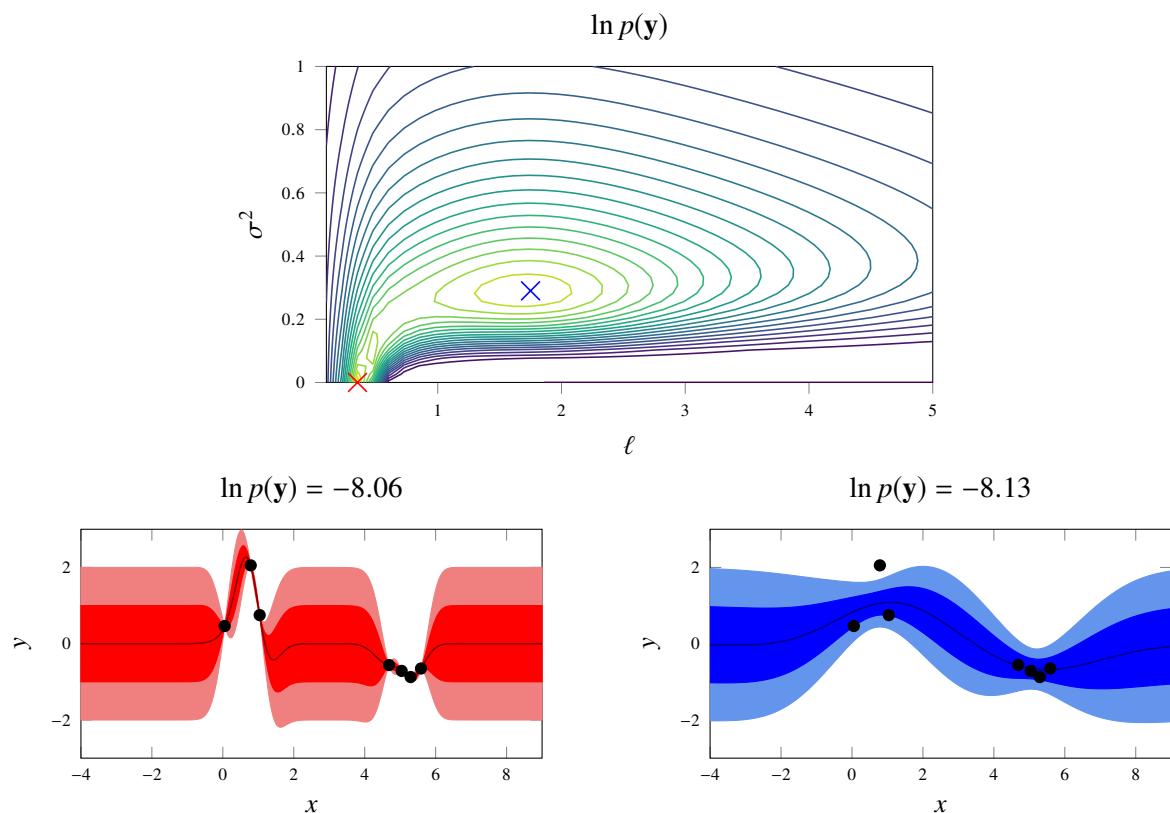


Figure 9.11: The landscape of $p(\mathbf{y})$ may have multiple local maxima. In this case there is one local maximum at the blue cross, with relatively large noise variance and length scale. There is also another local maximum, which also is the global one, at the red cross, with much less noise variance and a shorter lengthscale. There is also a third local maximum in between (not shown). It is not uncommon that the different maximum points corresponds to different “interpretations” of the data, as in this case. As a machine learning engineer it is important to be aware that this can happen; the red one does indeed optimize the marginal likelihood, but the blue one can also be practically useful.

9.A The multivariate Gaussian distribution

This appendix contains some results on the multivariate Gaussian distribution that are useful for Bayesian linear regression and the Gaussian process. Figure 9.12 summarizes how they relate to each other.

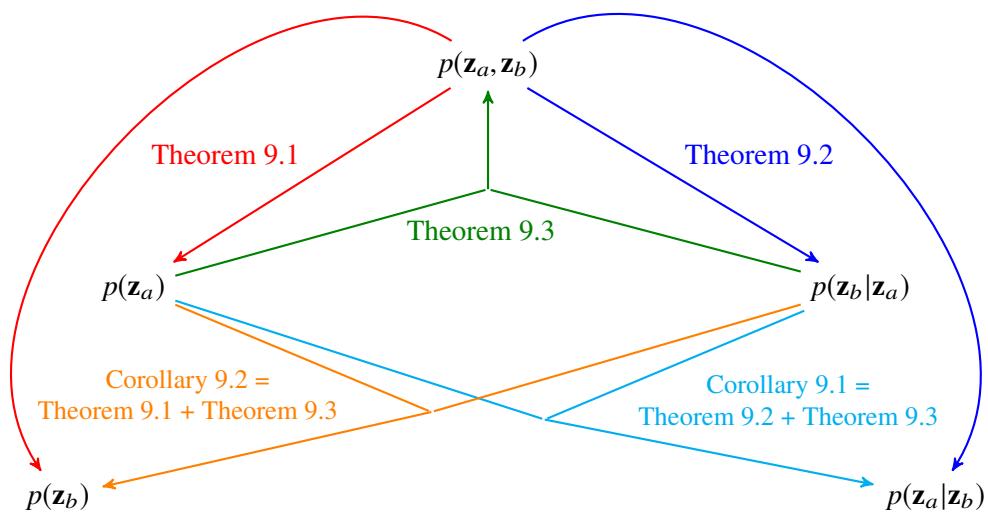


Figure 9.12: A graphical summary of how Theorem 9.1, 9.2, 9.3, Corollary 9.1 and Corollary 9.2 relate to each other. In all results, \mathbf{z}_a and \mathbf{z}_b are dependent multivariate Gaussian random variables.

Theorem 9.1 (Marginalization) Partition the Gaussian random vector $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ according to

$$\mathbf{z} = \begin{pmatrix} \mathbf{z}_a \\ \mathbf{z}_b \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}. \quad (9.29)$$

The marginal distribution $p(\mathbf{z}_a)$ is then given by

$$p(\mathbf{z}_a) = \mathcal{N}(\mathbf{z}_a; \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}). \quad (9.30)$$

Theorem 9.2 (Conditioning) Partition the Gaussian random vector $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ according to

$$\mathbf{z} = \begin{pmatrix} \mathbf{z}_a \\ \mathbf{z}_b \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}. \quad (9.31)$$

The conditional distribution $p(\mathbf{z}_a | \mathbf{z}_b)$ is then given by

$$p(\mathbf{z}_a | \mathbf{z}_b) = \mathcal{N}\left(\mathbf{z}_a; \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b}\right), \quad (9.32)$$

where

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{z}_b - \boldsymbol{\mu}_b), \quad (9.33)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}. \quad (9.34)$$

Theorem 9.3 (Affine transformation) Assume that \mathbf{z}_a , as well as \mathbf{z}_b conditioned on \mathbf{z}_a , are Gaussian distributed according to

$$p(\mathbf{z}_a) = \mathcal{N}(\mathbf{z}_a; \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_a), \quad (9.35)$$

$$p(\mathbf{z}_b | \mathbf{z}_a) = \mathcal{N}(\mathbf{z}_b; \mathbf{A}\mathbf{z}_a, \boldsymbol{\Sigma}_{b|a}). \quad (9.36)$$

Then the joint distribution of \mathbf{z}_a and \mathbf{z}_b is

$$p(\mathbf{z}_a, \mathbf{z}_b) = \mathcal{N}\left(\begin{bmatrix} \mathbf{z}_a \\ \mathbf{z}_b \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}_a \\ \mathbf{A}\boldsymbol{\mu}_a \end{bmatrix}, \mathbf{R}\right) \quad (9.37)$$

with

$$\mathbf{R} = \begin{bmatrix} \boldsymbol{\Sigma}_a & \boldsymbol{\Sigma}_a \mathbf{A}^\top \\ \mathbf{A} \boldsymbol{\Sigma}_a & \boldsymbol{\Sigma}_{b|a} + \mathbf{A} \boldsymbol{\Sigma}_a \mathbf{A}^\top \end{bmatrix} \quad (9.38)$$

Corollary 9.1 (Affine transformation – conditional) Assume that \mathbf{z}_a , as well as \mathbf{z}_b conditioned on \mathbf{z}_a , are Gaussian distributed according to

$$p(\mathbf{z}_a) = \mathcal{N}(\mathbf{z}_a; \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_a), \quad (9.39)$$

$$p(\mathbf{z}_b | \mathbf{z}_a) = \mathcal{N}(\mathbf{z}_b; \mathbf{A}\mathbf{z}_a + \mathbf{b}, \boldsymbol{\Sigma}_{b|a}). \quad (9.40)$$

Then the conditional distribution of \mathbf{z}_a given \mathbf{z}_b is

$$p(\mathbf{z}_a | \mathbf{z}_b) = \mathcal{N}(\mathbf{z}_a; \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b}), \quad (9.41)$$

with

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\Sigma}_{a|b} \left(\boldsymbol{\Sigma}_a^{-1} \boldsymbol{\mu}_a + \mathbf{A}^\top \boldsymbol{\Sigma}_{b|a}^{-1} (\mathbf{z}_b - \mathbf{b}) \right), \quad (9.42)$$

$$\boldsymbol{\Sigma}_{a|b} = \left(\boldsymbol{\Sigma}_a^{-1} + \mathbf{A}^\top \boldsymbol{\Sigma}_{b|a}^{-1} \mathbf{A} \right)^{-1}. \quad (9.43)$$

Corollary 9.2 (Affine transformation – Marginalization) Assume that \mathbf{z}_a , as well as \mathbf{z}_b conditioned on \mathbf{z}_a , are Gaussian distributed according to

$$p(\mathbf{z}_a) = \mathcal{N}(\mathbf{z}_a; \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_a), \quad (9.44)$$

$$p(\mathbf{z}_b | \mathbf{z}_a) = \mathcal{N}(\mathbf{z}_b; \mathbf{A}\mathbf{z}_a + \mathbf{b}, \boldsymbol{\Sigma}_{b|a}). \quad (9.45)$$

Then the marginal distribution of \mathbf{z}_b is then given by

$$p(\mathbf{z}_b) = \mathcal{N}(\mathbf{z}_b; \boldsymbol{\mu}_b, \boldsymbol{\Sigma}_b), \quad (9.46)$$

where

$$\boldsymbol{\mu}_b = \mathbf{A}\boldsymbol{\mu}_a + \mathbf{b}, \quad (9.47)$$

$$\boldsymbol{\Sigma}_b = \boldsymbol{\Sigma}_{b|a} + \mathbf{A} \boldsymbol{\Sigma}_a \mathbf{A}^\top. \quad (9.48)$$

10 User aspects of machine learning

Dealing with supervised machine learning problems in practice is to a great extent an engineering discipline where many practical issues have to be considered and where the scarce resource is in the end the available work-hours. To use this resource efficiently, we need to have a well-structured procedure for how to develop and improve the model. Multiple actions can potentially be taken. How do we know which action to take and is it really worth spending the time implementing them? Is it for example worth spending an extra week collecting and labeling more training data or should we do something else? These issues will be addressed in this chapter. Note that the layout of this chapter is thematic and does not represent the sequential order the different issues should be addressed.

10.1 Defining the machine learning problem

Solving a machine learning problem in practice is an iterative process. We train the model, evaluate the model, and from there suggest an action for improvement and train the model again, and so on. To do this efficiently, we need to be able to tell whether a new model is an improvement over the previous model or not. One way to evaluate the model after each iteration would be to put it in production (for example running a traffic-sign classifier in a self-driving car for a few hours). Besides the obvious safety issues, this would be very time inefficient as well as inaccurate since it could still be hard to tell whether the proposed change was an actual improvement or not.

A better strategy is to automate this evaluation procedure without the need to put the model in production each time we want to evaluate its performance. We do this by putting aside a *validation* dataset and *test* dataset and evaluate the performance using a single number *evaluation metric*. This will define the machine learning problem that we are solving.

Training, validation and test data

In Chapter 4 we introduced the strategy of splitting the available data into training data, validation data and test data.

- **Training data** is used for training the model.
- **Hold-out validation data** is used for comparing different model structures, choosing hyperparameters of the model, feature selection, and so on.
- **Test data** is used to evaluate the performance of the final model.

If the amount of available data is small, it is possible to perform k -fold cross-validation instead of putting aside hold-out validation data, the idea of how to use it in the iterative procedure is unchanged. To get a final estimate of the performance, test data is always needed.

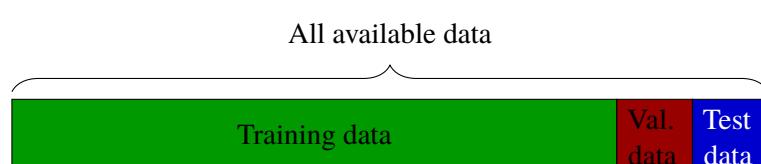


Figure 10.1: Splitting the data into training data, hold-out validation data and test data.

In the iterative procedure, the hold-out validation data (or k -fold cross-validation) is used to judge if the new model is an improvement over the previous model. This step is called *validation*. During the validation stage, we can also choose to train several new models and not just one new model. For example, if we have a neural network model and are interested in the number of hidden units to use in a certain layer, we can train several models, each with a different choice of hidden units. Afterwards, we pick the one that performs best on the validation data. In the next iteration we can use validation to choose the learning rate (by training multiple models using different learning rates), and so on.

Eventually, we will effectively have used the validation data to compare many models. Depending on the size of the validation data, we might risk picking a model that does particularly well on the validation data in comparison to completely unseen data. To detect this and to get a fair estimate of the actual performance of a model, we use the test data, which has neither been used during training nor validation. If the performance on the validation data is substantially better than the performance on the test data, we have overfitted on the validation data. The easiest solution in that case would be to extend the size of the validation data, if possible. Make sure not to use the test data too often. If you start making major decisions based on your test data, the test data will become your validation data and you can no longer trust that the performance on the test data is an objective indication of the actual performance of your model.

It is important that both the validation data and the test data always come from the same data distribution, namely the data distribution that we are expecting to see when we put the model into production. If they do not stem from the same distribution, we are validating and improving our model towards something that is not represented in the test data and hence "aiming for the wrong target".

Eventually, we should also evaluate the performance of the model in production. Also after deploying the model, we should monitor the performance. If we realize that the model is doing systematically worse in production than on the test data, the test data and validation data have to be updated such they actually represent what we expect to see at runtime. Preferably, also the training data should come from the same data distribution as the test and validation data, but this requirement can be relaxed if we have good reasons to do so, more about this in Section 10.2.

Size of validation and test data

How much data should we set aside as hold-out validation data and test data, or should we even avoid setting aside hold-out validation data and use k -fold cross-validation instead? This depends on how much data we have available, which performance difference we plan to detect, and how many models we plan to compare. For example, if we have a classification model with a 99.8% accuracy and want to know if a new model is even better, a validation dataset of 100 samples will not be able to tell that difference. Also, if we plan to compare many (say, hundreds or more) different hyperparameter values and model structures using 100 validation data points, we will most likely overfit to that validation data.

If we have, say, 500 data points, a reasonable split would be 60%-20%-20% (i.e. 300-100-100 samples) for training-validation-test. With such a small validation dataset, we cannot afford to compare several hyperparameter values and model structures or detecting an improvement in accuracy of 0.1%. In this situation, we are probably better off using k -fold cross-validation to decrease the risk of overfitting to the validation data. Be aware, however, that the risk of overfitting to the data still exists even with k -fold cross-validation. We also still need to set aside test data if we want a final unbiased estimate of the performance.

Many machine learning problems have substantially larger datasets. Assume we have a dataset of 1 000 000 data points. In this scenario it would be enough to use a split of 98%-1%-1%, i.e. leaving 10 000 data points for validation and test, respectively, unless we really care about the very last decimals in performance. Here, k -fold cross-validation is of less use, since having all 99% = 98% + 1% (training+validation) available for training would make a small difference in comparison to using "only" 98%. Also the price for training k models (instead of only one) with this amount of data would be much higher.

Another advantage of having a hold-out validation dataset is that we can allow the training data to come from a slightly different distribution than the validation and test dataset, for example if that would enable

us to find a much larger training dataset. We will discuss this more in Section 10.2.

Single number evaluation metric

In Section 4.5 we introduced additional metrics, besides the misclassification rate such as precision, recall and F1-score, for evaluating binary classifiers. There is no unique answer to which metric is the most appropriate for your. What metric to pick is rather a part of the problem definition. To improve you model quickly and in a more automated fashion, it is advisable to agree on a single number evaluation metric, especially if you are a larger team of engineers working on the problem. Such a metric, together with the validation and test data, defines the problem, and without a proper problem definition it is difficult to solve that problem and to improve its solution. If we during the development process realize that the metric does not favor the properties we want a good model to have, we can always change that metric later.

The single number evaluation metric together with the validation data defines the supervised machine learning problem. Having an efficient procedure in place where we can evaluate the model on the hold-out validation data (or by k -fold cross-validation) using the metric, allows us to speed up the iterations since we quickly can see if a proposed change to the model improves the performance or not. This is important in order to manage an efficient workflow of trying out and accepting or rejecting new models.

10.2 Improving a machine learning model

As already mentioned, solving a machine learning problem is an iterative procedure where we train, evaluate, and suggest action for improvement, for instance by changing some hyperparameters or trying another methods. How do we start this iterative procedure?

Try simple things first

A good strategy is to *try simple things first*, for example k -NN or linear/logistic regression. Trying a simple thing first allows us to start early with the iterative procedure of finding a good model. This is important, since it might reveal important aspects of the problem formulation that we need to re-think before it makes sense to proceed with more complicated models. It also reduces the risk of ending up with a too complicated model when a much simpler model would be just as good (or even better).

There are many actions that could be taken to improve the model, for example changing the type of model, increasing/decreasing model complexity, changing input variables, collecting more data, starting correcting mislabeled data (if there are any) etc. What should we do next? Two possible strategies for guiding us to meaningful actions to improve the solution are by trading *training error and generalization gap*, or by applying *error analysis*.

Training error vs generalization gap

With the notation from Chapter 4, E_{train} is the performance of the model on training data and $E_{\text{hold-out}}$ the performance on hold-out validation data. In the validation step, we are interested in changing the model such that $E_{\text{hold-out}}$ is minimized. We can write the hold-out validation error as a sum of the training error and the generalization gap as

$$E_{\text{hold-out}} = E_{\text{train}} + \underbrace{(E_{\text{hold-out}} - E_{\text{train}})}_{\approx \text{generalization gap}}. \quad (10.1)$$

In words, the generalization gap is the difference between the validation error $E_{\text{hold-out}}$ and the training error E_{train} .¹

¹This can be related to (4.11), if approximating $\bar{E}_{\text{train}} \approx E_{\text{train}}$ and $\bar{E}_{\text{new}} \approx E_{\text{hold-out}}$. If using k -fold cross validation, $\bar{E}_{\text{new}} \approx E_{k\text{-fold}}$ would be an equally good approximation when computing the generalization gap.

Once the validation step is completed ($E_{\text{hold-out}}$ or $E_{k\text{-fold}}$ computed), we can easily compute the training error E_{train} and the generalization gap $E_{\text{hold-out}} - E_{\text{train}}$. By computing these quantities, we can actually get good guidance for what changes we may consider for the next iteration.

As we discussed in Chapter 4, if the training error is small and the generalization gap is big (E_{train} small, $E_{\text{hold-out}}$ big), we have typically overfitted the model. The opposite situation, big training error and small generalization gap (both E_{train} and $E_{\text{hold-out}}$ big), typically indicates underfitting.

If we want to reduce the generalization gap $E_{\text{hold-out}} - E_{\text{train}}$ (reduce overfitting), the following actions can be explored:

- *Use a less flexible model.* If we have a very flexible model we might start overfitting to the training data, i.e. that E_{train} is much smaller than $E_{\text{hold-out}}$. If we use a less flexible model, we also reduce this gap.
- *Use more regularization.* Using more regularization will reduce the flexibility of the model, and hence also reduce the generalization gap.
- *Use bagging*, or use more ensemble members if we already are using it. Bagging is a method for reducing the variance of the model, which typically also means that we reduce the generalization gap.
- *Collect more training data.* If we collect more training data, the model is less prone to overfit that extended training dataset and is forced to only focus on aspects which generalizes to the validation data.
- *Early stopping* For models that are trained iteratively we can stop the training before reaching the minimum. One good practice is to monitor $E_{\text{hold-out}}$ during training and stop if it starts increasing.

If we want to reduce the training error E_{train} (reduce underfitting), the following actions can be considered:

- *Use a more flexible model* that is able to fit the training data better. This can be to change a hyperparameter in the model we are considering, for example decreasing k in k -NN or changing the model to a more flexible one, for example change the linear regression model to deep neural network.
- *Extend the set of input variables.* If we suspect that there are more input variables that carry information, we might want to extend the data with these input variables.
- *Use less regularization.* Can of course only be applied if regularization is used at all.
- *Train the model longer* For models that are trained iteratively we can reduce E_{train} by training longer.

We summarize the above discussion in Figure 10.2. Fortunately, evaluating E_{train} and $E_{\text{hold-out}}$ is cheap. We only have to evaluate the model on the training data and the validation data, respectively. Yet, it gives us good advice on what actions to take next. Besides suggesting what action to explore next, this procedure also tells us what *not* to do: If $E_{\text{train}} \gg E_{\text{hold-out}} - E_{\text{train}}$, collecting more training data will most likely not help. Furthermore, if $E_{\text{train}} \ll E_{\text{hold-out}} - E_{\text{train}}$ a more flexible model will most likely not help.

Error analysis

Another strategy to identify actions that can improve the model is to perform *error analysis*. Below we only describe error analysis for classification problems, but the same strategy can be applied to regression problems as well.

In error analysis we manually look at a subset, say 100 data points, of the validation data that the model classified incorrectly. Such an analysis does not take much time, but might give valuable clues as to what type of data the model is struggling with and how much improvement we can expect by fixing these issues. We illustrate the procedure with an example.

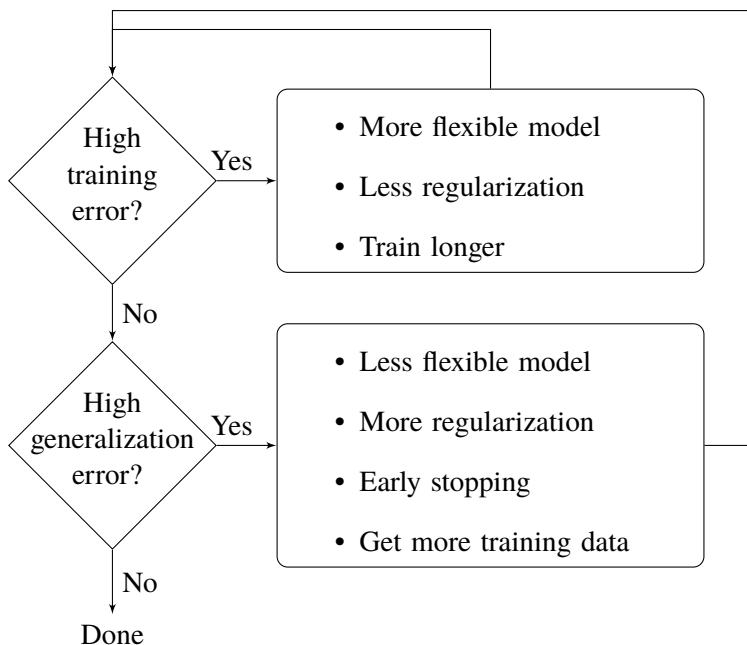


Figure 10.2: The iterative procedure of improving the model based on the decomposition into on the training error and generalization gap .

Example 10.1: Error analysis applied to vehicle detection

Consider a classification problem of detecting cars, bicycles and pedestrians in an image. The model takes an image as input and one of the four classes **car**, **bike**, **pedestrian**, or **other** as output. Assume that the model has a classification accuracy of 90% on validation data.

When looking at a subset of 100 images that were misclassified in the validation data, we make the following observations:

- All 10 images of class **pedestrian** that were incorrectly classified as **bike** contained a baby carrier.
- 30 images were substantially tilted.
- 15 images were mislabeled.

From this observation we can conclude:

- If we launch a project for improving the model to classify pedestrians with baby carriers as **pedestrian** and not incorrectly as **bike**, an improvement of at most a $\sim 1\%$ (a tenth of the 10% classification error rate) can be expected.
- If we improve the performance on tilted images, an improvement of at most $\sim 3\%$ can be expected.
- If we correct all mislabeled data, an improvement of at most $\sim 1.5\%$ can be expected.

Following the example, we get an indication on what improvement we can expect by tackling these three issues. These numbers should be considered as the maximal possible improvement. To prioritize which aspect to focus on, we should also consider what strategies are available for improving them, how much progress we expect to make applying these strategies and how much effort we have to invest fixing these issues.

For example, to improve the performance on tilted images we could try to extend the training data by augmenting it with more tilted images. This strategy could be investigated without too much extra effort by augmenting the training data with tilted versions of the training data point that we already have. Since, this could be applied fairly quickly and have a maximal performance increase of 3%, it seems to be a good thing to try out.

To improve the performance on the images with baby carriers, one approach would be to collect more

images of pedestrians with baby carriers. This obviously requires some more manual work and can be questioned if it is worth the effort since it would only give performance improvement of at most 1%.

Regarding the mislabeled data, the obvious actions to take to improve on this issue is to manually go through the data and correct these labels. In the example above we may say it is not quite worth the effort of getting an improvement of 1.5%. However, assume that we have improved the model with other actions to an accuracy of 98.0% on validation data and that still 1.5% of the total error is due to mislabeled data, this issue is now quite relevant to address if we want to improve the model further. Remember, the purpose of the validation data is to choose between different models. This purpose is degraded when the majority of the reported error on validation is due to incorrectly labeled data rather than the actual performance of the model.

There are three levels of ambitions for correcting the labels:

1. Go through the data points in the validation/test data that were misclassified by the best algorithm and correct these labels.
2. Go through all data points in the validation/test data and correct the labels.
3. Go through all data points, including the training data and correct the labels.

Approach 1 is in general not recommended since there could be mislabeled data points that the algorithm also classified to belong to that mislabeled class. These data points would then not be corrected, resulting in a too optimistic estimate of the performance on validation and test data. Therefore, it is safer to go through all data points in the validation and test data as suggested in approach 2. The advantage of approach 1, in comparison to approach 2, is the less amount of work it requires. If we have a model with 98% accuracy, there are 50 times less data to process in comparison to approach 2. Also, note that correcting labels in test and validation data and test data only does not necessarily increase the performance of model in production, but it will give us a more fair estimate of the actual performance of the model.

An even more ambitious approach is to go through all mislabeled data in the training data as well, as suggested in approach 3. This is also substantially more labor intensive. Assume, for example, that we have made a 98% – 1% – 1% split of training-validation-test data. Then it is yet another factor 50 more time consuming in comparison with approach 2. Unless the mislabeling is systematic, correcting the labels in the training data does not necessarily pay off.

Applying the data cleaning to validation and test data only, as suggested in approach 2, will result in the training data coming from a slightly different distribution than the validation and test data. However, if we are eager to correct the mislabeled data in the training data as well, a good recommendation would still be to start correcting validation and test data only and then use the techniques in the following section to see how much extra performance we can expect by cleaning data in the training data as well before launching that substantially more labor intensive data cleaning project.

Mismatched training and validation/test data

As already pointed out in Chapter 4, we should strive for letting the training data come from the same distribution as the validation and test data. However, there are situations where we, for different reasons, can accept the training data to come from a slightly different distribution than the validation and test data. One reason was presented in the previous section where we choose to correct mislabeled data in validation and test data, but not necessarily invest the time to do the same correction to the training data.

Another reason is that we might have access to another substantially larger data set which comes from a slightly different distribution than the data we care about, but similar enough that the advantage of having a larger training data overcomes the disadvantage of that data mismatch. This scenario is further elaborated in Section 10.3.

If we have a mismatch between training data and validation/test data, that mismatch contributes to yet another error source of the final validation error $E_{\text{hold-out}}$ that we care about. We want to estimate the magnitude of that error source. This can be done by revising the training-validation-test data split. From

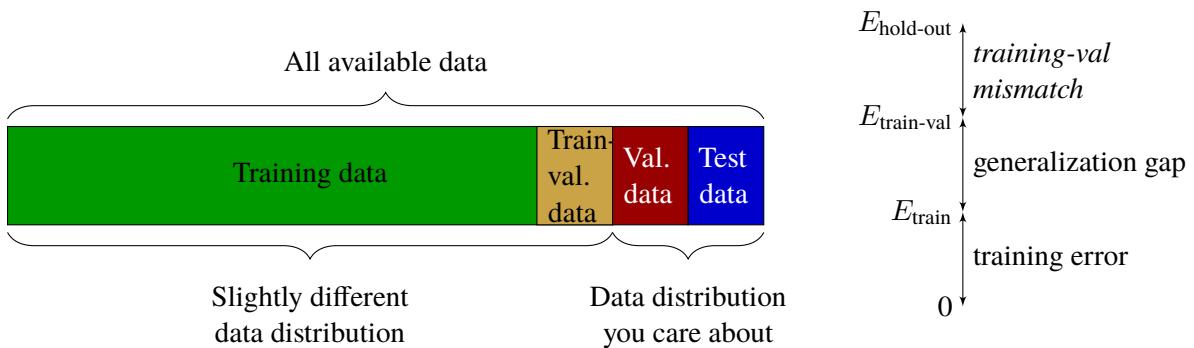


Figure 10.3: Revising the training-validation-test data split by carving out a separate *train-validation* dataset, see Figure 10.3. That dataset is neither used for training nor for validation. However, we do evaluate the performance of our model on that dataset as well. As before, the remaining part of the training data is used for training, the validation data is used for comparing different model structures, and test data is used for evaluating the final performance of the model.

This revised data split also allows us to revise the decomposition in (10.1)

$$E_{\text{hold-out}} = E_{\text{train}} + \underbrace{(E_{\text{train-val}} - E_{\text{train}})}_{\approx \text{generalization gap}} + \underbrace{(E_{\text{hold-out}} - E_{\text{train-val}})}_{\approx \text{Train-val mismatch}}, \quad (10.2)$$

where $E_{\text{train-val}}$ is the performance of the model on the new training-validation data and where, as before, $E_{\text{hold-out}}$ and E_{train} are the performances on validation and training data, respectively. With this new decomposition, the term $E_{\text{train-val}} - E_{\text{train}}$ is an approximation of the generalization gap, i.e. how well the model generalizes to unseen data *of the same distribution* as the training data, whereas the term $E_{\text{hold-out}} - E_{\text{train-val}}$ is the error related to the training-validation data mismatch. If the term $E_{\text{hold-out}} - E_{\text{train-val}}$ is small in comparison to the other two terms, it seems likely that the training-validation data mismatch is not a big problem and that it is better to focus on techniques reducing the other training error and the generalization gap as we talked about earlier. On the other hand, if $E_{\text{hold-out}} - E_{\text{train-val}}$ is significant, the data mismatch does have an impact and it might be worth investing time reducing that term. For example, if the mismatch is caused by the fact that we only corrected labels in the validation and test data, we might want to consider correcting labels for the training data as well.

10.3 What if we cannot collect more data?

We have seen in Section 10.2 that collecting more data is a good strategy to reduce the generalization gap and hence reduce overfitting. However, collecting labeled data is usually expensive and sometimes not even possible. What can we do if we cannot afford to collect more data but still want to benefit from the advantages that a larger dataset would give? In this section a few approaches are presented.

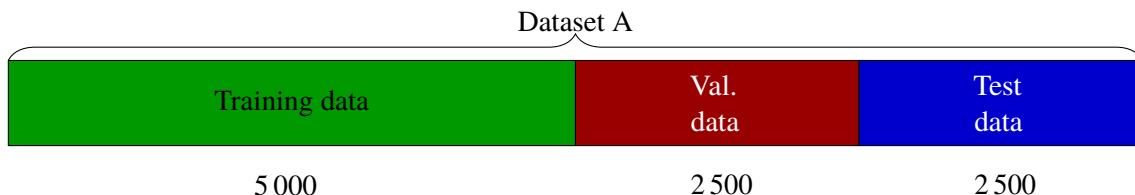
Extending the training data with slightly different data

As already mentioned, there are situations where we can accept the training data to come from a slightly different distribution than the validation and test data. One reason to accept this is if we then would have access to a substantially larger training data set.

Consider a problem with 10 000 data points representing the data that we also would expect to get when the model is deployed in production. We call this dataset A. We also have another dataset with 200 000 data points that come from a slightly different distribution, but which is similar enough that we think

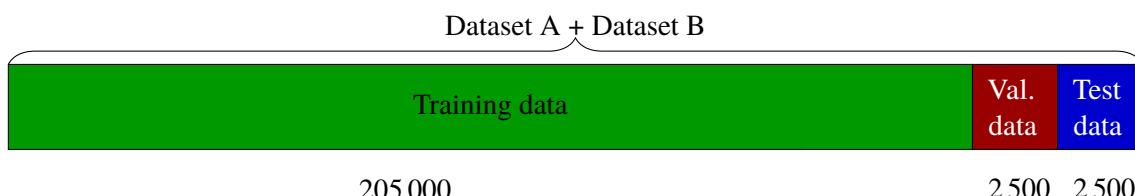
exploiting information from that data can improve the model. We call this dataset B. Some options to proceed would be the following:

- **Option 1** Use only dataset A and split it into training, validation, and test data.



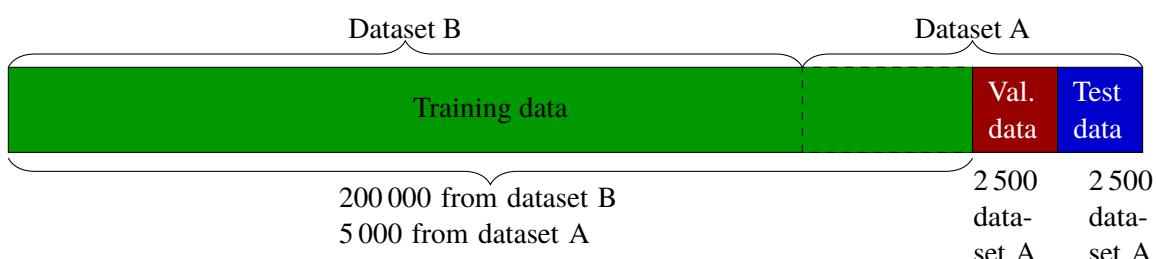
The advantage of this option is that we only train, validate, and evaluate on dataset A, which is also the type of data that we want our model to perform well on. The disadvantage is that we have quite few data points and we do not exploit potentially useful information in the larger dataset B.

- **Option 2** Use both dataset A and dataset B. Randomly shuffle the data and split it in training, validation and test data.



The advantage over option 1 is that we have a lot more data available for training. However, the disadvantage is that we mainly evaluate the model on data from dataset B, whereas we want our model to perform well on data from dataset A.

- **Option 3** Use both dataset A and dataset B. Use data points from dataset A for validation data and test data and some in the training data. Dataset B only goes into the training data.



Similar to option 2, the advantage is that we have more training data in comparison to option 1 and in contrast to option 2 we now evaluate the model on data from dataset A, which is the data we want our model to perform well on. However, one disadvantage is that the training data no longer have the same distribution as the validation and test data.

From these three options we would strongly recommend option 3. We do exploit the information available in the much larger dataset B, but evaluate only on the data where we want the model to perform well on (dataset A). The main disadvantage with option 3 is that the training data no longer come from the same distribution as the validation data and test data. In order to quantify how big impact this mismatch has on the final performance, the techniques described in Section 10.2 can be used. To push the model to do better on data from dataset A during training, we can also consider giving data from dataset A higher weight in the cost function than data from dataset B.

Artificially extending the training data set

Data augmentation is another approach to extend the dataset without the need to collect more data. In data augmentation we construct new data points duplicating the existing data by applying invariant

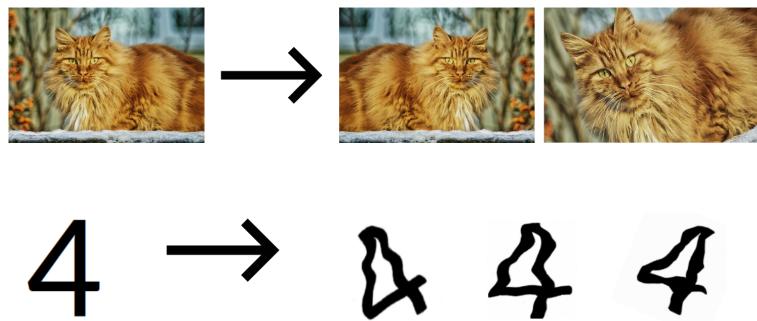


Figure 10.4: Different examples of data augmentation applied to images. Above: An image of a cat has been reproduced by tilting and vertical flipping. Below: An image of a digit has been reproduced by tilting and blurring.

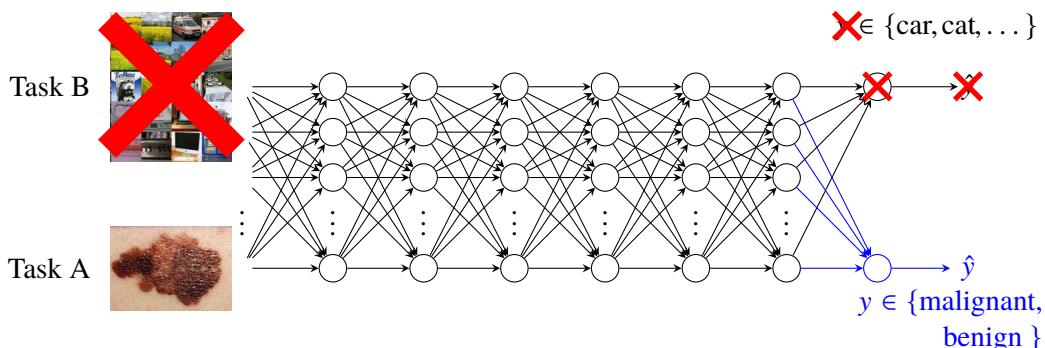


Figure 10.5: In transfer learning we reuse models that have been trained on a different task than the one we are interested in. Here we reuse a model which has been trained on images displaying all sorts of classes, such as cars, cats, computers and later train only the last few layers on the skin cancer data which is the task we are interested in.

transformations. This is especially common for images where such invariant transformations can be scaling, rotation, and vertical flipping of the images. For example, if we vertically flip an image of a cat, it still displays a cat, see the examples Figure 10.4.

Transfer learning

Yet another technique to effectively use more data than the dataset we have is by applying transfer learning. In transfer learning we use the knowledge from a model that has been trained on a different task with a different dataset and then apply that model in solving a different, but slightly related, problem.

Transfer learning is especially common for sequential model structures such as the neural network models introduced in Chapter 6. Consider an application where we want to detect whether a certain type of skin cancer is malignant or benign, and for this task we have 100 000 labeled images of skin cancer. We call this task A. Instead of training the full neural network from scratch on this data, we can reuse an already pre-trained network from another image classification task (task B), which preferably has been trained on a much larger dataset, not necessarily containing images even resembling skin cancer tumors. By using the weights from the model trained in task B and only train the last few layers on the data in task A, we can get a better model than if the whole model had been trained on just dataset A. The procedure is also displayed in Figure 10.5. The intuition is that the layers closer to the input accomplish tasks that are generic for all types of images, such as extracting lines, edges and corners in the image, whereas the layers closer to the output are more specific to the particular problem.

In order for transfer learning to be applicable, we need the two tasks to have the same type of input (in the example above, images of the same dimension). Further, for transfer learning to be an attractive option, the task that we transfer from should have been trained on substantially more data than the task we transfer to.

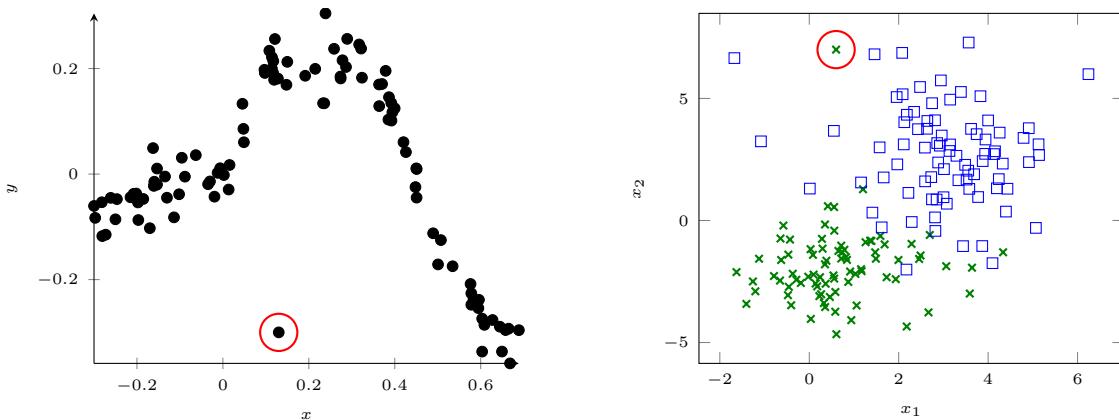


Figure 10.6: Two typical examples of outliers (marked with red circle) in regression (left) and classification (right), respectively.

10.4 Practical data issues

Besides the amount and distribution of data, a machine learning engineer may also face other data issues. In this section we will discuss some of the most common ones; outliers, missing data and if some features can be removed.

Outliers

In some applications, a common issue is *outliers*, meaning data points whose outputs do not follow the overall pattern. Two typical examples of outliers are sketched in Figure 10.6. Even though the situation in Figure 10.6 looks simple, it can be quite hard to find outliers when the data has more dimensions and is harder to visualize. The error analysis discussed in Section 10.2, which amounts to inspecting misclassified samples in the validation data, is a systematic way to discover outliers.

When facing a problem with outliers, the first question to ask is whether the outliers are meant to be captured by the model or not. Do the encircled data points in Figure 10.6 describe an interesting phenomena that we would like to predict, or are they irrelevant noise (possibly originating from a poor data collection process)? The answer to this question depends on the actual problem and ambition. Since outliers by definition (no matter their origin) do not follow the overall pattern, they are typically hard to predict.

If the outliers are not of any interest, the first is to consult the data provider and identify the reason for the outliers and if something could be changed in the data collection to avoid these outliers, for example replacing a malfunction sensor. If the outliers are unavoidable, there are basically two approaches one could take. The first approach is to simply delete (or replace) the outliers in the data. Unfortunately this means that one has to first find the outliers, which can be hard, but sometimes some thresholding and manual inspection (that is, look at all data points whose output value is smaller/larger than some value) can help. Once the outliers are removed from the data, one can proceed as usual. The second approach is to instead make sure that the learning algorithm is robust against outliers, for example by using a robust loss function such as absolute error instead of squared error loss (see Chapter 5 for more details). Making a model more robust amounts, to some extent, to making it less flexible. However, robustness amounts to making the model less flexible in a particular way, namely by putting less emphasis on the data points whose predictions are severely wrong.

If the outliers are of interest to the prediction, they are not really an issue, but rather a challenge. We have to use a model that is flexible enough to capture the behavior (small bias). This has to be done with care, since very flexible models have a high risk of overfitting also to noise. If it turns out that the outliers in a classification problem indeed are interesting and in fact are from an underrepresented class, we are rather facing an imbalanced problem, see Section 4.5.

Missing data

A common practical issue is that certain values are sporadically missing in the data. Throughout this book so far the data has always consisted of complete input-output pairs $\{\mathbf{x}_i, y_i\}_{i=1}^n$, and *missing data* refers to the situation where some (or a few) values from either the input \mathbf{x}_i or the output y_i , for some i , is missing. If the output y_i is missing, we can also refer to it as *unlabeled* data. It is a common practice to denote missing data in a computer with NaN (not a number), but less obvious codings also exists, such as 0. Reasons for missing data could for instance be a malfunctioning sensor or similar issues at data collection time, or that certain values for some reason have been discarded during the data processing.

As for outliers, a sensible first option is to figure out the reason for the missing data. By going back to the data provider this issue could potentially be fixed and the missing data recovered. If this is not possible, there is no universal solution for how to handle missing data. There is, however, some common practice which can serve as a guideline. First of all, if the output y_i would be missing, the data point is useless for supervised machine learning,² and can be discarded. In the following, we assume that the missing values are only in the input \mathbf{x}_i .

The easiest way to handle missing data is to discard the entire data points (“rows in \mathbf{X} ”) where data is missing. That is, if some feature is missing in \mathbf{x}_i , the entire input sample \mathbf{x}_i and its corresponding output value y_i is discarded from the data, and we are left with a smaller dataset. If the dataset that remains after this procedure still contains enough data, this approach can work well. However, if this would lead to a too small dataset, it is of course problematic. More subtle, but also important, is the situation when the data is missing in a systematic fashion, for example that missing data is more common for a certain class. In such a situation, discarding data points with missing data would lead to a mismatch between reality and training data, which may degrade the performance of the learned model further.

If missing data is common, but only for certain features, another easy option is to not use those features (“column of \mathbf{X} ”) which are suffering from missing data. It depends on the situation whether this is a fruitful approach or not.

Instead of discarding the missing data, it is possible to impute (fill in) the missing values using some heuristics. Say, for example, that the j th feature x_j is missing from sample \mathbf{x}_i . A simple imputation strategy would be to take the mean or median of x_j for all other samples (where it is not missing), or the mean or median of x_j for all samples of the same class (if it is a classification problem). It is also possible to come up with more complicated imputation strategies, but each imputation strategy implies some assumptions about the problem. Those assumptions might or might not be fulfilled, and it is hard to guarantee that imputation will help the performance in the end. A poor imputation can even degrade the performance compared to just discarding the missing data.

Some methods are actually able to handle missing data to some extent (which we have not discussed in this book), but under rather restrictive assumptions. Such an example is that the data is “completely missing at random”, meaning that which data that is missing is completely uncorrelated to what value it, and other features and the output, would have had, had it not been missing. Assumptions like these are very strong and rarely met in practice, and the performance can be severely degraded if those assumptions are not fulfilled.

Feature selection

When working with a supervised machine learning problem, the question whether all available input variables/features contribute to the performance is often relevant. Removing the right feature is indeed a type of regularization that possibly can reduce overfit and improve the performance, and the data collection might be simplified if a certain variable do not even have to be collected. Selecting between the available features is an important task for the machine learning engineer.

The connection between regularization and feature selection becomes clear by considering L^1 regularization. Since the main feature of L^1 regularization is that the learned parameter vector $\hat{\boldsymbol{\theta}}$ is sparse, it

²The “partly labeled data” problem is a semi-supervised problem, which is introduced in Chapter 11 but not covered in depth by this book.

effectively removes the influence of certain features. If using a model where L^1 regularization is possible, we can study $\hat{\theta}$ to see which features we simply can remove from the data set. However, if we can not or prefer not to use L^1 regularization, we have to use a more manual approach to feature selection.

Remember that our overall goal is to obtain a small new data error E_{new} , which we for most methods estimate using cross-validation. We can therefore always use cross-validation to tell whether we gain, or lose, by including a certain feature in \mathbf{x} . Depending on the amount of data, evaluating all possible combinations of removed features might not be a good idea, either due to computational aspects or the risk of overfit. There are, however, some rule of thumbs that possibly can give us some guidance for which features we should investigate closer whether they contribute to the performance or not.

To get a feeling for the different features, we can look at the correlation between each feature and the output, and thereby get a clue about which features that might be more informative about the output. If there is little correlation between a feature and the output, it is possibly a useless feature and we could investigate further if we can remove them. However, looking only at one feature at a time can be misleading and there exists cases where this would lead to the wrong conclusion, for example the case in Example 8.1.

Another approach is to explore if there are redundant features, with the reasoning that having two features that (essentially) contain the same information will lead to an increased variance compared to having only one feature with the same information. Based on this argument one may look at the pairwise correlation between the features, and investigate removing features that have high correlation to other features. This approach is somewhat related to PCA (Chapter 11).

10.5 Can I trust my machine learning model?

Supervised machine learning presents a powerful family of all-purpose general black-box methods, and has demonstrated impressive performance in many applications. The main argument for supervised machine learning is, frankly, that it works well empirically. However, depending on the requirements by the applications, supervised machine learning also has a potential shortcoming, in that it relies on “repeating patterns seen in training data” rather than “deduction from a set of carefully written rules”.

Understanding why a certain prediction was made

In some applications there might be an interest to “understand” why a certain prediction was made by a supervised machine learning model, for example in medicine or law. Unfortunately, the underlying design philosophy in machine learning is to deliver good predictions rather than explaining them.

With simpler model, like the ones in ??-Chapter 3, it can to some degree be possible for an engineer to inspect the learned model and explain the “reasoning” behind it for a non-expert. For the more complicated models it can be a rather hard task.

There are however methods on the research forefront, and the situation may look different in the future. A related topic is that of so-called adversarial examples, which essentially amounts to finding an input \mathbf{x}' which is as close as possible to \mathbf{x} but still gives another prediction. In the image classification setting, it can for example be the problem of having a picture of a car being predicted as a dog by only changing a few pixel values.

Worst case guarantees

In the view of this book, a supervised machine learning model is good if it attains a small E_{new} . It is, however, important to remember that E_{new} is a *statistical* claim, under the assumption that the training and/or test data resembles the reality which the model will face once it is put in production. And even if that non-trivial assumption is satisfied, there are no claims about how bad the model will predict in the worst individual cases. This is indeed a shortcoming of supervised machine learning, and potentially also a show-stopper for some applications.

Simpler and more interpretable models, like logistic regression and trees for example, can be inspected manually in order to deduce the “worst case” that could happen. By looking at the leaf nodes in a

regression tree, as an example, it is possible to give an interval to which all predictions will belong. With more complicated models, like random forests and deep learning, it is very hard to give any worst case guarantees about how the wrong the model can be in its predictions when faced with some particular input. However, an extensive testing scheme might reveal some of the potential issues.

10.6 Further reading

User aspects of machine learning is a fairly under-explored area both in academic research publications and in standard textbooks on machine learning. Two exceptions are Ng (2019) and Burkov (2020) from which parts of this chapter has been inspired. Connected to data augmentation, see further in Shorten and Khoshgoftaar (2019) for a review of different data augmenting techniques for images.

11 Generative models and learning from unlabeled data

The models introduced so far in this book are so-called *discriminative* models, meaning (in the context of supervised machine learning) that they are designed to learn from data with the aim of predicting new outputs. We will in the first half of this chapter introduce models of another paradigm, namely so-called *generative* models. Generative models are indeed also learned from data, but their scope is wider. In contrast to discriminative models, a generative model is for instance able to simulate more data, and predictions are merely obtained as a by-product. Generative modeling is therefore a natural way to take us beyond supervised learning, which we will do in the second half of this chapter.

A generative model aims to describe the distribution $p(\mathbf{x}, y)$, that is, how the data (both its input and its output) is generated. Perhaps we should write $p(\mathbf{x}, y | \theta)$ to emphasize that generative models also contain some parameters that we will learn from data, but to ease the notation we settle for $p(\mathbf{x}, y)$. To make a prediction with a generative model, the expression for $p(y | \mathbf{x})$ has to be derived from $p(\mathbf{x}, y)$ using probability theory. We will make this idea concrete by considering the rather simple, yet useful, generative Gaussian mixture model. As any generative model also the Gaussian mixture model can be used for different purposes. When used for prediction it results in a method traditionally called linear or quadratic discriminant analysis (LDA or QDA, respectively). We will thereafter see how the Gaussian mixture model also can be used for semi-supervised learning (where labels y are partly missing) and unsupervised learning (where no labels are present at all; there are only \mathbf{x} and no y) as well. In the latter case, the Gaussian mixture model can be used for solving the so-called clustering problem where similar \mathbf{x} are to be grouped together.

Generative models bridge the gap between supervised and unsupervised machine learning, but not all methods for unsupervised learning comes from generative models. We therefore finally also review two other popular unsupervised methods (not derived from generative models), namely k -means and principal component analysis (PCA). k -means is another clustering method, whereas the purpose of PCA is to infer which dimensions of \mathbf{x} that are most informative, so-called dimensionality reduction.

11.1 The Gaussian mixture model and the LDA & QDA classifiers

We will now introduce a generative model, the Gaussian mixture model, from which we will derive several methods for different purposes. The Gaussian mixture model attempts to model $p(\mathbf{x}, y)$, that is, the *joint* distribution between inputs \mathbf{x} and outputs y . (The discriminative models in previous chapters only attempt to model the the conditional distribution $p(y | \mathbf{x})$, a less ambitious problem since $p(y | \mathbf{x})$ can be derived from $p(\mathbf{x}, y)$ but not vice versa.) For the Gaussian mixture model, we assume that \mathbf{x} is a numerical and y a categorical variable.

The Gaussian mixture model

The Gaussian mixture model makes use of the factorization

$$p(\mathbf{x}, y) = p(\mathbf{x} | y)p(y), \quad (11.1a)$$

where $p(\mathbf{x} | y)$ is *assumed* to be a Gaussian distribution

$$p(\mathbf{x} | y) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) \quad (11.1b)$$

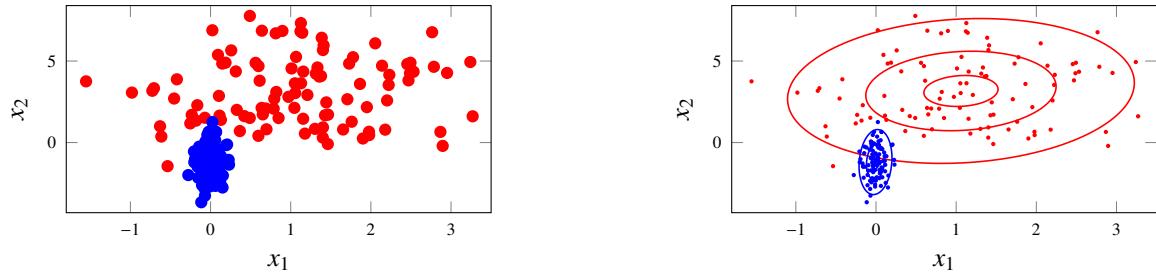


Figure 11.1: The Gaussian mixture model is a generative model, and we think about the input variables \mathbf{x} as random and *assume* that they have a certain distribution. The Gaussian mixture model assumes that $p(\mathbf{x} | y)$ has a Gaussian distribution for each y . In this figure \mathbf{x} is two-dimensional and there are two classes y (red and blue). The left panel shows some data with this nature. The right panel shows, for each value of y , the contour lines of the Gaussians $p(\mathbf{x} | y)$ that are learned from the data using (11.3).

where μ_y and Σ_y depends on y . Since y is categorical, and thereby takes values $1, \dots, M$, the distribution $p(y)$ is modeled as a categorical distribution with M parameters $\{\pi_m\}_{m=1}^M$ as

$$\begin{aligned} p(y = 1) &= \pi_1, \\ &\vdots \\ p(y = M) &= \pi_M. \end{aligned} \tag{11.1c}$$

In words the model (11.1) assumes a categorical distribution for y and, for each possible value of y , a Gaussian distribution for \mathbf{x} . If we look at \mathbf{x} , as we do in Figure 11.1, the model corresponds to a mixture of Gaussians (one component for each value of y), and hence its name. Altogether (11.1) is a generative model which describes an *assumption* on how data (\mathbf{x}, y) is generated.

The Gaussian mixture model (11.1) will lead us to classifiers (for supervised learning) which are easy to learn (requires no numerical optimization, in contrast to logistic regression) and useful in practice also when the data do not obey the Gaussian assumption (11.1b) perfectly. The classifiers are for historical reasons called linear and quadratic discriminant analysis (LDA¹ and QDA, respectively). It will also lead us also to a clustering method for unsupervised learning as well as being useful also for the in-between problem when the output label y is missing for some of the data points.

Supervised learning of the Gaussian mixture model

Like any machine learning model the Gaussian mixture model (11.1) is learned from training data. The unknown parameters to be learned are $\theta = \{\mu_m, \Sigma_m, \pi_m\}_{m=1}^M$. We start by the supervised case, meaning that the training data does contain inputs \mathbf{x} and outputs y as $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ (which has been the case for all other methods in this book so far). We will, however, later see how we can learn the Gaussian mixture model also when the output y is missing.

Mathematically we learn the Gaussian mixture model by maximizing the likelihood

$$\widehat{\theta} = \arg \max_{\theta} p(\underbrace{\{\mathbf{x}_i, y_i\}_{i=1}^n}_{\mathcal{T}} | \theta). \tag{11.2}$$

Alternatively it is possible to learn the Gaussian mixture model also using the Bayesian approach, although we do not pursue that any further here. The optimization problem (11.2) turns out to have the closed form solution

$$\widehat{\pi}_m = \frac{n_m}{n}, \tag{11.3a}$$

¹Note to be confused with Latent Dirichlet Allocation, also abbreviated LDA, which is a completely different method.

where n_m is the number of training data samples in class m . (Consequently, all n_m must sum to n , and thereby $\sum_m \widehat{\pi}_m = 1$.) Furthermore, the mean vector $\boldsymbol{\mu}_m$ of each class is learned as

$$\widehat{\boldsymbol{\mu}}_m = \frac{1}{n_m} \sum_{i:y_i=m} \mathbf{x}_i, \quad (11.3b)$$

the empirical mean among all training samples of class m , and the covariance matrix $\boldsymbol{\Sigma}_m$ of each class $m = 1, \dots, M$, is learned as

$$\widehat{\boldsymbol{\Sigma}}_m = \frac{1}{n_m} \sum_{i:y_i=m} (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_m)(\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_m)^\top. \quad (11.3c)$$

Note that we can compute the learned parameters $\widehat{\boldsymbol{\theta}}$ no matter if the data comes from a distribution where $p(\mathbf{x} | y)$ is Gaussian or not. In fact (11.3b)-(11.3c) learns a Gaussian distribution for \mathbf{x} for each class such that the mean and covariance fits the data (so-called moment-matching).

Predicting output labels for new inputs: QDA

We have so far described the generative Gaussian mixture model $p(\mathbf{x}, y)$, where \mathbf{x} is numerical and y categorical, and how to learn the unknown parameters in $p(\mathbf{x}, y)$ from training data. We will now see how this can be used as a classifier for supervised machine learning.

The key insight for using a generative model $p(\mathbf{x}, y)$ to make predictions is to realize that a prediction “only” amounts to compute $p(y | \mathbf{x})$. From probability theory we have

$$p(y | \mathbf{x}) = \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} = \frac{p(\mathbf{x}, y)}{\sum_{m=1}^M p(\mathbf{x}, m)}. \quad (11.4)$$

The left hand side $p(y | \mathbf{x})$ is the prediction whereas all expressions on the right hand side are defined by the generative Gaussian mixture model (11.1). We therefore get the classifier

$$p(y = m | \mathbf{x}_*) = \frac{\widehat{\pi}_m \mathcal{N}(\mathbf{x}_* | \widehat{\boldsymbol{\mu}}_m, \widehat{\boldsymbol{\Sigma}}_m)}{\sum_{j=1}^M \widehat{\pi}_j \mathcal{N}(\mathbf{x}_* | \widehat{\boldsymbol{\mu}}_j, \widehat{\boldsymbol{\Sigma}}_j)}, \quad (11.5)$$

which historically has been called quadratic discriminant analysis (QDA). As usual, we can obtain “hard” predictions \widehat{y}_* by selecting the class which is predicted to be the most probable,

$$\widehat{y}_* = \arg \max_m p(y = m | \mathbf{x}_*), \quad (11.6)$$

and compute corresponding decision boundaries. It turns out that the decision boundary for QDA is always a quadratic function (hence its name). We summarize this by Algorithm 18 and Figure 11.3, and in Figure 11.2 we show the decision boundary when the Gaussian mixture model from Figure 11.1 is turned into a QDA classifier.

It is possible to make the restriction that the covariance matrix is equal for all classes, $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \dots = \boldsymbol{\Sigma}_M = \boldsymbol{\Sigma}$ in (11.1b). With that restriction (11.3c) is replaced with

$$\widehat{\boldsymbol{\Sigma}} = \frac{1}{n - M} \sum_{m=1}^M \sum_{i:y_i=m} (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_m)(\mathbf{x}_i - \widehat{\boldsymbol{\mu}}_m)^\top. \quad (11.7)$$

Using (11.7) in (11.5) leads to the so-called linear discriminant analysis (LDA), a classifier with linear decision boundaries. Note that LDA is obtained by replacing (11.3c) with (11.7) in Algorithm 18. We compare LDA and QDA in Figure 11.4 by applying both of them to the music classification problem from Example 2.1.

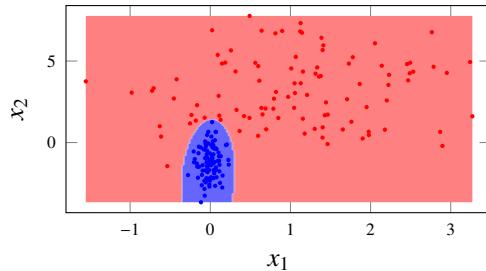


Figure 11.2: The decision boundary for the QDA classifier (obtained by (11.5) and (11.6)) corresponding to the learned Gaussian mixture model in the right panel of Figure 11.1.

Algorithm 18: Quadratic Discriminant Analysis, QDA

Data: Training data $\{\mathbf{x}_i, y_i\}_{i=1}^n$ (with output classes $k = 1, \dots, M$) and test input \mathbf{x}_*
Result: Predicted test output \hat{y}_*

Learn

- 1 **for** $m = 1, \dots, M$ **do**
 - 2 |_ Compute $\hat{\pi}_m$ (11.3a), $\hat{\mu}_m$ (11.3b) and $\hat{\Sigma}_m$ (11.3c)
 - 3 **end**

 - 4 **Predict**
 - 5 **for** $m = 1, \dots, M$ **do**
 - 6 |_ Compute $p(y = m | \mathbf{x}_*)$ (11.5)
 - 7 **end**
 - 7 Find largest $p(y = m | \mathbf{x}_*)$ and set \hat{y}_* to that m
-

Time to reflect 11.1: In the Gaussian mixture model it was assumed that $p(\mathbf{x} | y)$ is Gaussian. When applying LDA or QDA “out of the box” for a supervised problem, is there any check that the Gaussian assumption actually holds? If yes, what? If no, is that a problem?

We have now derived a classifier, QDA, from a generative model. In practice the QDA classifier can be employed just like any discriminative classifier. It can be argued that a generative model contains more assumptions than a discriminative model, and if the assumptions are fulfilled we could possibly expect QDA to be slightly more data efficient (requiring fewer data points to reach a certain performance) than a discriminative model. In most practical cases that makes, however, not a major difference. The difference between using a generative and a discriminative model will, however, be larger when we next look at the semi-supervised problem.

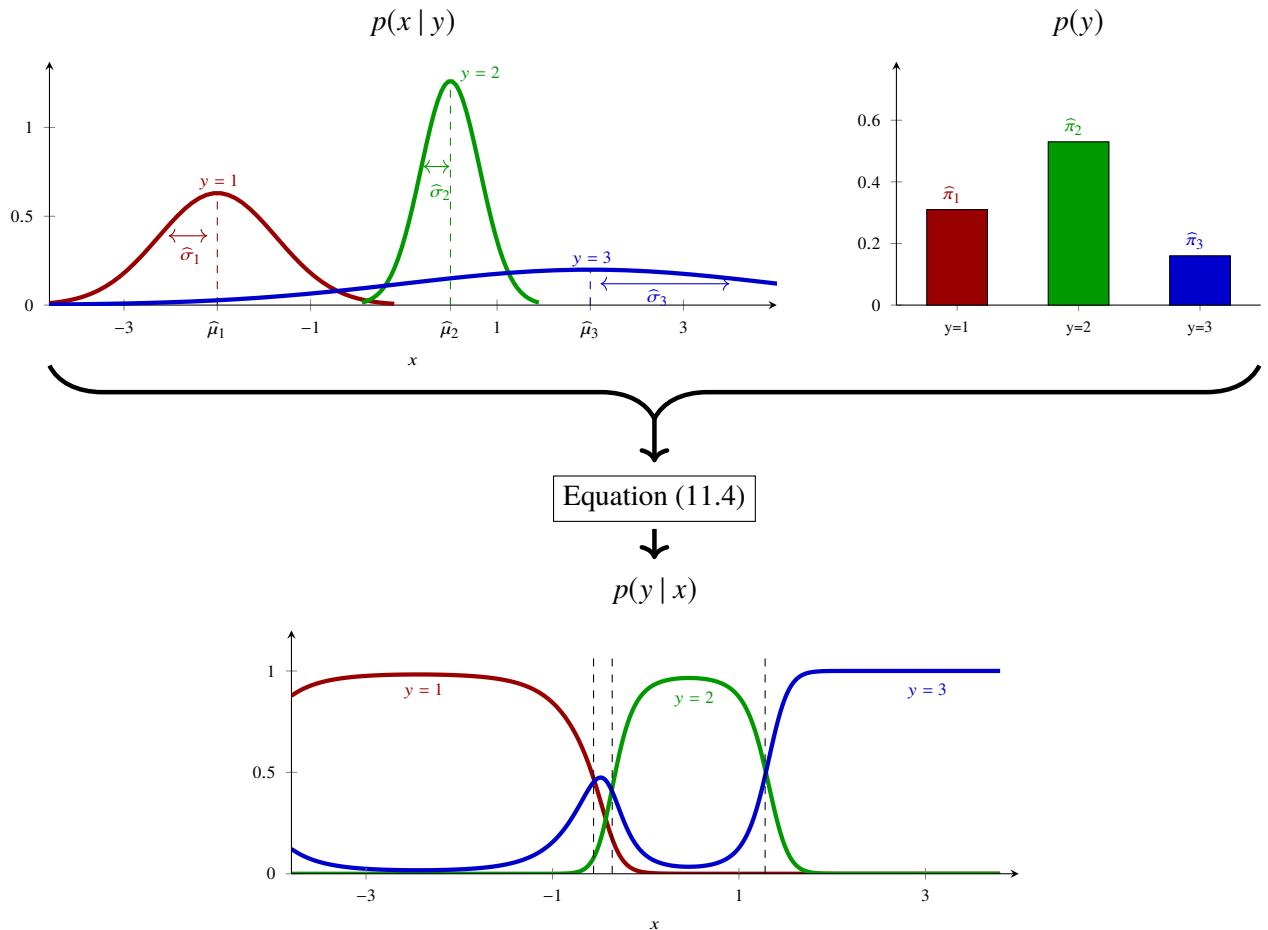
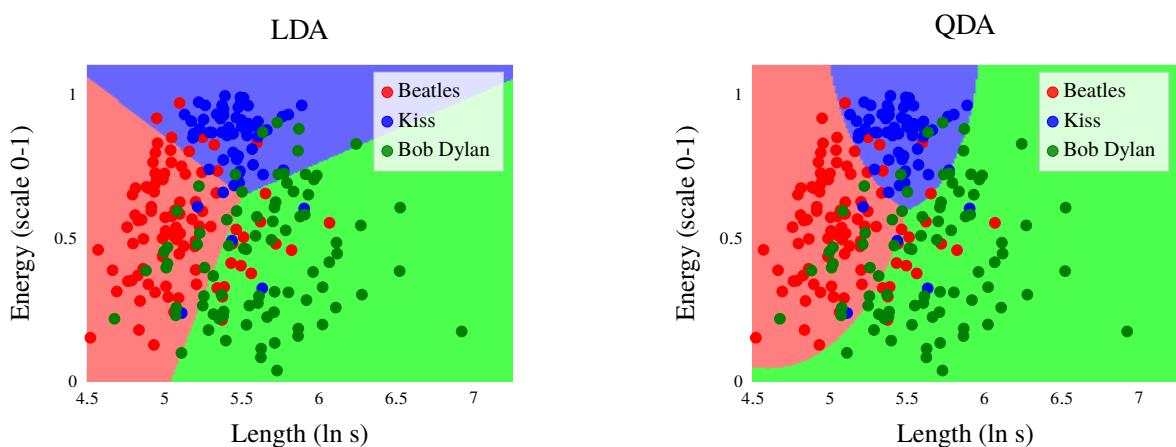


Figure 11.3: An illustration of QDA for $M = 3$ classes, with dimension $p = 1$ of the input x . In the top is the generative Gaussian mixture model shown. To the left is the Gaussian model of $p(\mathbf{x} | k)$, parameterized by $\hat{\mu}_m$ and Σ_m . To the right is the model of $p(y)$ shown, parameterized by $\hat{\pi}_m$. All parameters are learned from training data, not shown in the figure. (Since $p = 1$, we only have a scalar variance $\hat{\sigma}_m^2$, instead of a covariance matrix Σ_m). By (11.4) is the generative model “warped” into $p(y = m | x)$, shown in the bottom. The decision boundaries are shown as vertical dotted lines in the bottom plot.



(a) Decision boundaries for the music classification problem for an LDA classifier.

(b) Decision boundaries for the music classification problem for a QDA classifier.

Figure 11.4: We apply an LDA and QDA classifier to the music classification problem from Example 2.1 and plot the decision boundaries. Note that the LDA classifier gives linear decision boundaries, whereas the QDA classifier has decision boundaries with quadratic shapes.

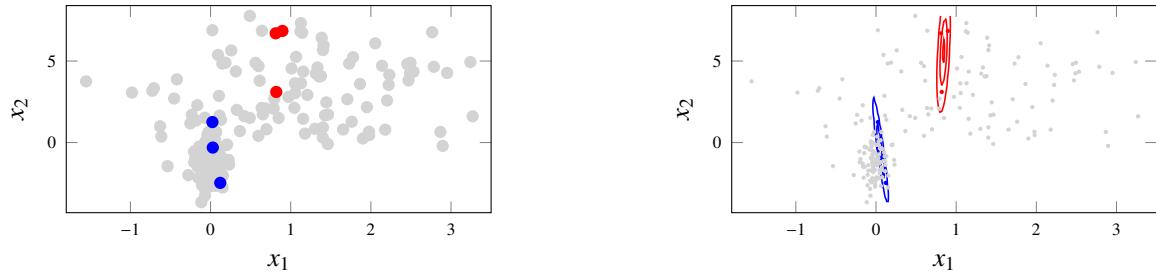


Figure 11.5: We consider the same situation as in Figure 11.1, except for the fact that we have “lost” the output y_i for most of the data points. The problem is now semi-supervised. The unlabeled data points $\{\mathbf{x}_i\}_{i=1}^{n_u}$ are illustrated in the left panel as gray dots, whereas the $\{\mathbf{x}_i, y_i\}_{i=1}^{n_c} = 6$ complete data points $\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}$ are red or blue. In the right panel we have learned a Gaussian mixture model using only the complete data points, as if the problem was supervised with only n_c data points. Clearly the unlabeled data points have made the problem harder, compared to Figure 11.1. We will, however, continue this story in Figure 11.6 where we use this as an initialization to a semi-supervised procedure.

11.2 The Gaussian mixture model when some or all labels are missing

We have so far discussed how the Gaussian mixture model can be learned in the supervised setting, that is, from training data that contains input as well as output values. We will now have a look at the so-called *semi-supervised* problem where some output values y_i are missing in the training data. The input values \mathbf{x}_i for which the corresponding y_i are missing are called *unlabeled* data points. As before we denote the total number of training data points as n , out of which now only n_c are complete input-output pairs $\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}$ and the remaining n_u unlabeled data points $\{\mathbf{x}_i\}_{i=1}^{n_u}$. All in all we have the training data $\mathcal{T} = \{\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}, \{\mathbf{x}_i\}_{i=1}^{n_u}\}$, where $n = n_c + n_u$. Semi-supervised problems are common to encounter, since supervised problems in practice often appears to rather be semi-supervised because of various data issues that gives corrupted or missing output values y_i .

Remember that a generative model is a model of $p(\mathbf{x}, \mathbf{y})$, which can be factorized as $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y} | \mathbf{x})$. Since $p(\mathbf{x})$ in this sense is a part of the generative model, it seems intuitively plausible that also the unlabeled data points $\{\mathbf{x}_i\}_{i=1}^{n_u}$ can be useful when learning a generative model. This is in contrast to discriminative models, which only cares about how \mathbf{y} depends on \mathbf{x} , and unlabeled data points are harder to make use of.

We will later also consider the *unsupervised* problem, where all data points are unlabeled.

Semi-supervised learning of the Gaussian mixture model

The simplest solution to the semi-supervised problem would be to discard the n_u unlabeled data points, and thereby turn the problem into a standard supervised machine learning problem. That is indeed a pragmatic solution if the number of unlabeled data points n_u is small compared to the total number of data points n . However, if the unlabeled data points are a significant part of all available data \mathcal{T} , we might be left with rather few complete data points $\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}$. We illustrate this with Figure 11.5, which depicts a semi-supervised problem where we have learned a (poor) Gaussian mixture model by only using the few n_c complete data points. From before we also know that the more data the better, which also is an incentive for a more careful way to handle the semi-supervised problem.

Conceptually we would like to solve the maximum likelihood problem again (as for the supervised problem, (11.2)), but leaving out the missing outputs. That is, we would like to solve

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} p(\underbrace{\{\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}, \{\mathbf{x}_i\}_{i=1}^{n_u}\}}_{\mathcal{T}} | \boldsymbol{\theta}). \quad (11.8)$$

Unfortunately this problem has no closed-form solution for the Gaussian mixture model since, if we are to spell $p(\{\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}, \{\mathbf{x}_i\}_{i=1}^{n_u}\} | \boldsymbol{\theta})$ out, it contains integrals over the missing y_i .

Leaving the seemingly hard problem (11.8) aside, a pragmatic approach to the semi-supervised problem would be to use the Gaussian mixture model to predict the missing output values $\{y_i\}_{i=1}^{n_u}$, in the iterative fashion

- (i) Learn the Gaussian mixture model from the n_c complete available input-output pairs $\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}$
- (ii) Use the Gaussian mixture model to predict (as a QDA classifier) the missing outputs to $\{\mathbf{x}_i\}_{i=1}^{n_u}$
- (iii) Update the Gaussian mixture model using also the predicted outputs from step (ii)

and then repeat step (ii) and (iii) until convergence.

It is far from obvious whether such a procedure makes any sense, and even less so whether it will converge. As it turns out, however, it can be rigorously shown that this is indeed a very sensible idea, and it does converge. If we only pay attention to some important details, the suggested procedure is actually the so-called expectation-maximization (EM) algorithm applied to (11.8).

The EM algorithm is a widely used computational statistics tool for solving maximum likelihood problems which contain latent (unobserved) variables, which can be shown to converge to a local optimum of (11.8). By considering the missing output values $\{y_i\}_{i=1}^{n_u}$ as being such latent variables, the EM algorithm fits perfectly and boils down to the suggested iteration with a few important details: From step (ii) we should return the predicted class probabilities $p(y = m | \mathbf{x})$ (and not the class prediction $\hat{y}(\mathbf{x}_*)$), and in step (iii) we make use of the predicted class probabilities by introducing the notation

$$w_i(m) = \begin{cases} p(y = m | \mathbf{x}_i) & \text{if } y_i \text{ is missing} \\ 1 & \text{if } y_i = m \\ 0 & \text{otherwise} \end{cases} \quad (11.9a)$$

and learn the parameters as

$$\hat{\pi}_m = \frac{1}{n} \sum_{i=1}^n w_i(m), \quad (11.9b)$$

$$\hat{\mu}_m = \frac{1}{\sum_{i=1}^n w_i(m)} \sum_{i=1}^n \mathbf{x}_i w_i(m), \quad (11.9c)$$

$$\hat{\Sigma}_m = \frac{1}{\sum_{i=1}^n w_i(m)} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu}_m)(\mathbf{x}_i - \hat{\mu}_m)^T w_i(m). \quad (11.9d)$$

Note that (11.9) is a generalization of the supervised case, and that (11.3) is a special case of (11.9) when no labels y_i are missing. We also summarize this as Algorithm 19, and illustrate it in Figure 11.6 by applying it to the situation introduced in Figure 11.5.

Algorithm 19: Semi-supervised learning of the Gaussian mixture model

- Data:** Incomplete training data $\mathcal{T} = \{\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}, \{\mathbf{x}_i\}_{i=1}^{n_u}\}$ (with output classes $k = 1, \dots, M$)
- 1 Compute, for each $m = 1, \dots, M$, $\hat{\pi}_m$ (11.3a), $\hat{\mu}_m$ (11.3b) and $\hat{\Sigma}_m$ (11.3c) using only the complete data $\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}$
 - 2 **repeat**
 - 3 | Compute, for each \mathbf{x}_i in $\{\mathbf{x}_i\}_{i=1}^{n_u}$, the prediction $p(y | \mathbf{x}_i)$.
 - 4 | Compute, for each $m = 1, \dots, M$, $\hat{\pi}_m$ (11.9b), $\hat{\mu}_m$ (11.9c) and $\hat{\Sigma}_m$ (11.9d)
 - 5 **until** convergence
-

We have now devised a way to handle semi-supervised classification problems using the Gaussian mixture model, and thereby extended the QDA classifier such that it can be used also in the semi-supervised setting when some output values y_i are missing from the training data.

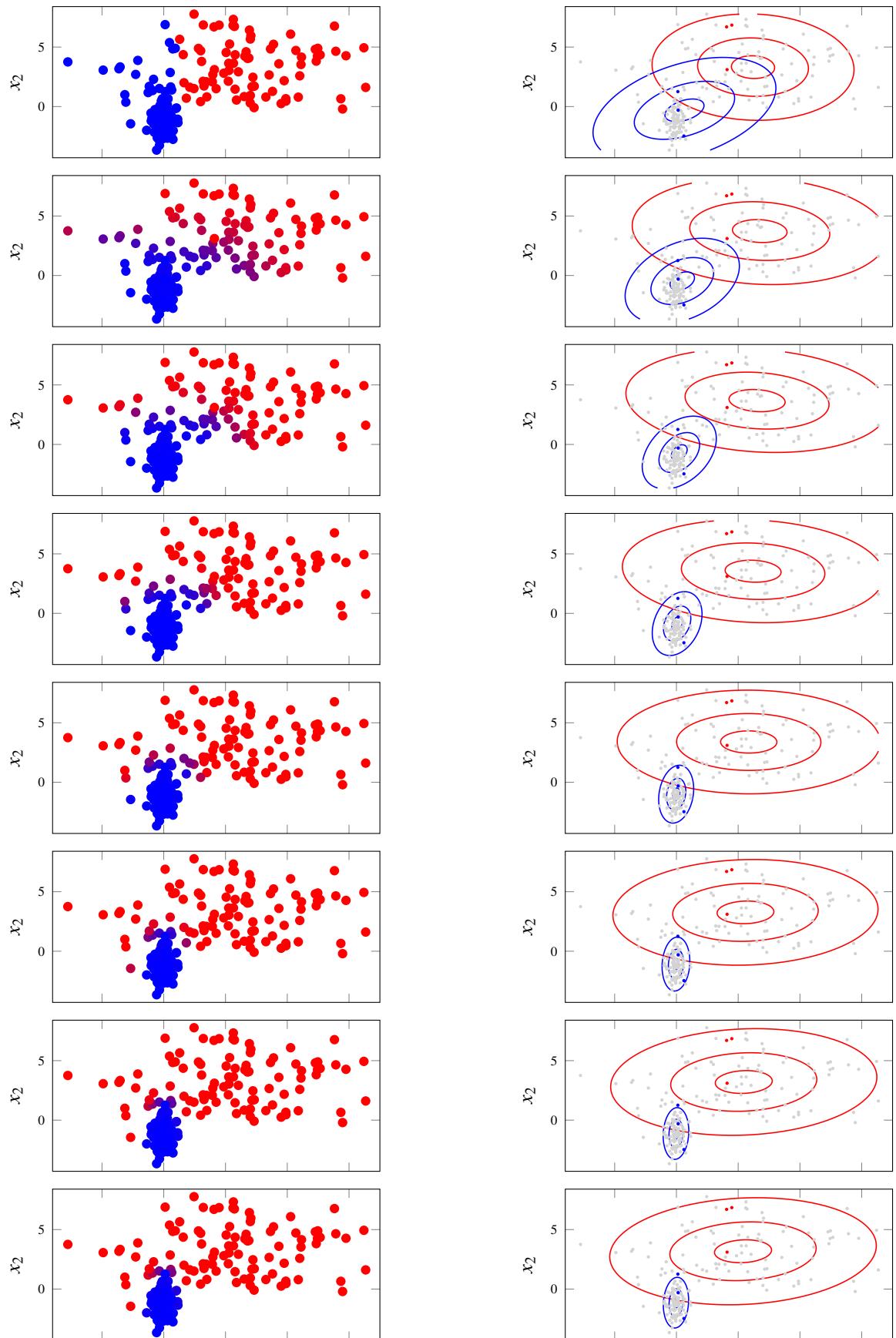


Figure 11.6: The iterative Algorithm 19 applied to the problem from Figure 11.5. For each iteration, the left panel shows the predicted class probabilities from the previously learned model (using color coding; purple is in the middle between red and blue). The previous model to the top row was shown in the right panel of Figure 11.5. The new model learned using (11.9) and the predictions in the left panel is shown in the right panel for each row.

It has perhaps yet not been clear why we have chosen to introduce the semi-supervised problem in connection to generative models. We could indeed think of using any discriminative model (instead of the Gaussian mixture model) for iteratively predicting the missing y_i and updating the model using the predictions. That could sometimes work to some extent, in particular if the fraction of unlabeled data points is small, but the procedure becomes much more robust and able to handle a larger proportion unlabeled datapoints when using a generative model. The reason for this is, however, somewhat subtle. Let us first remind ourselves about the difference between a discriminative and a generative model.

A discriminative model is a model for $p(y | \mathbf{x})$ only, which encodes information like “if \mathbf{x} , then y is ...”. Consequently the discriminative model does *not* contain any model of \mathbf{x} whatsoever, and makes no assumptions about the connection from y to \mathbf{x} . This is in contrast to a generative model which models $p(\mathbf{y}, \mathbf{x})$, from which $p(y | \mathbf{x})$ as well as $p(\mathbf{x} | y)$ follows. The presence of $p(\mathbf{x} | y)$ in a generative model is what makes it more useful for semi-supervised learning.

Consider the semi-supervised problem in Figure 11.5, and let us think about using a discriminative model, say a neural network, instead of a Gaussian mixture model in the iterative fashion of Algorithm 19 and Figure 11.6. We would first learn the model (the neural network) using only the complete datapoints $\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}$. That would give us a decision boundary that separates the blue and red data points well. We would thereafter use that model for predicting the output for the unlabeled (gray) data points. The next step would be to train the model (the neural network) anew using also the predicted outputs. However, a discriminative model (like a neural network) would in this step only care about fitting the training data (including the predicted ones) as close as possible, and the learned model would in practice look very similar to the model learned in the previous iteration. The procedure would therefore in practice get stuck after the first iteration, and the discriminative model would not converge to something much better than a model trained using only the complete datapoints $\{\mathbf{x}_i, y_i\}_{i=1}^{n_c}$.

The generative model, on the other hand, contains more assumptions than a discriminative model, which helps when handling the semi-supervised problem. The generative model implies a model $p(\mathbf{x} | y)$, which for the Gaussian mixture model encodes the information that all inputs \mathbf{x} that belongs to a certain output y should all have the same Gaussian distribution and thereby live in a “cluster”. This assumption is the key enabler for Algorithm 19 to work in practice, also in such a challenging situation as in Figure 11.5-11.6 where the vast majority of the data points are unlabeled. The downside of using a generative model is that the result can be rather misleading if the assumptions are not fulfilled.

Unsupervised learning of the Gaussian mixture model

In the supervised setting, all datapoints are complete (both \mathbf{x} and y), and in the semi-supervised problem some datapoints are complete and some datapoints are unlabeled (only \mathbf{x}). In the *unsupervised* problem, all datapoints are unlabeled, $\mathcal{T} = \{\mathbf{x}_i\}_{i=1}^n$. Whereas the overall goal in supervised and semi-supervised machine learning is to make as good predictions as possible, the overall goal of unsupervised machine learning is in general less pronounced. In the context of generative models, however, unsupervised machine learning amounts to learning $p(\mathbf{x})$.

The Gaussian mixture model (11.1) is a joint model for \mathbf{x} and y , as $p(\mathbf{x}, y)$. To obtain a model only for \mathbf{x} we can marginalize out y as $p(\mathbf{x}) = \sum_y p(\mathbf{x}, y)$ from it. The marginalization implies that we consider y as being a latent random variable, that is, a random variable which exists in the model but is not observed in the data. In practice we still learn $p(\mathbf{x}, y)$ (11.1), but from data only containing inputs $\{\mathbf{x}_i\}_{i=1}^n$.

Luckily we already have a tool for learning the Gaussian mixture model from unlabeled data. Algorithm 19, which we devised for the semi-supervised case, also works for completely unlabeled data $\{\mathbf{x}_i\}_{i=1}^n$, if we only replace the initialization (line 1) with some pragmatic choice of initial $\widehat{\pi}_m, \widehat{\mu}_m$ and $\widehat{\Sigma}_m$. It can be shown that Algorithm 19 also is the EM algorithm applied to solving the unsupervised maximum likelihood problem

$$\widehat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} p(\underbrace{\{\mathbf{x}_i\}_{i=1}^n}_{\mathcal{T}} | \boldsymbol{\theta}). \quad (11.10)$$

If we have unlabeled data and the number of classes M we can use Algorithm 19 to learn a Gaussian mixture model, as we illustrate by Figure 11.7.

There are a few important details when learning the Gaussian mixture model in the unsupervised setting which deserve some attention. First of all, even though the data is unlabeled, the number of classes (or rather the number of Gaussian components in the mixture) M still has to be known. Furthermore, since there are only unlabeled datapoints, the indexation of the M Gaussian components becomes arbitrary. In Figure 11.7 it means that the colors (red and blue) are interchangeable, and the way it happened to turn out is only an effect of the initialization.

The initialization is also an important detail to carefully consider, since Algorithm 19 (which builds on the EM algorithm) is only guaranteed to find local optima. With an unfortunate initialization Algorithm 19 may therefore end up in a poor local optima of (11.10). A pragmatic remedy is to run Algorithm 19 with different (random) initializations. Another issue with Algorithm 19 is that it might find a local optima which corresponds to a degenerate solution where some Gaussian component have a singular covariance matrix that causes numerical problems.

We have so far only said that unsupervised learning of a generative model corresponds to learning $p(\mathbf{x})$, but it is perhaps not clear why $p(\mathbf{x})$ can be of interest. As the name suggests, a generative model can be used to generate more data, by drawing samples from $p(\mathbf{x})$. That is true also for the Gaussian mixture model. The Gaussian mixture model in particular can be used also for solving the so-called clustering problem. By learning a Gaussian mixture model and then use it as a QDA classifier (even though it was learned from unlabeled data), we get a method for clustering (classifying) “similar” \mathbf{x}_i into M different clusters (classes). We will next have a look at another common clustering methods, which can be understood as a simplified version of the Gaussian mixture model and Algorithm 19.

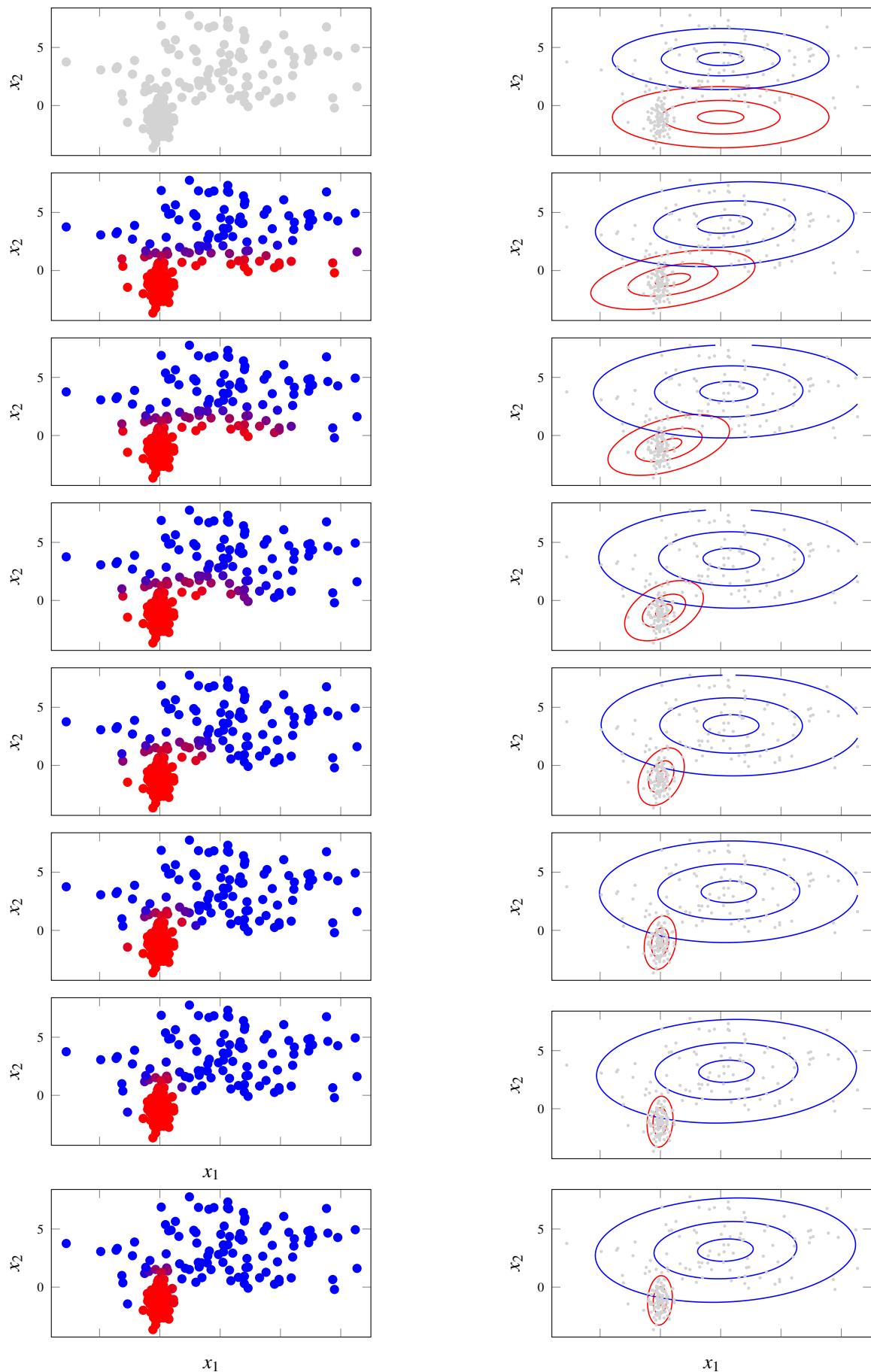


Figure 11.7: The Algorithm 19 applied to an unsupervised problem, where all training data is unlabeled. The only difference, in practice, to Figure 11.6 is the initialization (in the upper row), which here is done arbitrarily instead of using the complete datapoints.

11.3 More unsupervised methods: k -means and PCA

Whereas supervised machine learning deals with the problem of predicting $\hat{y}(\mathbf{x}_*)$ by learning a model from input-output data $\{\mathbf{x}_i, y_i\}_{i=1}^n$, *unsupervised machine learning* is more of an umbrella term for methods that somehow attempt to extract information from unlabeled data $\{\mathbf{x}_i\}_{i=1}^n$. We have introduced generative models, illustrated by the Gaussian mixture model, which is a general concept that can be used for supervised as well as unsupervised machine learning. We will now have a look at two more methods for unsupervised machine learning, for solving the clustering and dimensionality reduction problem respectively. The clustering problem amounts to grouping similar data points \mathbf{x}_i , whereas dimensionality reduction amounts to finding a lower dimensional representation of the data onto that still contain most of the information.

k -means clustering

The clustering problem amounts to group similar data points \mathbf{x}_i . There are several alternative ways of mathematically defining as well as solving the clustering problem. In k -means it is assumed that all elements of \mathbf{x} are numerical and that similarity is defined in terms of the Euclidian distance. We can define the k -means objective if we first introduce k sets (clusters) R_1, R_2, \dots, R_k . With the condition that each data point \mathbf{x}_i in $\{\mathbf{x}_i\}_{i=1}^n$ should be a member of exactly one cluster, k -means clustering amounts to selecting the clusters such that the sum of pairwise squared Euclidian distances within each cluster are minimized,

$$\arg \min_{R_1, R_2, \dots, R_k} \sum_{j=1}^k \frac{1}{|R_j|} \sum_{\mathbf{x}, \mathbf{x}' \in R_j} \|\mathbf{x} - \mathbf{x}'\|_2^2, \quad (11.11)$$

where $|R_j|$ is the number of data points in cluster R_j . The intention of (11.11) is to select the clusters such that all points within each cluster are as similar as possible. Solving (11.11) can be shown to be equivalent to selecting the clusters such that the distance to the cluster center, summed over all data points, is minimized,

$$\arg \min_{R_1, R_2, \dots, R_k} \sum_{j=1}^k \sum_{\mathbf{x} \in R_j} \|\mathbf{x} - \boldsymbol{\mu}_j\|_2^2. \quad (11.12)$$

Here $\boldsymbol{\mu}_j$ is the center (average) of all data points $\mathbf{x} \in R_j$. Unfortunately both (11.11) and (11.12) are combinatorial problems, meaning that we can not expect to solve them exactly if the number of data points n is large.

A common way to heuristically solve (11.11)/(11.12) is to

- (i) Set the cluster centers $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k$ to some initial values
- (ii) Determine which cluster R_j each \mathbf{x}_i belongs to, that is, compute which $\boldsymbol{\mu}_j$ is closest to each \mathbf{x}_i
- (iii) Update each cluster center $\boldsymbol{\mu}_j$ as the average of all \mathbf{x}_i that belongs to R_j

and then iterate step (ii) and (iii) until convergence. This procedure, which is an instance of the so-called Lloyds algorithm and often simply called “the k -means algorithm”, is guaranteed to converge (since we can verify that the objective in (11.12) decreases for each iteration) but it is not guaranteed to find the global optima. In practice it is common to run it multiple times, each with a different initialization in step (i), and pick the result of the run for which the objective in (11.11)/(11.12) attains the smallest value. As an illustration of k -means, we apply it to the input data from the music classification problem in Example 2.1.

Indeed k -means has some similarities to k -NN, in particular its heavy use of the Euclidian distance to define similarities in the input space. This implies that also k -means is very sensitive to the normalization of the input values. It is, however, important not to confuse the two methods. Whereas k -NN is a method for supervised machine learning, k -means is a method for the unsupervised clustering problem. Note, in particular, that k has a different meaning in the two methods.

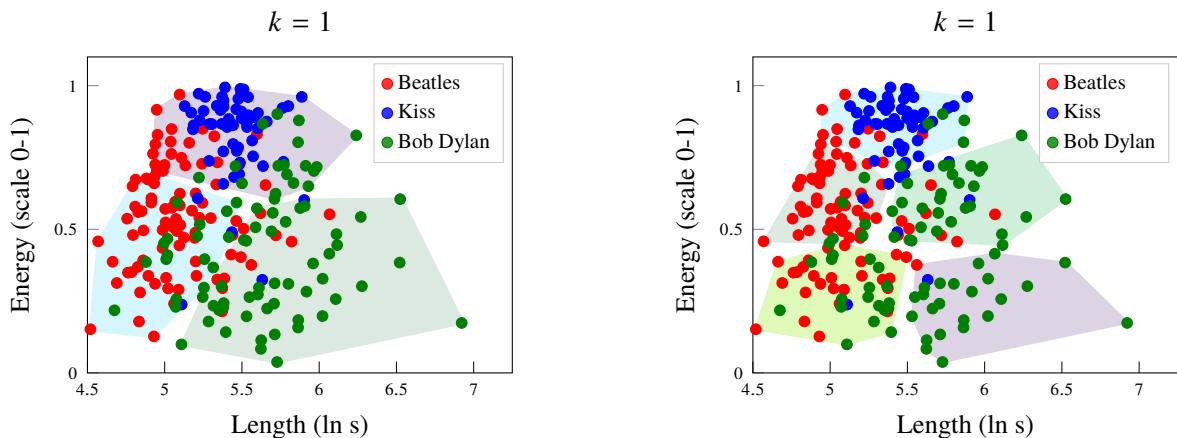


Figure 11.8: k -means applied to the music classification Example 2.1. In this example we do have input as well as output data, so the purpose of applying a clustering algorithm to its input only is rather unclear, but if nothing else we do it for the sake of curiosity (we could perhaps imagine that we have lost the artist labels, and are attempting to recover them). We try $k = 3$ (left) and $k = 5$ (right). It is worth to note that it so happens that for $k = 3$ there is almost one artist per cluster.

It is actually possible to see k -means almost as a special case of the unsupervised learning of the Gaussian mixture model with $k = M$. In the view of a Gaussian mixture model, all covariance matrices in k -means are set to the identity matrix, and furthermore is \hat{y} (instead of $p(y | \mathbf{x})$) used in Algorithm 19.

Selecting k

In k -means it is assumed that k is known. However, when doing clustering in practice it is common that k is a design variable left for the user to choose. Unfortunately we cannot choose k using cross-validation, since clustering is an unsupervised method and the concept of new data error E_{new} (which cross-validation builds upon, see Chapter 4) is not applicable to clustering. One heuristics to decide k instead is to compare the objective in (11.11), the sum of pairwise squared Euclidian distances within each cluster, for various values k .

Selecting k based on (11.11) is however not a fully automatic procedure, but there is some interpretation left to the user. By scrutinizing the equivalent formulation (11.12) we can see that it is minimized (zero) for the meaningless solution with $k = n$ clusters and one data point per cluster. We can furthermore see that if comparing (11.12) with k and $k + 1$ clusters, the objective will always take a smaller value for $k + 1$ clusters. We can therefore not blindly search for the minimum of (11.11) (which always must be found at $k = n$). Instead we should look for the point when the decrease in the objective becomes smaller, that is, selecting k such that the gain of going from $k - 1$ to k clusters is significantly larger than the gain of going from k to $k + 1$ clusters. For a dataset which has very distinct clusters, a graph of the objective in (11.11) looks visually like an elbow, and this method for selecting k in k -means is therefore sometimes called the *elbow method*. We illustrate it in Figure 11.9.

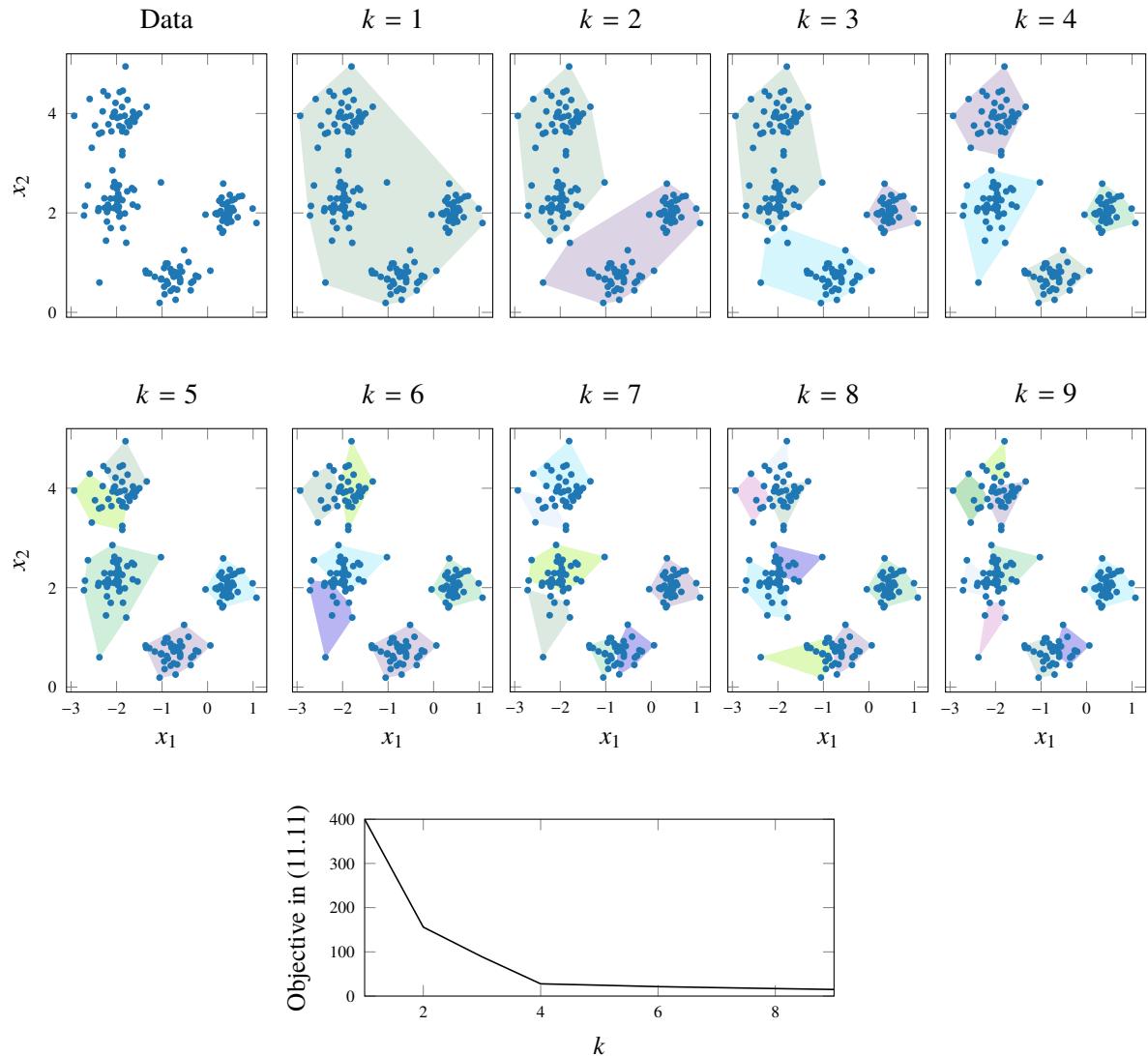


Figure 11.9: For selecting k in k -means we can use the so-called elbow method, which amounts to try different values of k (the upper panels) and record what the objective in (11.11) (the bottom panel). To select k , we look for a ‘bend’ in the bottom panel. In an ideal case there is a very distinct kink, but for this particular data we could either draw the conclusion $k = 2$ or $k = 4$ and it is up to the user to decide. Note that in this example the data has only 2 dimensions, and we can therefore visualize also the clusters themselves and visually compare them. If there are more than two dimensions of the data, however, we have to select k based only on the ‘elbow plot’ in the bottom panel.

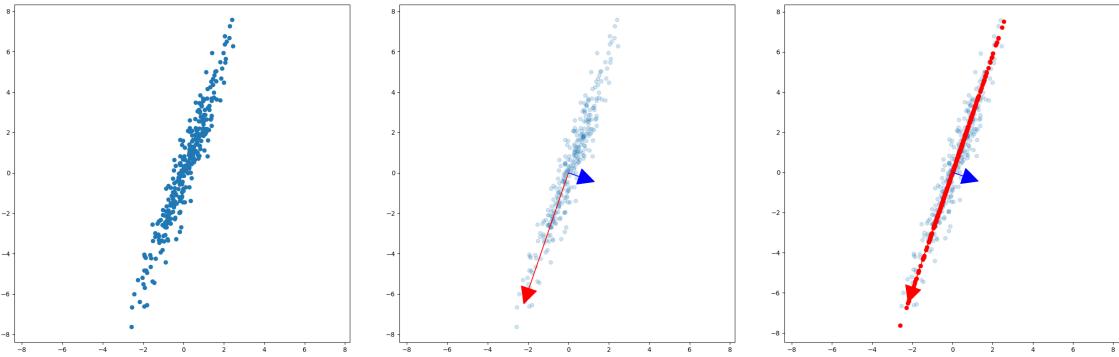


Figure 11.10: The idea of PCA is to find the dimensions along which the data varies the most. These dimensions are parameterized by the so-called principal components. In the left panel is some data $\{\mathbf{x}_i\}_{i=1}^n$ shown. The middle panel also shows the first (red) and second (blue) principal component. The right panel shows the data when projected onto only the first principal component (red dots). The red dots are now a one-dimensional representation $\{\mathbf{t}_i\}_{i=1}^n$ of the original two-dimensional blue dots $\{\mathbf{x}_i\}_{i=1}^n$.

Principal component analysis

We will now introduce principal component analysis (PCA), a method for *dimensionality reduction*. We start with a dataset $\{\mathbf{x}_i\}_{i=1}^n$, where \mathbf{x} is a p -dimensional vector $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_p]^T$ of numerical variables. The underlying assumption of dimensionality reduction is that there is some amount of redundancy in the p -dimensional representation of the data, and that it lives on (or close to) some manifold that allows a more compact representation. Dimensionality reduction therefore amounts to creating another dataset $\{\mathbf{t}_i\}_{i=1}^n$ where \mathbf{t} is a q -dimensional vector $\mathbf{t} = [t_1 \ t_2 \ \dots \ t_q]^T$. Importantly the new dataset has a smaller dimension, $q < p$, but with as much “information” as possible preserved from the original dataset. The purpose of doing dimensionality reduction is to get a more compact representation of the data, and thereby (hopefully) also filter out noise and only keep the most dominant patterns. In *linear* dimensionality reduction we also impose the condition that \mathbf{t} has to be a linear transformation of \mathbf{x} , meaning that it must hold that

$$\underbrace{\mathbf{t}}_{q \times 1} = \underbrace{\mathbf{U}}_{q \times p} \underbrace{\mathbf{x}}_{p \times 1}. \quad (11.13)$$

for some \mathbf{U} . Designing a linear dimensionality reduction method therefore “only” amounts to prescribing a way of choosing the transformation matrix \mathbf{U} .

In PCA, which probably is the most common method for linear dimensionality reduction, q is a user choice and the transformation \mathbf{U} is selected such that the variance of the transformed data $\{\mathbf{t}_i\}_{i=1}^n$ is maximized. The matrix \mathbf{U} is built up of q so-called principal components \mathbf{u} , which are p -dimensional vectors,

$$\mathbf{U} = \begin{bmatrix} \mathbf{u}_1^T \\ \mathbf{u}_2^T \\ \vdots \\ \mathbf{u}_q^T \end{bmatrix}. \quad (11.14)$$

The idea is illustrated in Figure 11.10.

Let us start with how to compute the first principal component \mathbf{u}_1 . Since \mathbf{u}_1 prescribes a direction onto which \mathbf{x} will be projected, it is convenient to let it be normalized, $\mathbf{u}_1^T \mathbf{u}_1 = 1$. We are searching for \mathbf{u}_1 such that the variance of $t_1 = \mathbf{u}_1^T \mathbf{x}$ is as large as possible. Strictly speaking we can not maximize the variance of t_1 , since we do not know the distribution of \mathbf{x} but only have n samples of it $\{\mathbf{x}_i\}_{i=1}^n$. Instead we therefore maximize the sample variance of t_1 . If we denote the mean of $\{\mathbf{x}_i\}_{i=1}^n$ as $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$, the

sample variance of $t_1 = \mathbf{u}_1^\top \mathbf{x}$ is by definition

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (\mathbf{u}_1^\top \mathbf{x}_i - \mathbf{u}_1^\top \bar{\mathbf{x}})^2 &= \frac{1}{n} \sum_{i=1}^n (\mathbf{u}_1^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}_1 - 2\mathbf{u}_1^\top \mathbf{x}_i \bar{\mathbf{x}}^\top \mathbf{u}_1 + \mathbf{u}_1^\top \bar{\mathbf{x}} \bar{\mathbf{x}}^\top \mathbf{u}_1) \\ &= \mathbf{u}_1^\top \underbrace{\left(\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top \right)}_{\mathbf{S}} \mathbf{u}_1 = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1. \end{aligned} \quad (11.15)$$

where \mathbf{S} is the sample covariance matrix for $\{\mathbf{x}_i\}_{i=1}^n$. Maximizing the sample variance of t_1 can thus be done by solving the optimization problem

$$\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 \quad (11.16)$$

subject to the constraint that $\mathbf{u}_1^\top \mathbf{u}_1 = 1$. Note that the constraint is important for the optimization problem to make sense. By introducing a Lagrange multiplier λ_1 we can rewrite the above optimization problem as the unconstrained problem

$$\arg \max_{\mathbf{u}_1, \lambda_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1(1 - \mathbf{u}_1^\top \mathbf{u}_1), \quad (11.17)$$

which is a quadratic problem. By taking the derivative with respect to \mathbf{u}_1 and set to zero, we know that the solution must fulfill

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1. \quad (11.18)$$

We recognize this equation as the eigenequation for \mathbf{S} , implying that the solution \mathbf{u}_1 to (11.16) is an eigenvector to \mathbf{S} with eigenvalue λ_1 . Furthermore, by multiplying (11.18) from the left with \mathbf{u}_1^\top we obtain an expression for the sample variance of $t_1 = \mathbf{u}_1^\top \mathbf{x}$ (11.15) as

$$\mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 = \mathbf{u}_1^\top \lambda_1 \mathbf{u}_1 = \lambda_1. \quad (11.19)$$

Since we seek to maximize the sample variance of t_1 , we conclude that the solution \mathbf{u}_1 to (11.16) must be the eigenvector to \mathbf{S} which has the largest eigenvalue λ_1 .

The number q of principal components in \mathbf{U} is a user choice, and if q is chosen larger than 1 we have to find (at least) a second one, \mathbf{u}_2 . We will not give a complete derivation of the second principal component \mathbf{u}_2 , but considering that it should be chosen such that the sample variance of t_2 is maximized but at the same time not “overlap” with t_1 , and the fact that all $p \times p$ symmetric matrices have a set of p orthogonal eigenvectors, it is perhaps not hard to accept that the second principal component \mathbf{u}_2 is the eigenvector to \mathbf{S} with the second largest eigenvalue λ_2 . In fact we have the general pattern that the j th principal component \mathbf{u}_j is the eigenvector to \mathbf{S} with the j th largest eigenvalue λ_j . If some eigenvalues to \mathbf{S} are not distinct, there is a corresponding degree of freedom in choosing the principal components. There is also a strong connection between PCA and singular value decomposition which we do not discuss any further here.

It can be a good practice to normalize the data (2.2) before applying PCA, in particular if the different variables have different scales. It is, however, more of an art than a science, in particular since the overall goal of dimensionality reduction is defined in rather vague terms. If all variables are measured in the same unit, for example, one can argue that PCA is more informative if not using normalization.