

Chapter 5

Machine Learning Basics

Deep learning is a specific kind of machine learning. To understand it well, one must have a solid understanding of the basic principles of machine learning. This chapter provides a brief course in the most important general concepts that are applied throughout the rest of the book. Novice readers or those with a wider perspective are encouraged to consider machine learning texts for more comprehensive coverage of the fundamentals, such as Murphy (2006). If you are already familiar with machine learning basics, feel free to skip ahead to section 5.11. That section covers some perspectives on traditional learning techniques that have strongly influenced the development of deep learning algorithms.

We begin with a definition of what a learning algorithm is and an example: the linear regression algorithm. We then proceed to describe the challenge of fitting the training data differs from the challenge of finding models that generalize to new data. Most machine learning algorithms have parameters called *hyperparameters*, which must be determined outside the learning process.

itself; we discuss how to set these using additional data. Machine learning is essentially a form of applied statistics with increased emphasis on using computers to statistically estimate complicated functions and a decision on proving confidence intervals around these functions; we therefore focus on two central approaches to statistics: frequentist estimators and Bayesian methods.

Most machine learning algorithms can be divided into the categories of supervised learning and unsupervised learning; we describe these categories and provide examples of simple learning algorithms from each category. Most machine learning algorithms are based on an optimization algorithm called stochastic gradient descent.

descent. We describe how to combine various algorithm components: an optimization algorithm, a cost function, a model, and a data representation, to construct a machine learning algorithm. Finally, in section 5.11, we describe some factors that have limited the ability of traditional machine learning algorithms to learn complex functions. These challenges have motivated the development of deep learning, which we describe in section 5.12. We will see that deep learning has overcome these obstacles.

5.1 Learning Algorithms

A machine learning algorithm is an algorithm that is able to learn from experience. But what do we mean by learning? Mitchell (1997) provides a succinct definition: “A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .” One can imagine a computer program that experiences E , tasks T , and performance measures P , and we do not need to understand exactly what each of these means to construct machine learning algorithms. In this book to formally define what may be used for each of these elements. In the following sections, we provide intuitive descriptions and examples of different kinds of tasks, performance measures, and experiences that can be used to construct machine learning algorithms.

5.1.1 The Task, T

Machine learning enables us to tackle tasks that are too difficult or impossible for fixed programs written and designed by human beings. From a philosophical point of view, machine learning is interesting because it requires an understanding of it entails developing our understanding of the nature of intelligence.

In this relatively formal definition of the word “task,” the program itself is not the task. Learning is our means of attaining the ability to perform the task.

task. For example, if we want a robot to be able to walk, then we could program the robot to learn to walk, or we could attempt to write a program that specifies how to walk manually.

Machine learning tasks are usually described in terms of how a learning system should process an **example**. An example is a collection of features that have been quantitatively measured from some object or event that the machine learning system to process. We typically represent an example as a vector $\mathbf{x} \in \mathbb{R}^n$ where each entry x_i of the vector is another feature. For example, the features of an image are usually the values of the pixels in the image.

Many kinds of tasks can be solved with machine learning. Some common machine learning tasks include the following:

- **Classification:** In this type of task, the computer program is asked to determine which of k categories some input belongs to. To solve this task, the algorithm is usually asked to produce a function $f : \mathbb{R}^n \rightarrow \{1, 2, \dots, k\}$. Given an input \mathbf{x} , the model assigns an input described by vector \mathbf{x} to one of the k categories identified by numeric code y . There are other variants of this task, for example, where f outputs a probability distribution over the categories. An example of a classification task is object recognition, where the input \mathbf{x} is an image (usually described as a set of pixel brightness values) and the output is a numeric code identifying the object in the image. For example, the Willow Garage PR2 robot is able to act as a waiter that can identify different kinds of drinks and deliver them to people on command (Koch et al., 2009; Koch and Culbreth-Ward, 2010). Modern object recognition is best accomplished using deep learning (Krizhevsky et al., 2012; Ioffe and Szegedy, 2015). In fact, object recognition is the same basic technology that enables computers to recognize faces (Taigman et al., 2014), which can be used to automatically tag people in photo collections and for computers to interact more naturally with their users.
- **Classification with missing inputs:** Classification becomes more challenging if the computer program is not guaranteed that every input variable will always be provided. To solve the classification task, the learning algorithm only has to define a *single* function mapping each input to a categorical output. When some of the inputs may be missing, rather than providing a single classification function, the learning algorithm must learn a *set* of functions. Each function corresponds to classifying a different subset of its inputs missing. This kind of situation arises frequently in medical diagnosis, because many kinds of medical tests are non-invasive. One way to efficiently define such a large set of functions is to use a decision tree.

Invasive. One way to efficiently define such a large set of learn a probability distribution over all the relevant variables classification task by marginalizing out the missing variables, we can now obtain all 2^n different classification functions for each possible set of missing inputs, but the computer to learn only a single function describing the joint probability. See Goodfellow *et al.* (2013b) for an example of a deep probability applied to such a task in this way. Many of the other tasks discussed in this section can also be generalized to work with missing inputs. A neural network with missing inputs is just one example of what machine learning

- **Regression:** In this type of task, the computer program is asked to output a numerical value given some input. To solve this task, the learning algorithm is asked to output a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. This type of task is similar to classification, except that the format of output is different. An example of a regression task is the prediction of the expected claim amount that an insured person will make (used to set insurance premiums), or the prediction of future prices of securities. These kinds of predictions are used in algorithmic trading.
 - **Transcription:** In this type of task, the machine learning algorithm is asked to observe a relatively unstructured representation of some information and transcribe the information into discrete textual form. For example, in optical character recognition, the computer program is shown an image containing an image of text and is asked to return this text as a sequence of characters (e.g., in ASCII or Unicode format). Address extraction and street view uses deep learning to process address numbers in this way (Krause *et al.*, 2014d). Another example is speech recognition, where the computer program is provided an audio waveform and emits a sequence of word ID codes describing the words that were spoken in the waveform. Deep learning is a crucial component of modern speech recognition systems used at major companies, including Microsoft, IBM and Google (Hinton *et al.*, 2012b).
 - **Machine translation:** In a machine translation task, the input consists of a sequence of symbols in some language, and the computer program must convert this into a sequence of symbols in another language. This is commonly applied to natural languages, such as translating from English to French. Deep learning has recently begun to have an important impact on this kind of task (Sutskever *et al.*, 2014; Bahdanau *et al.*, 2014).
- : Structured output tasks involve any of the following
- **Structured output**

output is a vector (or other data structure containing multiple values) that captures important relationships between the different elements. This category and subsumes the transcription and translation tasks mentioned above, as well as many other tasks. One example is parsing a natural language sentence into a tree that describes its grammatical structure by tagging nodes of the trees as being verbs, nouns, adverbs, etc. See Collobert (2011) for an example of deep learning applied to this task. Another example is pixel-wise segmentation of images. In this task, a computer program assigns every pixel in an image to a specific category.

For example, deep learning can be used to annotate the location of objects in aerial photographs (Mnih and Hinton, 2010). The output does not mirror the structure of the input as closely as in these annotation tasks. For example, in image captioning, the computer program takes an image and outputs a natural language sentence describing the image (Vinyals *et al.*, 2014a,b; Mao *et al.*, 2015; Vinyals *et al.*, 2015b; Donahue *et al.*, 2015; Karpathy and Li, 2015; Fang *et al.*, 2015; Xu *et al.*, 2015). These tasks are called *structured output tasks* because the program must produce multiple values that are all tightly interrelated. For example, the words produced by an image captioning program must form a valid sentence.

- **Anomaly detection:** In this type of task, the computer program looks through a set of events or objects and flags some of them as normal or atypical. An example of an anomaly detection task is credit card fraud detection. By modeling your purchasing habits, a credit card company can detect misuse of your cards. If a thief steals your credit card information, the thief’s purchases will often come from a different geographic distribution over purchase types than your own. The credit card company can prevent fraud by placing a hold on an account as soon as it has been used for an uncharacteristic purchase. See Chandola *et al.* (2019) for a survey of anomaly detection methods.
- **Synthesis and sampling:** In this type of task, the machine learning algorithm is asked to generate new examples that are similar to the training data. Synthesis and sampling via machine learning is common for media applications when generating large volumes of content that would be expensive, boring, or require too much time. For example, games can automatically generate textures for large objects in a scene rather than requiring an artist to manually label each pixel (Lam et al., 2019). In some cases, we want the sampling or synthesis procedure to produce a specific kind of output given the input. For example, in a speech recognition system, we might want the synthesis procedure to generate a waveform that sounds like a particular person speaking.

task, we provide a written sentence and ask the program to generate a waveform containing a spoken version of that sentence. This is a structured output task, but with the added qualification that there is only one single correct output for each input, and we explicitly desire some degree of variation in the output, in order for the output to seem more realistic.

- **Imputation of missing values:** In this type of task, the machine learning algorithm is given a new example $\mathbf{x} \in \mathbb{R}^n$, but with some missing values.

missing. The algorithm must provide a prediction of the value of the entries.

- **Denoising:** In this type of task, the machine learning algorithm takes as input a *corrupted example* $\tilde{\mathbf{x}} \in \mathbb{R}^n$ obtained by an unknown corruption process from a *clean example* $\mathbf{x} \in \mathbb{R}^n$. The learner must predict the original \mathbf{x} from its corrupted version $\tilde{\mathbf{x}}$, or more generally predict the conditional probability distribution $p(\mathbf{x} \mid \tilde{\mathbf{x}})$.
- **Density estimation or probability mass function estimation:** In this density estimation problem, the machine learning algorithm is asked to learn a function $p_{\text{model}} : \mathbb{R}^n \rightarrow \mathbb{R}$, where $p_{\text{model}}(\mathbf{x})$ can be interpreted as a density function (if \mathbf{x} is continuous) or a probability mass function (if \mathbf{x} is discrete) on the space that the examples were drawn from. To learn this well (we will specify exactly what that means when we discuss learning measures P), the algorithm needs to learn the structure of the underlying distribution. It must know where examples cluster tightly and where they are sparse. Most of the tasks described above require the learning algorithm to at least implicitly capture the structure of the probability distribution. Density estimation enables us to explicitly capture that distribution, so that we can then perform computations on that distribution to solve related tasks as well. For example, if we have performed density estimation and learned a probability distribution $p(\mathbf{x})$, we can use that distribution to solve a missing value imputation task. If a value x_i is missing, and all other values, denoted \mathbf{x}_{-i} , are given, then we know the distribution of x_i given \mathbf{x}_{-i} by $p(x_i \mid \mathbf{x}_{-i})$. In practice, density estimation does not always solve all these related tasks, because in many cases the required computations on $p(\mathbf{x})$ are computationally intractable.

Of course, many other tasks and types of tasks are possible. The ones we list here are intended only to provide examples of what machine learning can do.

do, not to define a rigid taxonomy of tasks.

5.1.2 The Performance Measure, P

To evaluate the abilities of a machine learning algorithm, we need a quantitative measure of its performance. Usually this performance is specific to the task T being carried out by the system.

For tasks such as classification, classification with missing information, or regression, we often measure the **accuracy** of the model. Accuracy is

proportion of examples for which the model produces the correct output. We can also obtain equivalent information by measuring the **error rate**, which is the proportion of examples for which the model produces an incorrect output. We can compute the error rate as the expected 0-1 loss. The 0-1 loss on a particular example is 0 if it is correctly classified and 1 if it is not. For tasks such as denoising, it does not make sense to measure accuracy, error rate, or any other metric based on 0-1 loss. Instead, we must use a different performance metric that gives a continuous-valued score for each example. The most common metric for such tasks is the average log-probability the model assigns to some example.

Usually we are interested in how well the machine learning algorithm performs on data that it has not seen before, since this determines how well it will perform when deployed in the real world. We therefore evaluate these performance metrics on a **test set** of data that is separate from the data used for training the machine learning system.

The choice of performance measure may seem straightforward in some cases, but it is often difficult to choose a performance measure that correctly rewards the desired behavior of the system.

In some cases, this is because it is difficult to decide what should be measured. For example, when performing a transcription task, should we measure the accuracy of the system at transcribing entire sequences, or should we use a metric like a performance measure that gives partial credit for getting some elements of a sequence correct? When performing a regression task, should we reward the system more if it frequently makes medium-sized mistakes or if it rarely makes very large mistakes? These kinds of design choices depend on the context of the application.

In other cases, we know what quantity we would ideally like to measure, but measuring it is impractical. For example, this arises frequently in density estimation. Many of the best probabilistic models represent probability distributions only implicitly. Computing the actual probability value at a specific point in space in many such models is intractable. In these cases,

must design an alternative criterion that still corresponds to the desired criterion, or design a good approximation to the desired criterion.

5.1.3 The Experience, E

Machine learning algorithms can be broadly categorized as **unsupervised** by what kind of experience they are allowed to have during the learning process.

Most of the learning algorithms in this book can be understood as

to experience an entire **dataset**. A dataset is a collection of many examples, as defined in section 5.1.1. Sometimes we call examples **data points**.

One of the oldest datasets studied by statisticians and machine learning researchers is the Iris dataset (Fisher, 1936). It is a collection of measurements of different parts of 150 iris plants. Each individual plant corresponds to one example. The features within each example are the measurements of the plant: the sepal length, sepal width, petal length and petal width. The dataset also records which species each plant belonged to. Three species of iris are represented in the dataset.

Unsupervised learning algorithms experience a dataset containing unlabeled features, then learn useful properties of the structure of this dataset. In contrast to deep learning, we usually want to learn the entire probability distribution over the dataset, whether explicitly, as in density estimation, or implicitly, as in generative models. Some unsupervised learning algorithms perform other roles, like clustering, which consists of dividing the dataset into clusters of similar examples.

Supervised learning algorithms experience a dataset containing labeled features, but each example is also associated with a **label** or **target**. For example, the Iris dataset is annotated with the species of each iris plant. A supervised learning algorithm can study the Iris dataset and learn to classify iris plants into different species based on their measurements.

Roughly speaking, unsupervised learning involves observing several examples of a random vector \mathbf{x} and attempting to implicitly or explicitly learn the underlying probability distribution $p(\mathbf{x})$, or some interesting properties of that distribution. Supervised learning involves observing several examples of a random vector \mathbf{x} along with an associated value or vector \mathbf{y} , then learning to predict \mathbf{y} from \mathbf{x} by estimating $p(\mathbf{y} \mid \mathbf{x})$. The term **supervised learning** originates from the fact that the target \mathbf{y} being provided by an instructor or teacher who shows the learning system what to do. In unsupervised learning, there is no explicit target provided.

teacher, and the algorithm must learn to make sense of the data with

Unsupervised learning and supervised learning are not formally
The lines between them are often blurred. Many machine learning t
be used to perform both tasks. For example, the chain rule of pr
that for a vector $\mathbf{x} \in \mathbb{R}^n$, the joint distribution can be decomposed

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}).$$

This decomposition means that we can solve the ostensibly unsuperv

modeling $p(\mathbf{x})$ by splitting it into n supervised learning problems. A model can solve the supervised learning problem of learning $p(y | \mathbf{x})$ by using unsupervised learning technologies to learn the joint distribution and inferring

$$p(y | \mathbf{x}) = \frac{p(\mathbf{x}, y)}{\sum_{y'} p(\mathbf{x}, y')}.$$

Though unsupervised learning and supervised learning are not conceptually or distinct concepts, they do help roughly categorize some of the thousands of machine learning algorithms. Traditionally, people refer to regression and classification and structured output problems as supervised learning. Density estimation and support of other tasks is usually considered unsupervised learning.

Other variants of the learning paradigm are possible. For example, in semi-supervised learning, some examples include a supervision target and others do not. In multi-instance learning, an entire collection of examples is labeled as containing or not containing an example of a class, but the individual examples within the collection are not labeled. For a recent example of multi-instance learning with deep models, see [Kotzias et al. \(2015\)](#).

Some machine learning algorithms do not just experience a fixed dataset. For example, **reinforcement learning** algorithms interact with an environment and there is a feedback loop between the learning system and its experience. Reinforcement learning algorithms are beyond the scope of this book. Please see [Sutton and Barto \(1998\)](#) or [Bertsekas and Tsitsiklis \(1996\)](#) for information about reinforcement learning and [Mnih et al. \(2013\)](#) for the deep learning approach to reinforcement learning.

Most machine learning algorithms simply experience a dataset. A dataset can be described in many ways. In all cases, a dataset is a collection of examples which are in turn collections of features.

One common way of describing a dataset is with a **design matrix**. A design matrix is a matrix containing a different example in each row. Each column of the matrix corresponds to a different feature. For instance, the Iris dataset

150 examples with four features for each example. This means we can represent the dataset with a design matrix $\mathbf{X} \in \mathbb{R}^{150 \times 4}$, where $X_{i,1}$ is the plant i , $X_{i,2}$ is the sepal width of plant i , etc. We describe most algorithms in this book in terms of how they operate on design matrices.

Of course, to describe a dataset as a design matrix, it must be possible to describe each example as a vector, and each of these vectors must be of the same length. This is not always possible. For example, if you have a collection of photographs of different objects with different widths and heights, then different photographs will contain different numbers of pixels, so not all the photographs may be described as vectors.

length of vector. In Section 9.7 and chapter 10, we describe how to types of such heterogeneous data. In cases like these, rather than dataset as a matrix with m rows, we describe it as a set containing $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$. This notation does not imply that any two $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ have the same size.

In the case of supervised learning, the example contains a label as well as a collection of features. For example, if we want to use a learned model to perform object recognition from photographs, we need to specify what objects appear in each of the photos. We might do this with a numerical label, such as 0 signifying a person, 1 signifying a car, 2 signifying a cat, and so forth. When working with a dataset containing a design matrix of feature observations, we also provide a vector of labels \mathbf{y} , with y_i providing the label for example i .

Of course, sometimes the label may be more than just a single word. For example, if we want to train a speech recognition system to transcribe sentences, then the label for each example sentence is a sequence of words.

Just as there is no formal definition of supervised and unsupervised learning, there is no rigid taxonomy of datasets or experiences. The structures described here cover most cases, but it is always possible to design new ones for novel situations.

5.1.4 Example: Linear Regression

Our definition of a machine learning algorithm as an algorithm of improving a computer program's performance at some task is somewhat abstract. To make this more concrete, we present a simple machine learning algorithm: **linear regression**. We will use this example repeatedly as we introduce more machine learning concepts to understand the algorithm's behavior.

As the name implies, linear regression solves a regression problem. In other words, the goal is to build a system that can take a vector

predict the value of a scalar $y \in \mathbb{R}$ as its output. The output of linear models is a linear function of the input. Let \hat{y} be the value that our model predicts. We define the output to be

$$\hat{y} = \mathbf{w}^\top \mathbf{x},$$

where $\mathbf{w} \in \mathbb{R}^n$ is a vector of **parameters**.

Parameters are values that control the behavior of the system. In the equation above, \mathbf{w} is the coefficient that we multiply by feature x_i before summing up the results from all the features. We can think of \mathbf{w} as a set of **weights** that

each feature affects the prediction. If a feature x_i receives a positive weight, then increasing the value of that feature increases the value of our prediction. If a feature receives a negative weight, then increasing the value of that feature decreases the value of our prediction. If a feature's weight is large, then it has a large effect on the prediction. If a feature's weight is small, then it has a small effect on the prediction.

We thus have a definition of our task T : to predict y from \mathbf{x} . We can write this as $\hat{y} = \mathbf{w}^\top \mathbf{x}$. Next we need a definition of our performance measure, which we will call the mean squared error.

Suppose that we have a design matrix of m example inputs \mathbf{X} that we can use for training, only for evaluating how well the model performs. We also have a vector of regression targets providing the correct value of y for each of the m examples. Because this dataset will only be used for evaluation, we refer to it as the test set. We refer to the design matrix of inputs as $\mathbf{X}^{(\text{test})}$ and the vector of regression targets as $\mathbf{y}^{(\text{test})}$.

One way of measuring the performance of the model is to compute the **squared error** of the model on the test set. If $\hat{\mathbf{y}}^{(\text{test})}$ gives the predicted values of y by our model on the test set, then the mean squared error is given by

$$\text{MSE}_{\text{test}} = \frac{1}{m} \sum_i (\hat{\mathbf{y}}^{(\text{test})} - \mathbf{y}^{(\text{test})})_i^2.$$

Intuitively, one can see that this error measure decreases to 0 when $\hat{\mathbf{y}}^{(\text{test})} = \mathbf{y}^{(\text{test})}$. We can also see that

$$\text{MSE}_{\text{test}} = \frac{1}{m} \|\hat{\mathbf{y}}^{(\text{test})} - \mathbf{y}^{(\text{test})}\|_2^2,$$

so the error increases whenever the Euclidean distance between the predicted values and the targets increases.

To make a machine learning algorithm, we need to design an optimization problem that minimizes the mean squared error.

will improve the weights ω in a way that reduces MSE_{test} when it is allowed to gain experience by observing a training set ($X^{(\text{train})}$) in an intuitive way of doing this (which we justify later, in section 5) minimize the mean squared error on the training set, MSE_{train} .

To minimize MSE_{train} , we can simply solve for where its gradient

$$\begin{aligned}\nabla_{\boldsymbol{w}} MSE_{\text{train}} &= 0 \\ \Rightarrow \nabla_{\boldsymbol{w}} \frac{1}{m} \|\hat{\boldsymbol{y}}^{(\text{train})} - \boldsymbol{y}^{(\text{train})}\|_2^2 &= 0\end{aligned}$$

106

$$\begin{aligned}
 & \Rightarrow \frac{1}{m} \nabla_{\mathbf{w}} \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = 0 \\
 & \Rightarrow \nabla_{\mathbf{w}} (\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})})^\top (\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}) \\
 & \Rightarrow \nabla_{\mathbf{w}} (\mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})\top} \mathbf{y}^{(\text{train})}) \\
 & \Rightarrow 2\mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2\mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} = 0 \\
 & \Rightarrow \mathbf{w} = (\mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})})^{-1} \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})}
 \end{aligned}$$

The system of equations whose solution is given by equation 5.12 is called the **normal equations**. Evaluating equation 5.12 constitutes a learning algorithm. For an example of the linear regression learning algorithm see figure 5.1.

It is worth noting that the term **linear regression** is often used to refer to a slightly more sophisticated model with one additional parameter, the term b . In this model

$$\hat{y} = \mathbf{w}^\top \mathbf{x} + b,$$

so the mapping from parameters to predictions is still a linear function, but the mapping from features to predictions is now an affine function. This is shown in figure 5.1.

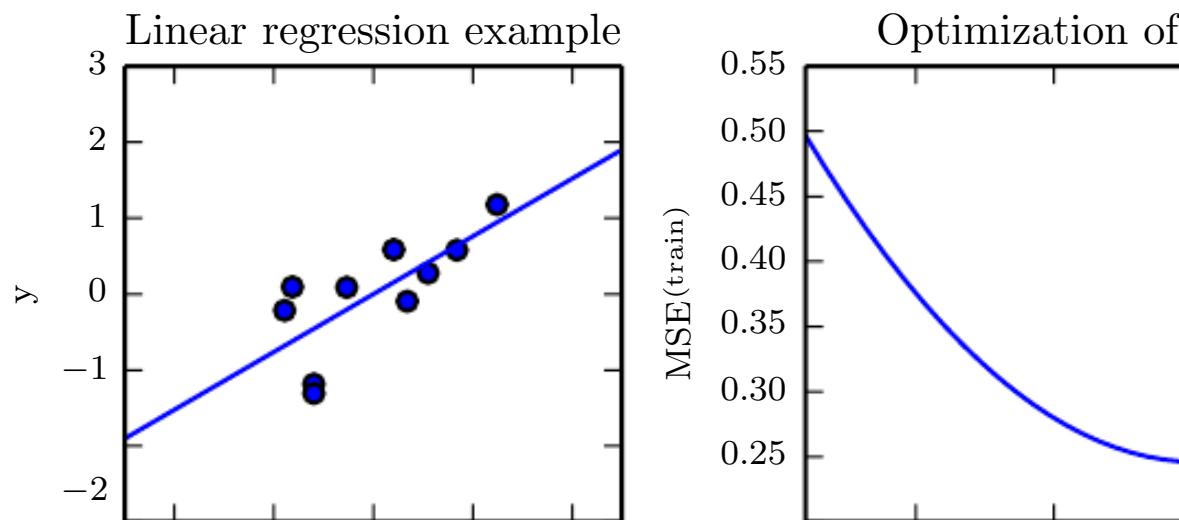




Figure 5.1: A linear regression problem, with a training set consisting of each containing one feature. Because there is only one feature, the contains only a single parameter to learn, w_1 . (*Left*)Observe that linear to set w_1 such that the line $y = w_1 x$ comes as close as possible to passing training points. (*Right*)The plotted point indicates the value of w_1 found equations, which we can see minimizes the mean squared error on the t

affine functions means that the plot of the model’s predictions is a straight line, but it need not pass through the origin. Instead of adding the intercept term b , one can continue to use the model with only weights but augment the input vector with an extra entry that is always set to 1. The weight corresponding to this extra entry plays the role of the bias parameter. We frequently use the term **affine function** referring to affine functions throughout this book.

The intercept term b is often called the **bias** parameter of the linear transformation. This terminology derives from the point of view that this transformation is biased toward being b in the absence of any information. This is different from the idea of a statistical bias, in which a statistical learning algorithm’s expected estimate of a quantity is not equal to the true value.

Linear regression is of course an extremely simple and limited learning algorithm. But it provides an example of how a learning algorithm can work. In the next few sections we describe some of the basic principles underlying learning algorithms and demonstrate how these principles can be used to build more complex learning algorithms.

5.2 Capacity, Overfitting and Underfitting

The central challenge in machine learning is that our algorithm must perform well on *new, previously unseen* inputs—not just those on which it was trained. The ability to perform well on previously unobserved inputs is called **generalization**.

Typically, when training a machine learning model, we have access to a training set; we can compute some error measure on the training set, called the **training error**; and we reduce this training error. So far, what we have described is an optimization problem. What separates machine learning from other optimization problems is that we want the **generalization error**, also called the **test error**, to be small. The generalization error is defined as the expected value of the error

well. The generalization error is defined as the expected value of new input. Here the expectation is taken across different possible from the distribution of inputs we expect the system to encounter

We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected from the training set.

In our linear regression example, we trained the model by minimizing the training error,

$$\frac{1}{m^{(\text{train})}} \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2,$$

but we actually care about the test error, $\frac{1}{m^{(\text{test})}} \|\mathbf{X}^{(\text{test})}\mathbf{w} - \mathbf{y}^{(\text{test})}\|$

How can we affect performance on the test set when we can only affect the training set? The field of **statistical learning theory** provides some answers. If the training and the test set are collected arbitrarily, there is indeed nothing we can do. If we are allowed to make some assumptions about how the training and test set are collected, then we can make some progress.

The training and test data are generated by a probability distribution over datasets called the **data-generating process**. We typically make several assumptions known collectively as the **i.i.d. assumptions**. These assumptions state that the examples in each dataset are **independent** from each other and that the training set and test set are **identically distributed**, drawn from the same probability distribution as each other. This assumption enables us to model the data-generating process with a probability distribution over all possible datasets. The same distribution is then used to generate every train example and every test example. We call that shared underlying distribution the **data-generating distribution**, denoted p_{data} . This probabilistic framework and the i.i.d. assumptions enable us to mathematically study the relationship between training error and test error.

One immediate connection we can observe between training error and test error is that the expected training error of a randomly selected model is equal to the expected test error of that model. Suppose we have a probability distribution $p(\mathbf{x}, y)$ and we sample from it repeatedly to generate the training set and test set. For some fixed value \mathbf{w} , the expected training set error is exactly equal to the expected test set error, because both expectations are formed over the same dataset sampling process. The only difference between the two computations is the name we assign to the dataset we sample.

Of course, when we use a machine learning algorithm, we first choose parameters ahead of time, then sample both datasets. We sample the training set and use it to choose the parameters to reduce training set error, then sample the test set to evaluate the test set error.

test set. Under this process, the expected test error is greater than the expected value of training error. The factors determining how well a learning algorithm will perform are its ability to

1. Make the training error small.
2. Make the gap between training and test error small.

These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model

obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.

We can control whether a model is more likely to overfit or underfit by adjusting its **capacity**. Informally, a model's capacity is its ability to fit a wide variety of functions. Models with low capacity may struggle to fit the training data well. Models with high capacity can overfit by memorizing properties of the training data, which may not serve them well on the test set.

One way to control the capacity of a learning algorithm is by specifying the **hypothesis space**, the set of functions that the learning algorithm will select as being the solution. For example, the linear regression algorithm uses the set of all linear functions of its input as its hypothesis space. We can increase the capacity of linear regression to include polynomials, rather than just linear functions, by increasing the size of the hypothesis space. Doing so increases the model's capacity.

A polynomial of degree 1 gives us the linear regression model we are already familiar with, with the prediction

$$\hat{y} = b + wx.$$

By introducing x^2 as another feature provided to the linear regression algorithm, we can learn a model that is quadratic as a function of x :

$$\hat{y} = b + w_1x + w_2x^2.$$

Though this model implements a quadratic function of its *input*, it is still a linear function of the *parameters*, so we can still use the normal equations to train the model in closed form. We can continue to add more additional features, for example, to obtain a polynomial of degree n :

$$\hat{y} = b + \sum_{i=1}^n w_i x^i.$$

Machine learning algorithms will generally perform best when is appropriate for the true complexity of the task they need to p amount of training data they are provided with. Models with insu are unable to solve complex tasks. Models with high capacity can tasks, but when their capacity is higher than needed to solve the pr may overfit.

Figure 5.2 shows this principle in action. We compare a lin and degree-9 predictor attempting to fit a problem where the t

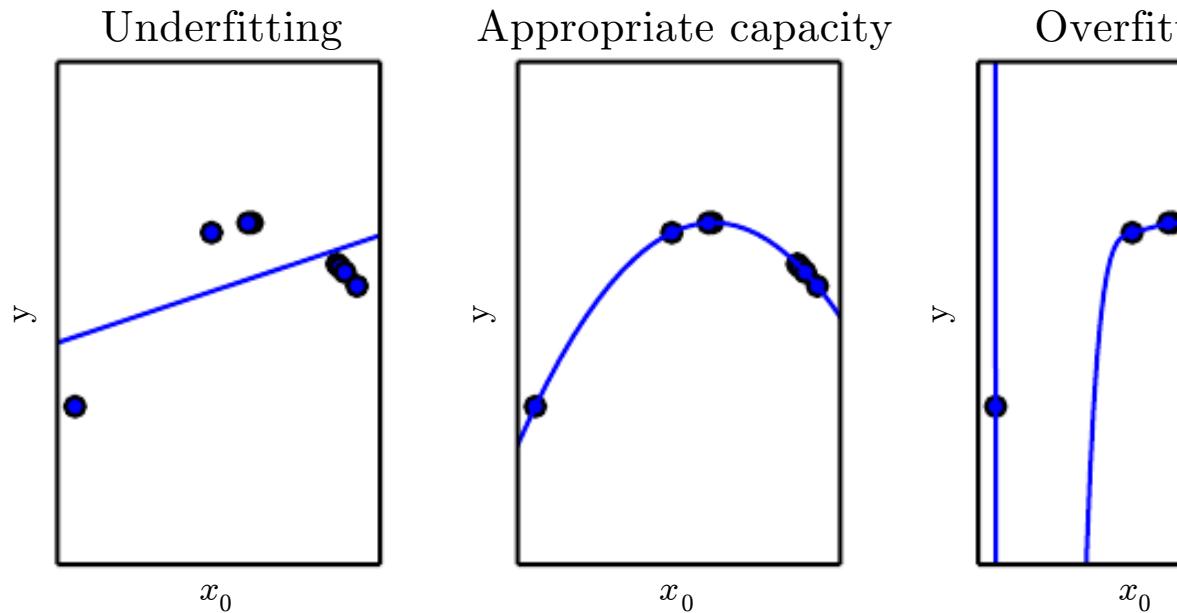


Figure 5.2: We fit three models to this example training set. The training data was generated synthetically, by randomly sampling x values and choosing y values by evaluating a quadratic function. (Left) A linear function fit to the data underfits—it cannot capture the curvature that is present in the data. A quadratic function fit to the data generalizes well to unseen points. It does not have a significant amount of overfitting or underfitting. (Right) A polynomial fit to the data suffers from overfitting. Here we used the Moore-Penrose pseudoinverse to solve the underdetermined normal equations. The solution passes through all the training points exactly, but we have not been lucky enough for it to extract the underlying quadratic function. It also increases sharply on the left side of the data and decreases sharply on the right side, where the true underlying function decreases in this area.

function is quadratic. The linear function is unable to capture the curvature of the true underlying problem, so it underfits. The degree-9 predictor is capable of representing the correct function, but it is also capable of representing many other functions that pass exactly through the training points. It has more parameters than training examples. We have little chance of finding a solution that generalizes well when so many wildly different solutions are possible.

this example, the quadratic model is perfectly matched to the training data, so it generalizes well to new data.

So far we have described only one way of changing a model's capacity: changing the number of input features it has, and simultaneously changing the values of the parameters associated with those features. There are in fact many other ways of changing a model's capacity. Capacity is not determined only by the choice of features; the learning algorithm can also change the model's capacity when varying the parameters in order to reduce a training objective function. This is what we mean by the **representational capacity** of the model. In many cases, fine-tuning the parameters of a model can increase its representational capacity.

function within this family is a difficult optimization problem. In fact, the learning algorithm does not actually find the best function, but merely significantly reduces the training error. These additional limitations, along with imperfection of the optimization algorithm, mean that the learned function's **effective capacity** may be less than the representational capacity of the function family.

Our modern ideas about improving the generalization of machine learning models are refinements of thought dating back to philosophers at least as far back as Ptolemy. Many early scholars invoke a principle of parsimony that has become widely known as **Occam's razor** (c. 1287–1347). This principle states that among competing hypotheses that explain known observations equally well, one should choose the “simplest” one. This idea was formalized and made mathematically precise in the twentieth century by the founders of statistical learning theory ([Vapnik and Chervonenkis, 1971](#); [Vapnik, 1982](#); [Blumer *et al.*, 1989](#); [Vapnik, 1995](#)).

Statistical learning theory provides various means of quantifying the capacity of a model. Among these, the most well known is the **Vapnik-Chervonenkis** (VC) dimension. The VC dimension measures the capacity of a binary classifier. The VC dimension is defined as being the largest possible value of m such that there exists a training set of m different \mathbf{x} points that the classifier can misclassify.

Quantifying the capacity of the model enables statistical learning theory to make quantitative predictions. The most important results in statistical learning theory show that the discrepancy between training error and generalization error is bounded from above by a quantity that grows as the model capacity increases and shrinks as the number of training examples increases ([Vapnik and Chervonenkis, 1971](#); [Vapnik, 1982](#); [Blumer *et al.*, 1989](#); [Vapnik, 1995](#)). These bounds provide an intellectual justification that machine learning algorithms can work well in practice. However, these bounds are rarely used in practice when working with deep learning algorithms. This is in part because the bounds are often quite loose and in part because it is very difficult to determine the capacity of deep learning algorithms. The

determining the capacity of a deep learning model is especially difficult because the effective capacity is limited by the capabilities of the optimization algorithm and we have little theoretical understanding of the general nonconvex optimization problems involved in deep learning.

We must remember that while simpler functions are more likely to have a small gap between training and test error, we must use a sufficiently complex hypothesis to achieve low training error. Typically, the training error decreases until it asymptotes to the minimum possible error and the capacity increases (assuming the error measure has a minimum value).

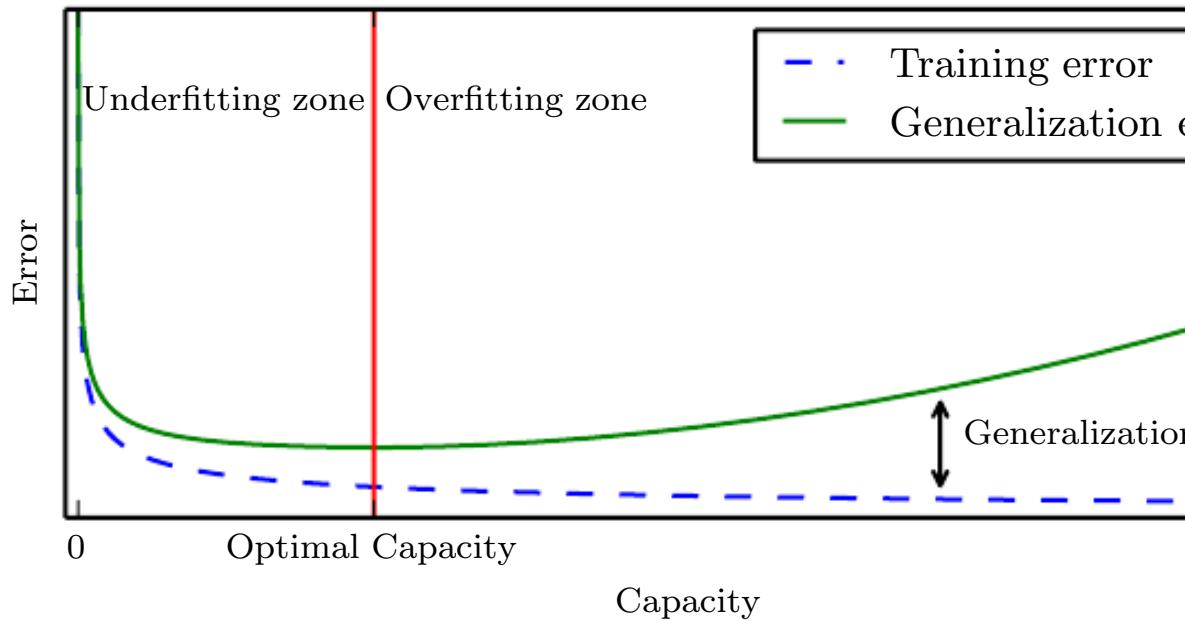


Figure 5.3: Typical relationship between capacity and error. Training error and generalization error behave differently. At the left end of the graph, training error and generalization error are both high. This is the **underfitting regime**. As we increase capacity, training error decreases, but the gap between training and generalization error increases. When the size of this gap outweighs the decrease in training error, and we enter the **overfitting regime**, where capacity is too large, above the **optimal capacity**.

generalization error has a U-shaped curve as a function of model capacity, as illustrated in figure 5.3.

To reach the most extreme case of arbitrarily high capacity, we introduce the concept of **nonparametric models**. So far, we have seen only parametric models, such as linear regression. Parametric models learn a function represented by a parameter vector whose size is finite and fixed before any data is observed. Nonparametric models have no such limitation.

Sometimes, nonparametric models are just theoretical abstractions, such as an algorithm that searches over all possible probability distributions that can be implemented in practice. However, we can also design practical nonparametric models by making their complexity a function of the training set size.

of such an algorithm is **nearest neighbor regression**. Unlike linear regression which has a fixed-length vector of weights, the nearest neighbor regression simply stores the \mathbf{X} and \mathbf{y} from the training set. When asked to predict the value of a point \mathbf{x} , the model looks up the nearest entry in the training set and returns its associated regression target. In other words, $\hat{y} = y_i$ where $i = \arg \min_j d(\mathbf{x}, \mathbf{X}_{j,:})$. The algorithm can also be generalized to distance metrics other than Euclidean such as learned distance metrics (Goldberger *et al.*, 2005). If the algorithm is allowed to break ties by averaging the y_i values for all $\mathbf{X}_{i,:}$ that are equidistant from \mathbf{x} , then this algorithm is able to achieve the minimum possible training error.

might be greater than zero, if two identical inputs are associated with different outputs) on any regression dataset.

Finally, we can also create a nonparametric learning algorithm by nesting a parametric learning algorithm inside another algorithm that increases the number of parameters as needed. For example, we could imagine an outer algorithm that changes the degree of the polynomial learned by linear regression based on a polynomial expansion of the input.

The ideal model is an oracle that simply knows the true probability distribution that generates the data. Even such a model will still incur some error. In practice, there are several reasons for this. First, even in supervised learning, the mapping from \mathbf{x} to y may be inherently noisy or y may be a deterministic function that involves other variables not included in \mathbf{x} . The error incurred by an oracle making predictions according to its own distribution $p(\mathbf{x}, y)$ is called the **Bayes error**.

Training and generalization error vary as the size of the training set increases. Expected generalization error can never increase as the number of training examples increases. For nonparametric models, more data yield better generalization, until the best possible error is achieved. Any fixed parametric model with finite optimal capacity will asymptote to an error value that exceeds the Bayes error. See figure 5.4 for an illustration. Note that it is possible for the model with optimal capacity and yet still have a large gap between training and generalization errors. In this situation, we may be able to reduce this gap by getting more training examples.

5.2.1 The No Free Lunch Theorem

Learning theory claims that a machine learning algorithm can generalize beyond a finite training set of examples. This seems to contradict some basic logic. Inductive reasoning, or inferring general rules from a limited set of examples, is a well-known method of inference.

is not logically valid. To logically infer a rule describing every member of a set, one must have information about every member of that set.

In part, machine learning avoids this problem by offering only *probabilistic* rules rather than the entirely certain rules used in purely logical reasoning. Machine learning promises to find rules that are *probably* correct about members of the set they concern.

Unfortunately, even this does not resolve the entire problem. The **no free lunch theorem** for machine learning (Wolpert, 1996) states that, over all possible data-generating distributions, every classification algorithm has the same performance.

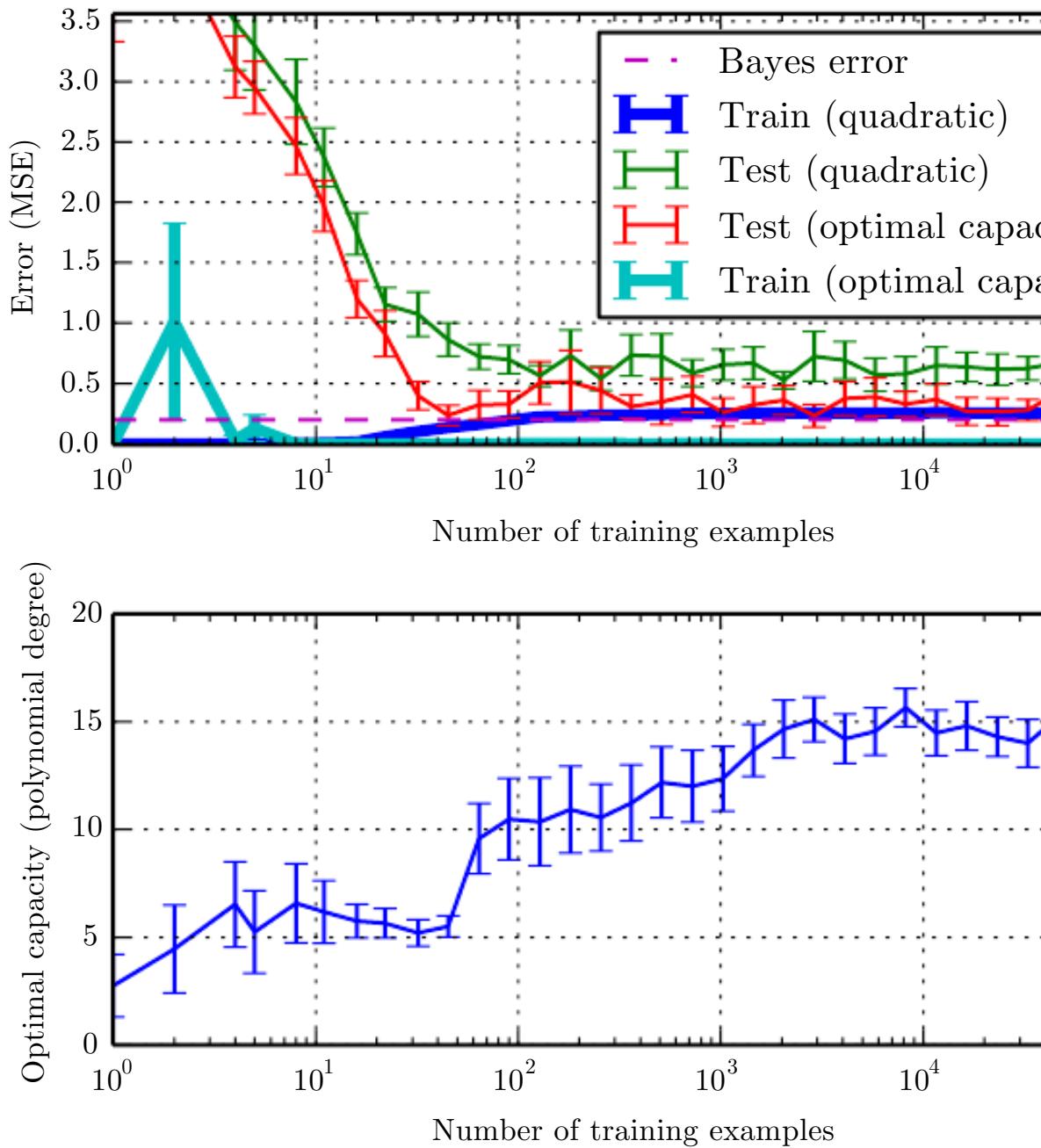


Figure 5.4: The effect of the training dataset size on the train and test error and on the optimal model capacity. We constructed a synthetic regression problem by adding a moderate amount of noise to a degree-5 polynomial, generated and then generated several different sizes of training set. For each size, we generated different training sets in order to plot error bars showing 95 percent confidence intervals.

(Top) The MSE on the training and test set for two different models: a quadratic model and a model with degree chosen to minimize the test error. Both are fit in

the quadratic model, the training error increases as the size of the training set increases. This is because larger datasets are harder to fit. Simultaneously, the test error decreases because fewer incorrect hypotheses are consistent with the training data. If the model does not have enough capacity to solve the task, so its test error will remain at a high value. The test error at optimal capacity asymptotes to the Bayes error. (Top) As the training set size increases, the training error can fall below the Bayes error, due to the ability of the trained model to memorize specific instances of the training set. As the training size increases, the training error of any fixed-capacity model (here, the quadratic model) asymptotically approaches at least the Bayes error. (Bottom) As the training set size increases, the degree of the model (shown here as the degree of the optimal polynomial regressor) increases until the capacity plateaus after reaching sufficient complexity to solve the task.

same error rate when classifying previously unobserved points. In some sense, no machine learning algorithm is universally any better than any other. The most sophisticated algorithm we can conceive of has the same performance (over all possible tasks) as merely predicting that every point belongs to the same class.

Fortunately, these results hold only when we average over *all* possible generating distributions. If we make assumptions about the kind of distributions we encounter in real-world applications, then we can find algorithms that perform well on these distributions.

This means that the goal of machine learning research is not to find the best learning algorithm or the absolute best learning algorithm. Instead, we must understand what kinds of distributions are relevant to the “real world” agent experiences, and what kinds of machine learning algorithms learn from data drawn from the kinds of data-generating distributions we care about.

5.2.2 Regularization

The no free lunch theorem implies that we must design our machine learning algorithms to perform well on a specific task. We do so by building specific preferences into the learning algorithm. When these preferences are taken into account, the learning problems that we ask the algorithm to solve, it performs well.

So far, the only method of modifying a learning algorithm that we have concretely is to increase or decrease the model’s representational capacity, or removing functions from the hypothesis space of solutions the learner is able to choose from. We gave the specific example of increasing the degree of a polynomial for a regression problem. The view we have had so far is oversimplified.

The behavior of our algorithm is strongly affected not just by how we make the set of functions allowed in its hypothesis space, but by the way we weight those functions.

of those functions. The learning algorithm we have studied so far, linear regression, has a hypothesis space consisting of the set of linear functions of inputs. Linear functions can be useful for problems where the relationship between inputs and outputs truly is close to linear. They are less useful for problems where the relationship is not linear. For example, linear regression would not be a good choice if we tried to use it to predict $\sin(x)$ from x . We can thus control the flexibility of our algorithms by choosing what kind of functions we allow them to learn from, as well as by controlling the amount of these functions.

We can also give a learning algorithm a preference for one solution over another.

in its hypothesis space. This means that both functions are eligible and preferred. The unpreferred solution will be chosen only if it fits the training data significantly better than the preferred solution.

For example, we can modify the training criterion for linear regression to include **weight decay**. To perform linear regression with weight decay, we can minimize the cost function $J(\mathbf{w})$ comprising both the mean squared error on the training and testing data, and a term that expresses a preference for the weights to have smaller squared L^2 norm:

$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^\top \mathbf{w},$$

where λ is a value chosen ahead of time that controls the strength of the preference for smaller weights. When $\lambda = 0$, we impose no preference, and larger weights are allowed. As λ increases, we prefer smaller weights to become smaller. Minimizing $J(\mathbf{w})$ results in a choice of weights that make a tradeoff between fitting the training data and being small. It can result in solutions that have a smaller slope, or that put weight on fewer features. As an example of how we can control a model's tendency to overfit by expressing a preference via weight decay, we can train a high-degree polynomial regression model on the same data with different values of λ . See figure 5.5 for the results.

More generally, we can regularize a model that learns a function by expressing a preference for one function over another by adding a penalty called a **regularizer** to the cost function. In the case of weight decay, the regularizer is $\Omega(\mathbf{w}) = \mathbf{w}^\top \mathbf{w}$. In chapter 7, we will see that many other types of regularizers are possible.

Expressing preferences for one function over another is a more general way of controlling a model's capacity than including or excluding members of a hypothesis space. We can think of excluding a function from a hypothesis space as expressing an infinitely strong preference against that function.

In our weight decay example, we expressed our preference for smaller weights by defining a preference for functions defined with smaller weights explicitly, via an extra term in the cost function that we minimize. There are many other ways of expressing preferences between different solutions, both implicitly and explicitly. Together, these different ways of expressing preferences are called **regularization**.

are known as **regularization**. *Regularization is any modification learning algorithm that is intended to reduce its generalization error.* Regularization is one of the central concerns of the field of learning, rivaled in its importance only by optimization.

The no free lunch theorem has made it clear that there is no single best learning algorithm, and, in particular, no best form of regularization. Instead, we must choose a form of regularization that is well suited to the specific problem or task that we want to solve. The philosophy of deep learning in general and of neural networks in particular is that a wide range of tasks (such as all the intellectual challenges of the world) can be solved by a single architecture if it is given enough training data and computation power.

people can do) may all be solved effectively using very general-purpose regularization.

5.3 Hyperparameters and Validation Sets

Most machine learning algorithms have hyperparameters, settings we use to control the algorithm's behavior. The values of hyperparameters are adapted by the learning algorithm itself (though we can design a procedure in which one learning algorithm learns the best hyperparameters for another learning algorithm).

The polynomial regression example in figure 5.2 has a single hyperparameter: the degree of the polynomial, which acts as a **capacity** hyperparameter. A λ value used to control the strength of weight decay is another hyperparameter.

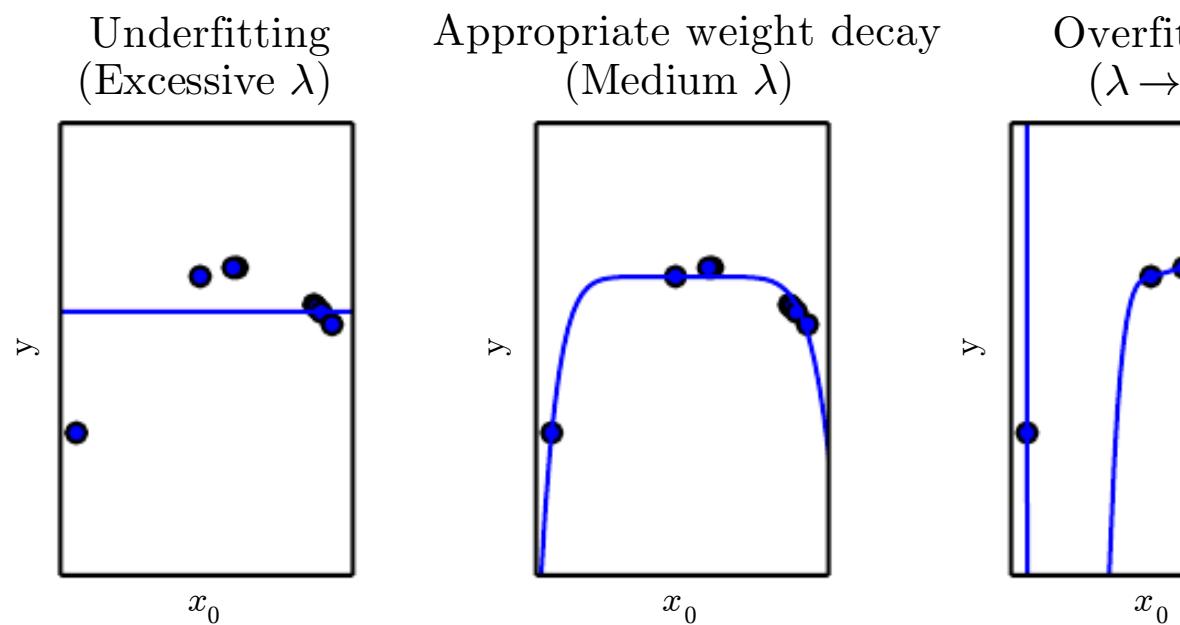


Figure 5.5: We fit a high-degree polynomial regression model to our example from figure 5.2. The true function is quadratic, but here we use only mod-

We vary the amount of weight decay to prevent these high-degree models. (*Left*) With very large λ , we can force the model to learn a function that is flat. This underfits because it can only represent a constant function. With a medium value of λ , the learning algorithm recovers a curve with the right overall shape, even though the model is capable of representing functions with much more complex shapes. (*Right*) With weight decay approaching zero (i.e., using the pseudoinverse to solve the underdetermined problem with minimal regularization), a degree-9 polynomial overfits significantly, as we saw in figure 5.2.

Sometimes a setting is chosen to be a hyperparameter that the algorithm does not learn because the setting is difficult to optimize. Not all hyperparameters are learned; some must be set by hand. For example, the setting must be a hyperparameter because it is not appropriate to learn the value of a hyperparameter on the training set. This applies to all hyperparameters that control model capacity. If learned on the training set, such hyperparameters will always choose the maximum possible model capacity, resulting in overfitting (see figure 5.3). For example, we can always fit the training set perfectly with a higher-degree polynomial and a weight decay setting of $\lambda = 0$ than with a lower-degree polynomial and a positive weight decay setting.

To solve this problem, we need a **validation set** of examples that the algorithm does not observe.

Earlier we discussed how a held-out test set, composed of examples drawn from the same distribution as the training set, can be used to estimate the generalization error of a learner, after the learning process has completed. It is important to note that the test examples are not used in any way to make choices about the model or its hyperparameters. For this reason, no example from the test set should appear in the validation set. Therefore, we always construct the validation set before splitting the *training data*. Specifically, we split the training data into two disjoint subsets. One of these subsets is used to learn the parameters. The other subset is called the validation set, used to estimate the generalization error during optimization, allowing for the hyperparameters to be updated accordingly. The subset used to learn the parameters is still typically called the training set, even though this may be confused with the larger pool of data used for the learning process. The subset of data used to guide the selection of hyperparameters is called the validation set. Typically, one uses about 80 percent of the data for training and 20 percent for validation. Since the validation set is not used to “train” the hyperparameters, the validation set error will under-estimate the true generalization error, though typically by a smaller amount than the test set error does. After all hyperparameter optimization is complete, the generalization error may be estimated using the test set.

In practice, when the same test set has been used repeatedly to compare the performance of different algorithms over many years, and especially when we consider all the attempts from the scientific community at beating the record for state-of-the-art performance on that test set, we end up having optimistic expectations about the true performance of a trained system. Thankfully, the community can move on to new (and usually more ambitious and larger) benchmarks.

5.3.1 Cross-Validation

Dividing the dataset into a fixed training set and a fixed test set can be problematic if it results in the test set being small. A small test set implies statistical uncertainty around the estimated average test error, making it difficult to claim that algorithm A works better than algorithm B on the given task.

When the dataset has hundreds of thousands of examples or more, this can be a serious issue. When the dataset is too small, alternative procedures must be used to use all the examples in the estimation of the mean test error, which leads to increased computational cost. These procedures are based on the idea of performing the training and testing computation on different randomly chosen subsets of the original dataset. The most common of these is the k -fold cross-validation procedure, shown in algorithm 5.1, in which a partition of the data is obtained by splitting it into k nonoverlapping subsets. The test error may then be estimated by taking the average test error across k trials. On trial i , the i -th subset of data is used as the test set, and the rest of the data is used as the training set. One problem is that no unbiased estimators of the variance of such estimators exist (Bengio and Grandvalet, 2004), but approximation methods are used.

5.4 Estimators, Bias and Variance

The field of statistics gives us many tools to achieve the machine learning goal of solving a task not only on the training set but also to generalize to new data. Concepts such as parameter estimation, bias and variance are used to characterize notions of generalization, underfitting and overfitting.

5.4.1 Point Estimation

Point estimation is the attempt to provide the single “best” predicted quantity of interest. In general the quantity of interest can be a single number or a vector of parameters in some parametric model, such as the linear regression example in section 5.1.4, but it can also be a whole function.

To distinguish estimates of parameters from their true value, we will be to denote a point estimate of a parameter θ by $\hat{\theta}$.

Let $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ be a set of m independent and identically distributed observations.

Algorithm 5.1 The k -fold cross-validation algorithm. It can be used to estimate the generalization error of a learning algorithm A when the given dataset \mathbb{D} is too small for a simple train/test or train/valid split to yield accurate estimates. The generalization error, because the mean of a loss L on a small test set has a high variance. The dataset \mathbb{D} contains as elements the abstract expression $\mathbf{z}^{(i)}$ of the i -th example, which could stand for an (input,target) pair (\mathbf{x}, y) in the case of supervised learning, or for just an input $\mathbf{z}^{(i)} = \mathbf{x}$ in the case of unsupervised learning. The algorithm returns the vector of estimated errors e , one for each example in \mathbb{D} , whose mean is the estimated generalization error. The individual examples can be used to compute a confidence interval around the mean (equation 5.47). Though these confidence intervals are no longer available after the use of cross-validation, it is still common practice to use them to compare two algorithms. We say that algorithm A is better than algorithm B only if the confidence interval of the generalization error of algorithm A lies below and does not intersect the confidence interval of algorithm B .

Define KFoldXV(\mathbb{D}, A, L, k):

Require: \mathbb{D} , the given dataset, with elements $\mathbf{z}^{(i)}$

Require: A , the learning algorithm, seen as a function that takes an input and outputs a learned function

Require: L , the loss function, seen as a function from a learned function f to a scalar $\in \mathbb{R}$ that maps an example $\mathbf{z}^{(i)} \in \mathbb{D}$ to a scalar $\in \mathbb{R}$

Require: k , the number of folds

Split \mathbb{D} into k mutually exclusive subsets \mathbb{D}_i , whose union is \mathbb{D}

for i from 1 to k **do**

$f_i = A(\mathbb{D} \setminus \mathbb{D}_i)$

for $\mathbf{z}^{(j)}$ in \mathbb{D}_i **do**

$e_j = L(f_i, \mathbf{z}^{(j)})$

end for

end for

Return e

(i.i.d.) data points. A **point estimator** or **statistic** is any function

$$\hat{\theta}_m = g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}).$$

The definition does not require that g return a value that is close to the true parameter value or even that the range of g be the same as the set of allowable values. This definition of a point estimator is very general and would enable the estimator great flexibility. While almost any function thus qualifies

a good estimator is a function whose output is close to the true unknown value that generated the training data.

For now, we take the frequentist perspective on statistics. That is, we assume that the true parameter value θ is fixed but unknown, while the estimator $\hat{\theta}$ is a function of the data. Since the data is drawn from a random process, the estimator $\hat{\theta}$ is a function of the data is random. Therefore $\hat{\theta}$ is a random variable.

Point estimation can also refer to the estimation of the relationship between input and target variables. We refer to these types of point estimators as function estimators.

Function Estimation Sometimes we are interested in performing function estimation (or function approximation). Here, we are trying to predict a target variable y given an input vector x . We assume that there is a function $f(x)$ that describes the approximate relationship between y and x . For example, we might assume that $y = f(x) + \epsilon$, where ϵ stands for the part of y that is not predictable by x . In function estimation, we are interested in approximating f with a function \hat{f} . Function estimation is really just the same as estimating a parameter θ ; the function estimator \hat{f} is simply a point estimator in function estimation. The linear regression example (discussed in section 5.1.4) and the polynomial regression example (discussed in section 5.2) both illustrate scenarios that may be considered function estimation, as either estimating a parameter w or estimating a function \hat{f} that maps x to y .

We now review the most commonly studied properties of point estimators and discuss what they tell us about these estimators.

5.4.2 Bias

The bias of an estimator is defined as

$$\text{bias}(\hat{\theta}_m) = \mathbb{E}(\hat{\theta}_m) - \theta$$

$$\text{bias}(\boldsymbol{\theta}^m) = \mathbb{E}(\boldsymbol{\theta}^m) - \boldsymbol{\theta},$$

where the expectation is over the data (seen as samples from a random variable) and $\boldsymbol{\theta}$ is the true underlying value of $\boldsymbol{\theta}$ used to define the data-generating distribution. An estimator $\hat{\boldsymbol{\theta}}_m$ is said to be **unbiased** if $\text{bias}(\hat{\boldsymbol{\theta}}_m) = \mathbf{0}$, that is, $\mathbb{E}(\hat{\boldsymbol{\theta}}_m) = \boldsymbol{\theta}$. An estimator $\hat{\boldsymbol{\theta}}_m$ is said to be **asymptotically unbiased** if $\lim_{m \rightarrow \infty} \text{bias}(\hat{\boldsymbol{\theta}}_m) = \mathbf{0}$, which implies that $\lim_{m \rightarrow \infty} \mathbb{E}(\hat{\boldsymbol{\theta}}_m) = \boldsymbol{\theta}$.

Example: Bernoulli Distribution Consider a set of samples that are independently and identically distributed according to a Bernoulli distribution.

bution with mean θ :

$$P(x^{(i)}; \theta) = \theta^{x^{(i)}} (1 - \theta)^{(1-x^{(i)})}.$$

A common estimator for the θ parameter of this distribution is the training samples:

$$\hat{\theta}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}.$$

To determine whether this estimator is biased, we can substitute it into equation 5.20:

$$\begin{aligned} \text{bias}(\hat{\theta}_m) &= \mathbb{E}[\hat{\theta}_m] - \theta \\ &= \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m x^{(i)}\right] - \theta \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}[x^{(i)}] - \theta \\ &= \frac{1}{m} \sum_{i=1}^m \sum_{x^{(i)}=0}^1 \left(x^{(i)} \theta^{x^{(i)}} (1 - \theta)^{(1-x^{(i)})}\right) - \theta \\ &= \frac{1}{m} \sum_{i=1}^m (\theta) - \theta \\ &= \theta - \theta = 0 \end{aligned}$$

Since $\text{bias}(\hat{\theta}) = 0$, we say that our estimator $\hat{\theta}$ is unbiased.

Example: Gaussian Distribution Estimator of the Mean
 a set of samples $x^{(1)}, \dots, x^{(m)}$ that are independently and identically distributed according to a Gaussian distribution $p(x^{(i)}) = \mathcal{N}(x^{(i)}; \mu, \sigma^2)$, where

Recall that the Gaussian probability density function is given by

$$p(x^{(i)}; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(x^{(i)} - \mu)^2}{\sigma^2}\right).$$

A common estimator of the Gaussian mean parameter is known **mean**:

$$\hat{\mu}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

123

To determine the bias of the sample mean, we are again interested in its expectation:

$$\begin{aligned}
 \text{bias}(\hat{\mu}_m) &= \mathbb{E}[\hat{\mu}_m] - \mu \\
 &= \mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m x^{(i)}\right] - \mu \\
 &= \left(\frac{1}{m} \sum_{i=1}^m \mathbb{E}[x^{(i)}]\right) - \mu \\
 &= \left(\frac{1}{m} \sum_{i=1}^m \mu\right) - \mu \\
 &= \mu - \mu = 0
 \end{aligned}$$

Thus we find that the sample mean is an unbiased estimator of the parameter.

Example: Estimators of the Variance of a Gaussian Distribution
 In this example, we compare two different estimators of the variance parameter of a Gaussian distribution. We are interested in knowing if either estimator is unbiased.

The first estimator of σ^2 we consider is known as the **sample variance**:

$$\hat{\sigma}_m^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} - \hat{\mu}_m^2,$$

where $\hat{\mu}_m$ is the sample mean. More formally, we are interested in

$$\text{bias}(\hat{\sigma}_m^2) = \mathbb{E}[\hat{\sigma}_m^2] - \sigma^2.$$

We begin by evaluating the term $\mathbb{E}[\hat{\sigma}_m^2]$:

$$\begin{aligned}\mathbb{E}[\hat{\sigma}_m^2] &= \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \left(x^{(i)} - \hat{\mu}_m \right)^2 \right] \\ &= \frac{m-1}{m} \sigma^2\end{aligned}$$

Returning to equation 5.37, we conclude that the bias of $\hat{\sigma}_m^2$ is $-\sigma^2$. The sample variance is a biased estimator.

The **unbiased sample variance** estimator

$$\tilde{\sigma}_m^2 = \frac{1}{m-1} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2$$

provides an alternative approach. As the name suggests this estimator is unbiased. That is, we find that $\mathbb{E}[\tilde{\sigma}_m^2] = \sigma^2$:

$$\begin{aligned}\mathbb{E}[\tilde{\sigma}_m^2] &= \mathbb{E} \left[\frac{1}{m-1} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2 \right] \\ &= \frac{m}{m-1} \mathbb{E}[\hat{\sigma}_m^2] \\ &= \frac{m}{m-1} \left(\frac{m-1}{m} \sigma^2 \right) \\ &= \sigma^2.\end{aligned}$$

We have two estimators: one is biased, and the other is not. While unbiased estimators are clearly desirable, they are not always the “best” estimators. In fact, we will see we often use biased estimators that possess other important properties.

5.4.3 Variance and Standard Error

Another property of the estimator that we might want to consider is its variance. Just as we expect the mean to vary as a function of the data sample, we expect the expectation of the estimator to determine its bias, we can compute its variance. The **variance** of an estimator is simply the variance of its expectation:

$$\text{Var}(\hat{\theta})$$

where the random variable is the training set. Alternately, the square root of the variance is called the **standard error**.

variance is called the **standard error**, denoted $\text{SE}(\theta)$.

The variance, or the standard error, of an estimator provides a measure of how much we would expect the estimate we compute from data to vary as we were to resample the dataset from the underlying data-generating process. While we might like an estimator to exhibit low bias, we would also like it to have low variance.

When we compute any statistic using a finite number of samples, our estimate of the true underlying parameter is uncertain, in the sense that if we had obtained other samples from the same distribution and their statistics

been different. The expected degree of variation in any estimator is the error that we want to quantify.

The standard error of the mean is given by

$$\text{SE}(\hat{\mu}_m) = \sqrt{\text{Var}\left[\frac{1}{m} \sum_{i=1}^m x^{(i)}\right]} = \frac{\sigma}{\sqrt{m}},$$

where σ^2 is the true variance of the samples x^i . The standard deviation is estimated by using an estimate of σ . Unfortunately, neither the sample variance nor the square root of the unbiased estimator provide an unbiased estimate of the standard deviation. Both appear to underestimate the true standard deviation but are still used in practice because the square root of the unbiased estimator of the variance is less of an approximation. For large m , the approximation is quite reasonable.

The standard error of the mean is very useful in machine learning. We often estimate the generalization error by computing the sample error on the test set. The number of examples in the test set affects the accuracy of this estimate. Taking advantage of the central limit theorem tells us that the mean will be approximately distributed with a normal distribution. Thus, we can use the standard error to compute the probability that the test error falls in any chosen interval. For example, the 95 percent confidence interval on the mean $\hat{\mu}_m$ is

$$(\hat{\mu}_m - 1.96\text{SE}(\hat{\mu}_m), \hat{\mu}_m + 1.96\text{SE}(\hat{\mu}_m)),$$

under the normal distribution with mean $\hat{\mu}_m$ and variance $\text{SE}(\hat{\mu}_m)^2$. In machine learning experiments, it is common to say that algorithm A is better than algorithm B if the upper bound of the 95 percent confidence interval for the error of algorithm A is less than the lower bound of the 95 percent confidence interval for the error of algorithm B .

of algorithm B .

Example: Bernoulli Distribution We once again consider a $\{x^{(1)}, \dots, x^{(m)}\}$ drawn independently and identically from a Bernoulli distribution (recall $P(x^{(i)}; \theta) = \theta^{x^{(i)}}(1 - \theta)^{1-x^{(i)}}$). This time we are interested in calculating the variance of the estimator $\hat{\theta}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$.

$$\text{Var}(\hat{\theta}_m) = \text{Var}\left(\frac{1}{m} \sum_{i=1}^m x^{(i)}\right)$$

126

$$\begin{aligned} &= \frac{1}{m^2} \sum_{i=1}^m \text{Var}(x^{(i)}) \\ &= \frac{1}{m^2} \sum_{i=1}^m \theta(1 - \theta) \\ &= \frac{1}{m^2} m\theta(1 - \theta) \\ &= \frac{1}{m} \theta(1 - \theta) \end{aligned}$$

The variance of the estimator decreases as a function of m , the number of observations in the dataset. This is a common property of popular estimators, which we will return to when we discuss consistency (see section 5.4.5).

5.4.4 Trading off Bias and Variance to Minimize Mean Squared Error

Bias and variance measure two different sources of error in an estimator. Bias measures the expected deviation from the true value of the function. Variance on the other hand, provides a measure of the deviation from the estimator value that any particular sampling of the data is likely to produce.

What happens when we are given a choice between two estimators, one with more bias and one with more variance? How do we choose between them? For example, imagine that we are interested in approximating the function shown in figure 5.2 and we are only offered the choice between a model with low bias but high variance and one that suffers from large variance. How do we choose between the two?

The most common way to negotiate this trade-off is to use cross-validation. Empirically, cross-validation is highly successful on many real-world problems. Alternatively, we can also compare the **mean squared error** (MSE) of the two

$$\text{MSE} = \mathbb{E}[(\hat{\theta}_m - \theta)^2]$$

$$= \text{Bias}(\hat{\theta}^m)^2 + \text{Var}(\hat{\theta}^m)$$

The MSE measures the overall expected deviation—in a squared sense—between the estimator and the true value of the parameter θ . As shown in equation 5.54, evaluating the MSE incorporates both the bias and the variance. Desirable estimators are those with small MSE and these are often called unbiased estimators. It is important to manage to keep both their bias and variance somewhat in check.

The relationship between bias and variance is tightly linked to the machine learning concepts of capacity, underfitting and overfitting. When the model has low capacity, it underfits the training data. When the model has high capacity, it overfits the training data.

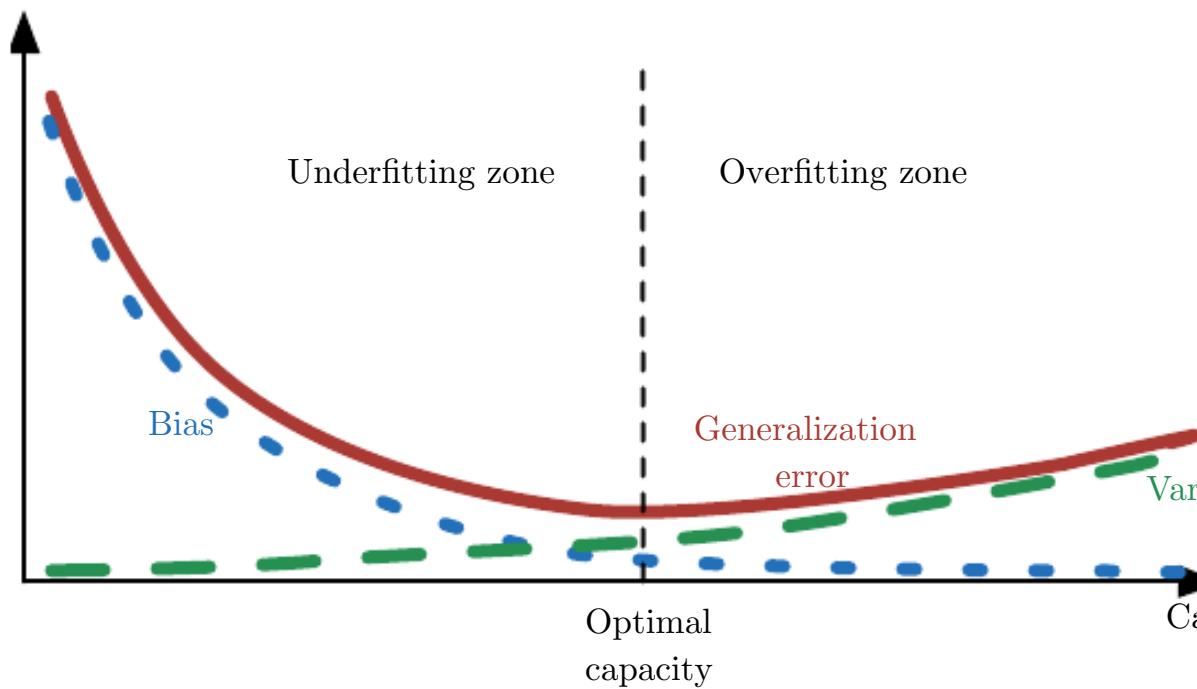


Figure 5.6: As capacity increases (x -axis), bias (dotted) tends to decrease (dashed) tends to increase, yielding another U-shaped curve for generalization error. If we vary capacity along one axis, there is an optimal capacity, when the capacity is below this optimum and overfitting when it is above. This is similar to the relationship between capacity, underfitting, and overfitting section 5.2 and figure 5.3.

error is measured by the MSE (where bias and variance are meaningful components of generalization error), increasing capacity tends to increase variance rather than bias. This is illustrated in figure 5.6, where we see again the U-shaped nature of generalization error as a function of capacity.

5.4.5 Consistency

So far we have discussed the properties of various estimators for a fixed size. Usually, we are also concerned with the behavior of an estimator as the amount of training data grows. In particular, we usually wish that

of data points m in our dataset increases, our point estimates converge to the true value of the corresponding parameters. More formally, we would like to say that

$$\text{plim}_{m \rightarrow \infty} \hat{\theta}_m = \theta.$$

The symbol plim indicates convergence in probability, meaning that $P(|\hat{\theta}_m - \theta| > \epsilon) \rightarrow 0$ as $m \rightarrow \infty$. The condition described by equation (1) is known as **consistency**. It is sometimes referred to as weak consistency, and it is also called strong consistency referring to the **almost sure** convergence of $\hat{\theta}_m$.

sure convergence of a sequence of random variables $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$ occurs when $p(\lim_{m \rightarrow \infty} \mathbf{x}^{(m)} = \mathbf{x}) = 1$.

Consistency ensures that the bias induced by the estimator decreases as the number of data examples grows. However, the reverse is not true: unbiasedness does not imply consistency. For example, consider the mean parameter μ of a normal distribution $\mathcal{N}(x; \mu, \sigma^2)$, with a dataset of m samples: $\{x^{(1)}, \dots, x^{(m)}\}$. We could use the first sample $x^{(1)}$ as an unbiased estimator: $\hat{\theta} = x^{(1)}$. In that case, $\mathbb{E}(\hat{\theta}_m) = \theta$, so it is unbiased no matter how many data points are seen. This, of course, does not mean that the estimate is asymptotically unbiased. However, this is not the case for the first sample as it is *not* the case that $\hat{\theta}_m \rightarrow \theta$ as $m \rightarrow \infty$.

5.5 Maximum Likelihood Estimation

We have seen some definitions of common estimators and analyzed their properties. But where did these estimators come from? Rather than guessing which function might make a good estimator and then analyzing its bias and variance, we would like to have some principle from which we can derive specific estimators that are good estimators for different models.

The most common such principle is the maximum likelihood principle.

Consider a set of m examples $\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ drawn independently from the true but unknown data-generating distribution $p_{\text{data}}(\mathbf{x})$.

Let $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$ be a parametric family of probability distributions over the same space indexed by $\boldsymbol{\theta}$. In other words, $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$ maps any dataset \mathbf{x} to a real number estimating the true probability $p_{\text{data}}(\mathbf{x})$.

The maximum likelihood estimator for $\boldsymbol{\theta}$ is then defined as

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}),$$

$$= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}).$$

This product over many probabilities can be inconvenient for computation.
For example, it is prone to numerical underflow. To obtain a more numerically stable but equivalent optimization problem, we observe that taking the logarithm of the likelihood does not change its arg max but does conveniently transform the product into a sum.

into a sum:

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}).$$

Because the $\arg \max$ does not change when we rescale the cost function, we can divide by m to obtain a version of the criterion that is expressed as an average with respect to the empirical distribution \hat{p}_{data} defined by the training set:

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}).$$

One way to interpret maximum likelihood estimation is to view it as a measure of the dissimilarity between the empirical distribution \hat{p}_{data} , defined by the training set and the model distribution, with the degree of dissimilarity being measured by the KL divergence. The KL divergence is given by

$$D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

The term on the left is a function only of the data-generating process and the model. This means when we train the model to minimize the KL divergence, we need only minimize

$$-\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})],$$

which is of course the same as the maximization in equation 5.59.

Minimizing this KL divergence corresponds exactly to minimizing the cross-entropy between the distributions. Many authors use the term “cross-entropy loss” to identify specifically the negative log-likelihood of a Bernoulli or softmax distribution, but that is a misnomer. Any loss consisting of a negative log-likelihood plus the cross-entropy between the empirical distribution defined by the training set and the probability distribution defined by the model. For example, mean squared error is a loss function that measures the cross-entropy between the empirical distribution and a Gaussian model.

We can thus see maximum likelihood as an attempt to make the model’s probability distribution match the empirical distribution \hat{p}_{data} . Ideally, we would like to make the model’s distribution equal to the true data-generating distribution p_{data} , but we have no direct access to this distribution.

distribution.

While the optimal θ is the same regardless of whether we are maximizing likelihood or minimizing the KL divergence, the values of the objective are different. In software, we often phrase both as minimizing a loss function. Maximum likelihood thus becomes minimization of the negative log-likelihood (NLL), or equivalently, minimization of the cross-entropy. The maximum likelihood as minimum KL divergence becomes helpful because the KL divergence has a known minimum value of zero. Note that the negative log-likelihood can actually become negative when x is real-valued.

5.5.1 Conditional Log-Likelihood and Mean Squared

The maximum likelihood estimator can readily be generalized to estimate conditional probability $P(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta})$ in order to predict \mathbf{y} given \mathbf{x} . This is the most common situation because it forms the basis for most supervised learning algorithms. If \mathbf{X} represents all our inputs and \mathbf{Y} all our observed targets, then the maximum likelihood estimator is

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{Y} \mid \mathbf{X}; \boldsymbol{\theta}).$$

If the examples are assumed to be i.i.d., then this can be decomposed into

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log P(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}; \boldsymbol{\theta}).$$

Example: Linear Regression as Maximum Likelihood Linear regression, introduced in section 5.1.4, may be justified as a maximum likelihood estimation problem. Previously, we motivated linear regression as an algorithm that learns a function \hat{y} from input \mathbf{x} and produce an output value \hat{y} . The mapping from \mathbf{x} to \hat{y} was justified by minimizing mean squared error, a criterion that we introduced more generally in section 5.1.4. We now revisit linear regression from the point of view of maximum likelihood estimation. Instead of producing a single prediction \hat{y} , we now think of the algorithm as producing a conditional distribution $p(y \mid \mathbf{x})$. We can imagine that, given an infinitely large training set, we might see several training examples with the same input value \mathbf{x} but different values of y . The goal of the learning algorithm is to fit the distribution $p(y \mid \mathbf{x})$ to all those different y values that are associated with \mathbf{x} . To derive the same linear regression algorithm we obtain the conditional distribution by defining $p(y \mid \mathbf{x}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \mathbf{w}), \sigma^2)$. The function $\hat{y}(\mathbf{x}; \mathbf{w})$ gives the mean of the Gaussian. In this example, we assume that the variance is some constant σ^2 chosen by the user. We will see that this choice of

form of $p(y | \boldsymbol{x})$ causes the maximum likelihood estimation proceed same learning algorithm as we developed before. Since the example to be i.i.d., the conditional log-likelihood (equation 5.63) is given

$$\begin{aligned} & \sum_{i=1}^m \log p(y^{(i)} | \boldsymbol{x}^{(i)}; \boldsymbol{\theta}) \\ &= -m \log \sigma - \frac{m}{2} \log(2\pi) - \sum_{i=1}^m \frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2}, \end{aligned}$$

where $\hat{y}^{(i)}$ is the output of the linear regression on the i -th input \mathbf{x} number of the training examples. Comparing the log-likelihood squared error,

$$\text{MSE}_{\text{train}} = \frac{1}{m} \sum_{i=1}^m \|\hat{y}^{(i)} - y^{(i)}\|^2,$$

we immediately see that maximizing the log-likelihood with respect to the same estimate of the parameters \mathbf{w} as does minimizing the mean squared error. The two criteria have different values but the same location of the minimum, which justifies the use of the MSE as a maximum likelihood estimation procedure. As we will see, the maximum likelihood estimator has several desirable properties.

5.5.2 Properties of Maximum Likelihood

The main appeal of the maximum likelihood estimator is that it can be the best estimator asymptotically, as the number of examples m increases, provided that we know the true distribution of the data and its rate of convergence as m increases.

Under appropriate conditions, the maximum likelihood estimator has the property of consistency (see section 5.4.5), meaning that as the number of examples approaches infinity, the maximum likelihood estimate converges to the true value of the parameter. These conditions are:

- The true distribution p_{data} must lie within the model family. Otherwise, no estimator can recover p_{data} .
- The true distribution p_{data} must correspond to exactly one value of θ . Otherwise, maximum likelihood can recover the correct p_{data} but not necessarily the correct value of θ to determine which value of θ was used by the data-generating process.

There are other inductive principles besides the maximum likelihood principle, many of which share the property of being consistent estimators.

estimators can differ, however, in their **statistical efficiency**, meaning that a consistent estimator may obtain lower generalization error for a fixed number of samples m , or equivalently, may require fewer examples to obtain a fixed level of generalization error.

Statistical efficiency is typically studied in the **parametric case** (e.g., linear regression), where our goal is to estimate the value of a parameter (e.g., the slope of a line, for which it is known that it is possible to identify the true parameter), not the value of a function. A common way to measure how close we are to the true parameter is by the expected squared prediction error, computing the squared difference between the estimated and true values of the function.

values, where the expectation is over m training samples from the distribution. That parametric mean squared error decreases as m for m large, the Cramér-Rao lower bound (Rao, 1945; Cramér, 1946) no consistent estimator has a lower MSE than the maximum likelihood estimator.

For these reasons (consistency and efficiency), maximum likelihood has been considered the preferred estimator to use for machine learning. While the number of examples is small enough to yield overfitting behavior, regularization techniques such as weight decay may be used to obtain a biased version of maximum likelihood that has less variance when training data is limited.

5.6 Bayesian Statistics

So far we have discussed **frequentist statistics** and approaches based on fitting a single value of θ , then making all predictions thereafter based on this point estimate. Another approach is to consider all possible values of θ and their associated probability of prediction. The latter is the domain of **Bayesian statistics**.

As discussed in section 5.4.1, the frequentist perspective is that the true parameter value θ is fixed but unknown, while the point estimate $\hat{\theta}$ is a random variable on account of it being a function of the dataset (which is subject to noise).

The Bayesian perspective on statistics is quite different. The key difference is that probabilities are assigned to parameters based on their probability to reflect degrees of certainty in states of knowledge. The true value of θ is directly observed and so is not random. On the other hand, the true value of θ is unknown or uncertain and thus is represented as a random variable θ .

Before observing the data, we represent our knowledge of θ via a **prior probability distribution**, $p(\theta)$ (sometimes referred to as simply $p(\theta)$). Generally, the machine learning practitioner selects a prior distribution that is quite broad (i.e., with high entropy) to reflect a high degree of uncertainty about the true value of θ before observing any data. For example, one might assume a uniform prior distribution for θ if no information is available.

$\boldsymbol{\theta}$ lies in some finite range or volume, with a uniform distribution instead reflect a preference for “simpler” solutions (such as small coefficients, or a function that is closer to being constant).

Now consider that we have a set of data samples $\{x^{(1)}, \dots, x^{(m)}\}$, recover the effect of data on our belief about $\boldsymbol{\theta}$ by combining the $p(x^{(1)}, \dots, x^{(m)} | \boldsymbol{\theta})$ with the prior via Bayes’ rule:

$$p(\boldsymbol{\theta} | x^{(1)}, \dots, x^{(m)}) = \frac{p(x^{(1)}, \dots, x^{(m)} | \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(x^{(1)}, \dots, x^{(m)})}$$

In the scenarios where Bayesian estimation is typically used, the prior is often a relatively uniform or Gaussian distribution with high entropy, and the arrival of the data usually causes the posterior to lose entropy and concentrate on a few highly likely values of the parameters.

Relative to maximum likelihood estimation, Bayesian estimation has two important differences. First, unlike the maximum likelihood approach which makes predictions using a point estimate of $\boldsymbol{\theta}$, the Bayesian approach is to make predictions using a full distribution over $\boldsymbol{\theta}$. For example, after observing m data samples, the predicted distribution over the next data sample, $x^{(m+1)}$, is given by

$$p(x^{(m+1)} | x^{(1)}, \dots, x^{(m)}) = \int p(x^{(m+1)} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | x^{(1)}, \dots, x^{(m)}) d\boldsymbol{\theta}$$

Here each value of $\boldsymbol{\theta}$ with positive probability density contributes to the prediction of the next example, with the contribution weighted by the posterior probability of $\boldsymbol{\theta}$. After having observed $\{x^{(1)}, \dots, x^{(m)}\}$, if we are still quite uncertain about the true value of $\boldsymbol{\theta}$, then this uncertainty is incorporated directly into any prediction we might make.

In section 5.4, we discussed how the frequentist approach addressed the question of uncertainty in a given point estimate of $\boldsymbol{\theta}$ by evaluating its variance. The Bayesian answer to the question of uncertainty in an estimator is to simply integrate over it, averaging the uncertainty in the estimator with the uncertainty in the estimator. This integral is of course just another application of the laws of probability, making the Bayesian approach simple to justify. The frequentist machinery for constructing an estimator is based on the idea that the decision to summarize all knowledge contained in the dataset will lead to a good estimate.

The second important difference between the Bayesian approach and the maximum likelihood approach is due to the contribution of the prior distribution.

prior distribution. The prior has an influence by shifting probability towards regions of the parameter space that are preferred a priori. The prior often expresses a preference for models that are simpler or more parsimonious. Critics of the Bayesian approach identify the prior as a source of subjective judgment affecting the predictions.

Bayesian methods typically generalize much better when limited labeled training data is available but typically suffer from high computational cost when the number of training examples is large.

Example: Bayesian Linear Regression Here we consider the Bayesian approach to learning the linear regression parameters. In linear regression, we learn a linear mapping from an input vector $\mathbf{x} \in \mathbb{R}^n$ to predict a scalar $y \in \mathbb{R}$. The prediction is parametrized by the vector $\mathbf{w} \in \mathbb{R}^n$

$$\hat{y} = \mathbf{w}^\top \mathbf{x}.$$

Given a set of m training samples $(\mathbf{X}^{(\text{train})}, \mathbf{y}^{(\text{train})})$, we can express the prediction \hat{y} over the entire training set as

$$\hat{\mathbf{y}}^{(\text{train})} = \mathbf{X}^{(\text{train})} \mathbf{w}.$$

Expressed as a Gaussian conditional distribution on $\mathbf{y}^{(\text{train})}$, we have

$$p(\mathbf{y}^{(\text{train})} | \mathbf{X}^{(\text{train})}, \mathbf{w}) = \mathcal{N}(\mathbf{y}^{(\text{train})}; \mathbf{X}^{(\text{train})} \mathbf{w}, \mathbf{I})$$

$$\propto \exp\left(-\frac{1}{2}(\mathbf{y}^{(\text{train})} - \mathbf{X}^{(\text{train})} \mathbf{w})^\top (\mathbf{y}^{(\text{train})} - \mathbf{X}^{(\text{train})} \mathbf{w})\right)$$

where we follow the standard MSE formulation in assuming that the variance on y is one. In what follows, to reduce the notational burden, we will denote the training data $(\mathbf{X}^{(\text{train})}, \mathbf{y}^{(\text{train})})$ as simply (\mathbf{X}, \mathbf{y}) .

To determine the posterior distribution over the model parameters, we first need to specify a prior distribution. The prior should reflect our beliefs about the value of these parameters. While it is sometimes difficult to express our prior beliefs in terms of the parameters of the model, we typically assume a fairly broad distribution, expressing a high degree of uncertainty about θ . For real-valued parameters it is common to use a Gaussian distribution,

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}^0, \boldsymbol{\Lambda}_0) \propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}^0)^\top \boldsymbol{\Lambda}_0^{-1}(\mathbf{w} - \boldsymbol{\mu}^0)\right)$$

where μ^0 and Λ^0 are the prior distribution mean vector and covariance matrix respectively.

With the prior thus specified, we can now proceed in determining the posterior *distribution* over the model parameters:

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) \propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w})p(\mathbf{w})$$

¹Unless there is a reason to use a particular covariance structure, we typically use a diagonal covariance matrix $\Lambda_0 = \text{diag}(\boldsymbol{\lambda}_0)$.

$$\propto \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^\top(\mathbf{y} - \mathbf{X}\mathbf{w})\right) \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_0)^\top(\mathbf{w} - \boldsymbol{\mu}_0)\right)$$

$$\propto \exp\left(-\frac{1}{2}\left(-2\mathbf{y}^\top\mathbf{X}\mathbf{w} + \mathbf{w}^\top\mathbf{X}^\top\mathbf{X}\mathbf{w} + \mathbf{w}^\top\boldsymbol{\Lambda}_0^{-1}\mathbf{w} - 2\mathbf{w}^\top\boldsymbol{\mu}_0\right)\right)$$

We now define $\boldsymbol{\Lambda}_m = (\mathbf{X}^\top\mathbf{X} + \boldsymbol{\Lambda}_0^{-1})^{-1}$ and $\boldsymbol{\mu}_m = \boldsymbol{\Lambda}_m(\mathbf{X}^\top\mathbf{y})$. Using these new variables, we find that the posterior may be rewritten as a Gaussian distribution:

$$\begin{aligned} p(\mathbf{w} | \mathbf{X}, \mathbf{y}) &\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_m)^\top\boldsymbol{\Lambda}_m^{-1}(\mathbf{w} - \boldsymbol{\mu}_m) + \frac{1}{2}\boldsymbol{\mu}_m^\top\boldsymbol{\Lambda}_m^{-1}\boldsymbol{\mu}_m\right) \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_m)^\top\boldsymbol{\Lambda}_m^{-1}(\mathbf{w} - \boldsymbol{\mu}_m)\right). \end{aligned}$$

All terms that do not include the parameter vector \mathbf{w} have been omitted. These terms are implied by the fact that the distribution must be normalized to unity. Equation 3.23 shows how to normalize a multivariate Gaussian distribution.

Examining this posterior distribution enables us to gain some insight into the effect of Bayesian inference. In most situations, we set $\boldsymbol{\mu}_0$ to $\mathbf{0}$. If we do this, then $\boldsymbol{\mu}_m$ gives the same estimate of \mathbf{w} as does frequentist linear regression. There is also a weight decay penalty of $\alpha\mathbf{w}^\top\mathbf{w}$. One difference is that the Bayesian estimate is undefined if α is set to zero—we are not allowed to begin the Bayesian process with an infinitely wide prior on \mathbf{w} . The more important difference is that the Bayesian estimate provides a covariance matrix, showing how different values of \mathbf{w} are, rather than providing only the estimate.

5.6.1 Maximum a Posteriori (MAP) Estimation

While the most principled approach is to make predictions using the posterior distribution over the parameter θ , it is still often desired to have a single point estimate. One common reason for desiring a point estimate is that most operations involving the Bayesian posterior for most interesting problems are intractable, and a point estimate offers a tractable approximation. Instead of simply returning to the maximum likelihood estimate, we can still obtain the benefit of the Bayesian approach by allowing the prior to influence the choice of the point estimate. One rational way to do this is to choose the **a posteriori** (MAP) point estimate. The MAP estimate chooses the value

maximal posterior probability (or maximal probability density in the case of continuous $\boldsymbol{\theta}$):

$$\boldsymbol{\theta}_{\text{MAP}} = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathbf{x}) = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{x} \mid \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$$

We recognize, on the righthand side, $\log p(\mathbf{x} \mid \boldsymbol{\theta})$, that is, the likelihood term, and $\log p(\boldsymbol{\theta})$, corresponding to the prior distribution.

As an example, consider a linear regression model with a Gaussian prior on the weights \mathbf{w} . If this prior is given by $\mathcal{N}(\mathbf{w}; \mathbf{0}, \frac{1}{\lambda} \mathbf{I}^2)$, then the logarithm in equation 5.79 is proportional to the familiar $\lambda \mathbf{w}^\top \mathbf{w}$ weight decay term that does not depend on \mathbf{w} and does not affect the learning rule. Bayesian inference with a Gaussian prior on the weights thus corresponds to weight decay.

As with full Bayesian inference, MAP Bayesian inference has the advantage of leveraging information that is brought by the prior and cannot be obtained from the training data. This additional information helps to reduce the variance of the MAP point estimate (in comparison to the ML estimate). However, this comes at the price of increased bias.

Many regularized estimation strategies, such as maximum likelihood regularized with weight decay, can be interpreted as making the transition to Bayesian inference. This view applies when the regularization is implemented by adding an extra term to the objective function that corresponds to the regularization penalties. All regularization penalties correspond to MAP Bayesian inference, although some regularizer terms may not be the logarithm of a probability distribution. Other regularization terms depend on the data, which of course a probability distribution is not allowed to do.

MAP Bayesian inference provides a straightforward way to design yet interpretable regularization terms. For example, a more complex term can be derived by using a mixture of Gaussians, rather than a single

distribution, as the prior (Nowlan and Hinton, 1992).

5.7 Supervised Learning Algorithms

Recall from section 5.1.3 that supervised learning algorithms are, roughly speaking, learning algorithms that learn to associate some input with some output. A training set of examples of inputs \mathbf{x} and outputs \mathbf{y} . In many cases, the outputs \mathbf{y} may be difficult to collect automatically and must be provided by a human.

“supervisor,” but the term still applies even when the training samples are collected automatically.

5.7.1 Probabilistic Supervised Learning

Most supervised learning algorithms in this book are based on specifying a probability distribution $p(y | \mathbf{x})$. We can do this simply by using maximum likelihood estimation to find the best parameter vector $\boldsymbol{\theta}$ for a particular family of distributions $p(y | \mathbf{x}; \boldsymbol{\theta})$.

We have already seen that linear regression corresponds to the distribution

$$p(y | \mathbf{x}; \boldsymbol{\theta}) = \mathcal{N}(y; \boldsymbol{\theta}^\top \mathbf{x}, \mathbf{I}).$$

We can generalize linear regression to the classification scenario by using a different family of probability distributions. If we have two classes, class 1 and class 0, then we need only specify the probability of one of these classes, because the probability of the other class must be determined by the fact that the probabilities must add up to 1.

The normal distribution over real-valued numbers that we used for linear regression is parametrized in terms of a mean. Any value we supply for the mean is valid. A distribution over a binary variable is slightly more constrained, because its mean must always be between 0 and 1. One way to solve this problem is to use the logistic sigmoid function to squash the output of the linear function into the interval $(0, 1)$ and interpret that value as a probability:

$$p(y = 1 | \mathbf{x}; \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \mathbf{x}).$$

This approach is known as **logistic regression** (a somewhat misleading name, because we use the model for classification rather than regression).

In the case of linear regression, we were able to find the optimal parameters by minimizing the sum of squared errors. In the case of logistic regression, we can find the optimal parameters by maximizing the log-likelihood function, which is defined as the sum of the log-probabilities of the observed data points:

solving the normal equations. Logistic regression is somewhat more is no closed-form solution for its optimal weights. Instead, we n them by maximizing the log-likelihood. We can do this by minimizi log-likelihood using gradient descent.

This same strategy can be applied to essentially any supervised le by writing down a parametric family of conditional probability dis the right kind of input and output variables.

5.7.2 Support Vector Machines

One of the most influential approaches to supervised learning is the machine (Boser *et al.*, 1992; Cortes and Vapnik, 1995). This model is similar to logistic regression in that it is driven by a linear function $\mathbf{w}^\top \mathbf{x} + b$. Unlike logistic regression, the support vector machine does not provide probability outputs. Instead, it outputs a class identity. The SVM predicts that the positive class if $\mathbf{w}^\top \mathbf{x} + b > 0$, and the negative class if $\mathbf{w}^\top \mathbf{x} + b \leq 0$.

One key innovation associated with support vector machines is the **kernel trick**. The kernel trick consists of observing that many machine learning models can be written exclusively in terms of dot products between example vectors. For example, it can be shown that the linear function used by the support vector machine can be re-written as

$$\mathbf{w}^\top \mathbf{x} + b = b + \sum_{i=1}^m \alpha_i \mathbf{x}^\top \mathbf{x}^{(i)},$$

where $\mathbf{x}^{(i)}$ is a training example, and $\boldsymbol{\alpha}$ is a vector of coefficients. Using the kernel trick, we can learn a linear function this way enables us to replace \mathbf{x} with the output of a feature mapping function $\phi(\mathbf{x})$ and the dot product with a function $k(\mathbf{x}, \mathbf{x}^{(i)}) = \phi(\mathbf{x})^\top \phi(\mathbf{x}^{(i)})$. This function $k(\cdot, \cdot)$ is called a **kernel**. The \cdot operator represents an inner product analogous to the dot product. For some feature spaces, we may not use literally the vector inner product. For example, in some infinite dimensional spaces, we need to use other kinds of inner products, such as example, inner products based on integration rather than summation. The mathematical development of these kinds of inner products is beyond the scope of this book.

After replacing dot products with kernel evaluations, we can make the function $f(\mathbf{x}) = b + \sum_{i=1}^m \alpha_i k(\mathbf{x}, \mathbf{x}^{(i)})$ nonlinear with respect to \mathbf{x} , but the relationship is still linear with respect to the kernel evaluations $\mathbf{x}^\top \mathbf{x}^{(i)}$.

$$f(\mathbf{x}) = b + \sum_{i=1}^m \alpha_i k(\mathbf{x}, \mathbf{x}^{(i)}).$$

This function is nonlinear with respect to \mathbf{x} , but the relationship is still linear with respect to the kernel evaluations $\mathbf{x}^\top \mathbf{x}^{(i)}$.

and $f(\mathbf{x})$ is linear. Also, the relationship between α and $f(\mathbf{x})$ kernel-based function is exactly equivalent to preprocessing the data $\phi(\mathbf{x})$ to all inputs, then learning a linear model in the new transformed space.

The kernel trick is powerful for two reasons. First, it enables us to learn functions that are nonlinear as a function of \mathbf{x} using convex optimization techniques that are guaranteed to converge efficiently. This is possible because we can learn and optimize only α , that is, the optimization algorithm can view the function as being linear in a different space. Second, the kernel function $K(\mathbf{x}, \mathbf{x}')$ is often much easier to compute than the transformed feature vector $\phi(\mathbf{x})$.

admits an implementation that is significantly more computational than naively constructing two $\phi(\mathbf{x})$ vectors and explicitly taking their dot product.

In some cases, $\phi(\mathbf{x})$ can even be infinite dimensional, which incurs an infinite computational cost for the naive, explicit approach. For example, $k(\mathbf{x}, \mathbf{x}') = \min(x_i, x'_i)$ is a nonlinear, tractable function of \mathbf{x} even when $\phi(\mathbf{x})$ is an example of an infinite-dimensional feature space with a tractable basis. We can construct a feature mapping $\phi(x)$ over the nonnegative integers x such that $\phi(x) = [x, x^2, x^3, \dots]$. In this case, $k(\mathbf{x}, \mathbf{x}')$ is the dot product of this mapping returns a vector containing x ones followed by infinitely many zeros. We can write a kernel function $k(x, x^{(i)}) = \min(x, x^{(i)})$ that is exactly equivalent to the corresponding infinite-dimensional dot product.

The most commonly used kernel is the **Gaussian kernel**,

$$k(\mathbf{u}, \mathbf{v}) = \mathcal{N}(\mathbf{u} - \mathbf{v}; 0, \sigma^2 \mathbf{I}),$$

where $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the standard normal density. This kernel is called the **radial basis function** (RBF) kernel, because its value decreases exponentially as the distance between \mathbf{u} and \mathbf{v} increases. It corresponds to a Gaussian probability distribution in \mathbf{v} space radiating outward from \mathbf{u} . The Gaussian kernel corresponds to a Gaussian process, which is a generalization of the linear model to an infinite-dimensional space, but the derivation of the Gaussian kernel is much more straightforward than in our example of the min kernel over the infinite-dimensional space.

We can think of the Gaussian kernel as performing a kind of **template matching**. A training example \mathbf{x} associated with training label y becomes a template for class y . When a test point \mathbf{x}' is near \mathbf{x} according to Euclidean distance, the Gaussian kernel has a large response, indicating that \mathbf{x}' is very similar to the template. The model then puts a large weight on the associated template. Overall, the prediction will combine many such training labels weighted by the similarity of the corresponding training examples.

Support vector machines are not the only algorithm that can benefit from using the kernel trick. Many other linear models can be enhanced in this way. A category of algorithms that employ the kernel trick is known as **kernel methods** or **kernel methods** (Williams and Rasmussen, 1996; Schölkopf et al., 1998).

A major drawback to kernel machines is that the cost of evaluating the function is linear in the number of training examples, because the function contributes a term $\alpha_i k(\mathbf{x}, \mathbf{x}^{(i)})$ to the decision function. Support vectors are able to mitigate this by learning an $\boldsymbol{\alpha}$ vector that contains only the values of α_i for the training examples that have nonzero α_i . These training examples are called **support vectors**.

Kernel machines also suffer from a high computational cost of evaluation if the dataset is large. We revisit this idea in section 5.9. Kernel

generic kernels struggle to generalize well. We explain why in section 5.6. The modern incarnation of deep learning was designed to overcome these limitations of kernel machines. The current deep learning renaissance began when LeCun et al. (2006) demonstrated that a neural network could outperform the RBF kernel machine on the MNIST benchmark.

5.7.3 Other Simple Supervised Learning Algorithms

We have already briefly encountered another nonprobabilistic supervised learning algorithm, nearest neighbor regression. More generally, *k*-nearest neighbor regression is a family of techniques that can be used for classification or regression. As a nonparametric learning algorithm, *k*-nearest neighbors is not restricted by the number of parameters. We usually think of the *k*-nearest neighbor algorithm as not having any parameters but rather implementing a simple averaging rule over the training data. In fact, there is not even really a training stage or learning step. Instead, at test time, when we want to produce an output y for a new input \mathbf{x} , we find the *k*-nearest neighbors to \mathbf{x} in the training data \mathbf{X} . We then average the corresponding y values in the training set. This works well for any kind of supervised learning where we can define an average operation. For example, in the case of classification, we can average over one-hot code vectors. If $c_i = 1$ for one value of i and $c_i = 0$ for all other values of i , we can then interpret the average of these one-hot codes as giving a probability distribution over classes. As a probabilistic learning algorithm, *k*-nearest neighbor can achieve very high capacity. For example, suppose we have a multiclass classification task and measure performance using the cross-entropy loss. In this setting, 1-nearest neighbor converges to double the Bayes error as the number of training examples approaches infinity. The error in excess of the Bayes error results from choosing a single neighbor by breaking ties between equally distant neighbors randomly. When there is infinite training data, a *k*-nearest neighbor algorithm will have infinitely many training set neighbors at distance zero. One way to fix this is to have the algorithm to use all these neighbors to vote, rather than randomizing the choice of a single neighbor.

of them, the procedure converges to the Bayes error rate. The k -nearest neighbors enables it to obtain high accuracy given a large training set. It does so at high computational cost, however, and it may generalize poorly given a small finite training set. One weakness of k -nearest neighbors is that it cannot learn that one feature is more discriminative than another. For example, imagine we have a regression task with $\mathbf{x} \in \mathbb{R}^{100}$ drawn from an isotropic Gaussian distribution, but only a single variable x_1 is relevant to the outcome. In this case, further that this feature simply encodes the output directly, that is, $y = x_1 + \epsilon$ for some noise ϵ . Nearest neighbor regression will not be able to detect this

The nearest neighbor of most points x will be determined by the last few features x_2 through x_{100} , not by the lone feature x_1 . Thus the order of training sets will essentially be random.

Another type of learning algorithm that also breaks the input space into regions and has separate parameters for each region is the **decision tree** (Breiman 1984) and its many variants. As shown in figure 5.7, each node in the decision tree is associated with a region in the input space, and internal nodes split the region into one subregion for each child of the node (typically using an axis-aligned cut). Space is thus subdivided into nonoverlapping regions, with a one-to-one correspondence between leaf nodes and input regions. Each leaf node outputs a single value, corresponding to every point in its input region to the same output. Decision trees can be trained with specialized algorithms that are beyond the scope of this book, but any learning algorithm can be considered nonparametric if it is allowed to have an arbitrary number of nodes. While decision trees are of arbitrary size, though decision trees are usually regularized with pruning rules that turn them into parametric models in practice. Decision trees are typically used, with axis-aligned splits and constant outputs within regions. They struggle to solve some problems that are easy even for logistic regression. For example, if we have a two-class problem, and the positive class of points satisfies $x_2 > x_1$, the decision boundary is not axis aligned. The decision tree must need to approximate the decision boundary with many nodes, implemented as a function that constantly walks back and forth across the true decision boundary with axis-aligned steps.

As we have seen, nearest neighbor predictors and decision trees have their limitations. Nonetheless, they are useful learning algorithms when computation resources are constrained. We can also build intuition for more sophisticated learning algorithms by thinking about the similarities and differences between these simpler algorithms and k -nearest neighbors or decision trees.

See Murphy (2012), Bishop (2006), Hastie *et al.* (2001) or other machine learning textbooks for more material on traditional supervised learning algorithms.

5.8 Unsupervised Learning Algorithms

Recall from section 5.1.3 that unsupervised algorithms are those that learn only “features” but not a supervision signal. The distinction between supervised and unsupervised algorithms is not formally and rigidly defined because there is no objective test for distinguishing whether a value is a feature or a target variable. In fact, there is no supervisor. Informally, unsupervised learning refers to most attempts to learn information from a distribution that do not require human labeling.

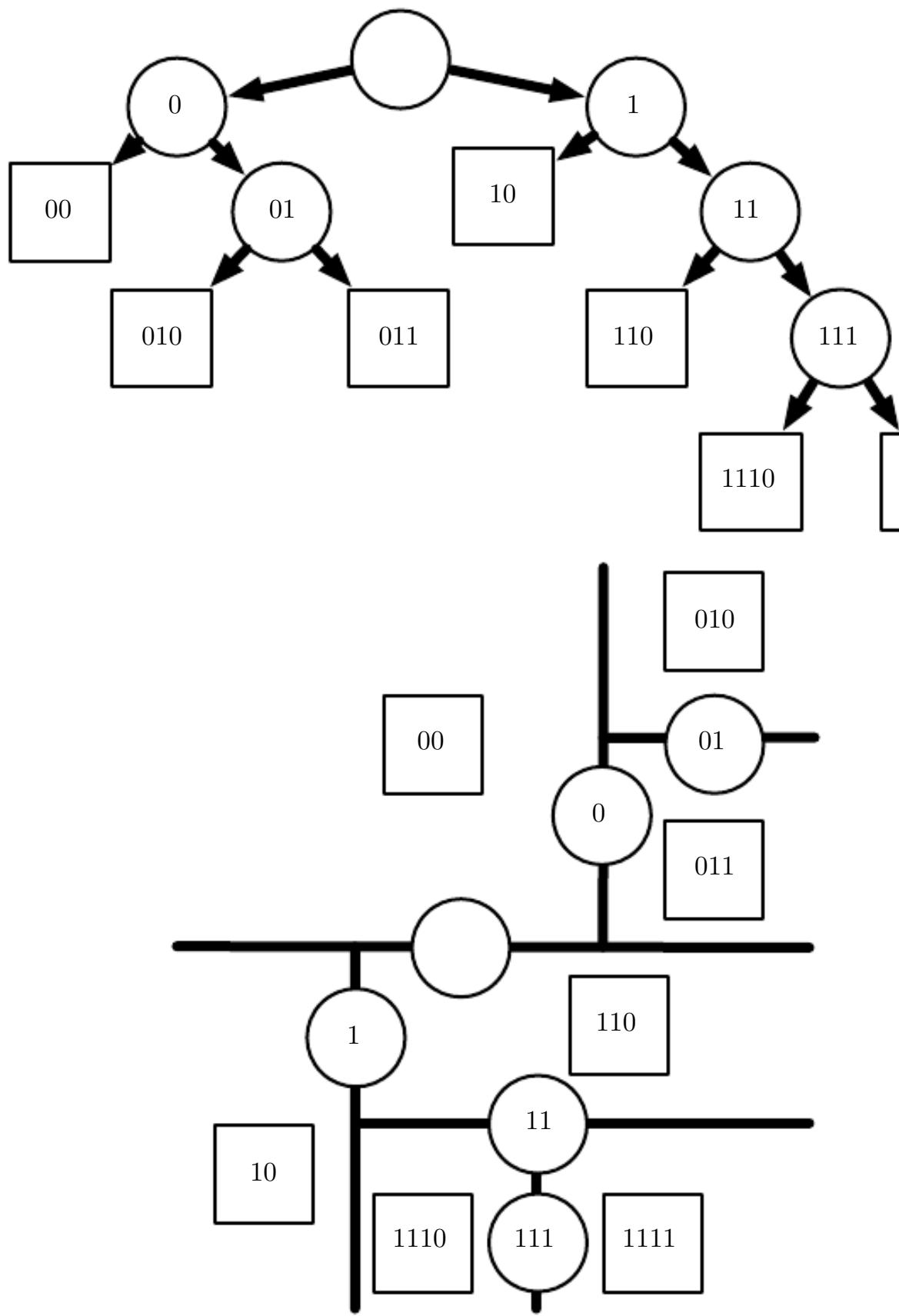


Figure 5.7: Diagrams describing how a decision tree works. (*Top*)Each

chooses to send the input example to the child node on the left (0) or to the right (1). Internal nodes are drawn as circles and leaf nodes as rectangles. A node is displayed with a binary string identifier corresponding to its position in the tree, obtained by appending a bit to its parent identifier (0 = choose left or right, 1 = choose bottom). *(Bottom)* The tree divides space into regions. The figure shows how a decision tree might divide \mathbb{R}^2 . The nodes of the tree are plotted along the dividing lines they use to categorize examples, and leaf nodes are drawn in the center of the region of examples they receive. The function is a piecewise-constant function, with one piece per leaf. Each leaf requires at least one training example to define, so it is not possible for the decision tree to have more leaves than the number of training examples.

examples. The term is usually associated with density estimation, drawing samples from a distribution, learning to denoise data from noise, finding a manifold that the data lies near, or clustering the data into related examples.

A classic unsupervised learning task is to find the “best” representation of the data. By “best” we can mean different things, but generally speaking we are looking for a representation that preserves as much information about x as possible while obeying some penalty or constraint aimed at keeping the representation more accessible than x itself.

There are multiple ways of defining a simpler representation. The most common include lower-dimensional representations, sparse representations, and independent representations. Low-dimensional representations attempt to compress as much information about x as possible in a smaller space. Sparse representations (Barlow, 1989; Olshausen and Field, 1996; Ghahramani, 1997) embed the dataset into a representation where most elements are mostly zeros for most inputs. The use of sparse representations typically increases the dimensionality of the representation, so that the process of becoming mostly zeros does not discard too much information. The overall structure of the representation that tends to distribute data points across the representation space. Independent representations attempt to capture the sources of variation underlying the data distribution such that the elements of the representation are statistically independent.

Of course these three criteria are certainly not mutually exclusive. Lower-dimensional representations often yield elements that have fewer dependencies than the original high-dimensional data. This is because one way to reduce the size of a representation is to find and remove redundancy, and removing more redundancy enables the dimensionality reduction to achieve more compression while discarding less information.

The notion of representation is one of the central themes of deep learning, and we will return to it many times throughout this book.

therefore one of the central themes in this book. In this section, we present simple examples of representation learning algorithms. Together, these algorithms show how to operationalize all three of the criteria above. The remaining chapters introduce additional representation learning algorithms and develop these criteria in different ways or introduce other criteria.

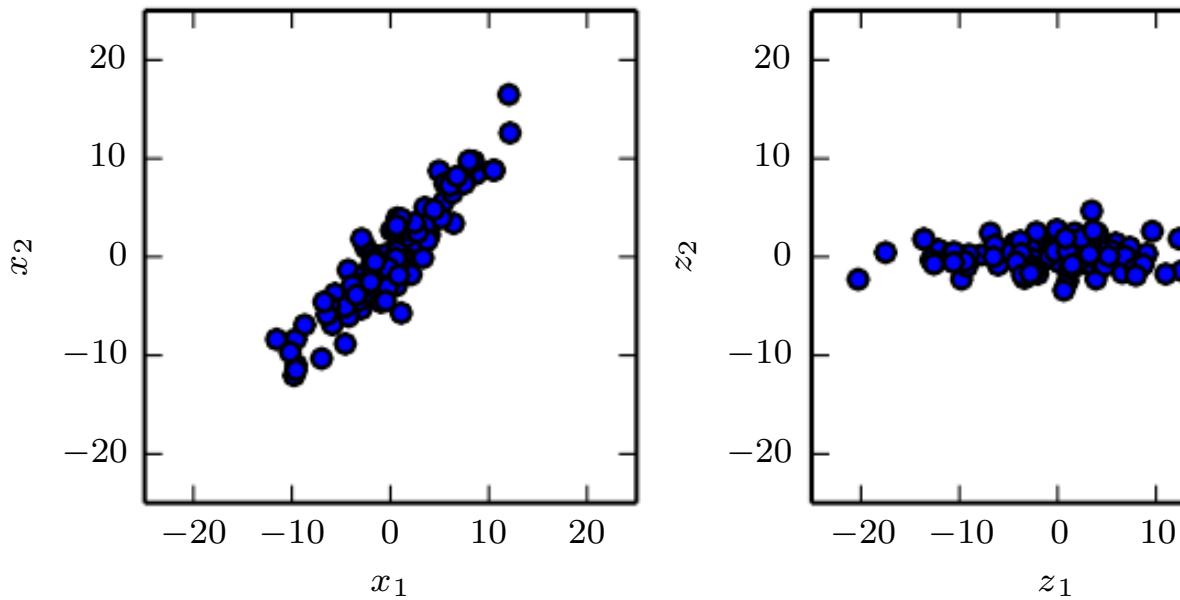


Figure 5.8: PCA learns a linear projection that aligns the direction of greatest variance with the axes of the new space. (*Left*) The original data consist of samples of \mathbf{x} . Variance might occur along directions that are not axis aligned. (*Right*) The data $\mathbf{z} = \mathbf{x}^\top \mathbf{W}$ now varies most along the axis z_1 . The direction of second most variance is now along z_2 .

5.8.1 Principal Components Analysis

In section 2.12, we saw that the principal components analysis algorithm is a means of compressing data. We can also view PCA as an unsupervised learning algorithm that learns a representation of data. This representation satisfies two of the criteria for a simple representation described above. It is a representation that has lower dimensionality than the original input, and it is a representation whose elements have no linear correlation with each other. This is a first step toward the criterion of learning representations whose elements are statistically independent. To achieve full independence, a representation learning algorithm must also remove the nonlinear relationships between variables.

PCA learns an orthogonal, linear transformation of the data that captures the maximum variance in the data.

input \mathbf{x} to a representation \mathbf{z} as shown in figure 5.8. In section 2.1 we could learn a one-dimensional representation that best reconstructs the data (in the sense of mean squared error) and that this representation corresponds to the first principal component of the data. Thus we can view PCA as a simple and effective dimensionality reduction method that preserves as much of the information in the data as possible (again, as measured by the reconstruction error). In the following, we will study how the PCA algorithm decorrelates the original data representation \mathbf{X} .

Let us consider the $m \times n$ design matrix \mathbf{X} . We will assume that

a mean of zero, $\mathbb{E}[\mathbf{x}] = \mathbf{0}$. If this is not the case, the data can easily be centered by subtracting the mean from all examples in a preprocessing step.

The unbiased sample covariance matrix associated with \mathbf{X} is given by

$$\text{Var}[\mathbf{x}] = \frac{1}{m-1} \mathbf{X}^\top \mathbf{X}.$$

PCA finds a representation (through linear transformation) $\mathbf{z} = \mathbf{W}\mathbf{x}$ such that $\text{Var}[\mathbf{z}]$ is diagonal.

In section 2.12, we saw that the principal components of a dataset are given by the eigenvectors of $\mathbf{X}^\top \mathbf{X}$. From this view,

$$\mathbf{X}^\top \mathbf{X} = \mathbf{W} \Lambda \mathbf{W}^\top.$$

In this section, we exploit an alternative derivation of the principal components. The principal components may also be obtained via singular value decomposition (SVD). Specifically, they are the right singular vectors of \mathbf{X} . To see this, consider the right singular vectors in the decomposition $\mathbf{X} = \mathbf{U} \Sigma \mathbf{W}^\top$. We can substitute the original eigenvector equation with \mathbf{W} as the eigenvector basis:

$$\mathbf{X}^\top \mathbf{X} = (\mathbf{U} \Sigma \mathbf{W}^\top)^\top \mathbf{U} \Sigma \mathbf{W}^\top = \mathbf{W} \Sigma^2 \mathbf{W}^\top.$$

The SVD is helpful to show that PCA results in a diagonal Variance matrix. Using the SVD of \mathbf{X} , we can express the variance of \mathbf{X} as:

$$\begin{aligned} \text{Var}[\mathbf{x}] &= \frac{1}{m-1} \mathbf{X}^\top \mathbf{X} \\ &= \frac{1}{m-1} (\mathbf{U} \Sigma \mathbf{W}^\top)^\top \mathbf{U} \Sigma \mathbf{W}^\top \\ &= \frac{1}{m-1} \mathbf{W} \Sigma^\top \mathbf{U}^\top \mathbf{U} \Sigma \mathbf{W}^\top \end{aligned}$$

$$= \frac{1}{m-1} \mathbf{W} \boldsymbol{\Sigma}^2 \mathbf{W}^\top,$$

where we use the fact that $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$ because the \mathbf{U} matrix of the decomposition is defined to be orthogonal. This shows that the covariance is diagonal as required:

$$\begin{aligned}\text{Var}[\mathbf{z}] &= \frac{1}{m-1} \mathbf{Z}^\top \mathbf{Z} \\ &= \frac{1}{m-1} \mathbf{W}^\top \mathbf{X}^\top \mathbf{X} \mathbf{W}\end{aligned}$$

$$\begin{aligned} &= \frac{1}{m-1} \mathbf{W}^\top \mathbf{W} \Sigma^2 \mathbf{W}^\top \mathbf{W} \\ &= \frac{1}{m-1} \Sigma^2, \end{aligned}$$

where this time we use the fact that $\mathbf{W}^\top \mathbf{W} = \mathbf{I}$, again from the SVD.

The above analysis shows that when we project the data \mathbf{x} to a transformation \mathbf{W} , the resulting representation has a diagonal covariance (as given by Σ^2), which immediately implies that the individual elements are mutually uncorrelated.

This ability of PCA to transform data into a representation where elements are mutually uncorrelated is a very important property of PCA. An example of a representation that attempts to *disentangle the unknown variation* underlying the data. In the case of PCA, this disentanglement is in the form of finding a rotation of the input space (described by \mathbf{W}) such that the principal axes of variance with the basis of the new representation \mathbf{z} are uncorrelated.

While correlation is an important category of dependency between the data, we are also interested in learning representations that disentangle more complicated forms of feature dependencies. For this, we will need more sophisticated methods, which can be done with a simple linear transformation.

5.8.2 k -means Clustering

Another example of a simple representation learning algorithm is k -means clustering. The k -means clustering algorithm divides the training set into k disjoint clusters of examples that are near each other. We can thus think of the algorithm as providing a k -dimensional one-hot code vector \mathbf{h} representing an example. If an example belongs to cluster i , then $h_i = 1$, and all other entries of the representation are zero.

zero.

The one-hot code provided by k -means clustering is an example of a sparse representation, because the majority of its entries are zero for every input \mathbf{x} . In this section, we develop other algorithms that learn more flexible sparse representations, where more than one entry can be nonzero for each input \mathbf{x} . One-hot codes are a special case of sparse representations that lose many of the benefits of sparsity. For example, they do not naturally convey the idea that all examples in the same cluster are similar to each other, and they do not confer the computational advantage that the entire set of examples in a cluster share the same representation.

may be captured by a single integer.

The k -means algorithm works by initializing k different centroids to different values, then alternating between two different steps until convergence. In one step, each training example is assigned to cluster i , where i is the index of the nearest centroid $\mu^{(i)}$. In the other step, each centroid $\mu^{(i)}$ is updated to be the mean of all training examples $x^{(j)}$ assigned to cluster i .

One difficulty pertaining to clustering is that the clustering problem is ill posed, in the sense that there is no single criterion that measures how well the clustering of the data corresponds to the real world. We can measure the quality of a clustering, such as the average Euclidean distance from a data point to the members of the cluster. This enables us to tell how well we can reconstruct the training data from the cluster assignments. We do not know, however, how well the cluster assignments correspond to properties of the real world. For example, there may be many different clusterings that all correspond well to some aspect of the real world. We may hope to find a clustering that relates to a particular aspect of the real world, but obtain a different, equally valid clustering that is not relevant to that aspect. For example, suppose that we run two clustering algorithms on a dataset containing images of red trucks, images of red cars, images of gray trucks, and images of gray cars. If we ask each clustering algorithm to find two clusters, one may find a cluster of cars and a cluster of trucks, while another may find a cluster of red vehicles and a cluster of gray vehicles. Suppose we also run a clustering algorithm that is allowed to determine the number of clusters. Then it may find the examples to four clusters, red cars, red trucks, gray cars, and gray trucks. This new clustering now at least captures information about both attributes, but has lost information about similarity. Red cars are in a different cluster than red trucks, just as they are in a different cluster from gray trucks. The clustering algorithm does not tell us that red cars are more similar to red trucks than they are to gray trucks. They are different from both things, but we know.

These issues illustrate some of the reasons that we may prefer a distributed representation to a one-hot representation. A distributed representation encodes two attributes for each vehicle—one representing its color and one representing whether it is a car or a truck. It is still not entirely clear why a distributed representation is useful (how can the learning algorithm know which of the two attributes we are interested in are color and car-versus-truck, manufacturer and age?), but having many attributes reduces the burden on the algorithm to guess which single attribute we care about, and gives us a way to measure similarity between objects in a fine-grained way by comparing

attributes instead of just testing whether one attribute matches.

5.9 Stochastic Gradient Descent

Nearly all of deep learning is powered by one very important algorithm: **gradient descent** (SGD). Stochastic gradient descent is an extended gradient descent algorithm introduced in section 4.3.

A recurring problem in machine learning is that large training sets are needed for good generalization, but large training sets are also more computationally expensive.

The cost function used by a machine learning algorithm often consists of a sum over training examples of some per-example loss function. For example, the negative conditional log-likelihood of the training data can be written as

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{\text{data}}} L(\mathbf{x}, y, \boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta}),$$

where L is the per-example loss $L(\mathbf{x}, y, \boldsymbol{\theta}) = -\log p(y | \mathbf{x}; \boldsymbol{\theta})$.

For these additive cost functions, gradient descent requires computing

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta}).$$

The computational cost of this operation is $O(m)$. As the training set grows to billions of examples, the time to take a single gradient step becomes increasingly long.

The insight of SGD is that the gradient is an expectation. The expectation can be approximately estimated using a small set of samples. Specifically, at each step of the algorithm, we can sample a **minibatch** of examples $\mathbb{B} = \{$

drawn uniformly from the training set. The minibatch size m' is typically to be a relatively small number of examples, ranging from one to thousands. Crucially, m' is usually held fixed as the training set size m grows, so that we can fit a model to a training set with billions of examples using updates computed on small batches of examples.

The estimate of the gradient is formed as

$$\mathbf{g} = \frac{1}{m'} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{m'} L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$

using examples from the minibatch \mathbb{B} . The stochastic gradient descent then follows the estimated gradient downhill:

$$\theta \leftarrow \theta - \epsilon g,$$

where ϵ is the learning rate.

Gradient descent in general has often been regarded as slow or unreliable. In the past, the application of gradient descent to nonconvex optimization was regarded as foolhardy or unprincipled. Today, we know that deep learning models described in part II work very well when trained with gradient descent. The optimization algorithm may not be guaranteed to always reach a global minimum in a reasonable amount of time, but it often finds a local minimum of the cost function quickly enough to be useful.

Stochastic gradient descent has many important uses outside of deep learning. It is the main way to train large linear models on large datasets. For a fixed model size, the cost per SGD update does not depend on the training set size m . In practice, we often use a larger model as the training set size m increases, but we are not forced to do so. The number of updates required for convergence usually increases with training set size. However, as m goes to infinity, the model will eventually converge to its best possible test error. From this point of view, one can argue that the asymptotic behavior of a model with SGD is $O(1)$ as a function of m .

Prior to the advent of deep learning, the main way to learn nonlinear functions was to use the kernel trick in combination with a linear model. Many machine learning algorithms require constructing an $m \times m$ matrix $G_{i,j} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$. For a fixed model size, this matrix has computational cost $O(m^2)$, which is clearly undesirable for datasets with billions of examples. In academia, starting in 2006, deep learning was initially interesting because it was able to generalize to new examples without explicitly learning a feature representation.

than competing algorithms when trained on medium-sized datasets of thousands of examples. Soon after, deep learning garnered additional attention from the industry because it provided a scalable way of training nonlinear models on large datasets.

Stochastic gradient descent and many enhancements to it are described in chapter 8.

5.10 Building a Machine Learning Algorithm

Nearly all deep learning algorithms can be described as particular combinations of a fairly simple recipe: combine a specification of a dataset, a cost function, an optimization procedure and a model.

For example, the linear regression algorithm combines a dataset \mathbf{X} and \mathbf{y} , the cost function

$$J(\mathbf{w}, b) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(y \mid \mathbf{x}),$$

the model specification $p_{\text{model}}(y \mid \mathbf{x}) = \mathcal{N}(y; \mathbf{x}^\top \mathbf{w} + b, 1)$, and, in this case, an optimization algorithm defined by solving for where the gradient of the cost function is zero, using the normal equations.

By realizing that we can replace any of these components mostly independently from the others, we can obtain a wide range of algorithms.

The cost function typically includes at least one term that causes the optimization process to perform statistical estimation. The most common cost function is the negative log-likelihood, so that minimizing the cost function can be interpreted as maximum likelihood estimation.

The cost function may also include additional terms, such as regularization terms. For example, we can add weight decay to the linear regression cost function to obtain

$$J(\mathbf{w}, b) = \lambda \|\mathbf{w}\|_2^2 - \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(y \mid \mathbf{x}).$$

This still allows closed form optimization.

If we change the model to be nonlinear, then most cost functions cannot be optimized in closed form. This requires us to choose an iterative optimization procedure, such as gradient descent.

The recipe for constructing a learning algorithm by combining model specification, cost function, optimization algorithm, and initialization is

optimization algorithms supports both supervised and unsupervised learning. A linear regression example shows how to support supervised learning while unsupervised learning can be supported by defining a dataset that contains only unlabeled data \mathcal{X} . An appropriate unsupervised cost and model. For example, we can find the principal component (PCA) vector by specifying that our loss function is

$$J(\mathbf{w}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \|\mathbf{x} - r(\mathbf{x}; \mathbf{w})\|_2^2$$

while our model is defined to have \mathbf{w} with norm one and reconstruct \mathbf{x} as $r(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \mathbf{w}$.

In some cases, the cost function may be a function that we can't evaluate, for computational reasons. In these cases, we can still minimize it using iterative numerical optimization, as long as we have a way of approximating its gradients.

Most machine learning algorithms make use of this recipe, though it may not be immediately obvious. If a machine learning algorithm seems especially hand designed, it can usually be understood as using a special-case of this recipe. Some models, such as decision trees and k -means, require special-case optimizers because their cost functions have flat regions that make them inappropriate for optimization by gradient-based optimizers. Recognizing that most machine learning algorithms can be described using this recipe helps to see the different algorithms as part of a taxonomy of methods for doing related tasks that work for similar reasons rather than as a long list of algorithms that each have separate justifications.

5.11 Challenges Motivating Deep Learning

The simple machine learning algorithms described in this chapter have been successful at solving a wide variety of important problems. They have not succeeded, however, at solving some of the central problems in AI, such as recognizing speech or recognizing images.

The development of deep learning was motivated in part by the desire to use traditional machine learning algorithms to generalize well on such AI tasks.

This section is about how the challenge of generalizing to new examples becomes exponentially more difficult when working with high-dimensional data. It also describes the mechanisms used to achieve generalization in traditional machine learning algorithms. These mechanisms are insufficient to learn complicated functions in high-dimensional spaces, and the computational costs of learning in such spaces also often impose high computational costs. Deep learning provides ways to overcome these and other obstacles.

5.11.1 The Curse of Dimensionality

Many machine learning problems become exceedingly difficult when the number of dimensions in the data is high. This phenomenon is known as the **curse of dimensionality**. Of particular concern is that the number of possible configurations of a set of variables increases exponentially as the number of dimensions increases.

The curse of dimensionality arises in many places in computer science and engineering, including machine learning.

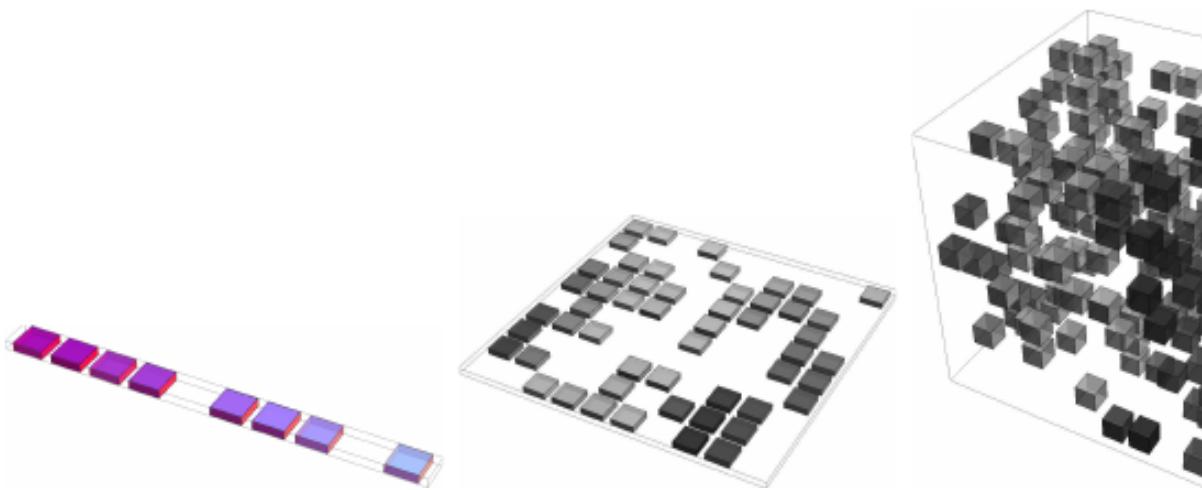


Figure 5.9: As the number of relevant dimensions of the data increases (right), the number of configurations of interest may grow exponentially. In the one-dimensional example, we have one variable for which we only care about discrete regions of interest. With enough examples falling within each of these regions (each region corresponds to a cell in the illustration), learning algorithms can easily generalize. A straightforward way to generalize is to estimate the value of the target variable within each region (and possibly interpolate between neighboring regions). (Left) With two dimensions, it is more difficult to distinguish 10 different values of each variable. We need to keep track of up to $10 \times 10 = 100$ regions, and we need at least that many training examples to cover all those regions. (Right) With three dimensions, this grows to $10^3 = 1000$ regions, and at least that many examples. For d dimensions and v values to be distinguished along each axis, we seem to need $O(v^d)$ regions and examples. This is an instantiation of the curse of dimensionality. Figure graciously provided by Nicolas Chapados.

One challenge posed by the curse of dimensionality is a statistical challenge. As illustrated in figure 5.9, a statistical challenge arises because the number of possible configurations of \mathbf{x} is much larger than the number of training examples. To understand the issue, let us consider that the input space is organized as a grid, as in the figure. We can describe low-dimensional space with a small number of grid cells that are mostly occupied by the data. When generalizing to a new point, we can usually tell what to do simply by inspecting the training points that lie in the same cell as the new input. For example, if estimating the density at some point \mathbf{x} , we can just return the number of training points

density at some point \mathbf{x} , we can just return the number of training examples in the same unit volume cell as \mathbf{x} , divided by the total number of training examples. If we wish to classify an example, we can return the most common class among the examples in the same cell. If we are doing regression, we can average the values observed over the examples in that cell. But what about the case where we have seen no example? Because in high-dimensional spaces, the number of possible configurations is huge, much larger than our number of examples, a new configuration \mathbf{x} will almost certainly have no training example associated with it. How could we possibly make any meaningful statement about these new configurations? Many traditional machine learning

algorithms simply assume that the output at a new point should be the same as the output at the nearest training point.

5.11.2 Local Constancy and Smoothness Regularization

To generalize well, machine learning algorithms need to be guided about what kind of function they should learn. We have seen these rated as explicit beliefs in the form of probability distributions over the model. More informally, we may also discuss prior beliefs as directly influencing the *function* itself and influencing the parameters only indirectly, as a relationship between the parameters and the function. Additionally, we may discuss prior beliefs as being expressed implicitly by choosing a class of functions. For example, if we are biased toward choosing a class of functions that are smooth, then these biases may not be expressed (or even be possible to express) as a probability distribution representing our degree of belief in various functions.

Among the most widely used of these implicit “priors” is the **prior**, or **local constancy prior**. This prior states that the function should not change very much within a small region.

Many simpler algorithms rely exclusively on this prior to generalize. As a result, they fail to scale to the statistical challenges involved in real-world level tasks. Throughout this book, we describe how deep learning models can incorporate additional (explicit and implicit) priors in order to reduce the error on sophisticated tasks. Here, we explain why the smoothness prior is insufficient for these tasks.

There are many different ways to implicitly or explicitly express the belief that the learned function should be smooth or locally constant. All of these methods are designed to encourage the learning process to learn a function that satisfies the condition

$$f^*(\mathbf{x}) \approx f^*(\mathbf{x} + \epsilon)$$

for most configurations \mathbf{x} and small change ϵ . In other words, if we have an answer for an input \mathbf{x} (for example, if \mathbf{x} is a labeled training example), the answer is probably good in the neighborhood of \mathbf{x} . If we have several answers in some neighborhood, we would combine them (by some form of weighted interpolation) to produce an answer that agrees with as many of the neighbors as possible.

An extreme example of the local constancy approach is the k -nearest neighbor family of learning algorithms. These predictors are literally constant over a region containing all the points \mathbf{x} that have the same set of k nearest neighbors.

the training set. For $k = 1$, the number of distinguishable regions equals than the number of training examples.

While the k -nearest neighbors algorithm copies the output from training examples, most kernel machines interpolate between training set outputs with nearby training examples. An important class of kernels is the **kernels**, where $k(\mathbf{u}, \mathbf{v})$ is large when $\mathbf{u} = \mathbf{v}$ and decreases as \mathbf{u} and \mathbf{v} are apart from each other. A local kernel can be thought of as a simple function that performs template matching, by measuring how closely a test example \mathbf{x} resembles each training example $\mathbf{x}^{(i)}$. Much of the modern motivation for learning is derived from studying the limitations of local template matching, and how deep models are able to succeed in cases where local template matching fails (Bengio *et al.*, 2006b).

Decision trees also suffer from the limitations of exclusively small regions of learning, because they break the input space into as many regions as possible, with one leaf per region. They have a separate parameter (or sometimes many parameters) associated with each region. If the target function requires at least n leaves to be represented accurately, then at least n training examples are required to fit the tree. A multiple of n is needed to achieve some level of confidence in the predicted output.

In general, to distinguish $O(k)$ regions in input space, all these regions require $O(k)$ training examples. Typically there are $O(k)$ parameters, with $O(k)$ associated with each of the $O(k)$ regions. The nearest neighbor scheme ensures that each training example can be used to define at most one region, as shown in figure 5.10.

Is there a way to represent a complex function that has many regions that cannot be distinguished than the number of training examples? Clearly, the smoothness of the underlying function will not allow a learner to do this. For example, imagine that the target function is a kind of checkerboard. The function contains many variations, but there is a simple structure to them.

happens when the number of training examples is substantially smaller than the number of black and white squares on the checkerboard. Based on the generalization and the smoothness or local constancy prior, the learner is guaranteed to correctly guess the color of a new point if it lay within the same square as a training example. There is no guarantee that the learner could correctly extend the checkerboard pattern to squares that do not contain training examples. With this prior information that an example tells us is the color of its square, and the goal of getting the colors of the entire checkerboard right is to cover each of the

least one example.

The smoothness assumption and the associated nonparametric algorithms work extremely well as long as there are enough examples for the algorithm to observe high points on most peaks and low points on most valleys of the true underlying function to be learned. This is generally true if the function to be learned is smooth enough and varies in few enough dimensions. In high dimensions, even a very smooth function can change smoothly in a different way along each dimension. If the function additionally behaves differently in various regions, it can become extremely complicated to describe it based on the training examples. If the function is complicated (we want to distinguish between the number of regions compared to the number of examples), is the function likely to generalize well?

The answer to both of these questions—whether it is possible to learn a complicated function efficiently, and whether it is possible for a complicated function to generalize well to new inputs—is yes. The key insight is that a function with k regions, such as $O(2^k)$, can be defined with $O(k)$ examples.

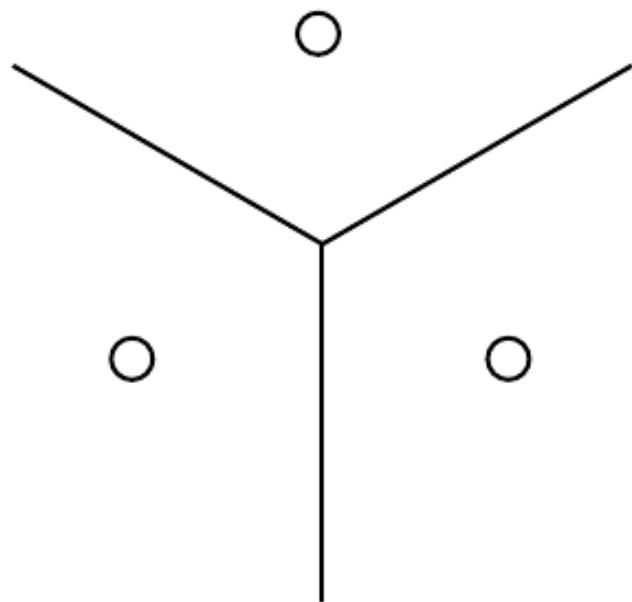


Figure 5.10: Illustration of how the nearest neighbor algorithm breaks up the space into regions.

into regions. An example (represented here by a circle) within each region is defined by the region boundary (represented here by the lines). The y value associated with each example defines what the output should be for all points within the corresponding region. Regions defined by nearest neighbor matching form a geometric pattern on a coordinate diagram. The number of these contiguous regions cannot grow faster than the number of training examples. While this figure illustrates the behavior of the nearest neighbor algorithm specifically, other machine learning algorithms that rely exclusively on local smoothness prior for generalization exhibit similar behaviors: each training example only informs the learner about how to generalize in some neighborhood surrounding that example.

introduce some dependencies between the regions through additional assumptions about the underlying data-generating distribution. In this way, we can generalize nonlocally (Bengio and Monperrus, 2005; Bengio *et al.*, 2007). Different deep learning algorithms provide implicit or explicit assumptions reasonable for a broad range of AI tasks in order to capture these dependencies.

Other approaches to machine learning often make stronger, task-specific assumptions. For example, we could easily solve the checkerboard task by assuming that the target function is periodic. Usually we do not make such strong, task-specific assumptions in neural networks so that they can apply to a much wider variety of structures. AI tasks have structure that is too complex to be limited to simple, manually specified properties such as periodicity. So we want learning algorithms that embody more general-purpose assumptions. The core idea in deep learning is that we assume that the data we are trying to predict is composed of factors, or features, potentially at multiple levels of hierarchy. Many other similarly generic assumptions can further improve learning algorithms. These apparently mild assumptions allow an exponential relationship between the number of examples and the number of representable functions to be distinguished. We describe these exponential gains more precisely in Sections 6.4.1, 15.4 and 15.5. The exponential advantages conferred by these distributed representations counter the exponential challenges posed by the curse of dimensionality.

5.11.3 Manifold Learning

An important concept underlying many ideas in machine learning is that of a manifold.

A **manifold** is a connected region. Mathematically, it is a set of points associated with a neighborhood around each point. From any given point on a manifold, the manifold locally appears to be a Euclidean space. In everyday life, we are familiar with the manifold of the world as a 2D plane, but it is in fact a complex manifold with many dimensions.

the surface of the world as a 2-D plane, but it is in fact a spherical 3-D space.

The concept of a neighborhood surrounding each point implies transformations that can be applied to move on the manifold from a neighboring one. In the example of the world's surface as a manifold, we can walk north, south, east, or west.

Although there is a formal mathematical meaning to the term "neighborhood" in machine learning it tends to be used more loosely to designate a set of points that can be approximated well by considering only a small neighborhood around a given point.

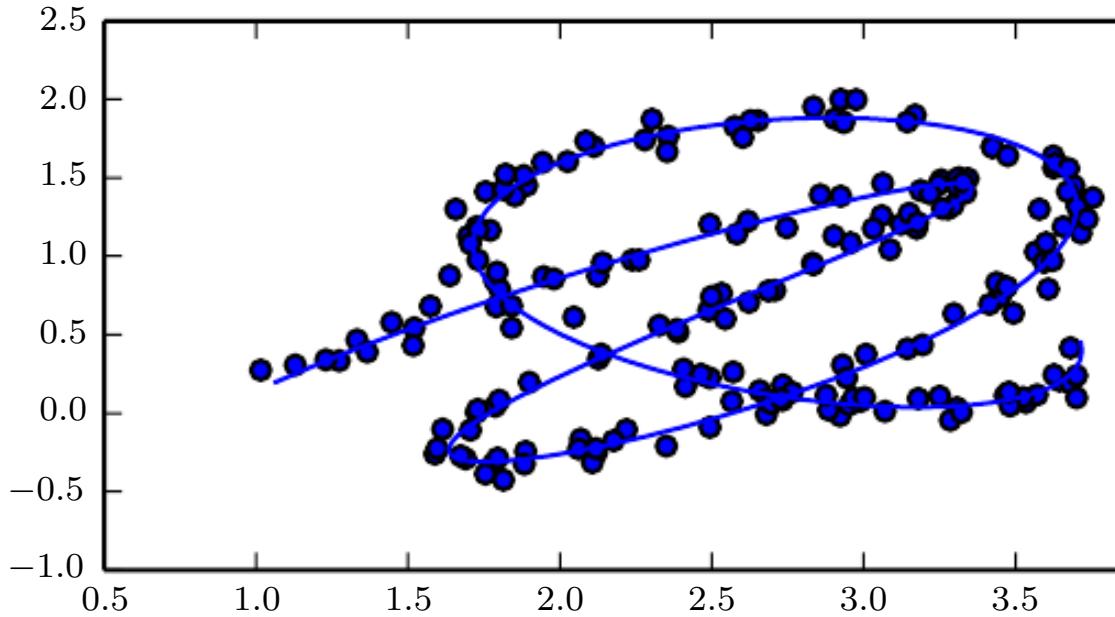


Figure 5.11: Data sampled from a distribution in a two-dimensional space concentrated near a one-dimensional manifold, like a twisted string. The set the underlying manifold that the learner should infer.

degrees of freedom, or dimensions, embedded in a higher-dimensional dimension corresponds to a local direction of variation. See figure example of training data lying near a one-dimensional manifold embedded in dimensional space. In the context of machine learning, we allow the of the manifold to vary from one point to another. This often happens if manifold intersects itself. For example, a figure eight is a manifold that has dimension in most places but two dimensions at the intersection area.

Many machine learning problems seem hopeless if we expect our learning algorithm to learn functions with interesting variations and intersections. **Manifold learning** algorithms surmount this obstacle by assuming that the space of \mathbb{R}^n consists of invalid inputs, and that interesting inputs occur on a collection of manifolds containing a small subset of points, with the variations in the output of the learned function occurring only at the points that lie on the manifold, or with interesting variations happening

that lie on the manifold, or with interesting variations happening move from one manifold to another. Manifold learning was introduced of continuous-valued data and in the unsupervised learning setting probability concentration idea can be generalized to both discrete supervised learning setting: the key assumption remains that probabilities are highly concentrated.

The assumption that the data lies along a low-dimensional manifold may not always be correct or useful. We argue that in the context of AI, especially those that involve processing images, sounds, or text, the manifolds

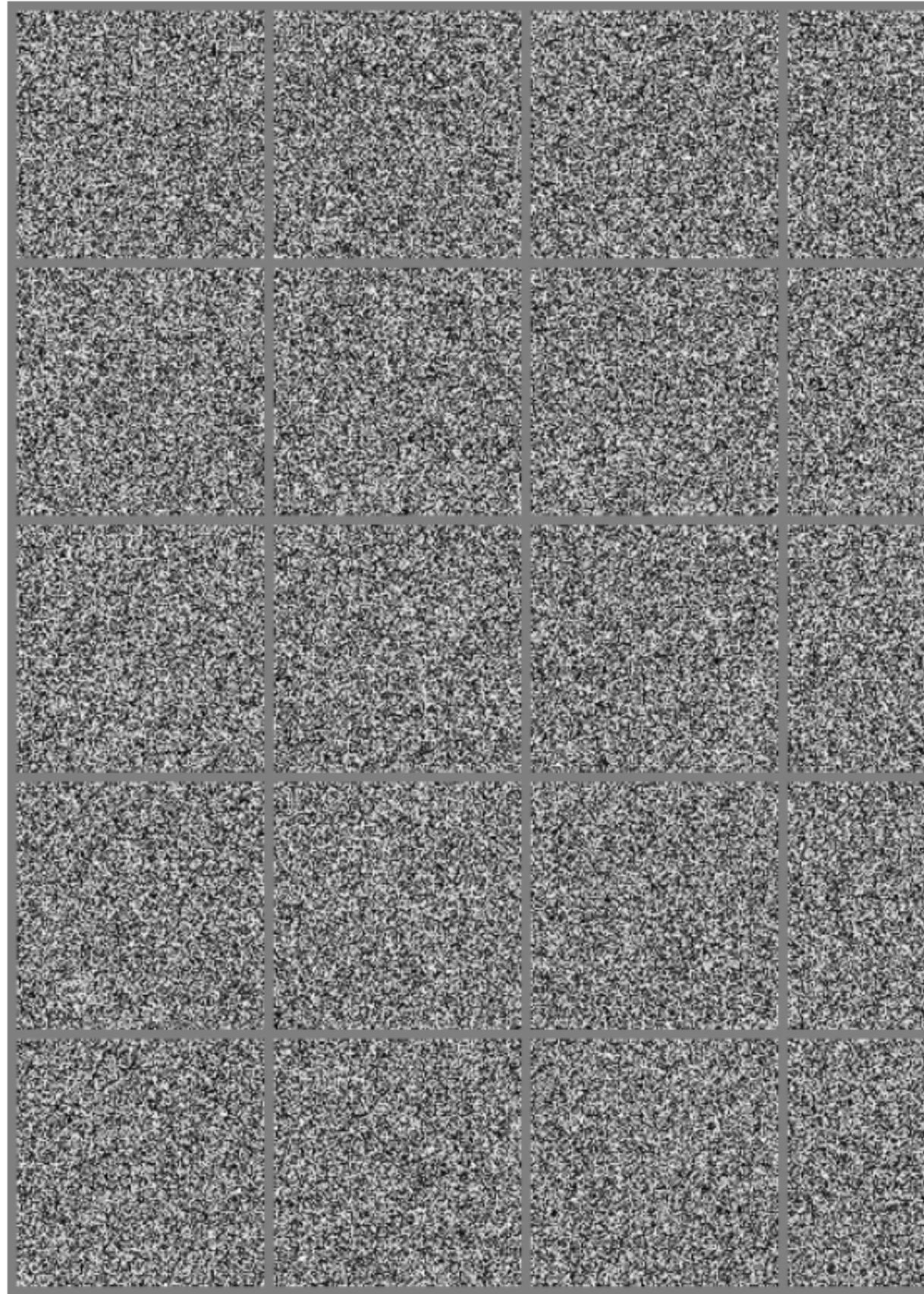


Figure 5.12: Sampling images uniformly at random (by randomly pi

according to a uniform distribution) gives rise to noisy images. Although the nonzero probability of generating an image of a face or of any other object encountered in AI applications, we never actually observe this happening suggests that the images encountered in AI applications occupy a negligible volume of image space.

at least approximately correct. The evidence in favor of this assumption consists of two categories of observations.

The first observation in favor of the **manifold hypothesis** is

bility distribution over images, text strings, and sounds that occurs highly concentrated. Uniform noise essentially never resembles strings from these domains. Figure 5.12 shows how, instead, uniformly distributed noise look like the patterns of static that appear on analog television sets. If you generate a sequence of letters at random, what is the probability that it is available. Similarly, if you generate a document by picking letters at random, what is the probability that you will get a meaningful English text? Almost zero, again, because most of the long sequences of letters correspond to a natural language sequence: the distribution of natural language sequences occupies a very little volume in the total space of sequences.

Of course, concentrated probability distributions are not sufficient to explain the manifold hypothesis. We must also assume that the data lies on a reasonably small number of manifolds. We must also assume that the examples we encounter are connected to each other by continuous paths, so that we can move from one example to another with each example surrounded by other highly similar examples that are close to it. We can then move along such paths by applying transformations to traverse the manifold. The second reason in favor of the manifold hypothesis is that we can imagine such neighborhoods of points as being generated by smooth continuous transformations, at least informally. In the case of images, we can imagine a manifold of images as being generated by a large number of many possible transformations that allow us to trace out a manifold in the space of all possible images. For example, we can gradually dim or brighten the lights, gradually rotate the camera, gradually change the colors of objects in the image, gradually alter the colors on the surfaces of objects, and so forth. Multiple manifolds are likely involved in most applications. For example, the manifold of human face images may not be connected to the manifold of images of cars or the manifold of images of faces of dogs.

These thought experiments convey some intuitive reasons supporting the manifold hypothesis. More rigorous experiments (Cayton, 2005; Narayan et al., 2010; Schölkopf *et al.*, 1998; Roweis and Saul, 2000; Tenenbaum *et al.*, 2000; Belkin and Niyogi, 2003; Donoho and Grimes, 2003; Weinberger and Saul, 2004) clearly support the hypothesis for a large class of datasets of images.

When the data lies on a low-dimensional manifold, it can be more easily learned by machine learning algorithms to represent the data in terms of coordinates in a low-dimensional space.

$$\mathbb{R}^n$$

manifold, rather than in terms of coordinates in . In everyday life we think of roads as 1-D manifolds embedded in 3-D space. We give direct addresses in terms of address numbers along these 1-D roads, not coordinates in 3-D space. Extracting these manifold coordinates is what holds the promise of improving many machine learning algorithms. This principle is applied in many contexts. Figure 5.13 shows the manifold structure of a dataset consisting of faces. By the end of this book, we will have learned the methods necessary to learn such a manifold structure. In figure 2.13 we show how a machine learning algorithm can successfully accomplish this.

This concludes part I, which has provided the basic concepts and machine learning that are employed throughout the remaining book. You are now prepared to embark on your study of deep learning.





Figure 5.13: Training examples from the QMUL Multiview Face Data ([2000](#)), for which the subjects were asked to move in such a way as to form a two-dimensional manifold corresponding to two angles of rotation. We would like our algorithms to be able to discover and disentangle such manifold coordinates. This figure illustrates such a feat.