

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. This chapter provides a brief course in the most important concepts that are applied throughout the rest of the book. Novice readers seeking a wider perspective are encouraged to consider machine learning from a more comprehensive coverage of the fundamentals, such as [M. Mitchell \(2006\)](#). If you are already familiar with machine learning, skip ahead to section [5.11](#). That section covers some perspectives on machine learning techniques that have strongly influenced the development of deep learning algorithms.

We begin with a definition of what a learning algorithm is, with the example: the linear regression algorithm. We then proceed to

descent. We describe how to combine various algorithms into an optimization algorithm, a cost function, a model, and a machine learning algorithm. Finally, in section 5.11, we discuss factors that have limited the ability of traditional machine learning. These challenges have motivated the development of deep learning to 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 definition: “A computer program is said to learn from experience E with respect to a class of tasks T and performance measure P , if its performance on tasks in T , measured by P , improves with experience E .” One can imagine a machine learning algorithm that takes experiences E , tasks T , and performance measures P , and produces a machine learning algorithm. In this book to formally define what may be used for each of these concepts, in the following sections, we provide intuitive descriptions of different kinds of tasks, performance measures, and experiences, and how to construct machine learning algorithms.

5.1.1 The Task, T

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

Many kinds of tasks can be solved with machine learning. Some of the most 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, a machine learning algorithm is usually asked to produce a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Given an input \mathbf{x} , $y = f(\mathbf{x})$, the model assigns an input described by \mathbf{x} to a category identified by numeric code y . There are other variations of this task, for example, where f outputs a probability distribution over the k categories. An example of a classification task is object recognition. Given an image (usually described as a set of pixel brightness values), the model's 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, recognizing different kinds of drinks and deliver them to people (Burgard and Roth, 2006; [fellow et al., 2010](#)). Modern object recognition is based on deep learning ([Krizhevsky et al., 2012](#); [Ioffe and Szegedy, 2015](#)). Face 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 effectively with users.
- **Classification with missing inputs:** Classification is challenging if the computer program is not guaranteed that its input vector will always be provided. To solve this task, a machine learning algorithm only has to define a *single* function that maps an input to a categorical output. When some of the inputs are missing, rather than providing a single classification function, the model must learn a *set* of functions. Each function corresponds to a specific input vector.

- **Regression:** In this type of task, the computer program is asked to output a numerical value given some input. To solve this task, the program 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 cost of a claim that an insured person will make (used to set insurance premiums) or the prediction of future prices of securities. These kinds of prediction are used in algorithmic trading.
- **Transcription:** In this type of task, the machine learning program is asked to observe a relatively unstructured representation of information and transcribe the information into discrete textual form. For example, in optical character recognition, the computer program is asked to take an image containing an image of text and is asked to return a sequence of characters (e.g., in ASCII or Unicode format). Microsoft's Bing View uses deep learning to process address numbers in images (Bengio *et al.*, 2014d). Another example is speech recognition: a machine learning program is provided an audio waveform and emits a sequence of word ID codes describing the words that were spoken. Microsoft's Cortana uses deep learning for speech recognition (Bengio *et al.*, 2012b). Deep learning is a crucial component of modern speech recognition systems, and is used at major companies, including Microsoft, IBM, and Google (Bengio *et al.*, 2012b).
- **Machine translation:** In a machine translation task, the input is a sequence of symbols in some language, and the output is a sequence of symbols in another language. This task must convert this into a sequence of symbols in another language. This is commonly applied to natural languages, such as translating English into French. Deep learning has recently begun to have an impact on machine translation (Bengio *et al.*, 2014c).

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

- **Anomaly detection:** In this type of task, the computer goes 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 distribution over purchase types than your own. The company can prevent fraud by placing a hold on an account as soon as it has been used for an uncharacteristic purchase. See Chameleon for a survey of anomaly detection methods.
- **Synthesis and sampling:** In this type of task, the computer algorithm is asked to generate new examples that are similar to the training data. Synthesis and sampling via machine learning are useful for media applications when generating large volumes of data would be expensive, boring, or require too much time. For example, video games can automatically generate textures for large

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

- **Denoising:** In this type of task, the machine learning algorithm is given as input a *corrupted example* $\tilde{\mathbf{x}} \in \mathbb{R}^n$ obtained by an unknown corruption of a *clean example* $\mathbf{x} \in \mathbb{R}^n$. The learner must predict \mathbf{x} from its corrupted version $\tilde{\mathbf{x}}$, or more generally predict the probability distribution $p(\mathbf{x} \mid \tilde{\mathbf{x}})$.
- **Density estimation or probability mass function estimation:** In a density estimation problem, the machine learning algorithm is given 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. In other words, well (we will specify exactly what that means when we discuss 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 learner to at least implicitly capture the structure of the probability distribution. Density estimation enables us to explicitly capture that distribution, so we can then perform computations on that distribution. This is useful for 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 the other values, denoted \mathbf{x}_{-i} , are given, then we know the distribution of x_i is given by $p(x_i \mid \mathbf{x}_{-i})$. In practice, density estimation does not solve all these related tasks, because in many cases the computations on $p(\mathbf{x})$ are computationally intractable.

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

Usually we are interested in how well the machine learning system performs on data that it has not seen before, since this determines how well it will be 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 learning system.

The choice of performance measure may seem straightforward, but it is often difficult to choose a performance measure that accurately reflects the desired behavior of the system.

In some cases, this is because it is difficult to decide what the desired behavior is. For example, when performing a transcription task, should we evaluate the performance of the system at transcribing entire sequences, or should we evaluate it at the word level? A performance measure that gives partial credit for getting parts of a sequence correct? When performing a regression task, should we evaluate the system more if it frequently makes medium-sized mistakes, or should we care more about very large mistakes? These kinds of design choices depend on the specific application.

In other cases, we know what quantity we would ideally like to optimize, but measuring it is impractical. For example, this arises frequently in the context of density estimation. Many of the best probabilistic models for

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

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

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

Supervised learning algorithms experience a dataset where each example is also associated with a **label** or **target**. In the Iris dataset is annotated with the species of each iris plant. A supervised algorithm can study the Iris dataset and learn to classify new examples into different species based on their measurements.

Roughly speaking, unsupervised learning involves observing a single example of a random vector \mathbf{x} and attempting to implicitly or explicitly estimate the probability distribution $p(\mathbf{x})$, or some interesting properties of this distribution. Supervised learning involves observing several examples of a random vector \mathbf{x} and an associated value or vector \mathbf{y} , then learning to predict \mathbf{y} from \mathbf{x} , or estimating $p(\mathbf{y} | \mathbf{x})$. The term **supervised learning** originates from the

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

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

Though unsupervised learning and supervised learning are not mutually exclusive or distinct concepts, they do help roughly categorize some of the machine learning algorithms. Traditionally, people refer to problems with a known and structured output problems as supervised learning. Problems without the support of other tasks is usually considered unsupervised learning.

Other variants of the learning paradigm are possible. In semi-supervised learning, some examples include a supervision signal and some 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 of 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 single episode. For example, **reinforcement learning** algorithms interact with an environment where there is a feedback loop between the learning system and the environment. These algorithms are beyond the scope of this book. Please see [Sutton and Barto \(2018\)](#) 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 sequence of examples that can be described in many ways. In all cases, a dataset is a collection of examples which are in turn collections of features.

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

In the case of supervised learning, the example contains not only a target value but also a collection of features. For example, if we want to use machine learning to perform object recognition from photographs, we need to know the label of the object that appears in each of the photos. We might do this with a target vector \mathbf{y} , where 0 signifying a person, 1 signifying a car, 2 signifying a cat, and so on. When working with a dataset containing a design matrix of features \mathbf{X} , we also provide a vector of labels \mathbf{y} , with y_i providing the label for the i th example.

Of course, sometimes the label may be more than just a single value. For example, if we want to train a speech recognition system to recognize 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 standard examples cover most cases, but it is always possible to design new ones.

5.1.4 Example: Linear Regression

Our definition of a machine learning algorithm as an algorithm that improves 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 algorithms to help understand the algorithm's behavior.

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 given \mathbf{x} . We can write this as $\hat{y} = \mathbf{w}^\top \mathbf{x}$. Next we need a definition of our performance measure.

Suppose that we have a design matrix of m examples \mathbf{X} and a vector of regression targets \mathbf{y} providing the correct values for each example. We use \mathbf{X} for training, only for evaluating how well the model predicts on new 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 targets as $\mathbf{y}^{(\text{test})}$.

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

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

Intuitively, one can see that this error measure decreases to zero as the model's predictions get closer to the targets. 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.

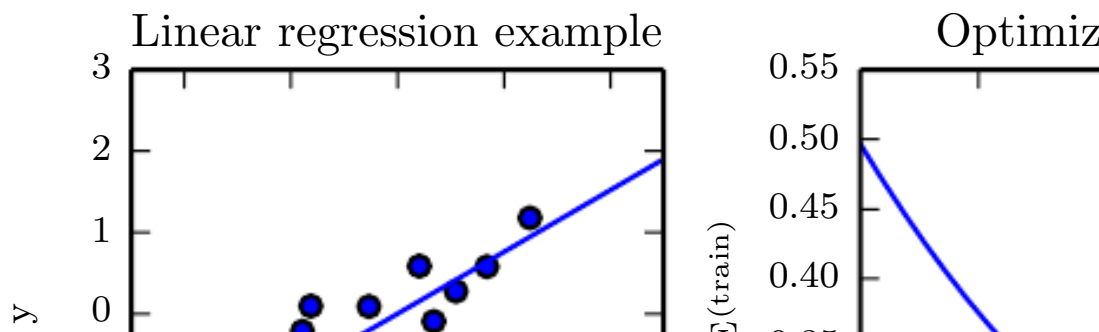
$$\begin{aligned}
 &\Rightarrow \frac{1}{m} \nabla_{\mathbf{w}} \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = 0 \\
 &\Rightarrow \nabla_{\mathbf{w}} \left(\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right)^\top \left(\mathbf{X}^{(\text{train})}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

The system of equations whose solution is given by equation 5.12 are called the **normal equations**. Evaluating equation 5.12 constitutes the **linear regression 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 a slightly more sophisticated model with one additional parameter b . In this model

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

so the mapping from parameters to predictions is still a linear mapping, the mapping from features to predictions is now an affine function.



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

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

Linear regression is of course an extremely simple and limited model, but it provides an example of how a learning algorithm can be designed. In the sections we describe some of the basic principles underlying the design and demonstrate how these principles can be used to design 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 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

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

How can we affect performance on the test set when we change the training set? The field of **statistical learning theory** provides a framework for understanding how the training and the test set are collected arbitrarily, there is no control. 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 probabilistic process. These datasets are called the **data-generating process**. We typically make several assumptions known collectively as the **i.i.d. assumptions**. The first assumption is that the examples in each dataset are **independent** from each other. The second assumption is that the training set and test set are **identically distributed**, meaning they are drawn from the same probability distribution as each other. This assumption allows us to model the data-generating process with a probability distribution. The same distribution is then used to generate every training and test example. We call that shared underlying distribution the **data distribution**, denoted p_{data} . This probabilistic framework and the i.i.d. assumptions enable us to mathematically study the relationship between training and test error.

One immediate connection we can observe between training 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 and test sets. For some fixed value \mathbf{w} , the expected training set error is equal to the expected test set error, because both expectations are taken over the same dataset sampling process. The only difference between the training and test error is the name we assign to the dataset we sample.

Of course, when we use a machine learning algorithm,

obtain a sufficiently low error value on the training set. On the other hand, if the gap between the training error and test error is too large, the model is overfitting.

We can control whether a model is more likely to overfit by controlling its **capacity**. Informally, a model's capacity is its ability to fit a wide range of functions. Models with low capacity may struggle to fit the training data, while models with high capacity can overfit by memorizing properties of the training data that do not serve them well on the test set.

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

A polynomial of degree 1 gives us the linear regression model, which we have already seen. The prediction function is already familiar, with the prediction

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

By introducing x^2 as another feature provided to the linear model, 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 x , it is still a linear function of the *parameters*, so we can still use the same algorithm to train the model in closed form. We can continue to add more features, for example, to obtain a polynomial of degree 2 or higher.

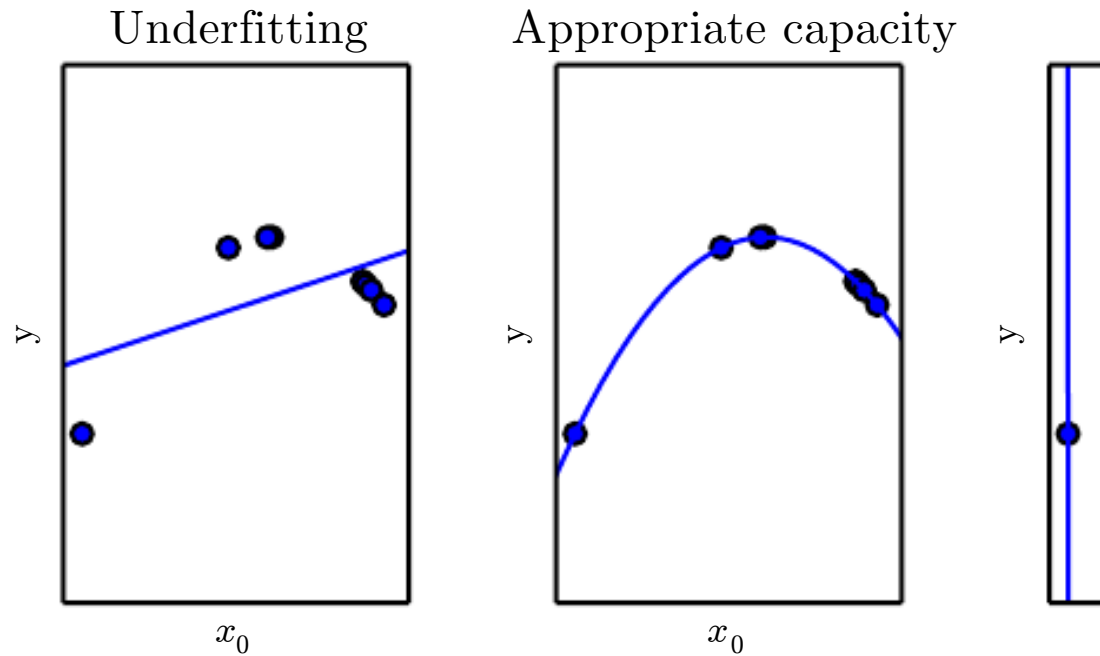


Figure 5.2: We fit three models to this example training set, which was generated synthetically, by randomly sampling x values and choosing y by evaluating a quadratic function. *(Left)* A linear function fit suffers from underfitting—it cannot capture the curvature that is present in the data. *(Middle)* A quadratic function fit to the data generalizes well to unseen points. *(Right)* A polynomial fit of degree 9 to the data suffers from overfitting. Here we used the Moore-Penrose pseudoinverse to solve the underdetermined normal equations. The solution passes through all training points exactly, but we have not been lucky enough for it to extrapolate well. It now has a deep valley between two training points that does not exist in the underlying function. It also increases sharply on the left side of the plot, while the function decreases in this area.

The underlying function is quadratic. The linear function is unable to capture the true underlying problem, so it underfits. The degree-9 polynomial is capable of representing the correct function, but it is also capable of overfitting to the training data.

function within this family is a difficult optimization problem. A learning algorithm does not actually find the best function, and this significantly reduces the training error. These additional limitations, due to the imperfection of the optimization algorithm, mean that the **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 philosophy, starting with Ptolemy. Many early scholars invoke a principle of parsimony, now widely known as **Occam's razor** (c. 1287–1347). This principle states that among competing hypotheses that explain known observations equally well, we choose the “simplest” one. This idea was formalized and developed in the twentieth century by the founders of statistical learning theory (Chervonenkis, 1971; Vapnik, 1982; Blumer *et al.*, 1989; Vapnik, 1995).

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

Quantifying the capacity of the model enables statisticians to make quantitative predictions. The most important results of statistical learning theory show that the discrepancy between training error and test error is bounded from above by a quantity that grows as the model capacity shrinks as the number of training examples increases (Vapnik, 1971; Vapnik, 1982; Blumer *et al.*, 1989; Vapnik, 1995). This provides an intellectual justification that machine learning algorithms are rarely used in practice when working with deep learning.

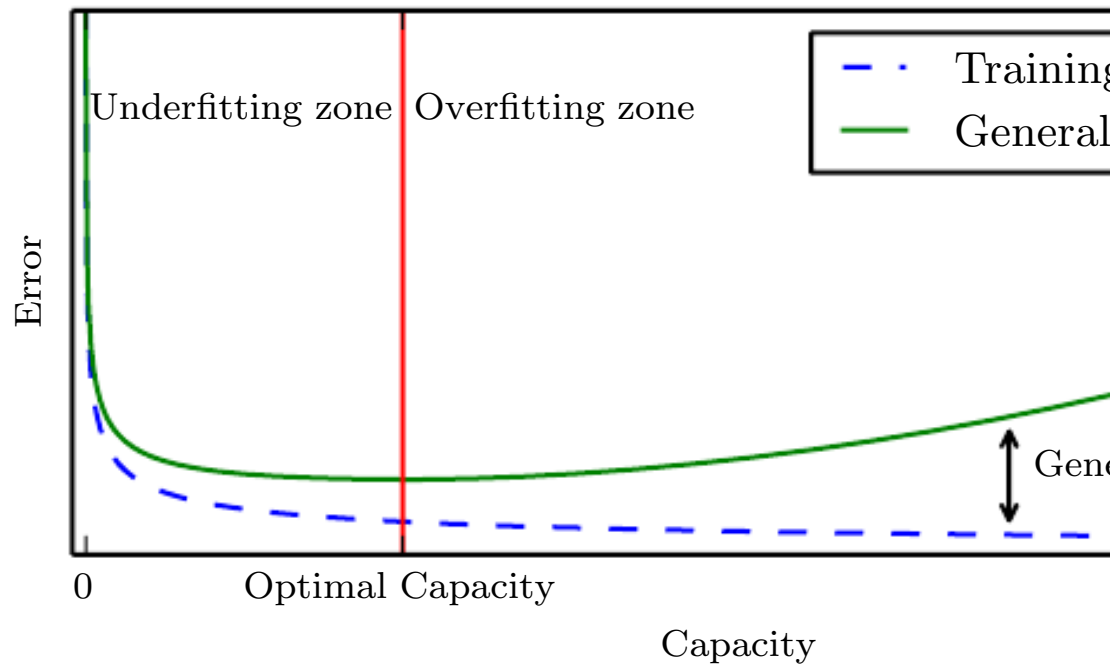


Figure 5.3: Typical relationship between capacity and error. Training and generalization errors behave differently. At the left end of the graph, training error and generalization error are both high. This is the **underfitting regime**. As we increase capacity, training error decreases, but the gap between training and generalization error increases. At the right end, the size of this gap outweighs the decrease in training error, and the generalization error increases. This is 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 parametric models, such as linear regression. Parametric models learn by a parameter vector whose size is finite and fixed before training. Nonparametric models have no such limitation.

Sometimes, nonparametric models are just theoretical

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

Finally, we can also create a nonparametric learning algorithm by embedding a parametric learning algorithm inside another algorithm that has a large number of parameters as needed. For example, we could imagine an algorithm that changes the degree of the polynomial learned by linear regression as a function of 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 supervised problems, because there may still be some noise in the data. In the context of supervised learning, the mapping from \mathbf{x} to y may be stochastic, or y may be a deterministic function that involves other variables not included in \mathbf{x} . The error incurred by an oracle making predictions from the true 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 results until the best possible error is achieved. Any fixed parametric model with an optimal capacity will asymptote to an error value that exists even with infinite data. See figure 5.4 for an illustration. Note that it is possible to have a 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 using more training examples.

5.2.1 The No Free Lunch Theorem

Learning theory claims that a machine learning algorithm cannot consistently outperform a random guesser across all possible datasets.

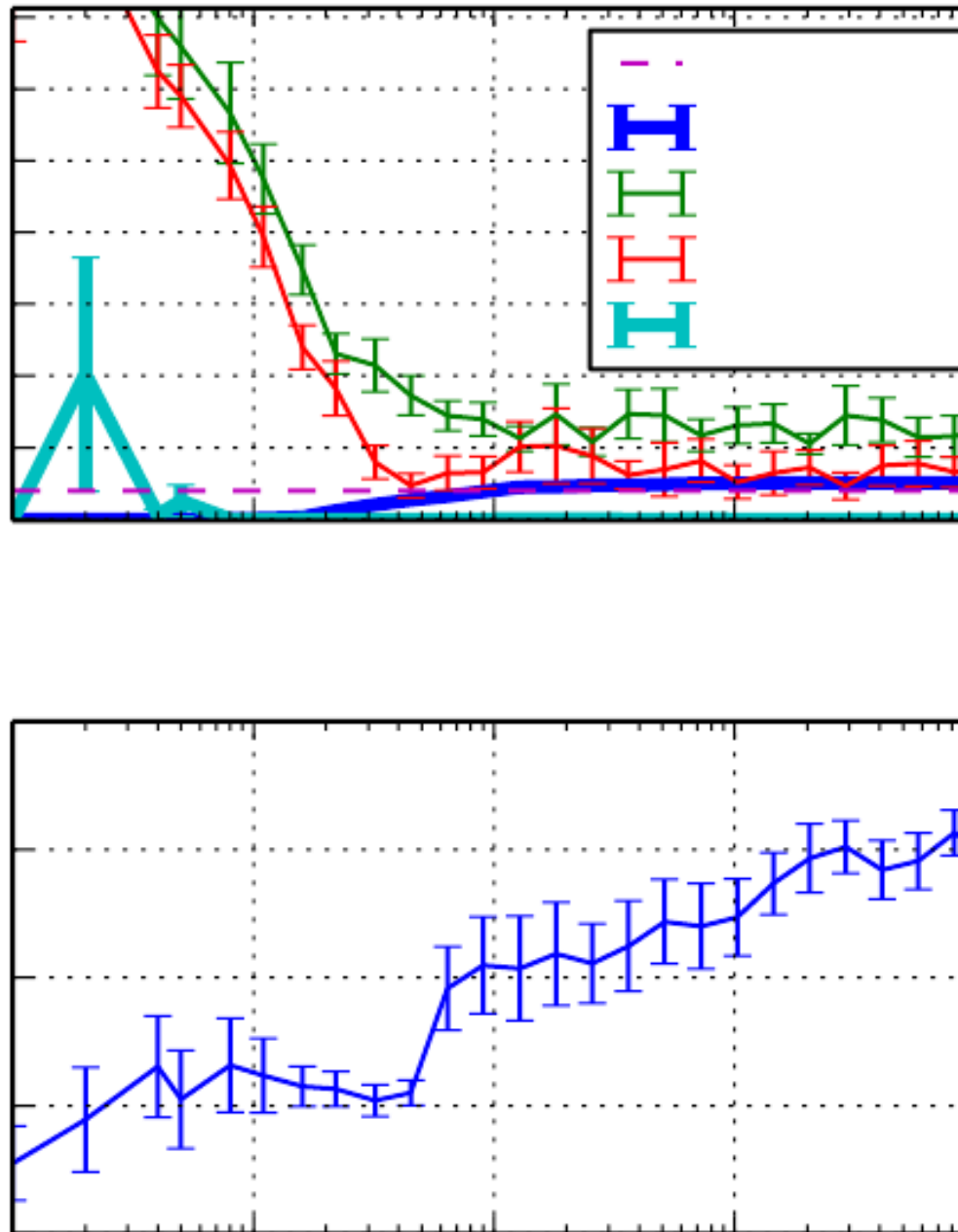


Figure 5.4: The effect of the training dataset size on the train and test error rates on the optimal model capacity. We constructed a synthetic regression problem by adding a moderate amount of noise to a degree-5 polynomial, generated several different sizes of training set. For each size, we generated 10 different training sets in order to plot error bars showing 95% confidence intervals.

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

Fortunately, these results hold only when we average over all possible data-generating distributions. If we make assumptions about the kinds 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 want to understand what kinds of distributions are relevant to the kinds of data an agent experiences, and what kinds of machine learning algorithms are good at learning from data drawn from the kinds of data-generating distributions.

5.2.2 Regularization

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

So far, the only method of modifying a learning algorithm to better match our preferences concretely is to increase or decrease the model's representational capacity by adding or removing functions from the hypothesis space of solutions the algorithm is able to choose from. We gave the specific example of increasing the degree of a polynomial for a regression problem. The point is that the current model so far is oversimplified.

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

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

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

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 the solution is to make the weights to become smaller. Minimizing $J(\mathbf{w})$ results in a tradeoff between fitting the training data and being conservative. Solutions that have a smaller slope, or that put weight on lower-degree terms, are preferred. As an example of how we can control a model's tendency to overfit via weight decay, we can train a high-degree polynomial model for different values of λ . See figure 5.5 for the results.

More generally, we can regularize a model that learns 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 other 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 functions from the hypothesis space. We can think of excluding a function from the hypothesis space as expressing an infinitely strong preference against that function.

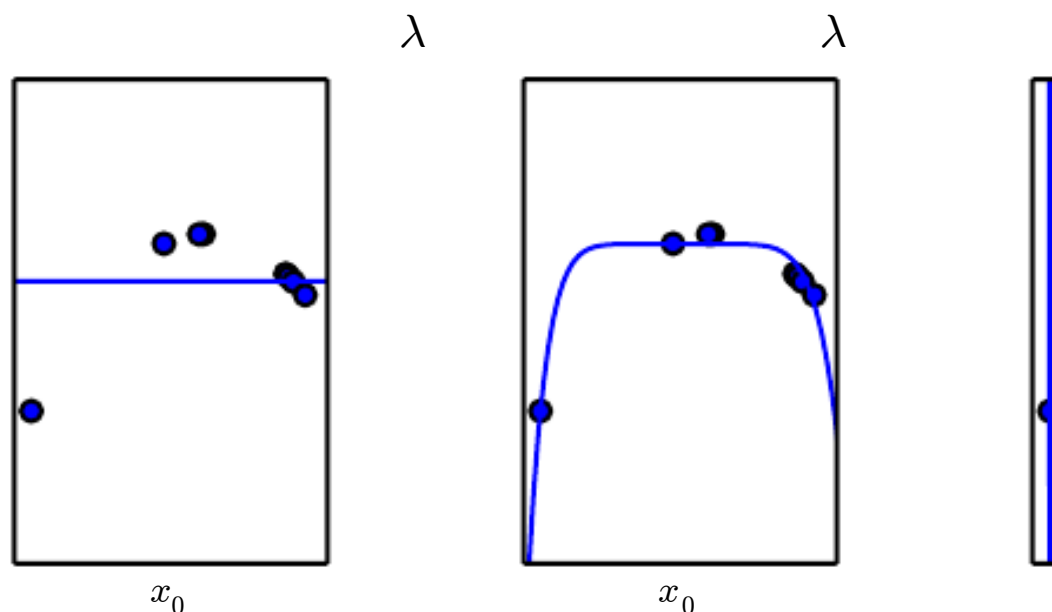
In our weight decay example, we expressed our preference

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

5.3 Hyperparameters and Validation Set

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

The polynomial regression example in figure 5.2 has a 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.



Sometimes a setting is chosen to be a hyperparameter. If an algorithm does not learn because the setting is difficult to optimize, the setting must be a hyperparameter because it is not applicable as a hyperparameter on the training set. This applies to all settings that control model capacity. If learned on the training set, such as λ , we always choose the maximum possible model capacity, resulting in overfitting (see figure 5.3). For example, we can always fit the training data with a higher-degree polynomial and a weight decay setting of $\lambda = 0$, or 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 from the same distribution as the training set, can be used to estimate the error of a learner, after the learning process has completed. In this case, the test examples are not used in any way to make choices about the learner's hyperparameters. For this reason, no example from the training set is in the validation set. Therefore, we always construct the validation set from the *training* data. Specifically, we split the training data into two subsets. One of these subsets is used to learn the parameters. The other is the validation set, used to estimate the generalization error during the learning process, allowing for the hyperparameters to be updated accordingly. The subset of data used to learn the parameters is still typically called the training set, though this may be confused with the larger pool of data used for the entire 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 used to “train” the hyperparameters, the validation set error v

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 high variance around the estimated average test error, making it difficult to compare algorithms. A works better than algorithm B on the given task.

When the dataset has hundreds of thousands of examples, this is not a serious issue. When the dataset is too small, alternative methods are needed to use all the examples in the estimation of the mean test error without increased computational cost. These procedures are based on repeating the training and testing computation on different random splits 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 split into k nonoverlapping subsets. The test error is 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 for training. One problem is that no unbiased estimators of the variance of the test error exist (Bengio and Grandvalet, 2004), but approximations are used.

5.4 Estimators, Bias and Variance

The field of statistics gives us many tools to achieve the 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 help to characterize notions of generalization, underfitting and overfitting.

Algorithm 5.1 The k -fold cross-validation algorithm. It computes the generalization error of a learning algorithm A when the dataset is small for a simple train/test or train/valid split to yield a reliable estimate of the generalization error, because the mean of a loss L on a small dataset has a high variance. The dataset \mathbb{D} contains as elements the absolute values of the i -th example), which could stand for an (input,target) pair in the case of supervised learning, or for just an input $\mathbf{z}^{(i)}$ in the case of unsupervised learning. The algorithm returns the vector of generalization errors, whose mean is the estimated generalization error. Individual examples can be used to compute a confidence interval around the mean (equation 5.47). Though these confidence intervals are not used after the use of cross-validation, it is still common practice to say that algorithm A is better than algorithm B only if the confidence interval of algorithm A lies below and does not intersect the confidence interval of algorithm B .

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

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 and 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

a good estimator is a function whose output is close to the generated the training data.

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

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

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

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

5.4.2 Bias

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 using m 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

To determine the bias of the sample mean, we are again 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
 In this example, we compare two different estimators of the variance of a Gaussian distribution. We are interested in knowing if either is unbiased.

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

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

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

The **unbiased sample variance** estimator

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

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

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

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

5.4.3 Variance and Standard Error

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

been different. The expected degree of variation in any 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 error 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 tend to underestimate the true standard deviation but are still consistent. The square root of the unbiased estimator of the variance is less biased. 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 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. We can use the standard error to compute the probability that the 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

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

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

5.4.4 Trading off Bias and Variance to Minimize Error

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

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

The most common way to negotiate this trade-off is to use cross-validation. Empirically cross-validation is highly successful on many real-world datasets.

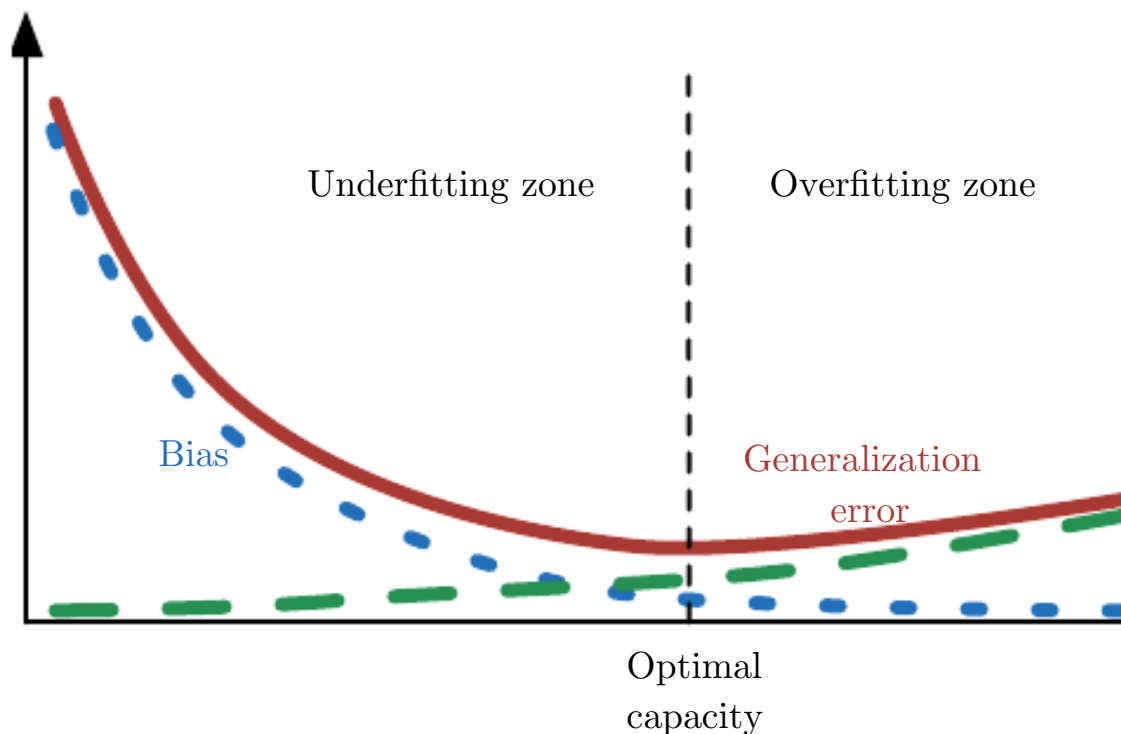


Figure 5.6: As capacity increases (x -axis), bias (dotted) tends to decrease, while variance (dashed) tends to increase, yielding another U-shaped curve for generalization error (solid red curve). If we vary capacity along one axis, there is an optimal capacity. When the capacity is below this optimum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

Generalization error is measured by the MSE (where bias and variance are components of generalization error), increasing capacity tends to increase variance and decrease bias. This is illustrated in figure 5.6, where we see again the relationship between generalization error as a function of capacity.

5.4.5 Consistency

So far we have discussed the computation of maximum likelihood estimates

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

Consistency ensures that the bias induced by the estimator goes to zero 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 m samples: $\{x^{(1)}, \dots, x^{(m)}\}$. We could use the first sample as an unbiased estimator: $\hat{\theta} = x^{(1)}$. In that case, $\mathbb{E}(\hat{\theta}_m)$ is unbiased no matter how many data points are seen. This does not mean that the estimate is asymptotically unbiased. However, this is not an estimator 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 an example of how to use them. But where did these estimators come from? Rather than just guessing, a function might make a good estimator and then analyzing its properties we would like to have some principle from which we can deduce which functions are good estimators for different models.

The most common such principle is the maximum likelihood estimation.

Consider a set of m examples $\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ drawn 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 to a real number estimating the true probability $p_{\text{data}}(\mathbf{x})$.

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

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 sum by dividing by m to obtain a version of the criterion that is expressed with respect to the empirical distribution \hat{p}_{data} defined by

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

One way to interpret maximum likelihood estimation is to view it as minimizing the dissimilarity between the empirical distribution \hat{p}_{data} , of the data set and the model distribution, with the degree of dissimilarity 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 (5.1).

Minimizing this KL divergence corresponds exactly to minimizing the cross-entropy between the distributions. Many authors use the term “cross-entropy” to identify specifically the negative log-likelihood of a Bernoulli distribution, but that is a misnomer. Any loss consisting of a negative log-likelihood is a cross-entropy between the empirical distribution defined by the data and the probability distribution defined by the model. For example, minimizing the cross-entropy between the empirical distribution and a Gaussian distribution is equivalent to maximizing the log-likelihood of the Gaussian model.

5.5.1 Conditional Log-Likelihood and Mean Squared Error

The maximum likelihood estimator can readily be generalized to a 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 machine learning algorithms. \mathbf{X} represents all our inputs and \mathbf{Y} all our observed target values. 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 as

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

Example: Linear Regression as Maximum Likelihood Estimation The linear regression introduced in section 5.1.4, may be justified as a maximum likelihood estimation. Previously, we motivated linear regression as an algorithm that takes an input \mathbf{x} and produce an output value \hat{y} . The mapping from input to output is learned by minimizing the mean squared error, a criterion that we introduced 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 view linear regression as producing a conditional distribution $p(y \mid \mathbf{x})$. We can think of this as 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 for each input \mathbf{x} . To derive the same linear regression algorithm we can assume that the data is generated by a linear model with Gaussian noise.

where $\hat{y}^{(i)}$ is the output of the linear regression on the i -th example, $y^{(i)}$ is the target value, and m is the number of the training examples. Comparing the log-likelihood and the 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 parameters \mathbf{w} is equivalent to minimizing the squared error. The two criteria have different values but the same location of the maximum. This justifies the use of the MSE as a maximum likelihood estimator. 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 is believed to be the best estimator asymptotically, as the number of examples approaches infinity. 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 estimator converges to the true value of the parameter. These conditions are:

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

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

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

5.6 Bayesian Statistics

So far we have discussed **frequentist statistics** and approached the problem of finding a single value of θ , then making all predictions thereafter using that estimate. Another approach is to consider all possible values of θ for a given prediction. The latter is the domain of **Bayesian statistics**.

As discussed in section 5.4.1, the frequentist perspective on statistics is that the parameter value θ is fixed but unknown, while the point estimate is variable on account of it being a function of the dataset (which is random).

The Bayesian perspective on statistics is quite different. In Bayesian statistics, probability is used to reflect degrees of certainty in states of knowledge. A parameter is directly observed and so is not random. On the other hand, a prediction is unknown or uncertain and thus is represented as a random variable.

Before observing the data, we represent our knowledge of the parameter as a **probability distribution**, $p(\theta)$ (sometimes referred to as the **prior**). Generally, the machine learning practitioner selects a prior distribution that is

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

Relative to maximum likelihood estimation, Bayesian has some important differences. First, unlike the maximum likelihood approach, predictions using a point estimate of θ , the Bayesian approach uses a full distribution over θ . For example, after observing $x^{(1)}, \dots, x^{(m)}$, the predicted distribution over the next data sample, $x^{(m+1)}$, is

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

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

In section 5.4, we discussed how the frequentist approach measures uncertainty in a given point estimate of θ by evaluating its variance. The estimator is an assessment of how the estimate might change across samplings of the observed data. The Bayesian answer to the question of how to deal with the uncertainty in the estimator is to simply integrate over the uncertainty. This integral is of course subject to the laws of probability, making the Bayesian approach simpler than the frequentist machinery for constructing an estimator based on a single decision to summarize all knowledge contained in the data into a point estimate.

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

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

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

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

Expressed as a Gaussian conditional distribution on $\mathbf{y}^{(\text{train})}$ given $\mathbf{X}^{(\text{train})}$ and \mathbf{w} ,

$$\begin{aligned} p(\mathbf{y}^{(\text{train})} | \mathbf{X}^{(\text{train})}, \mathbf{w}) &= \mathcal{N}(\mathbf{y}^{(\text{train})}; \mathbf{X}^{(\text{train})} \mathbf{w}, \mathbf{I}) \\ &\propto \exp \left(-\frac{1}{2} (\mathbf{y}^{(\text{train})} - \mathbf{X}^{(\text{train})} \mathbf{w})^\top (\mathbf{y}^{(\text{train})} - \mathbf{X}^{(\text{train})} \mathbf{w}) \right) \end{aligned}$$

where we follow the standard MSE formulation in assuming that the variance on y is one. In what follows, to reduce the notation, we write $(\mathbf{X}^{(\text{train})}, \mathbf{y}^{(\text{train})})$ 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 express 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 $\boldsymbol{\theta}$. For real-valued parameters it is common to use a Gaussian prior.

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

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

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

All terms that do not include the parameter vector \mathbf{w} have been moved to the right-hand side. These terms are implied by the fact that the distribution must be normalized. 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 the zero vector. Then $\boldsymbol{\mu}_m$ gives the same estimate of \mathbf{w} as does frequentist least squares. The weight decay penalty of $\alpha \mathbf{w}^\top \mathbf{w}$ is equivalent to a Gaussian prior on \mathbf{w} with covariance matrix $\boldsymbol{\Lambda}_0^{-1}$. One difference is that the prior is undefined if α is set to zero—we are not allowed to begin the learning 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.

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, $\log p(\mathbf{x} \mid \boldsymbol{\theta})$, the likelihood term, and $\log p(\boldsymbol{\theta})$, corresponding to the prior distribution.

As an example, consider a linear regression model with weights \mathbf{w} . If this prior is given by $\mathcal{N}(\mathbf{w}; \mathbf{0}, \frac{1}{\lambda} \mathbf{I}^2)$, then 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 likelihood. Bayesian inference with a Gaussian prior on the weights thus introduces weight decay.

As with full Bayesian inference, MAP Bayesian inference leverages information that is brought by the prior and the training data. This additional information helps to reduce the bias of the MAP point estimate (in comparison to the ML estimate) at the price of increased bias.

Many regularized estimation strategies, such as maximum a posteriori estimation, can be interpreted as making a connection to Bayesian inference. This view applies when the regularization is achieved by adding an extra term to the objective function that corresponds to the prior. In fact, all regularization penalties correspond to MAP Bayesian inference. However, 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.

“supervisor,” but the term still applies even when the training data is collected automatically.

5.7.1 Probabilistic Supervised Learning

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

We have already seen that linear regression corresponds to the normal distribution:

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

We can generalize linear regression to the classification problem by using a different family of probability distributions. If we have two classes, say class 0 and class 1, then we need only specify the probability of one class. The probability of class 1 determines the probability of class 0, because 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 is valid. A distribution over a binary variable is slightly more constrained: its mean must always be between 0 and 1. One way to solve this is to use the logistic sigmoid function to squash the output of the linear model into the interval (0, 1) and interpret that value as a probability:

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

This approach is known as **logistic regression** (a somewhat

5.7.2 Support Vector Machines

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

One key innovation associated with support vector machines is the **kernel trick**. The kernel trick consists of observing that many machine learning algorithms can be written exclusively in terms of dot products between examples. In this case, 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. By using this learning algorithm this way enables us to replace \mathbf{x} with the output of a feature function $\phi(\mathbf{x})$ and the dot product with a function $k(\mathbf{x}, \mathbf{x}^{(i)})$, 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 representation, but in some infinite dimensional spaces, we need to use other kinds of inner products. For example, inner products based on integration rather than summation. The development of these kinds of inner products is beyond the scope of this book.

After replacing dot products with kernel evaluations, we can write the support vector machine using the function

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

admits an implementation that is significantly more computationally efficient 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, incurring an infinite computational cost for the naive, explicit approach. However, $k(\mathbf{x}, \mathbf{x}')$ is a nonlinear, tractable function of \mathbf{x} even when $\phi(\mathbf{x})$ is infinite dimensional. An example of an infinite-dimensional feature space with a tractable kernel is the **Reproducing Kernel Hilbert Space (RKHS)**. We can construct a feature mapping $\phi(x)$ over the nonnegative integers \mathbb{N} such that this mapping returns a vector containing x ones followed by zeros. We can write a kernel function $k(x, x^{(i)}) = \min(x, x^{(i)})$ that corresponds 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 also known as the **radial basis function (RBF)** kernel, because its value at \mathbf{v} is like a wave in \mathbf{v} space radiating outward from \mathbf{u} . The Gaussian kernel is equivalent to the dot product in an infinite-dimensional space, but the derivation is more straightforward than in our example of the min kernel over the nonnegative integers.

We can think of the Gaussian kernel as performing a kind of **template matching**. A training example \mathbf{x} associated with training label y is 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 similar to the template. The model then puts a large weight on the associated training label y . 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 be used for classification. Many other linear models can be used for classification, such as logistic regression.

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

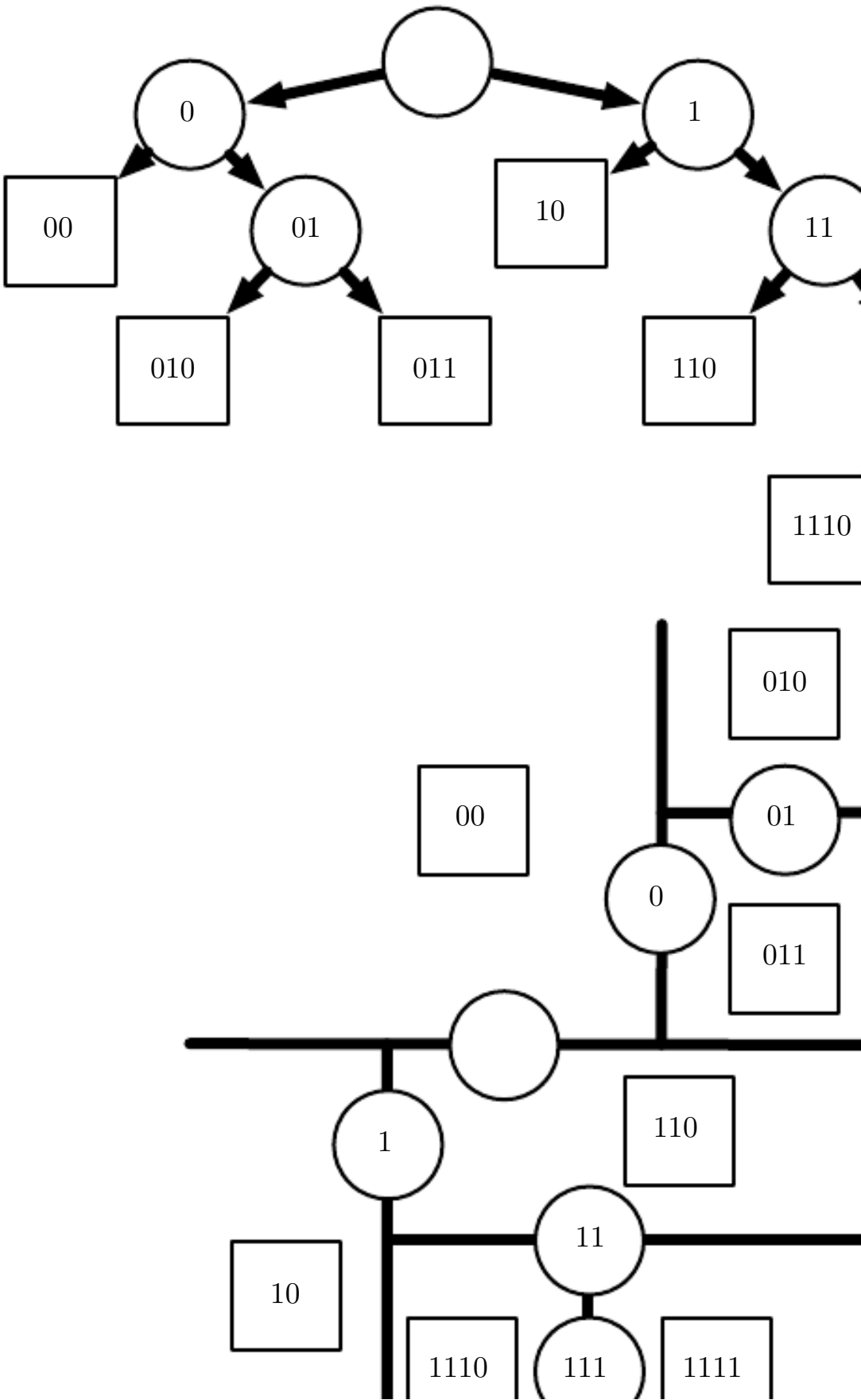
5.7.3 Other Simple Supervised Learning Algorithms

We have already briefly encountered another nonprobabilistic algorithm, nearest neighbor regression. More generally, a family of techniques that can be used for classification is nonparametric learning algorithm, k -nearest neighbors is number of parameters. We usually think of the k -nearest as not having any parameters but rather implementing a training data. In fact, there is not even really a training stage. Instead, at test time, when we want to produce an output y we find the k -nearest neighbors to \mathbf{x} in the training data \mathcal{X} and average of the corresponding y values in the training set. This is any kind of supervised learning where we can define an average. In the case of classification, we can average over one-hot codes and $c_i = 0$ for all other values of i . We can then interpret one-hot codes as giving a probability distribution over classes. In learning algorithm, k -nearest neighbor can achieve very high accuracy. Suppose we have a multiclass classification task and measure error using cross-entropy loss. In this setting, 1-nearest neighbor converges to double the error as the number of training examples approaches infinity. The error reduction results from choosing a single neighbor by breaking

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

Another type of learning algorithm that also breaks the iid assumption 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 a tree is associated with a region in the input space, and internal nodes split a region into one subregion for each child of the node (typically by a linear cut). Space is thus subdivided into nonoverlapping regions, establishing a correspondence between leaf nodes and input regions. Each leaf node maps every point in its input region to the same output. Decision trees are trained with specialized algorithms that are beyond the scope of this book. A learning algorithm can be considered nonparametric if it is not of arbitrary size, though decision trees are usually regularized in ways that turn them into parametric models in practice. Decision trees are typically used, with axis-aligned splits and constant outputs, to struggle to solve some problems that are easy even for linear models. For example, if we have a two-class problem, and the positive class is defined by $x_2 > x_1$, the decision boundary is not axis aligned. The only way to need to approximate the decision boundary with many nodes is to use a function that constantly walks back and forth across the space, which is not possible with axis-aligned steps.

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



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

A classic unsupervised learning task is to find the “best” representation of data. By “best” we can mean different things, but generally speaking, we look for a representation that preserves as much information about the data as possible, obeying some penalty or constraint aimed at keeping the representation more accessible than \mathbf{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 compress as much information about \mathbf{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 often involves increasing the dimensionality of the representation, so that the representation becoming mostly zeros does not discard too much information. The overall structure of the representation that tends to distribute the information across the representation space. Independent representations assume that the sources of variation underlying the data distribution sum to the total variation 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 suggests that to reduce the size of a representation is to find and remove redundancy. Finding and removing more redundancy enables the dimensionality

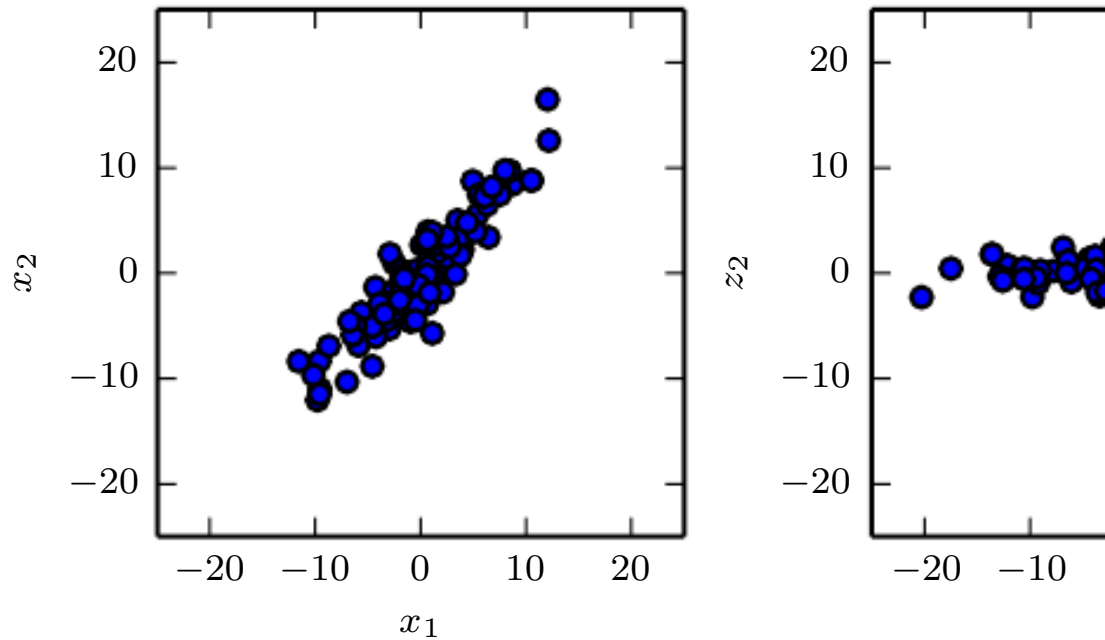


Figure 5.8: PCA learns a linear projection that aligns the direction of the axes of the new space. (Left) The original data consist of samples where the variance might occur along directions that are not axis aligned. (Right) The data $\mathbf{z} = \mathbf{x}^\top \mathbf{W}$ now varies most along the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

5.8.1 Principal Components Analysis

In section 2.12, we saw that the principal components analysis is a means of compressing data. We can also view PCA as an algorithm that learns a representation of data. This representation satisfies two of the criteria for a simple representation described in section 2.12: a representation that has lower dimensionality than the original data and a representation whose elements have no linear correlation. This is a first step toward the criterion of learning representations that are statistically independent. To achieve full independence, a representation must also be

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

The unbiased sample covariance matrix associated with the data is

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

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

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

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

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} , i.e., the right singular vectors in the decomposition $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$. We can rewrite the original eigenvector equation with \mathbf{W} as the eigenvectors of $\mathbf{X}^\top \mathbf{X}$ as

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

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

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

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

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

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

This ability of PCA to transform data into a representation where the features are mutually uncorrelated is a very important property of this representation. It is an example of a representation that attempts to *disentangle* the different sources of *variation* underlying the data. In the case of PCA, this decomposition is achieved by the form of finding a rotation of the input space (described by the principal axes of variance) with the basis of the new representation, \mathbf{z} .

While correlation is an important category of dependence in the data, we are also interested in learning representations that capture more complicated forms of feature dependencies. For this, we will see that more can be done with a simple linear transformation.

5.8.2 k -means Clustering

Another example of a simple representation learning algorithm is the k -means clustering algorithm. The k -means clustering algorithm divides the training set into

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. In one step, each training example is assigned to cluster i , the nearest centroid $\mu^{(i)}$. In the other step, each centroid is updated to the mean of all training examples $\mathbf{x}^{(j)}$ 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 guarantees that the clustering of the data corresponds to the real world. We can evaluate the quality of the clustering, such as the average Euclidean distance from each example to the members of the cluster. This enables us to tell how well the clustering reconstructs the training data from the cluster assignments, but it does not tell us how well the cluster assignments correspond to properties of the real world. There may be many different clusterings that all correspond to the real world. We may hope to find a clustering that reliably captures the real world, but we can also obtain a different, equally valid clustering that is not relevant to the real world. For example, suppose that we run two clustering algorithms on a set of 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, the first algorithm might find a cluster of cars and a cluster of trucks, while another algorithm might find a cluster of red vehicles and a cluster of gray vehicles. Suppose we also run a third clustering algorithm, which is allowed to determine the number of clusters. If we ask this algorithm to cluster the examples to four clusters, red cars, red trucks, gray cars, and gray trucks, the new clustering now at least captures information about both color and type, but it has lost information about similarity. Red cars are in a different cluster from gray cars, just as they are in a different cluster from gray trucks. A clustering algorithm does not tell us that red cars are more

attributes instead of just testing whether one attribute ma

5.9 Stochastic Gradient Descent

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

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

The cost function used by a machine learning algorithm is a sum over training examples of some per-example loss function. The negative conditional log-likelihood of the training data can

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

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

For these additive cost functions, gradient descent requires

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

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 long.

FIGURE 5.1.1: SGD: A single step of stochastic gradient descent.

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

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \mathbf{g},$$

where ϵ is the learning rate.

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

Stochastic gradient descent has many important uses in 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 is $O(1/m)$ where m is the training set size. In practice, we often use a larger model size as the dataset size 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 linear model. SGD has sampled every example in the training set. Increasing m only extends the amount of training time needed to reach the minimum error. From this point of view, one can argue that the asymptotic complexity of a model with SGD is $O(1)$ as a function of m .

Prior to the advent of deep learning, the main way to train linear models was to use the kernel trick in combination with a linear model. Kernel algorithms require constructing an $m \times m$ matrix $G_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$.

5.10 Building a Machine Learning Algorithm

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

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

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

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

By realizing that we can replace any of these components with something other than the others, we can obtain a wide range of algorithms.

The cost function typically includes at least one term that represents the process to perform statistical estimation. The most common is the negative log-likelihood, so that minimizing the cost function is equivalent to likelihood estimation.

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

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

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 optimization

In some cases, the cost function may be a function that is difficult to evaluate, for computational reasons. In these cases, we can minimize it using iterative numerical optimization, as long as we are approximating its gradients.

Most machine learning algorithms make use of this recipe, which can be immediately obvious. If a machine learning algorithm seems to be hand designed, it can usually be understood as using a special set of models, such as decision trees and k -means, require special cost functions, their cost functions have flat regions that make them inapproachable by gradient-based optimizers. Recognizing that most machine learning can be described using this recipe helps to see the different taxonomy of methods for doing related tasks that work for different tasks 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 solve a wide variety of important problems. They have not succeeded in solving the central problems in AI, such as recognizing speech or recognizing objects in images.

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

This section is about how the challenge of generalizing to new data is exponentially more difficult when working with high-dimensional data. The mechanisms used to achieve generalization in traditional machine learning are insufficient to learn complicated functions in high-dimensional spaces also often impose high computational costs. Deep learning

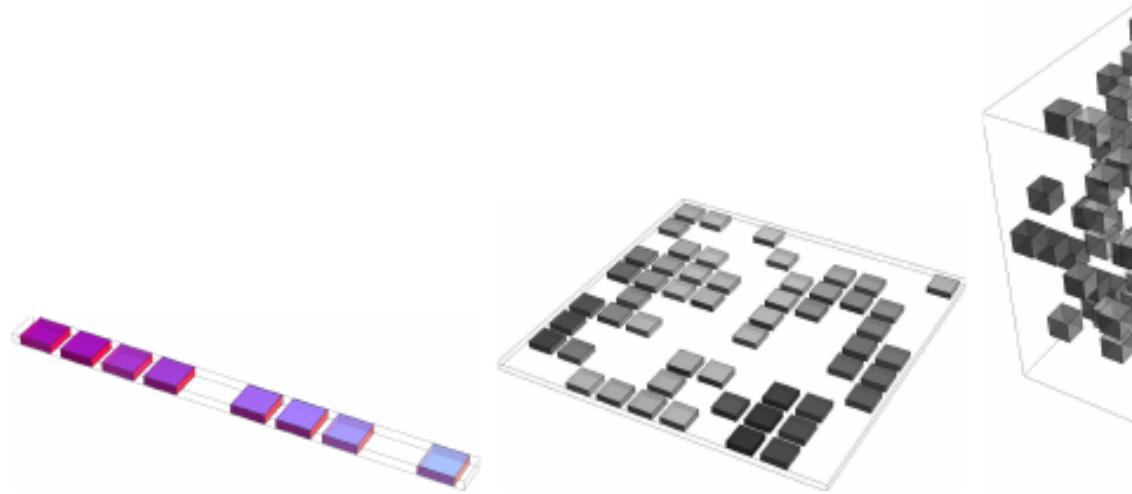


Figure 5.9: As the number of relevant dimensions of the data (left to right), the number of configurations of interest may grow exponentially. (Left) In a one-dimensional example, we have one variable for which we only need a few regions of interest. With enough examples falling within each of these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

One challenge posed by the curse of dimensionality is the sparsity of data. 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 examples. To understand the issue, let us consider that the input space is a 10x10x10 grid, as in the figure. We can describe low-dimensional spaces by the number of grid cells that are mostly occupied by the data. When gen-

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 informed about what kind of function they should learn. We have seen that this is represented as explicit beliefs in the form of probability distributions over the model. More informally, we may also discuss prior beliefs about the *function* itself and influencing the parameters only indirectly through the relationship between the parameters and the function. Additionally, we can discuss prior beliefs as being expressed implicitly by choices of algorithms that are biased toward choosing some class of functions over others. These biases may not be expressed (or even be possible to express) as a probability distribution representing our degree of belief in the space of functions.

Among the most widely used of these implicit “priors” is the **local constancy 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, and as a result, they fail to scale to the statistical challenges of high-dimensional tasks. Throughout this book, we describe how deeper models use additional (explicit and implicit) priors in order to reduce error on sophisticated tasks. Here, we explain why the local constancy prior is insufficient for these tasks.

There are many different ways to implicitly or explicitly encode the belief that the learned function should be smooth or locally constant. Many methods are designed to encourage the learning process to learn

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

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

Decision trees also suffer from the limitations of exclusive learning, because they break the input space into as many regions as leaves and use a separate parameter (or sometimes many parameters) of decision trees) in each region. If the target function requires at least n leaves to be represented accurately, then at least n regions are required to fit the tree. A multiple of n is needed to achieve high confidence in the predicted output.

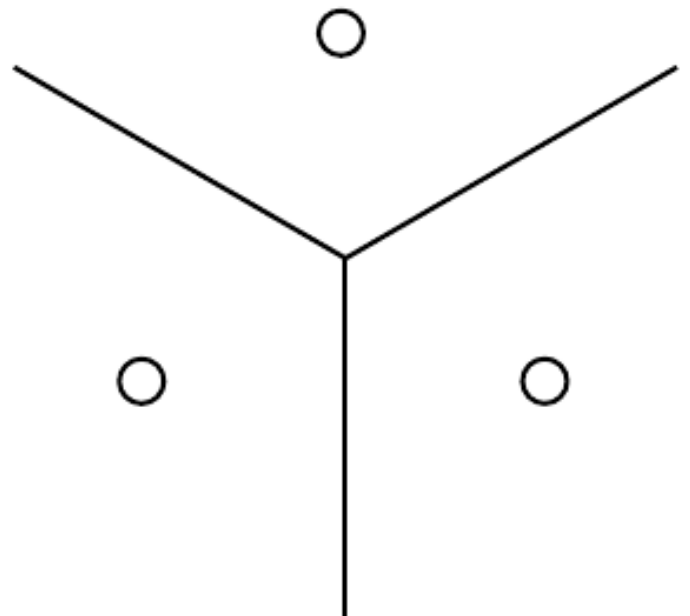
In general, to distinguish $O(k)$ regions in input space, at least $O(k)$ examples are needed. Typically there are $O(k)$ parameters, \mathbf{w}_i , associated with each of the $O(k)$ regions. The nearest neighbor model where each training example can be used to define at most one region is shown in figure 5.10.

Is there a way to represent a complex function that has more regions than can be distinguished than the number of training examples? (The smoothness of the underlying function will not allow a local model to succeed in this case.)

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 directions. In high dimensions, even a very smooth function can change in a very different way along each dimension. If the function additionally varies in various regions, it can become extremely complicated to learn from a small number of training examples. If the function is complicated (we want to know if the number of regions compared to the number of examples) can the algorithm 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 learned function to generalize well to new inputs—is yes. The key insight is that the number of regions, such as $O(2^k)$, can be defined with $O(k)$ parameters.



introduce some dependencies between the regions through a model that captures information about the underlying data-generating distribution. In this way, models can generalize nonlocally (Bengio and Monperrus, 2005; Bengio et al., 2007). Different deep learning algorithms provide implicit or explicit regularization, which is reasonable for a broad range of AI tasks in order to capture the underlying structure.

Other approaches to machine learning often make strong assumptions. For example, we could easily solve the checkerboard problem by making the assumption that the target function is periodic. Usually, deep learning makes strong, task-specific assumptions in neural networks so that they can be applied to a much wider variety of structures. AI tasks have structures that are too complex to be limited to simple, manually specified properties, so we want learning algorithms that embody more general assumptions. The core idea in deep learning is that we assume that the target function is the *composition of factors*, or features, potentially at multiple scales of complexity. Many other similarly generic assumptions can further guide the design of learning algorithms. These apparently mild assumptions allow an exponential relationship between the number of examples and the number of parameters to be distinguished. We describe these exponential gains more fully in sections 6.4.1, 15.4 and 15.5. The exponential advantages conferred by distributed representations counter the exponential challenge posed by the curse of dimensionality.

5.11.3 Manifold Learning

An important concept underlying many ideas in machine learning is the manifold.

A **manifold** is a connected region. Mathematically,

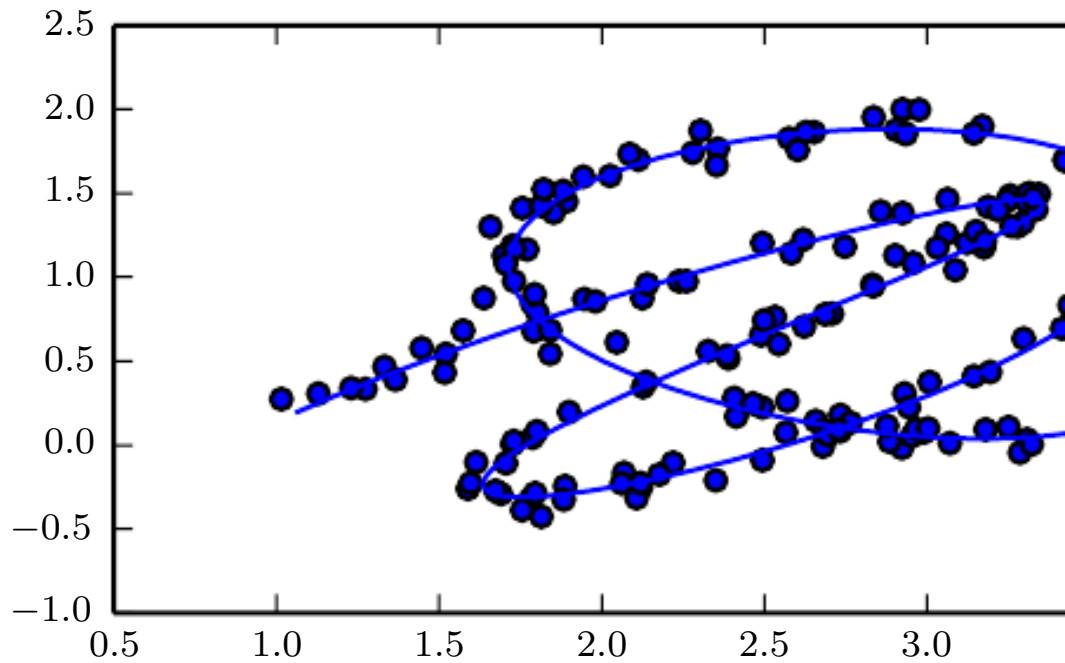
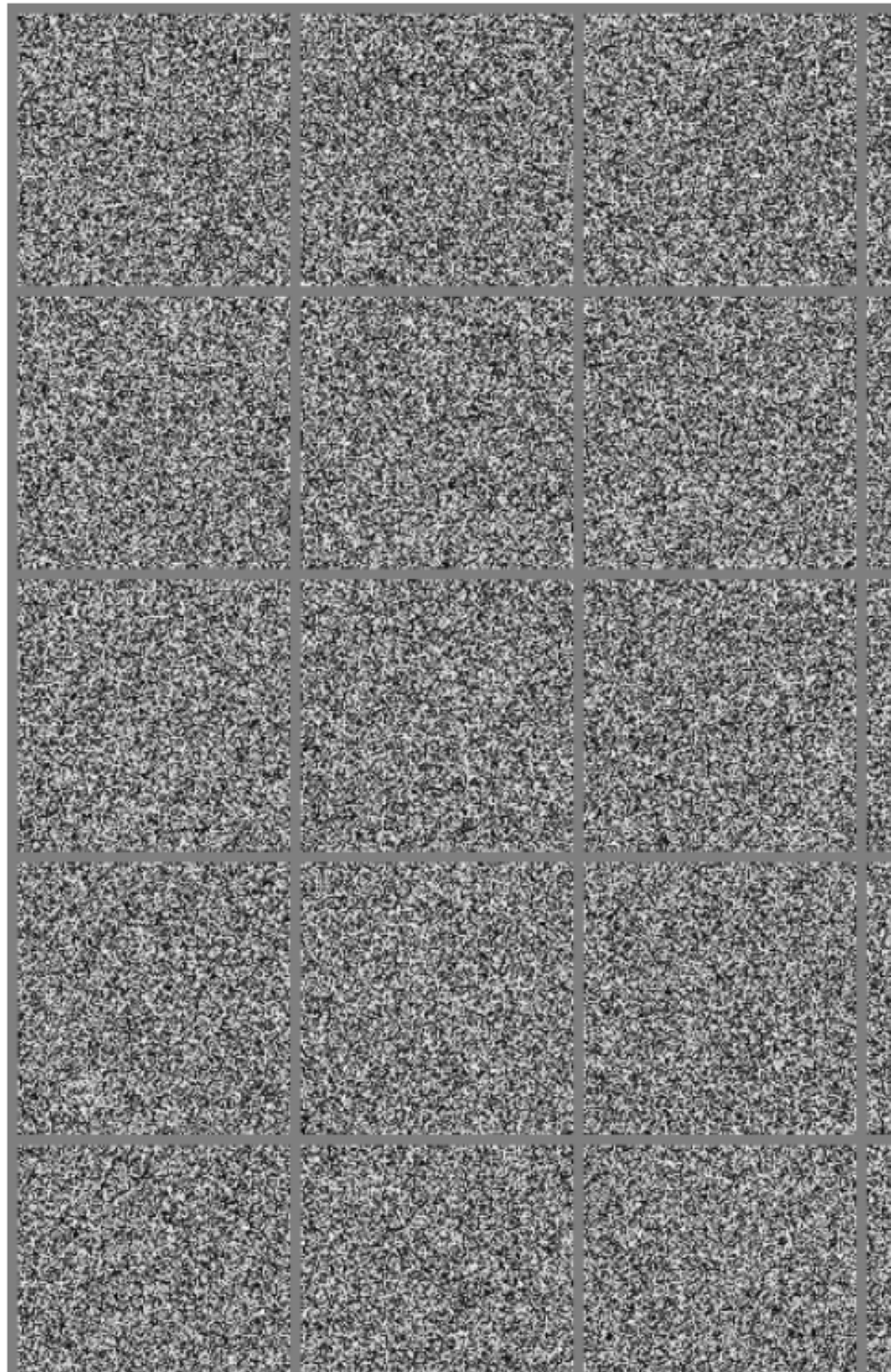


Figure 5.11: Data sampled from a distribution in a two-dimensional space concentrated near a one-dimensional manifold, like a twisted string. The goal is to infer the underlying manifold that the learner should infer.

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

Many machine learning problems seem hopeless if we try to use a standard learning algorithm to learn functions with interesting variation. **Manifold learning** algorithms surmount this obstacle by recognizing that the set of all possible inputs of \mathbb{R}^n consists of invalid inputs, and that interesting inputs



bility distribution over images, text strings, and sounds that is highly concentrated. Uniform noise essentially never resembles anything from these domains. Figure 5.12 shows how, instead, uniform noise looks like the patterns of static that appear on analog television when no signal is available. Similarly, if you generate a document by picking words at random, what is the probability that you will get a meaningful text? Almost zero, again, because most of the long sequences of words correspond to a natural language sequence: the distribution of natural language sequences occupies a very little volume in the total space of all possible sequences.

Of course, concentrated probability distributions are not unusual in nature: the data lies on a reasonably small number of manifolds. What is important is that the examples we encounter are connected to each other, and that with each example surrounded by other highly similar examples. We can move from one example to another by applying transformations to traverse the manifold. The main reason in favor of the manifold hypothesis is that we can imagine such transformations, at least informally. In the case of images, there are a lot of many possible transformations that allow us to trace out a path through the space: we can gradually dim or brighten the lights, gradually move objects in the image, gradually alter the colors on the surface, 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 other face images.

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

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

