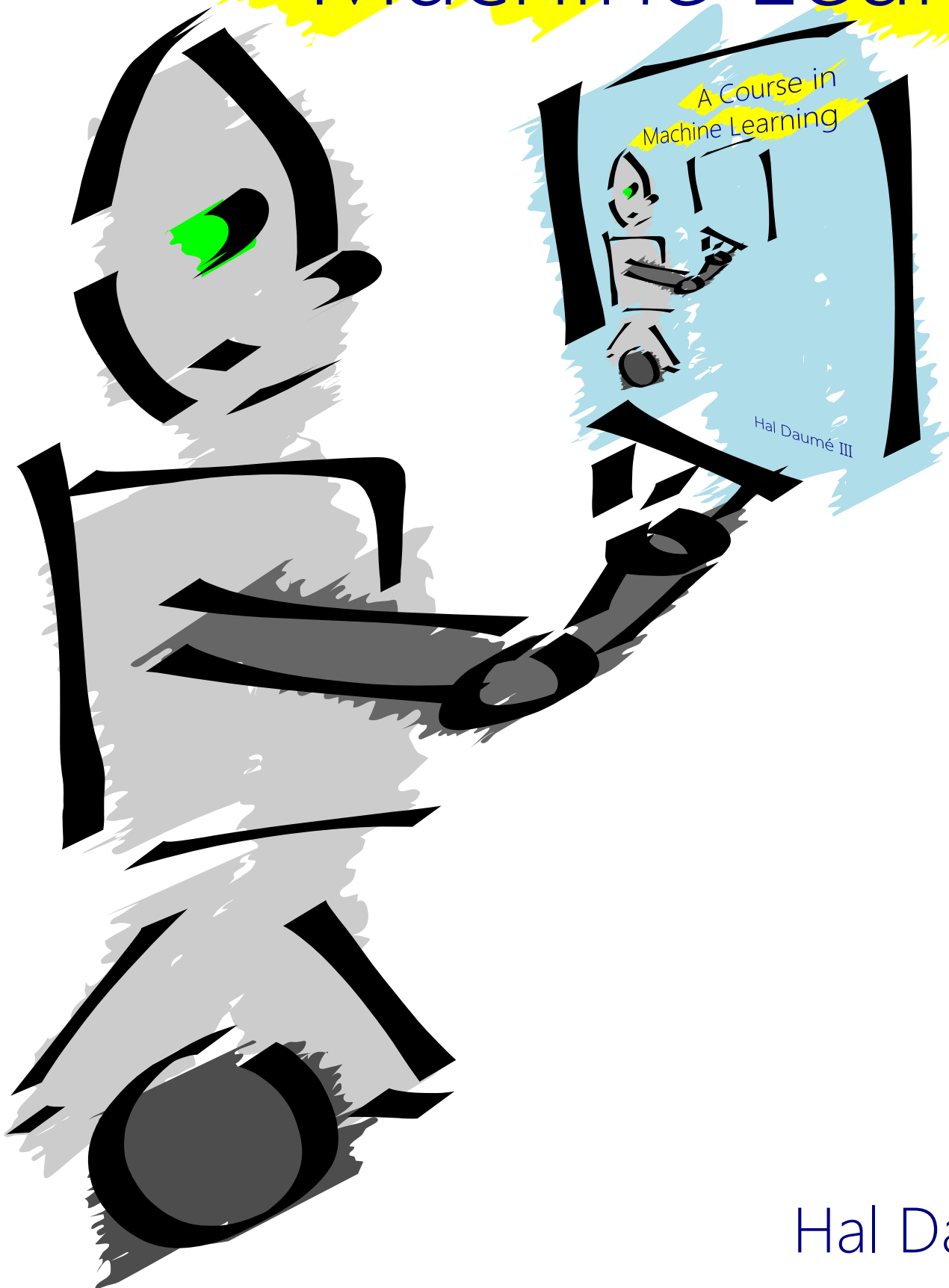


A Course in Machine Learning



Hal Daumé III

Copyright © 2013–2017 Hal Daumé III

SELF-PUBLISHED

<http://ciml.info/>

TODO...

Second printing, January 2017

For my students and teachers.

Often the same.

TABLE OF CONTENTS

	About this Book	6
1	Decision Trees	8
2	Limits of Learning	19
3	Geometry and Nearest Neighbors	29
4	The Perceptron	41
5	Practical Issues	55
6	Beyond Binary Classification	73
7	Linear Models	87
8	Bias and Fairness	104
9	Probabilistic Modeling	116
10	Neural Networks	129

11	Kernel Methods	141
12	Learning Theory	154
13	Ensemble Methods	164
14	Efficient Learning	171
15	Unsupervised Learning	178
16	Expectation Maximization	186
17	Structured Prediction	195
18	Imitation Learning	212
	Code and Datasets	222
	Bibliography	223
	Index	225

ABOUT THIS BOOK

MACHINE LEARNING IS A BROAD AND FASCINATING FIELD. Even today, machine learning technology runs a substantial part of your life, often without you knowing it. Any plausible approach to artificial intelligence must involve learning, at some level, if for no other reason than it's hard to call a system intelligent if it *cannot* learn. Machine learning is also fascinating in its own right for the philosophical questions it raises about what it means to learn and succeed at tasks.

Machine learning is also a very broad field, and attempting to cover everything would be a pedagogical disaster. It is also so quickly moving that any book that attempts to cover the latest developments will be outdated before it gets online. Thus, this book has two goals. First, to be a gentle introduction to what is a very deep field. Second, to provide readers with the skills necessary to pick up new technology as it is developed.

0.1 *How to Use this Book*

This book is designed to be read linearly, since it's goal is not to be a generic reference. That said, once you get through chapter 5, you can pretty much jump anywhere. When I teach a one-semester undergraduate course, I typically cover the chapter 1-13, sometimes skipping 7 or 9 or 10 or 12 depending on time and interest. For a graduate course for students with no prior machine learning background, I would very quickly blaze through 1-4, then cover the rest, augmented with some additional reading.

0.2 *Why Another Textbook?*

The purpose of this book is to provide a *gentle* and *pedagogically organized* introduction to the field. This is in contrast to most existing machine learning texts, which tend to organize things topically, rather

than pedagogically (an exception is Mitchell’s book¹, but unfortunately that is getting more and more outdated). This makes sense for researchers in the field, but less sense for learners. A second goal of this book is to provide a view of machine learning that focuses on ideas and models, not on math. It is not possible (or even advisable) to avoid math. But math should be there to *aid* understanding, not hinder it. Finally, this book attempts to have minimal dependencies, so that one can fairly easily pick and choose chapters to read. When dependencies exist, they are listed at the start of the chapter.

¹ Mitchell 1997

The *audience* of this book is anyone who knows differential calculus and discrete math, and can program reasonably well. (A little bit of linear algebra and probability will not hurt.) An undergraduate in their fourth or fifth semester should be fully capable of understanding this material. However, it should also be suitable for first year graduate students, perhaps at a slightly faster pace.

0.3 *Organization and Auxiliary Material*

There is an associated web page, <http://ciml.info/>, which contains an online copy of this book, as well as associated code and data. It also contains errata. Please submit bug reports on github: github.com/hal3/ciml.

0.4 *Acknowledgements*

Acknowledgements: I am indebted to many people for this book. My teachers, especially Rami Grossberg (from whom the title of this book was borrowed) and Stefan Schaal. Students who have taken machine learning from me over the past ten years, including those who suffered through the initial versions of the class before I figured out how to teach it. Especially Scott Alfeld, Josh de Bever, Cecily Heiner, Jeffrey Ferraro, Seth Juarez, John Moeller, JT Olds, Piyush Rai. People who have helped me edit, and who have submitted bug reports, including **TODO...**, but also check github for the latest list of contributors!

1 | DECISION TREES

The words printed here are concepts.
You must go through the experiences.

– Carl Frederick

Learning Objectives:

- Explain the difference between memorization and generalization.
- Implement a decision tree classifier.
- Take a concrete task and cast it as a learning problem, with a formal notion of input space, features, output space, generating distribution and loss function.

AT A BASIC LEVEL, machine learning is about predicting the future based on the past. For instance, you might wish to predict how much a user Alice will like a movie that she hasn't seen, based on her ratings of movies that she has seen. This prediction could be based on many factors of the movies: their category (drama, documentary, etc.), the language, the director and actors, the production company, etc. In general, this means making informed guesses about some unobserved property of some object, based on observed properties of that object.

The first question we'll ask is: what does it mean to learn? In order to develop learning machines, we must know what learning actually means, and how to determine success (or failure). You'll see this question answered in a very limited learning setting, which will be progressively loosened and adapted throughout the rest of this book. For concreteness, our focus will be on a very simple model of learning called a **decision tree**.

Dependencies: None.

1.1 What Does it Mean to Learn?

Alice has just begun taking a course on machine learning. She knows that at the end of the course, she will be expected to have “learned” all about this topic. A common way of gauging whether or not she has learned is for her teacher, Bob, to give her a exam. She has done well at learning if she does well on the exam.

But what makes a reasonable exam? If Bob spends the entire semester talking about machine learning, and then gives Alice an exam on History of Pottery, then Alice's performance on this exam will *not* be representative of her learning. On the other hand, if the exam only asks questions that Bob has answered exactly during lectures, then this is also a bad test of Alice's learning, especially if it's an “open notes” exam. What is desired is that Alice observes *specific* examples from the course, and then has to answer new, but related questions on the exam. This tests whether Alice has the ability to

generalize. Generalization is perhaps the most central concept in machine learning.

As a concrete example, consider a course recommendation system for undergraduate computer science students. We have a collection of students and a collection of courses. Each student has taken, and evaluated, a subset of the courses. The evaluation is simply a score from -2 (terrible) to $+2$ (awesome). The job of the recommender system is to **predict** how much a particular student (say, Alice) will like a particular course (say, Algorithms).

Given historical data from course ratings (i.e., the past) we are trying to predict unseen ratings (i.e., the future). Now, we could be unfair to this system as well. We could ask it whether Alice is likely to enjoy the History of Pottery course. This is unfair because the system has no idea what History of Pottery even is, and has no prior experience with this course. On the other hand, we could ask it how much Alice will like Artificial Intelligence, which she took last year and rated as $+2$ (awesome). We would expect the system to predict that she would really like it, but this isn't demonstrating that the system has learned: it's simply recalling its past experience. In the former case, we're expecting the system to generalize *beyond* its experience, which is unfair. In the latter case, we're not expecting it to generalize at all.

This general set up of predicting the future based on the past is at the core of most machine learning. The objects that our algorithm will make predictions about are **examples**. In the recommender system setting, an example would be some particular Student/Course pair (such as Alice/Algorithms). The desired prediction would be the rating that Alice would give to Algorithms.

To make this concrete, Figure 1.1 shows the general framework of **induction**. We are given **training data** on which our algorithm is expected to learn. This training data is the examples that Alice observes in her machine learning course, or the historical ratings data for the recommender system. Based on this training data, our learning algorithm induces a function f that will map a new example to a corresponding prediction. For example, our function might guess that $f(\text{Alice}/\text{Machine Learning})$ might be high because our training data said that Alice liked Artificial Intelligence. We want our algorithm to be able to make lots of predictions, so we refer to the collection of examples on which we will evaluate our algorithm as the **test set**. The test set is a closely guarded secret: it is the final exam on which our learning algorithm is being tested. If our algorithm gets to peek at it ahead of time, it's going to cheat and do better than it should.

The goal of inductive machine learning is to take some training data and use it to induce a function f . This function f will be evalu-

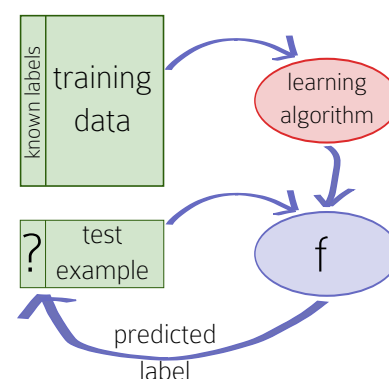


Figure 1.1: The general supervised approach to machine learning: a learning algorithm reads in training data and computes a learned function f . This function can then automatically label future test examples.

? Why is it bad if the learning algorithm gets to peek at the test data?

ated on the test data. The machine learning algorithm has succeeded if its performance on the test data is high.

1.2 Some Canonical Learning Problems

There are a large number of typical inductive learning problems. The primary difference between them is in what type of *thing* they're trying to predict. Here are some examples:

Regression: trying to predict a real value. For instance, predict the value of a stock tomorrow given its past performance. Or predict Alice's score on the machine learning final exam based on her homework scores.

Binary Classification: trying to predict a simple yes/no response. For instance, predict whether Alice will enjoy a course or not. Or predict whether a user review of the newest Apple product is positive or negative about the product.

Multiclass Classification: trying to put an example into one of a number of classes. For instance, predict whether a news story is about entertainment, sports, politics, religion, etc. Or predict whether a CS course is Systems, Theory, AI or Other.

Ranking: trying to put a set of objects in order of relevance. For instance, predicting what order to put web pages in, in response to a user query. Or predict Alice's ranked preferences over courses she hasn't taken.

The reason that it is convenient to break machine learning problems down by the type of object that they're trying to predict has to do with measuring error. Recall that our goal is to build a system that can make "good predictions." This begs the question: **what does it mean for a prediction to be "good?"** The different types of learning problems differ in how they define goodness. For instance, in regression, predicting a stock price that is off by \$0.05 is perhaps much better than being off by \$200.00. The same does not hold of multiclass classification. There, accidentally predicting "entertainment" instead of "sports" is no better or worse than predicting "politics."



For each of these types of canonical machine learning problems, come up with one or two concrete examples.

1.3 The Decision Tree Model of Learning

The **decision tree** is a classic and natural model of learning. It is closely related to the fundamental computer science notion of "divide and conquer." Although decision trees can be applied to many

learning problems, we will begin with the simplest case: binary classification.

Suppose that your goal is to predict whether some unknown user will enjoy some unknown course. You must simply answer “yes” or “no.” In order to make a guess, you’re allowed to ask binary questions about the user/course under consideration. For example:

You: Is the course under consideration in Systems?

Me: Yes

You: Has this student taken any other Systems courses?

Me: Yes

You: Has this student liked most previous Systems courses?

Me: No

You: *I predict this student will not like this course.*

The goal in learning is to figure out what questions to ask, in what order to ask them, and what answer to predict once you have asked enough questions.

The decision tree is so-called because we can write our set of questions and guesses in a tree format, such as that in Figure 1.2. In this figure, the questions are written in the internal tree nodes (rectangles) and the guesses are written in the leaves (ovals). Each non-terminal node has two children: the left child specifies what to do if the answer to the question is “no” and the right child specifies what to do if it is “yes.”

In order to learn, I will give you training data. This data consists of a set of user/course examples, paired with the correct answer for these examples (did the given user enjoy the given course?). From this, you must construct your questions. For concreteness, there is a small data set in Table 1 in the Appendix of this book. This training data consists of 20 course rating examples, with course ratings and answers to questions that you might ask about this pair. We will interpret ratings of 0, +1 and +2 as “liked” and ratings of −2 and −1 as “hated.”

In what follows, we will refer to the questions that you can ask as **features** and the responses to these questions as **feature values**. The rating is called the **label**. An example is just a set of feature values. And our training data is a set of examples, paired with labels.

There are a lot of logically possible trees that you could build, even over just this small number of features (the number is in the millions). It is computationally infeasible to consider all of these to try to choose the “best” one. Instead, we will build our decision tree *greedily*. We will begin by asking:

If I could only ask one question, what question would I ask?

You want to find a feature that is *most useful* in helping you guess whether this student will enjoy this course. A useful way to think

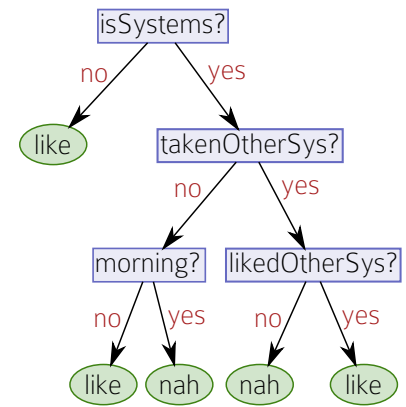


Figure 1.2: A decision tree for a course recommender system, from which the in-text “dialog” is drawn.

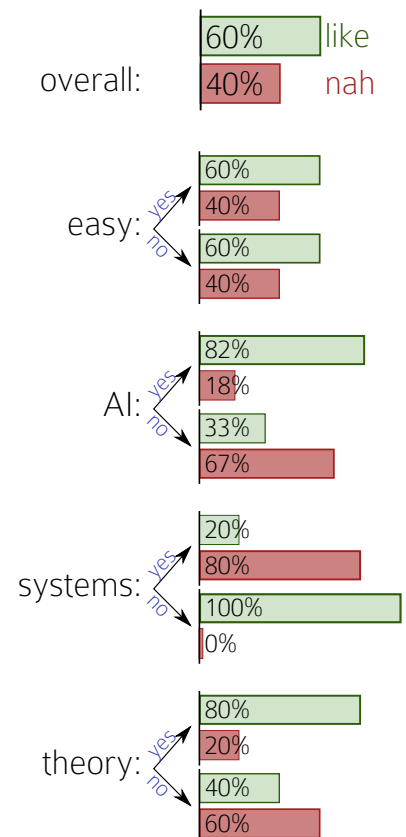


Figure 1.3: A histogram of labels for (a) the entire data set; (b-e) the examples in the data set for each value of the first four features.

about this is to look at the **histogram** of labels for each feature.¹ This is shown for the first four features in Figure 1.3. Each histogram shows the frequency of “like”/“hate” labels for each possible value of an associated feature. From this figure, you can see that asking the first feature is not useful: if the value is “no” then it’s hard to guess the label; similarly if the answer is “yes.” On the other hand, asking the second feature *is* useful: if the value is “no,” you can be pretty confident that this student will hate this course; if the answer is “yes,” you can be pretty confident that this student will like this course.

More formally, you will consider each feature in turn. You might consider the feature “Is this a System’s course?” This feature has two possible value: no and yes. Some of the training examples have an answer of “no” – let’s call that the “NO” set. Some of the training examples have an answer of “yes” – let’s call that the “YES” set. For each set (NO and YES) we will build a histogram over the labels. This is the second histogram in Figure 1.3. Now, suppose you were to ask this question on a random example and observe a value of “no.” Further suppose that you must *immediately* guess the label for this example. You will guess “like,” because that’s the more prevalent label in the NO set (actually, it’s the *only* label in the NO set). Alternatively, if you receive an answer of “yes,” you will guess “hate” because that is more prevalent in the YES set.

So, for this single feature, you know what you *would* guess if you had to. Now you can ask yourself: if I made that guess on the *training data*, how well would I have done? In particular, how many examples would I classify *correctly*? In the NO set (where you guessed “like”) you would classify all 10 of them correctly. In the YES set (where you guessed “hate”) you would classify 8 (out of 10) of them correctly. So overall you would classify 18 (out of 20) correctly. Thus, we’ll say that the *score* of the “Is this a System’s course?” question is 18/20.

You will then repeat this computation for each of the available features to us, compute the scores for each of them. When you must choose which feature consider first, you will want to choose the one with the highest score.

But this only lets you choose the *first* feature to ask about. This is the feature that goes at the *root* of the decision tree. How do we choose subsequent features? This is where the notion of divide and conquer comes in. You’ve already decided on your first feature: “Is this a Systems course?” You can now *partition* the data into two parts: the NO part and the YES part. The NO part is the subset of the data on which value for this feature is “no”; the YES half is the rest. This is the *divide* step.

¹ A colleague related the story of getting his 8-year old nephew to guess a number between 1 and 100. His nephew’s first four questions were: Is it bigger than 20? (YES) Is it even? (YES) Does it have a 7 in it? (NO) Is it 80? (NO). It took 20 more questions to get it, even though 10 should have been sufficient. At 8, the nephew hadn’t quite figured out how to divide and conquer. <http://blog.computationalcomplexity.org/2007/04/getting-8-year-old-interested-in.html>.



How many training examples would you classify correctly for each of the other three features from Figure 1.3?

Algorithm 1 `DECISIONTREETRAIN`(*data*, *remaining features*)

```

1: guess ← most frequent answer in data           // default answer for this data
2: if the labels in data are unambiguous then
3:   return LEAF(guess)                         // base case: no need to split further
4: else if remaining features is empty then
5:   return LEAF(guess)                         // base case: cannot split further
6: else                                           // we need to query more features
7:   for all  $f \in \text{remaining features}$  do
8:     NO ← the subset of data on which  $f=no$ 
9:     YES ← the subset of data on which  $f=yes$ 
10:    score[f] ← # of majority vote answers in NO
11:               + # of majority vote answers in YES
12:               // the accuracy we would get if we only queried on f
13:   end for
14:   f ← the feature with maximal score(f)
15:   NO ← the subset of data on which  $f=no$ 
16:   YES ← the subset of data on which  $f=yes$ 
17:   left ← DECISIONTREETRAIN(NO, remaining features \ {f})
18:   right ← DECISIONTREETRAIN(YES, remaining features \ {f})
19:   return NODE(f, left, right)

```

Algorithm 2 `DECISIONTREETEST`(*tree*, *test point*)

```

1: if tree is of the form LEAF(guess) then
2:   return guess
3: else if tree is of the form NODE(f, left, right) then
4:   if  $f = no$  in test point then
5:     return DECISIONTREETEST(left, test point)
6:   else
7:     return DECISIONTREETEST(right, test point)
8:   end if
9: end if

```

The *conquer* step is to recurse, and run the *same* routine (choosing the feature with the highest score) on the NO set (to get the left half of the tree) and then separately on the YES set (to get the right half of the tree).

At some point it will become useless to query on additional features. For instance, once you know that this is a Systems course, you *know* that everyone will hate it. So you can immediately predict “hate” without asking any additional questions. Similarly, at some point you might have already queried every available feature and still not whittled down to a single answer. In both cases, you will need to create a leaf node and guess the most prevalent answer in the current piece of the training data that you are looking at.

Putting this all together, we arrive at the algorithm shown in Algorithm 1.3.² This function, `DECISIONTREETRAIN` takes two argu-

² There are more nuanced algorithms for building decision trees, some of which are discussed in later chapters of this book. They primarily differ in how they compute the *score* function.

ments: our data, and the set of as-yet unused features. It has two base cases: either the data is unambiguous, or there are no remaining features. In either case, it returns a **LEAF** node containing the most likely guess at this point. Otherwise, it loops over all remaining features to find the one with the highest score. It then partitions the data into a NO/YES split based on the best feature. It constructs its left and right subtrees by recursing on itself. In each recursive call, it uses one of the partitions of the data, and removes the just-selected feature from consideration.

The corresponding *prediction* algorithm is shown in Algorithm 1.3. This function recurses down the decision tree, following the edges specified by the feature values in some *test point*. When it reaches a leaf, it returns the guess associated with that leaf.

? Is Algorithm 1.3 guaranteed to terminate?

1.4 Formalizing the Learning Problem

As you've seen, there are several issues that we must take into account when formalizing the notion of learning.

- The performance of the learning algorithm should be measured on unseen "test" data.
- The way in which we measure performance should depend on the problem we are trying to solve.
- There should be a strong relationship between the data that our algorithm sees at training time and the data it sees at test time.

In order to accomplish this, let's assume that someone gives us a **loss function**, $\ell(\cdot, \cdot)$, of two arguments. The job of ℓ is to tell us how "bad" a system's prediction is in comparison to the truth. In particular, if y is the truth and \hat{y} is the system's prediction, then $\ell(y, \hat{y})$ is a **measure of error**.

For three of the canonical tasks discussed above, we might use the following loss functions:

Regression: **squared loss** $\ell(y, \hat{y}) = (y - \hat{y})^2$
or **absolute loss** $\ell(y, \hat{y}) = |y - \hat{y}|$.

Binary Classification: **zero/one loss** $\ell(y, \hat{y}) = \begin{cases} 0 & \text{if } y = \hat{y} \\ 1 & \text{otherwise} \end{cases}$

Multiclass Classification: also zero/one loss.

This notation means that the loss is zero if the prediction is correct and is one otherwise.

Note that the loss function is something that *you* must decide on based on the goals of learning.

Now that we have defined our loss function, we need to consider where the data (training *and* test) comes from. The model that we

? Why might it be a bad idea to use zero/one loss to measure performance for a regression problem?

MATH REVIEW | EXPECTED VALUES

We write $\mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(y, f(x))]$ for the expected loss. **Expectation means “average.”** This is saying “if you drew a bunch of (x, y) pairs independently at random from \mathcal{D} , what would your *average* loss be? More formally, if \mathcal{D} is a discrete probability distribution, then this expectation can be expanded as:

$$\mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(y, f(x))] = \sum_{(x,y) \in \mathcal{D}} [\mathcal{D}(x, y) \ell(y, f(x))] \quad (1.1)$$

This is *exactly* the weighted average loss over the all (x, y) pairs in \mathcal{D} , weighted by their probability, $\mathcal{D}(x, y)$. If \mathcal{D} is a *finite discrete distribution*, for instance defined by a finite data set $\{(x_1, y_1), \dots, (x_N, y_N)\}$ that puts equal weight on each example (probability $1/N$), then we get:

$$\mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(y, f(x))] = \sum_{(x,y) \in \mathcal{D}} [\mathcal{D}(x, y) \ell(y, f(x))] \quad \text{definition of expectation} \quad (1.2)$$

$$= \sum_{n=1}^N [\mathcal{D}(x_n, y_n) \ell(y_n, f(x_n))] \quad \text{\textit{D} is discrete and finite} \quad (1.3)$$

$$= \sum_{n=1}^N \left[\frac{1}{N} \ell(y_n, f(x_n)) \right] \quad \text{definition of } \mathcal{D} \quad (1.4)$$

$$= \frac{1}{N} \sum_{n=1}^N [\ell(y_n, f(x_n))] \quad \text{rearranging terms} \quad (1.5)$$

Which is exactly the average loss on that dataset.

The most important thing to remember is that there are two equivalent ways to think about expectations: (1) The expectation of some function g is the *weighted average value* of g , where the weights are given by the underlying probability distribution. (2) The expectation of some function g is your *best guess of the value* of g if you were to draw a single item from the underlying probability distribution.

Figure 1.4:

will use is **the probabilistic model of learning**. Namely, there is a **probability distribution \mathcal{D} over input/output pairs**. This is often called the **data generating distribution**. If we write x for the input (the user/course pair) and y for the output (the rating), then \mathcal{D} is a distribution over (x, y) pairs.

A useful way to think about \mathcal{D} is that it gives *high probability* to *reasonable (x, y) pairs*, and *low probability* to *unreasonable (x, y) pairs*. A (x, y) pair can be unreasonable in two ways. First, x might be an unusual input. For example, a x related to an “Intro to Java” course might be highly probable; a x related to a “Geometric and Solid Modeling” course might be less probable. Second, y might be an unusual rating for the paired x . For instance, if Alice were to take AI 100 times (without remembering that she took it before!), she would give the course a +2 almost every time. Perhaps some

semesters she might give a slightly lower score, but it would be unlikely to see $x = \text{Alice/AI}$ paired with $y = -2$.

It is important to remember that we are not making *any* assumptions about what the distribution \mathcal{D} looks like. (For instance, we're not assuming it looks like a Gaussian or some other, common distribution.) We are also not assuming that we know what \mathcal{D} is. In fact, if you know *a priori* what your data generating distribution is, your learning problem becomes significantly easier. Perhaps the hardest thing about machine learning is that we *don't* know what \mathcal{D} is: all we get is a random sample from it. This random sample is our training data.

Our learning problem, then, is defined by two quantities:

1. The loss function ℓ , which captures our notion of what is *important* to learn.
2. The data generating distribution \mathcal{D} , which defines what sort of data we expect to see.

We are given access to **training data**, which is a random sample of input/output pairs drawn from \mathcal{D} . Based on this training data, we need to **induce** a function f that maps new inputs \hat{x} to corresponding prediction \hat{y} . The key property that f should obey is that it should do well (as measured by ℓ) on future examples that are *also* drawn from \mathcal{D} . Formally, it's **expected loss** ϵ over \mathcal{D} with respect to ℓ should be as small as possible:

$$\epsilon \triangleq \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(y, f(x))] = \sum_{(x,y)} \mathcal{D}(x,y) \ell(y, f(x)) \quad (1.6)$$

The difficulty in minimizing our **expected loss** from Eq (1.6) is that we *don't know what \mathcal{D} is!* All we have access to is some training data sampled from it! Suppose that we denote our training data set by D . The training data consists of N -many input/output pairs, $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$. Given a learned function f , we can compute our **training error**, $\hat{\epsilon}$:

$$\hat{\epsilon} \triangleq \frac{1}{N} \sum_{n=1}^N \ell(y_n, f(x_n)) \quad (1.7)$$

That is, our training error is simply our *average error* over the training data.

Of course, we can drive $\hat{\epsilon}$ to zero by simply memorizing our training data. But as Alice might find in memorizing past exams, this might not generalize well to a new exam!

This is the fundamental difficulty in machine learning: the thing we have access to is our training error, $\hat{\epsilon}$. But the thing we care about

Consider the following prediction task. Given a paragraph written about a course, we have to predict whether the paragraph is a *positive* or *negative* review of the course. (This is the sentiment analysis problem.) What is a reasonable loss function? How would you define the data generating distribution?

?

Verify by calculation that we can write our training error as $\mathbb{E}_{(x,y) \sim D} [\ell(y, f(x))]$, by thinking of D as a distribution that places probability $1/N$ to each example in D and probability 0 on everything else.

?

minimizing is our expected error ϵ . In order to get the expected error down, our learned function needs to **generalize** beyond the training data to some future data that it might not have seen yet!

So, putting it all together, we get a formal definition of induction machine learning: **Given (i) a loss function ℓ and (ii) a sample D from some unknown distribution \mathcal{D} , you must compute a function f that has low expected error ϵ over D with respect to ℓ .**

A very important comment is that we should *never* expect a machine learning algorithm to generalize beyond the data distribution it has seen at training time. In a famous—if possibly apocryphal—example from the 1970s, the US Government wanted to train a classifier to distinguish between US tanks and Russian tanks. They collected a training and test set, and managed to build a classifier with nearly 100% accuracy on that data. But when this classifier was run in the “real world”, it failed miserably. It had not, in fact, learned to distinguish between US tanks and Russian tanks, but rather just between clear photos and blurry photos. In this case, there was a *bias* in the training data (due to how the training data was collected) that caused the learning algorithm to learn something other than what we were hoping for. We will return to this issue in Chapter 8; for now, simply remember that the distribution D for training data *must match* the distribution D for the test data.

1.5 Chapter Summary and Outlook

At this point, you should be able to use decision trees to do machine learning. Someone will give you data. You’ll split it into training, development and test portions. Using the training and development data, you’ll find a good value for maximum depth that trades off between underfitting and overfitting. You’ll then run the resulting decision tree model on the test data to get an estimate of how well you are likely to do in the future.

You might think: why should I read the rest of this book? Aside from the fact that machine learning is just an awesome fun field to learn about, there’s a lot left to cover. In the next two chapters, you’ll learn about two models that have very different inductive biases than decision trees. You’ll also get to see a very useful way of thinking about learning: the geometric view of data. This will guide much of what follows. After that, you’ll learn how to solve problems more complicated than simple binary classification. (Machine learning people like binary classification a lot because it’s one of the simplest non-trivial problems that we can work on.) After that, things will diverge: you’ll learn about ways to think about learning as a formal optimization problem, ways to speed up learning, ways to learn

without labeled data (or with very little labeled data) and all sorts of other fun topics.

But throughout, we will focus on the view of machine learning that you've seen here. You select a model (and its associated inductive biases). You use data to find parameters of that model that work well on the training data. You use development data to avoid underfitting and overfitting. And you use test data (which you'll never look at or touch, right?) to estimate future model performance. Then you conquer the world.

1.6 Further Reading

In our discussion of decision trees, we used *misclassification rate* for selecting features. While simple and intuitive, misclassification rate has problems. There has been a significant amount of work that considers more advanced splitting criteria; the most popular is ID3, based on the mutual information quantity from information theory. We have also only considered a very simple mechanism for controlling inductive bias: limiting the depth of the decision tree. Again, there are more advanced “tree pruning” techniques that typically operate by growing deep trees and then pruning back some of the branches. These approaches have the advantage that different branches can have different depths, accounting for the fact that the amount of data that gets passed down each branch might differ dramatically³.

³ Quinlan 1986

2 | LIMITS OF LEARNING

Our lives sometimes depend on computers performing as predicted.

– Philip Emeagwali

Learning Objectives:

- Define “inductive bias” and recognize the role of inductive bias in learning.
- Illustrate how regularization trades off between underfitting and overfitting.
- Evaluate whether a use of test data is “cheating” or not.

MACHINE LEARNING IS A VERY GENERAL and useful framework, but it is not “magic” and will not always work. In order to better understand when it will and when it will not work, it is useful to formalize the learning problem more. This will also help us develop debugging strategies for learning algorithms.

2.1 Data Generating Distributions

Our underlying assumption for the majority of this book is that learning problems are characterized by some unknown probability distribution \mathcal{D} over input/output pairs $(x, y) \in \mathcal{X} \times \mathcal{Y}$. Suppose that someone *told* you what \mathcal{D} was. In particular, they gave you a Python function `COMPUTED` that took two inputs, x and y , and returned the probability of that x, y pair under \mathcal{D} . If you had access to such a function, classification becomes simple. We can define the **Bayes optimal classifier** as the classifier that, for any test input \hat{x} , simply returns the \hat{y} that maximizes `COMPUTED`(\hat{x}, \hat{y}), or, more formally:

$$f^{(\text{BO})}(\hat{x}) = \arg \max_{\hat{y} \in \mathcal{Y}} \mathcal{D}(\hat{x}, \hat{y}) \quad (2.1)$$

This classifier is optimal in one specific sense: of *all possible* classifiers, it achieves the smallest zero/one error.

Theorem 1 (Bayes Optimal Classifier). *The Bayes Optimal Classifier $f^{(\text{BO})}$ achieves minimal zero/one error of any deterministic classifier.*

This theorem assumes that you are comparing against *deterministic* classifiers. You can actually prove a stronger result that $f^{(\text{BO})}$ is optimal for randomized classifiers as well, but the proof is a bit messier. However, the intuition is the same: for a given x , $f^{(\text{BO})}$ chooses the label with highest probability, thus *minimizing* the probability that it makes an error.

Proof of Theorem 1. Consider some other classifier g that claims to be better than $f^{(\text{BO})}$. Then, there must be some x on which $g(x) \neq$

Dependencies: None.

$f^{(BO)}(x)$. Fix such an x . Now, the probability that $f^{(BO)}$ makes an error on this particular x is $1 - \mathcal{D}(x, f^{(BO)}(x))$ and the probability that g makes an error on this x is $1 - \mathcal{D}(x, g(x))$. But $f^{(BO)}$ was chosen in such a way to *maximize* $\mathcal{D}(x, f^{(BO)}(x))$, so this *must* be greater than $\mathcal{D}(x, g(x))$. Thus, the probability that $f^{(BO)}$ errs on this particular x is smaller than the probability that g errs on it. This applies to any x for which $f^{(BO)}(x) \neq g(x)$ and therefore $f^{(BO)}$ achieves smaller zero/one error than any g . \square

The **Bayes error rate** (or **Bayes optimal error rate**) is the error rate of the Bayes optimal classifier. It is the best error rate you can ever hope to achieve on this classification problem (under zero/one loss). The take-home message is that if someone gave you access to the data distribution, forming an *optimal* classifier would be trivial. Unfortunately, no one gave you this distribution, so we need to figure out ways of learning the mapping from x to y given only access to a training set *sampled from* \mathcal{D} , rather than \mathcal{D} itself.

2.2 Inductive Bias: What We Know Before the Data Arrives

In Figure 2.1 you'll find training data for a binary classification problem. The two labels are "A" and "B" and you can see four examples for each label. Below, in Figure 2.2, you will see some test data. These images are left unlabeled. Go through quickly and, based on the training data, label these images. (Really do it before you read further! I'll wait!)

Most likely you produced one of two labelings: either ABBA or AABB. Which of these solutions is right? The answer is that you cannot tell based on the training data. If you give this same example to 100 people, 60 – 70 of them come up with the ABBA prediction and 30 – 40 come up with the AABB prediction. Why? Presumably because the first group believes that the relevant distinction is between "bird" and "non-bird" while the second group believes that the relevant distinction is between "fly" and "no-fly."

This preference for one distinction (bird/non-bird) over another (fly/no-fly) is a bias that different human learners have. In the context of machine learning, it is called **inductive bias**: in the absence of data that narrow down the relevant concept, what type of solutions are we more likely to prefer? Two thirds of people seem to have an inductive bias in favor of bird/non-bird, and one third seem to have an inductive bias in favor of fly/no-fly.

Throughout this book you will learn about several approaches to machine learning. The decision tree model is the first such approach. These approaches differ primarily in the sort of inductive bias that

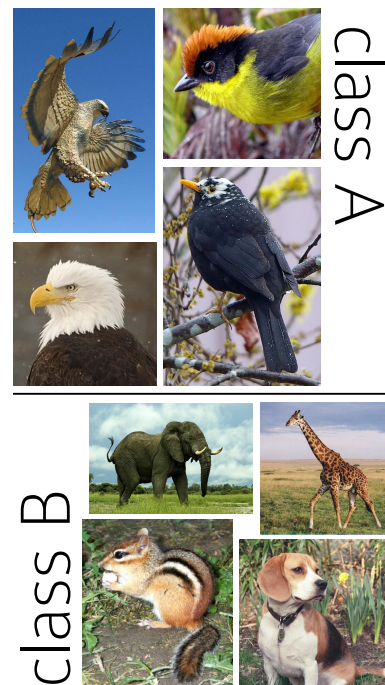


Figure 2.1: Training data for a binary classification problem.



Figure 2.2: Test data for the same classification problem.

? It is also possible that the correct classification on the test data is ABAB. This corresponds to the bias "is the background in focus." Somehow no one seems to come up with this classification rule.

they exhibit.

Consider a variant of the decision tree learning algorithm. In this variant, we will not allow the trees to grow beyond some pre-defined maximum depth, d . That is, once we have queried on d -many features, we cannot query on any more and must just make the best guess we can at that point. This variant is called a **shallow decision tree**.

The key question is: What is the inductive bias of shallow decision trees? Roughly, their bias is that decisions can be made by only looking at a small number of features. For instance, a shallow decision tree would be very good at learning a function like “students only like AI courses.” It would be very bad at learning a function like “if this student has liked an odd number of their past courses, they will like the next one; otherwise they will not.” This latter is the **parity** function, which requires you to inspect every feature to make a prediction. The inductive bias of a decision tree is that the sorts of things we want to learn to predict are more like the first example and less like the second example.

2.3 *Not Everything is Learnable*

Although machine learning works well—perhaps astonishingly well—in many cases, it is important to keep in mind that it is not magical. There are many reasons why a machine learning algorithm might fail on some learning task.

There could be **noise** in the training data. Noise can occur both at the feature level and at the label level. Some features might correspond to measurements taken by sensors. For instance, a robot might use a laser range finder to compute its distance to a wall. However, this sensor might fail and return an incorrect value. In a sentiment classification problem, someone might have a typo in their review of a course. These would lead to noise at the feature level. There might also be noise at the label level. A student might write a scathingly negative review of a course, but then accidentally click the wrong button for the course rating.

The features available for learning might simply be insufficient. For example, in a medical context, you might wish to diagnose whether a patient has cancer or not. You may be able to collect a large amount of data about this patient, such as gene expressions, X-rays, family histories, etc. But, even knowing all of this information exactly, it might still be impossible to judge for sure whether this patient has cancer or not. As a more contrived example, you might try to classify course reviews as positive or negative. But you may have erred when downloading the data and only gotten the first five char-

acters of each review. If you had the rest of the features you might be able to do well. But with this limited feature set, there's not much you can do.

Some examples may not have a single correct answer. You might be building a system for "safe web search," which removes offensive web pages from search results. To build this system, you would collect a set of web pages and ask people to classify them as "offensive" or not. However, what one person considers offensive might be completely reasonable for another person. It is common to consider this as a form of label noise. Nevertheless, since you, as the designer of the learning system, have some control over this problem, it is sometimes helpful to isolate it as a source of difficulty.

Finally, learning might fail because the inductive bias of the learning algorithm is too far away from the concept that is being learned. In the bird/non-bird data, you might think that if you had gotten a few more training examples, you might have been able to tell whether this was intended to be a bird/non-bird classification or a fly/no-fly classification. However, no one I've talked to has ever come up with the "background is in focus" classification. Even with many more training points, this is such an unusual distinction that it may be hard for anyone to figure out it. In this case, the inductive bias of the learner is simply too misaligned with the target classification to learn.

Note that the inductive bias source of error is fundamentally different than the other three sources of error. In the inductive bias case, it is the *particular* learning algorithm that you are using that cannot cope with the data. Maybe if you switched to a different learning algorithm, you would be able to learn well. For instance, Neptunians might have evolved to care greatly about whether backgrounds are in focus, and for them this would be an easy classification to learn. For the other three sources of error, it is not an issue to do with the particular learning algorithm. The error is a fundamental part of the learning problem.

2.4 Underfitting and Overfitting

As with many problems, it is useful to think about the *extreme cases* of learning algorithms. In particular, the extreme cases of decision trees. In one extreme, the tree is "empty" and we do not ask any questions at all. We simply immediately make a prediction. In the other extreme, the tree is "full." That is, every possible question is asked along every branch. In the full tree, there may be leaves with no associated training data. For these we must simply choose arbitrarily whether to say "yes" or "no."

Consider the course recommendation data from Table 1. Suppose we were to build an “empty” decision tree on this data. Such a decision tree will make the same prediction regardless of its input, because it is not allowed to ask any questions about its input. Since there are more “likes” than “hates” in the training data (12 versus 8), our empty decision tree will simply always predict “likes.” The training error, $\hat{\epsilon}$, is $8/20 = 40\%$.

On the other hand, we could build a “full” decision tree. Since each row in this data is unique, we can guarantee that any leaf in a full decision tree will have either 0 or 1 examples assigned to it (20 of the leaves will have one example; the rest will have none). For the leaves corresponding to training points, the full decision tree will always make the correct prediction. Given this, the training error, $\hat{\epsilon}$, is $0/20 = 0\%$.

Of course our goal is *not* to build a model that gets 0% error on the training data. This would be easy! Our goal is a model that will do well on *future, unseen* data. How well might we expect these two models to do on future data? The “empty” tree is likely to do not much better and not much worse on future data. We might expect that it would continue to get around 40% error.

Life is more complicated for the “full” decision tree. Certainly if it is given a test example that is identical to one of the training examples, it will do the right thing (assuming no noise). But for everything else, it will only get about 50% error. This means that even if every other test point happens to be identical to one of the training points, it would only get about 25% error. In practice, this is probably optimistic, and maybe only one in every 10 examples would match a training example, yielding a 35% error.

So, in one case (empty tree) we’ve achieved about 40% error and in the other case (full tree) we’ve achieved 35% error. This is not very promising! One would hope to do better! In fact, you might notice that if you simply queried on a *single* feature for this data, you would be able to get very low training error, but wouldn’t be forced to “guess” randomly.

This example illustrates the key concepts of **underfitting** and **overfitting**. Underfitting is when you had the opportunity to learn something but didn’t. A student who hasn’t studied much for an upcoming exam will be underfit to the exam, and consequently will not do well. This is also what the empty tree does. Overfitting is when you pay too much attention to idiosyncracies of the training data, and aren’t able to generalize well. Often this means that your model is fitting noise, rather than whatever it is supposed to fit. A student who memorizes answers to past exam questions without understanding them has overfit the training data. Like the full tree, this student

Convince yourself (either by proof or by simulation) that even in the case of imbalanced data – for instance data that is on average 80% positive and 20% negative – a predictor that guesses randomly (50/50 positive/negative) will get about 50% error.

?

?

Which feature is it, and what is its training error?

MATH REVIEW | LAW OF LARGE NUMBERS

Consider some random event, like spins of a roulette wheel, cars driving through an intersection, the outcome of an election, or pasta being appropriately *al dente*. We often want to make a conclusion about the entire population (the pot of pasta) based on a much smaller sample (biting a couple pieces of pasta). The law of large numbers tells us that under mild conditions this is an okay thing to do.

Formally, suppose that v_1, v_2, \dots, v_N are random variables (e.g., v_n measures if the n th spaghetti is *al dente*). Assume that these random variables are **independent** (i.e., v_2 and v_3 are uncorrelated—they weren't both taken from the same place in the pot) and **identically distributed** (they were all drawn from the same population—pot—that we wish to measure). We can compute the sample average $\bar{v} = \frac{1}{N} \sum_{n=1}^N v_n$ and under the **strong law of large numbers**, you can prove that $\bar{v} \rightarrow \mathbb{E}[v]$ as $N \rightarrow \infty$. Namely, the empirical sample average approaches the population average as the number of samples goes to infinity.

(Technical note: the notion of convergence here is almost sure convergence. In particular, the formal result is that $\Pr\left(\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N v_n = \mathbb{E}[v]\right) = 1$. Or, in words, with probability one the sample average reaches the population average.)

Figure 2.3:

also will not do well on the exam. A model that is neither overfit nor underfit is the one that is expected to do best in the future.

2.5 Separation of Training and Test Data

Suppose that, after graduating, you get a job working for a company that provides personalized recommendations for pottery. You go in and implement new algorithms based on what you learned in your machine learning class (you have learned the power of generalization!). All you need to do now is convince your boss that you have done a good job and deserve a raise!

How can you convince your boss that your fancy learning algorithms are really working?

Based on what we've talked about already with underfitting and overfitting, it is not enough to just tell your boss what your training error is. Noise notwithstanding, it is easy to get a training error of zero using a simple database query (or `grep`, if you prefer). Your boss will not fall for that.

The easiest approach is to *set aside* some of your available data as “test data” and use this to evaluate the performance of your learning algorithm. For instance, the pottery recommendation service that you work for might have collected 1000 examples of pottery ratings. You will select 800 of these as **training data** and set aside the final 200

as **test data**. You will run your learning algorithms *only* on the 800 training points. Only once you're done will you apply your learned model to the 200 test points, and report your **test error** on those 200 points to your boss.

The hope in this process is that however well you do on the 200 test points will be indicative of how well you are likely to do in the future. This is analogous to estimating support for a presidential candidate by asking a small (random!) sample of people for their opinions. Statistics (specifically, concentration bounds of which the “Central limit theorem” is a famous example) tells us that if the sample is large enough, it will be a good representative. The 80/20 split is not magic: it's simply fairly well established. Occasionally people use a 90/10 split instead, especially if they have a *lot* of data.

The cardinal rule of machine learning is: never touch your test data. Ever. If that's not clear enough:

Never ever touch your test data!

If there is only one thing you learn from this book, let it be that. Do not look at your test data. Even once. Even a tiny peek. Once you do that, it is not test data any more. Yes, perhaps your algorithm hasn't seen it. But you have. And you are likely a better learner than your learning algorithm. Consciously or otherwise, you might make decisions based on whatever you might have seen. Once you look at the test data, your model's performance on it is no longer indicative of it's performance on future unseen data. This is simply because future data is unseen, but your “test” data no longer is.



If you have more data at your disposal, why might a 90/10 split be preferable to an 80/20 split?

2.6 Models, Parameters and Hyperparameters

The general approach to machine learning, which captures many existing learning algorithms, is the **modeling** approach. The idea is that we come up with some formal model of our data. For instance, we might model the classification decision of a student/course pair as a decision tree. The choice of using a *tree* to represent this model is *our choice*. We also could have used an arithmetic circuit or a polynomial or some other function. The model tells us what sort of things we can learn, and also tells us what our inductive bias is.

For most models, there will be associated parameters. These are the things that we use the data to decide on. Parameters in a decision tree include: the specific questions we asked, the order in which we asked them, and the classification decisions at the leaves. The job of our decision tree learning algorithm `DECISIONTREETRAIN` is to take data and figure out a good set of parameters.

Many learning algorithms will have additional knobs that you can adjust. In most cases, these knobs amount to tuning the inductive bias of the algorithm. In the case of the decision tree, an obvious knob that one can tune is the **maximum depth** of the decision tree. That is, we could modify the `DECISIONTREE-TRAIN` function so that it *stops* recursing once it reaches some pre-defined maximum depth. By playing with this depth knob, we can adjust between underfitting (the empty tree, $\text{depth} = 0$) and overfitting (the full tree, $\text{depth} = \infty$).

Such a knob is called a **hyperparameter**. It is so called because it is a parameter that controls other parameters of the model. The exact definition of hyperparameter is hard to pin down: it's one of those things that are easier to identify than define. However, one of the key identifiers for hyperparameters (and the main reason that they cause consternation) is that they cannot be naively adjusted using the training data.

In `DECISIONTREE-TRAIN`, as in most machine learning, the learning algorithm is essentially trying to adjust the parameters of the model so as to minimize training error. This suggests an idea for choosing hyperparameters: choose them so that they minimize training error.

What is wrong with this suggestion? Suppose that you were to treat “maximum depth” as a hyperparameter and tried to tune it on your training data. To do this, maybe you simply build a collection of decision trees, $\text{tree}_0, \text{tree}_1, \text{tree}_2, \dots, \text{tree}_{100}$, where tree_d is a tree of maximum depth d . We then computed the training error of each of these trees and chose the “ideal” maximum depth as that which minimizes training error? Which one would it pick?

The answer is that it would pick $d = 100$. Or, in general, it would pick d as large as possible. Why? Because choosing a bigger d will *never hurt* on the training data. By making d larger, you are simply encouraging overfitting. But by evaluating on the training data, overfitting actually looks like a good idea!

An alternative idea would be to tune the maximum depth on test data. This is promising because test data performance is what we really want to optimize, so tuning this knob on the test data seems like a good idea. That is, it won't accidentally reward overfitting. Of course, it breaks our cardinal rule about test data: that you should never touch your test data. So that idea is immediately off the table.

However, our “test data” wasn't magic. We simply took our 1000 examples, called 800 of them “training” data and called the other 200 “test” data. So instead, let's do the following. Let's take our original 1000 data points, and select 700 of them as training data. From the remainder, take 100 as **development data**¹ and the remaining 200 as test data. The job of the development data is to allow us to tune

Go back to the `DECISIONTREE-TRAIN` algorithm and modify it so that it takes a maximum depth parameter. This should require adding two lines of code and modifying three others.

?

¹ Some people call this “**validation data**” or “**held-out data**.”

hyperparameters. The general approach is as follows:

1. Split your data into 70% training data, 10% development data and 20% test data.
2. For each possible setting of your hyperparameters:
 - (a) Train a model using that setting of hyperparameters on the training data.
 - (b) Compute this model's error rate on the development data.
3. From the above collection of models, choose the one that achieved the lowest error rate on development data.
4. Evaluate that model on the test data to estimate future test performance.

2.7 Real World Applications of Machine Learning

Figure 2.4 shows a typical sequence of decisions that must be made to deploy a machine learning approach in the real world. In the left column, you can see the generic decision being made. In the right column, an example of this decision for the particular case of advertising placement on a search engine we've built.

In this sequence, (1) we have some real world goal like increasing revenue for our search engine, and decide to try to increase revenue by (2) displaying better ads. We convert this task into a machine learning problem by (3) deciding to train a classifier to predict whether a user will click on an ad or not. In order to apply machine learning, we must collect some training data; in this case, (4) we collect data by logging user interactions with the current system. We must choose what to log; (5) we choose to log the ad being displayed, the query the user entered into our search engine, and binary value showing if they clicked or not.

In order to make these logs consumable by a machine learning algorithm, (6) we convert the data into input/output pairs: in this case, pairs of words from a bag-of-words representing the query and a bag-of-words representing the ad as input, and the click as a \pm label. We then (7) select a model family (e.g., depth 20 decision trees), and thereby an inductive bias, for instance depth ≤ 20 decision trees.

We're now ready to (8) select a specific subset of data to use as training data: in this case, data from April 2016. We split this into training and development and (9) learn a final decision tree, tuning the maximum depth on the development data. We can then use this decision tree to (10) make predictions on some held-out test data, in

? In step 3, you could either choose the model (trained on the 70% training data) that did the best on the development data. Or you could choose the hyperparameter settings that did best and *retrain* the model on the 80% union of training and development data. Is either of these options obviously better or worse?

1	real world goal	increase revenue
2	real world mechanism	better ad display
3	learning problem	classify click-through
4	data collection	interaction w/ current system
5	collected data	query, ad, click
6	data representation	bow ² , \pm click
7	select model family	decision trees, depth 20
8	select training data	subset from april'16
9	train model & hyperparams	final decision tree
10	predict on test data	subset from may'16
11	evaluate error	zero/one loss for \pm click
12	deploy!	(hope we achieve our goal)

Figure 2.4: A typical design process for a machine learning application.

this case from the following month. We can (11) measure the overall quality of our predictor as zero/one loss (classification error) on this test data and finally (12) deploy our system.

The important thing about this sequence of steps is that *in any one, things can go wrong*. That is, between any two rows of this table, we are *necessarily* accumulating some additional error against our original real world goal of increasing revenue. For example, in step 5, we decided on a representation that left out many possible variables we could have logged, like time of day or season of year. By leaving out those variables, we set an explicit upper bound on how well our learned system can do.

It is often an effective strategy to run an **oracle experiment**. In an oracle experiment, we assume that everything below some line can be solved perfectly, and measure how much impact that will have on a higher line. As an extreme example, before embarking on a machine learning approach to the ad display problem, we should measure something like: if our classifier were *perfect*, how much more money would we make? If the number is not very high, perhaps there is some better for our time.

Finally, although this sequence is denoted linearly, the entire process is highly interactive in practice. A large part of “debugging” machine learning (covered more extensively in Chapter 5) involves trying to figure out where in this sequence the biggest losses are and fixing that step. In general, it is often useful to *build the stupidest thing that could possibly work*, then look at how well it’s doing, and decide if and where to fix it.

2.8 Further Reading

TODO further reading

3 | GEOMETRY AND NEAREST NEIGHBORS

Our brains have evolved to get us out of the rain, find where the berries are, and keep us from getting killed. Our brains did not evolve to help us grasp really large numbers or to look at things in a hundred thousand dimensions. — Ronald Graham

Learning Objectives:

- Describe a data set as points in a high dimensional space.
- Explain the curse of dimensionality.
- Compute distances between points in high dimensional space.
- Implement a K -nearest neighbor model of learning.
- Draw decision boundaries.
- Implement the K -means algorithm for clustering.

YOU CAN THINK OF PREDICTION TASKS as mapping inputs (course reviews) to outputs (course ratings). As you learned in the previous chapter, decomposing an input into a collection of features (e.g., words that occur in the review) forms a useful abstraction for learning. Therefore, inputs are nothing more than lists of feature values. This suggests a **geometric view** of data, where we have one dimension for every feature. In this view, examples are points in a high-dimensional space.

Once we think of a data set as a collection of points in high dimensional space, we can start performing geometric operations on this data. For instance, suppose you need to predict whether Alice will like Algorithms. Perhaps we can try to find another student who is most “similar” to Alice, in terms of favorite courses. Say this student is Jeremy. If Jeremy liked Algorithms, then we might guess that Alice will as well. This is an example of a **nearest neighbor** model of learning. By inspecting this model, we’ll see a completely different set of answers to the key learning questions we discovered in Chapter 1.

Dependencies: Chapter 1

3.1 From Data to Feature Vectors

An example is just a collection of feature values about that example, for instance the data in Table 1 from the Appendix. To a person, these features have meaning. One feature might count how many times the reviewer wrote “excellent” in a course review. Another might count the number of exclamation points. A third might tell us if any text is underlined in the review.

To a machine, the **features** themselves have no meaning. Only the **feature values**, and how they vary across examples, mean something to the machine. From this perspective, you can think about an example as being represented by a **feature vector** consisting of one “dimension” for each feature, where each dimension is simply some real value.

Consider a review that said “excellent” three times, had one excla-

mation point and no underlined text. This could be represented by the feature vector $\langle 3, 1, 0 \rangle$. An almost identical review that happened to have underlined text would have the feature vector $\langle 3, 1, 1 \rangle$.

Note, here, that we have imposed the convention that for **binary features** (yes/no features), the corresponding feature values are 0 and 1, respectively. This was an arbitrary choice. We could have made them 0.92 and -16.1 if we wanted. But 0/1 is convenient and helps us interpret the feature values. When we discuss practical issues in Chapter 5, you will see other reasons why 0/1 is a good choice.

Figure 3.1 shows the data from Table 1 in three views. These three views are constructed by considering two features at a time in different pairs. In all cases, the plusses denote positive examples and the minuses denote negative examples. In some cases, the points fall on top of each other, which is why you cannot see 20 unique points in all figures.

The mapping from feature values to vectors is straightforward in the case of real valued features (trivial) and binary features (mapped to zero or one). It is less clear what to do with **categorical features**. For example, if our goal is to identify whether an object in an image is a tomato, blueberry, cucumber or cockroach, we might want to know its color: is it RED, BLUE, GREEN or BLACK?

One option would be to map RED to a value of 0, BLUE to a value of 1, GREEN to a value of 2 and BLACK to a value of 3. The problem with this mapping is that it turns an unordered set (the set of colors) into an ordered set (the set $\{0, 1, 2, 3\}$). In itself, this is not necessarily a bad thing. But when we go to *use* these features, we will measure examples based on their distances to each other. By doing this mapping, we are essentially saying that RED and BLUE are more similar (distance of 1) than RED and BLACK (distance of 3). This is probably not what we want to say!

A solution is to turn a categorical feature that can take four different values (say: RED, BLUE, GREEN and BLACK) into four binary features (say: IsItRed?, IsItBlue?, IsItGreen? and IsItBlack?). In general, if we start from a categorical feature that takes V values, we can map it to V -many binary indicator features.

With that, you should be able to take a data set and map each example to a feature vector through the following mapping:

- Real-valued features get copied directly.
- Binary features become 0 (for false) or 1 (for true).
- Categorical features with V possible values get mapped to V -many binary indicator features.

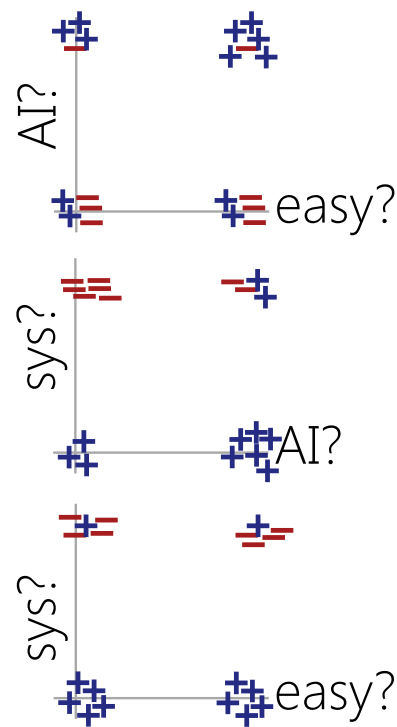


Figure 3.1: A figure showing projections of data in two dimension in three ways – see text. Top: horizontal axis corresponds to the first feature (easy) and the vertical axis corresponds to the second feature (AI?); Middle: horizontal is second feature and vertical is third (systems?); Bottom: horizontal is first and vertical is third. Truly, the data points would like exactly on $(0, 0)$ or $(1, 0)$, etc., but they have been perturbed slightly to show duplicates.

? Match the example ids from Table 1 with the points in Figure 3.1.

? The computer scientist in you might be saying: actually we could map it to $\log_2 V$ -many binary features! Is this a good idea or not?

After this mapping, you can think of a single example as a **vector** in a high-dimensional **feature space**. If you have D -many features (after expanding categorical features), then this **feature vector** will have D -many components. We will denote feature vectors as $\mathbf{x} = \langle x_1, x_2, \dots, x_D \rangle$, so that x_d denotes the value of the d th feature of \mathbf{x} . Since these are vectors with real-valued components in D -dimensions, we say that they belong to the space \mathbb{R}^D .

For $D = 2$, our feature vectors are just points in the plane, like in Figure 3.1. For $D = 3$ this is three dimensional space. For $D > 3$ it becomes quite hard to visualize. (You should resist the temptation to think of $D = 4$ as “time” – this will just make things confusing.) Unfortunately, for the sorts of problems you will encounter in machine learning, $D \approx 20$ is considered “low dimensional,” $D \approx 1000$ is “medium dimensional” and $D \approx 100000$ is “high dimensional.”

3.2 *K-Nearest Neighbors*

The biggest advantage to thinking of examples as vectors in a high dimensional space is that it allows us to apply geometric concepts to machine learning. For instance, one of the most basic things that one can do in a vector space is compute **distances**. In two-dimensional space, the distance between $\langle 2, 3 \rangle$ and $\langle 6, 1 \rangle$ is given by $\sqrt{(2-6)^2 + (3-1)^2} = \sqrt{18} \approx 4.24$. In general, in D -dimensional space, the **Euclidean distance** between vectors \mathbf{a} and \mathbf{b} is given by Eq (3.1) (see Figure 3.2 for geometric intuition in three dimensions):

$$d(\mathbf{a}, \mathbf{b}) = \left[\sum_{d=1}^D (a_d - b_d)^2 \right]^{\frac{1}{2}} \quad (3.1)$$

Now that you have access to distances between examples, you can start thinking about what it means to learn again. Consider Figure 3.3. We have a collection of training data consisting of positive examples and negative examples. There is a test point marked by a question mark. Your job is to guess the correct label for that point.

Most likely, you decided that the label of this test point is positive. One reason why you might have thought that is that you believe that the label for an example should be similar to the label of nearby points. This is an example of a new form of **inductive bias**.

The **nearest neighbor** classifier is build upon this insight. In comparison to decision trees, the algorithm is ridiculously simple. At training time, we simply store the entire training set. At test time, we get a test example $\hat{\mathbf{x}}$. To predict its label, we find the training example \mathbf{x} that is most similar to $\hat{\mathbf{x}}$. In particular, we find the training

Can you think of problems (perhaps ones already mentioned in this book!) that are low dimensional? That are medium dimensional? That are high dimensional?

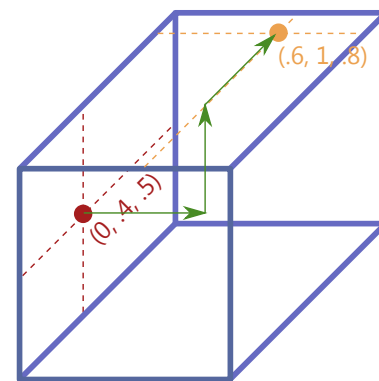


Figure 3.2: A figure showing Euclidean distance in three dimensions. The length of the green segments are 0.6, 0.6 and 0.3 respectively, in the x -, y -, and z -axes. The total distance between the red dot and the orange dot is therefore $\sqrt{0.6^2 + 0.6^2 + 0.3^2} = 0.9$.

Verify that d from Eq (3.1) gives the same result (4.24) for the previous computation.



example x that *minimizes* $d(x, \hat{x})$. Since x is a training example, it has a corresponding label, y . We predict that the label of \hat{x} is also y .

Despite its simplicity, this nearest neighbor classifier is incredibly effective. (Some might say *frustratingly* effective.) However, it is particularly prone to overfitting label noise. Consider the data in Figure 3.4. You would probably want to label the test point positive. Unfortunately, its nearest neighbor happens to be negative. Since the nearest neighbor algorithm only looks at the *single* nearest neighbor, it cannot consider the “preponderance of evidence” that this point should probably actually be a positive example. It will make an unnecessary error.

A solution to this problem is to consider more than just the single nearest neighbor when making a classification decision. We can consider the **K-nearest neighbors** and let them **vote** on the correct class for this test point. If you consider the 3-nearest neighbors of the test point in Figure 3.4, you will see that two of them are positive and one is negative. Through voting, positive would win.

The full algorithm for K-nearest neighbor classification is given in Algorithm 3.2. Note that there actually is no “training” phase for K-nearest neighbors. In this algorithm we have introduced five new conventions:

1. The training data is denoted by \mathbf{D} .
2. We assume that there are N -many training examples.
3. These examples are pairs $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$.
(Warning: do not confuse x_n , the n th training example, with x_d , the d th feature for example x .)
4. We use $[]$ to denote an empty list and $\oplus \cdot$ to append \cdot to that list.
5. Our prediction on \hat{x} is called \hat{y} .

The first step in this algorithm is to compute distances from the test point to all training points (lines 2-4). The data points are then sorted according to distance. We then apply a clever trick of *summing* the class labels for each of the K nearest neighbors (lines 6-10) and using the **SIGN** of this sum as our prediction.

The big question, of course, is how to choose K . As we’ve seen, with $K = 1$, we run the risk of overfitting. On the other hand, if K is large (for instance, $K = N$), then **KNN-PREDICT** will always predict the majority class. Clearly that is underfitting. So, K is a hyperparameter of the KNN algorithm that allows us to trade-off between overfitting (small value of K) and underfitting (large value of K).

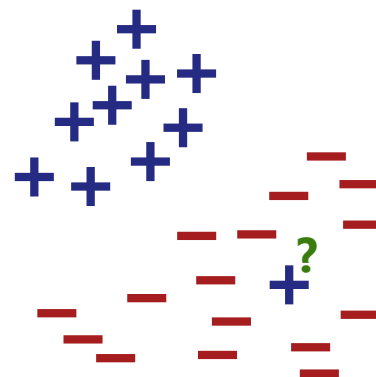


Figure 3.4: A figure showing an easy NN classification problem where the test point is a ? and should be positive, but its NN is actually a negative point that’s noisy.

? Why is it a good idea to use an odd number for K ?

? Why is the sign of the sum computed in lines 2-4 the same as the majority vote of the associated training examples?

? Why can’t you simply pick the value of K that does best on the training data? In other words, why do we have to treat it like a hyperparameter rather than just a parameter.

Algorithm 3 KNN-PREDICT(\mathbf{D}, K, \hat{x})

```

1:  $S \leftarrow []$ 
2: for  $n = 1$  to  $N$  do
3:    $S \leftarrow S \oplus \langle d(x_n, \hat{x}), n \rangle$            // store distance to training example  $n$ 
4: end for
5:  $S \leftarrow \text{SORT}(S)$                          // put lowest-distance objects first
6:  $\hat{y} \leftarrow 0$ 
7: for  $k = 1$  to  $K$  do
8:    $\langle \text{dist}, n \rangle \leftarrow S_k$                  //  $n$  this is the  $k$ th closest data point
9:    $\hat{y} \leftarrow \hat{y} + y_n$                      // vote according to the label for the  $n$ th training point
10: end for
11: return  $\text{SIGN}(\hat{y})$                            // return +1 if  $\hat{y} > 0$  and  $-1$  if  $\hat{y} < 0$ 

```

One aspect of **inductive bias** that we've seen for KNN is that it assumes that nearby points should have the same label. Another aspect, which is quite different from decision trees, is that all features are equally important! Recall that for decision trees, the key question was *which features are most useful for classification?* The whole learning algorithm for a decision tree hinged on finding a small set of good features. This is all thrown away in KNN classifiers: every feature is used, and they are all used the same amount. This means that if you have data with only a few relevant features and lots of irrelevant features, KNN is likely to do poorly.

A related issue with KNN is **feature scale**. Suppose that we are trying to classify whether some object is a ski or a snowboard (see Figure 3.5). We are given two features about this data: the width and height. As is standard in skiing, width is measured in millimeters and height is measured in centimeters. Since there are only two features, we can actually plot the entire training set; see Figure 3.6 where ski is the positive class. Based on this data, you might guess that a KNN classifier would do well.

Suppose, however, that our measurement of the width was computed in millimeters (instead of centimeters). This yields the data shown in Figure 3.7. Since the width values are now tiny, in comparison to the height values, a KNN classifier will effectively *ignore* the width values and classify almost purely based on height. The predicted class for the displayed test point had changed because of this feature scaling.

We will discuss feature scaling more in Chapter 5. For now, it is just important to keep in mind that KNN does not have the power to decide which features are important.

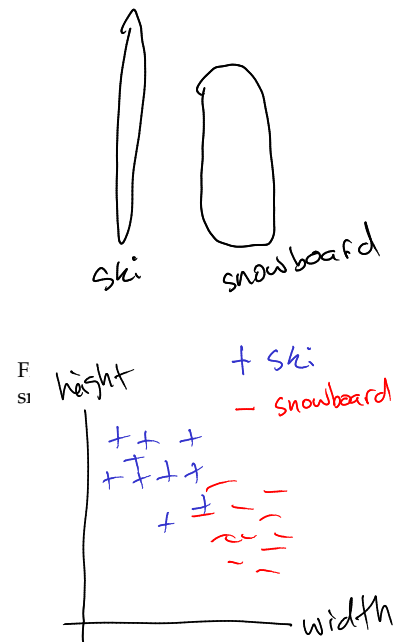


Figure 3.6: Classification data for ski vs snowboard in 2d

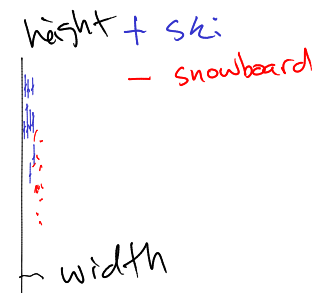


Figure 3.7: Classification data for ski vs snowboard in 2d, with width rescaled to mm.

MATH REVIEW | VECTOR ARITHMETIC AND VECTOR NORMS

A (real-valued) **vector** is just an array of real values, for instance $\mathbf{x} = \langle 1, 2.5, -6 \rangle$ is a three-dimensional vector. In general, if $\mathbf{x} = \langle x_1, x_2, \dots, x_D \rangle$, then x_d is its d th component. So $x_3 = -6$ in the previous example.

Vector sums are computed pointwise, and are only defined when dimensions match, so $\langle 1, 2.5, -6 \rangle + \langle 2, -2.5, 3 \rangle = \langle 3, 0, -3 \rangle$. In general, if $\mathbf{c} = \mathbf{a} + \mathbf{b}$ then $c_d = a_d + b_d$ for all d . Vector addition can be viewed geometrically as taking a vector \mathbf{a} , then tacking on \mathbf{b} to the end of it; the new end point is exactly \mathbf{c} .

Vectors can be **scaled** by real values; for instance $2\langle 1, 2.5, -6 \rangle = \langle 2, 5, -12 \rangle$; this is called scalar multiplication. In general, $a\mathbf{x} = \langle ax_1, ax_2, \dots, ax_D \rangle$.

The **norm** of a vector \mathbf{x} , written $\|\mathbf{x}\|$ is its length. Unless otherwise specified, this is its *Euclidean* length, namely: $\|\mathbf{x}\| = \sqrt{\sum_d x_d^2}$.

Figure 3.8:

3.3 Decision Boundaries

The standard way that we've been thinking about learning algorithms up to now is in the *query model*. Based on training data, you learn something. I then give you a query example and you have to guess its label.

An alternative, less passive, way to think about a learned model is to ask: what sort of test examples will it classify as positive, and what sort will it classify as negative. In Figure 3.9, we have a set of training data. The background of the image is colored blue in regions that *would* be classified as positive (if a query were issued there) and colored red in regions that *would* be classified as negative. This coloring is based on a 1-nearest neighbor classifier.

In Figure 3.9, there is a solid line separating the positive regions from the negative regions. This line is called the **decision boundary** for this classifier. It is the line with positive land on one side and negative land on the other side.

Decision boundaries are useful ways to visualize the **complexity** of a learned model. Intuitively, a learned model with a decision boundary that is really jagged (like the coastline of Norway) is really complex and prone to overfitting. A learned model with a decision boundary that is really simple (like the boundary between Arizona and Utah) is potentially underfit.

Now that you know about decision boundaries, it is natural to ask: what do decision boundaries for decision trees look like? In order

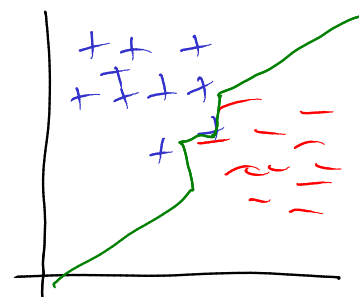
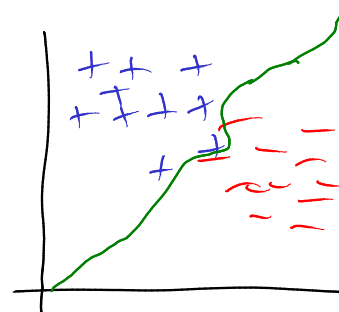


Figure 3.9: decision boundary for 1nn.

Figure 3.10: decision boundary for knn with $k=3$.

to answer this question, we have to be a bit more formal about how to build a decision tree on real-valued features. (Remember that the algorithm you learned in the previous chapter implicitly assumed *binary* feature values.) The idea is to allow the decision tree to ask questions of the form: “is the value of feature 5 greater than 0.2?” That is, for real-valued features, the decision tree nodes are parameterized by a feature and a threshold for that feature. An example decision tree for classifying skis versus snowboards is shown in Figure 3.11.

Now that a decision tree can handle feature vectors, we can talk about decision boundaries. By example, the decision boundary for the decision tree in Figure 3.11 is shown in Figure 3.12. In the figure, space is first split in half according to the first query along one axis. Then, depending on which half of the space you look at, it is either split again along the other axis, or simply classified.

Figure 3.12 is a good visualization of decision boundaries for decision trees in general. Their decision boundaries are axis-aligned cuts. The cuts must be axis-aligned because nodes can only query on a single feature at a time. In this case, since the decision tree was so shallow, the decision boundary was relatively simple.

3.4 K-Means Clustering

Up through this point, you have learned all about supervised learning (in particular, binary classification). As another example of the use of geometric intuitions and data, we are going to temporarily consider an **unsupervised learning** problem. In unsupervised learning, our data consists *only* of examples x_n and does *not* contain corresponding labels. Your job is to make sense of this data, even though no one has provided you with correct labels. The particular notion of “making sense of” that we will talk about now is the **clustering** task.

Consider the data shown in Figure 3.13. Since this is unsupervised learning and we do not have access to labels, the data points are simply drawn as black dots. Your job is to split this data set into three clusters. That is, you should label each data point as A, B or C in whatever way you want.

For this data set, it’s pretty clear what you should do. You probably labeled the upper-left set of points A, the upper-right set of points B and the bottom set of points C. Or perhaps you permuted these labels. But chances are your clusters were the same as mine.

The K-means clustering algorithm is a particularly simple and effective approach to producing clusters on data like you see in Figure 3.13. The idea is to represent each cluster by its cluster center. Given cluster centers, we can simply assign each point to its nearest

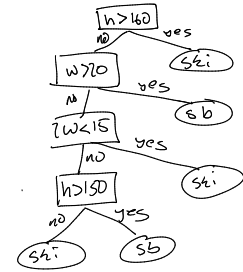


Figure 3.11: decision tree for ski vs. snowboard

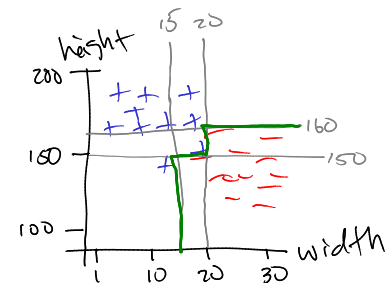


Figure 3.12: decision boundary for dt in previous figure

What sort of data might yield a very simple decision boundary with a decision tree and very complex decision boundary with 1-nearest neighbor? What about the other way around?

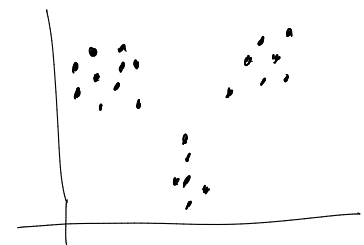


Figure 3.13: simple clustering data... clusters in UL, UR and BC.

center. Similarly, if we know the assignment of points to clusters, we can compute the centers. This introduces a chicken-and-egg problem. If we knew the clusters, we could compute the centers. If we knew the centers, we could compute the clusters. But we don't know either.

The general computer science answer to chicken-and-egg problems is **iteration**. We will start with a guess of the cluster centers. Based on that guess, we will assign each data point to its closest center. Given these new assignments, we can recompute the cluster centers. We repeat this process until clusters stop moving. The first few iterations of the K -means algorithm are shown in Figure 3.14. In this example, the clusters converge very quickly.

Algorithm 3.4 spells out the K -means clustering algorithm in detail. The cluster centers are initialized randomly. In line 6, data point x_n is compared against each cluster center μ_k . It is assigned to cluster k if k is the center with the smallest distance. (That is the “argmin” step.) The variable z_n stores the assignment (a value from 1 to K) of example n . In lines 8-12, the cluster centers are re-computed. First, X_k stores all examples that have been assigned to cluster k . The center of cluster k , μ_k is then computed as the mean of the points assigned to it. This process repeats until the centers converge.

An obvious question about this algorithm is: does it converge? A second question is: how long does it take to converge. The first question is actually easy to answer. Yes, it does. And in practice, it usually converges quite quickly (usually fewer than 20 iterations). In Chapter 15, we will actually *prove* that it converges. The question of how long it takes to converge is actually a really interesting question. Even though the K -means algorithm dates back to the mid 1950s, the best known convergence rates were *terrible* for a long time. Here, terrible means exponential in the number of data points. This was a sad situation because empirically we knew that it converged very quickly. New algorithm analysis techniques called “smoothed analysis” were invented in 2001 and have been used to show very fast convergence for K -means (among other algorithms). These techniques are well beyond the scope of this book (and this author!) but suffice it to say that K -means is fast in practice and is provably fast in theory.

It is important to note that although K -means is guaranteed to converge and guaranteed to converge quickly, it is *not* guaranteed to converge to the “right answer.” The key problem with unsupervised learning is that we have no way of knowing what the “right answer” is. Convergence to a bad solution is usually due to poor initialization.

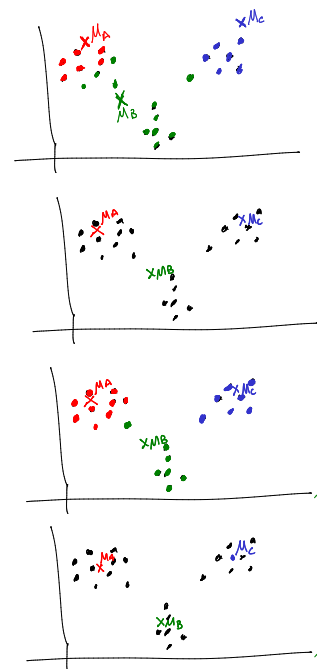


Figure 3.14: first few iterations of k -means running on previous data set



What is the difference between unsupervised and supervised learning that means that we know what the “right answer” is for supervised learning but not for unsupervised learning?

Algorithm 4 $K\text{-MEANS}(\mathbf{D}, K)$

```

1: for  $k = 1$  to  $K$  do
2:    $\mu_k \leftarrow$  some random location      // randomly initialize center for  $k$ th cluster
3: end for
4: repeat
5:   for  $n = 1$  to  $N$  do
6:      $z_n \leftarrow \operatorname{argmin}_k \|\mu_k - x_n\|$       // assign example  $n$  to closest center
7:   end for
8:   for  $k = 1$  to  $K$  do
9:      $X_k \leftarrow \{x_n : z_n = k\}$       // points assigned to cluster  $k$ 
10:     $\mu_k \leftarrow \operatorname{MEAN}(X_k)$       // re-estimate center of cluster  $k$ 
11:   end for
12: until  $\mu$ s stop changing
13: return  $z$       // return cluster assignments

```

3.5 Warning: High Dimensions are Scary

Visualizing one hundred dimensional space is incredibly difficult for humans. After huge amounts of training, some people have reported that they can visualize four dimensional space in their heads. But beyond that seems impossible.¹

In addition to being hard to visualize, there are at least two additional problems in high dimensions, both referred to as **the curse of dimensionality**. One is computational, the other is mathematical.

From a computational perspective, consider the following problem. For K -nearest neighbors, the speed of prediction is slow for a very large data set. At the very least you have to look at every training example every time you want to make a prediction. To speed things up you might want to create an *indexing* data structure. You can break the plane up into a grid like that shown in Figure 3.15. Now, when the test point comes in, you can quickly identify the grid cell in which it lies. Now, instead of considering *all* training points, you can limit yourself to training points *in that grid cell* (and perhaps the neighboring cells). This can potentially lead to huge computational savings.

In two dimensions, this procedure is effective. If we want to break space up into a grid whose cells are 0.2×0.2 , we can clearly do this with 25 grid cells in two dimensions (assuming the range of the features is 0 to 1 for simplicity). In three dimensions, we'll need $125 = 5 \times 5 \times 5$ grid cells. In four dimensions, we'll need 625. By the time we get to "low dimensional" data in 20 dimensions, we'll need 95,367,431,640,625 grid cells (that's 95 trillion, which is about 6 to 7 times the US national debt as of January 2011). So if you're in 20 dimensions, this gridding technique will only be useful if you have at least 95 trillion training examples.

¹ If you want to try to get an intuitive sense of what four dimensions looks like, I highly recommend the short 1884 book *Flatland: A Romance of Many Dimensions* by Edwin Abbott Abbott. You can even read it online at gutenberg.org/ebooks/201.

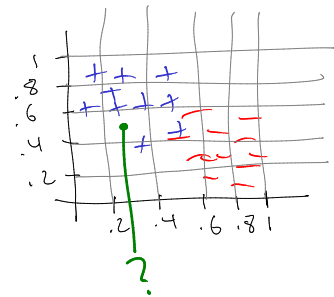


Figure 3.15: 2d knn with an overlaid grid, cell with test point highlighted

For “medium dimensional” data (approximately 1000) dimensions, the number of grid cells is a 9 followed by 698 numbers before the decimal point. For comparison, the number of atoms in the universe is approximately 1 followed by 80 zeros. So even if each atom yielded a googol training examples, we’d still have far fewer examples than grid cells. For “high dimensional” data (approximately 100000) dimensions, we have a 1 followed by just under 70,000 zeros. Far too big a number to even really comprehend.

Suffice it to say that for even moderately high dimensions, the amount of computation involved in these problems is enormous.

In addition to the computational difficulties of working in high dimensions, there are a large number of strange mathematical occurrences there. In particular, many of your intuitions that you’ve built up from working in two and three dimensions just do not carry over to high dimensions. We will consider two effects, but there are countless others. The first is that high dimensional spheres look more like porcupines than like balls.² The second is that distances between points in high dimensions are all approximately the same.

Let’s start in two dimensions as in Figure 3.16. We’ll start with four green spheres, each of radius one and each touching exactly two other green spheres. (Remember that in two dimensions a “sphere” is just a “circle.”) We’ll place a red sphere in the middle so that it touches all four green spheres. We can easily compute the radius of this small sphere. The pythagorean theorem says that $1^2 + 1^2 = (1 + r)^2$, so solving for r we get $r = \sqrt{2} - 1 \approx 0.41$. Thus, by calculation, the blue sphere lies entirely within the cube (cube = square) that contains the grey spheres. (Yes, this is also obvious from the picture, but perhaps you can see where this is going.)

Now we can do the same experiment in three dimensions, as shown in Figure 3.17. Again, we can use the pythagorean theorem to compute the radius of the blue sphere. Now, we get $1^2 + 1^2 + 1^2 = (1 + r)^2$, so $r = \sqrt{3} - 1 \approx 0.73$. This is still entirely enclosed in the cube of width four that holds all eight grey spheres.

At this point it becomes difficult to produce figures, so you’ll have to apply your imagination. In four dimensions, we would have 16 green spheres (called **hyperspheres**), each of radius one. They would still be inside a cube (called a **hypercube**) of width four. The blue hypersphere would have radius $r = \sqrt{4} - 1 = 1$. Continuing to five dimensions, the blue hypersphere embedded in 256 green hyperspheres would have radius $r = \sqrt{5} - 1 \approx 1.23$ and so on.

In general, in D -dimensional space, there will be 2^D green hyperspheres of radius one. Each green hypersphere will touch exactly n -many other hyperspheres. The blue hyperspheres in the middle will touch them all and will have radius $r = \sqrt{D} - 1$.

? How does the above analysis relate to the number of data points you would need to fill out a full decision tree with D -many features? What does this say about the importance of shallow trees?

² This result was related to me by Mark Reid, who heard about it from Marcus Hutter.

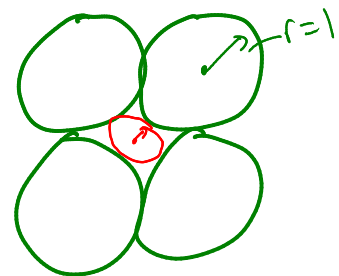


Figure 3.16: 2d spheres in spheres

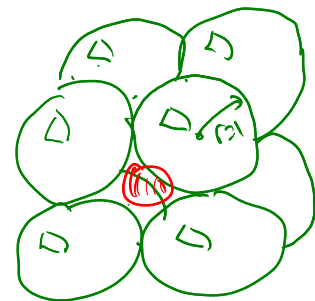


Figure 3.17: 3d spheres in spheres

Think about this for a moment. As the number of dimensions grows, the radius of the blue hypersphere *grows without bound!*. For example, in 9-dimensions the radius of the blue hypersphere is now $\sqrt{9} - 1 = 2$. But with a radius of two, the blue hypersphere is now “squeezing” between the green hypersphere and *touching* the edges of the hypercube. In 10 dimensional space, the radius is approximately 2.16 and it pokes outside the cube.

The second strange fact we will consider has to do with the distances between points in high dimensions. We start by considering random points in one dimension. That is, we generate a fake data set consisting of 100 random points between zero and one. We can do the same in two dimensions and in three dimensions. See Figure ?? for data distributed uniformly on the **unit hypercube** in different dimensions.

Now, pick two of these points at random and compute the distance between them. Repeat this process for all pairs of points and average the results. For the data shown in Figure ??, the average distance between points in one dimension is about 0.346; in two dimensions is about 0.518; and in three dimensions is 0.615. The fact that these *increase* as the dimension increases is not surprising. The furthest two points can be in a 1-dimensional hypercube (line) is 1; the furthest in a 2-dimensional hypercube (square) is $\sqrt{2}$ (opposite corners); the furthest in a 3-d hypercube is $\sqrt{3}$ and so on. In general, the furthest two points in a D -dimensional hypercube will be \sqrt{D} .

You can actually compute these values analytically. Write \mathcal{Uni}_D for the uniform distribution in D dimensions. The quantity we are interested in computing is:

$$\text{avgDist}(D) = \mathbb{E}_{a \sim \mathcal{Uni}_D} \left[\mathbb{E}_{b \sim \mathcal{Uni}_D} \left[\|a - b\| \right] \right] \quad (3.2)$$

We can actually compute this in closed form and arrive at $\text{avgDist}(D) = \sqrt{D}/3$. Because we know that the maximum distance between two points grows like \sqrt{D} , this says that the ratio between average distance and maximum distance converges to $1/3$.

What is more interesting, however, is the *variance* of the distribution of distances. You can show that in D dimensions, the variance is *constant* $1/\sqrt{18}$, *independent of* D . This means that when you look at (variance) divided-by (max distance), the variance behaves like $1/\sqrt{18D}$, which means that the effective variance continues to shrink as D grows³.

When I first saw and re-proved this result, I was skeptical, as I imagine you are. So I implemented it. In Figure 3.18 you can see the results. This presents a *histogram* of distances between random points in D dimensions for $D \in \{1, 2, 3, 10, 20, 100\}$. As you can see, all of these distances begin to concentrate around $0.4\sqrt{D}$, even for

³ Brin 1995

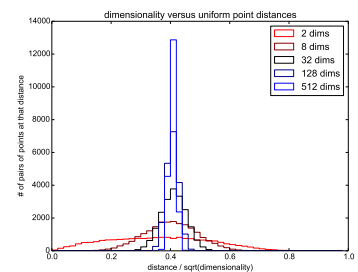


Figure 3.18: histogram of distances in $D=2,8,32,128,512$

“medium dimension” problems.

You should now be terrified: the only bit of information that KNN gets is distances. And you’ve just seen that in moderately high dimensions, all distances becomes equal. So then isn’t it the case that KNN simply cannot work?

The answer has to be no. The reason is that the data that we get is *not* uniformly distributed over the unit hypercube. We can see this by looking at two real-world data sets. The first is an image data set of hand-written digits (zero through nine); see Section ???. Although this data is originally in 256 dimensions (16 pixels by 16 pixels), we can artificially reduce the dimensionality of this data. In Figure 3.19 you can see the histogram of average distances between points in this data at a number of dimensions.

As you can see from these histograms, distances have *not* concentrated around a single value. This is very good news: it means that there is hope for learning algorithms to work! Nevertheless, the moral is that high dimensions are weird.

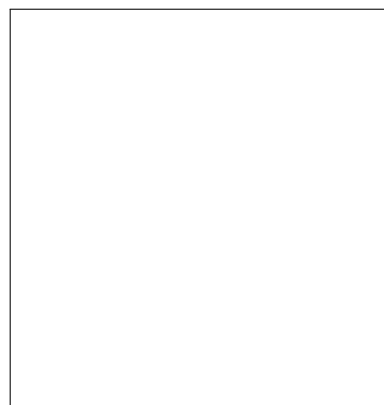


Figure 3.19: knn:mnist: histogram of distances in multiple D for mnist

3.6 Further Reading

TODO further reading

4 | THE PERCEPTRON

Algebra is nothing more than geometry, in words; geometry is nothing more than algebra, in pictures. — Sophie Germain

SO FAR, YOU'VE SEEN TWO TYPES of learning models: in decision trees, only a small number of features are used to make decisions; in nearest neighbor algorithms, all features are used equally. Neither of these extremes is always desirable. In some problems, we might want to use most of the features, but use some more than others.

In this chapter, we'll discuss the **perceptron** algorithm for learning **weights** for features. As we'll see, learning weights for features amounts to learning a **hyperplane** classifier: that is, basically a division of space into two halves by a straight line, where one half is "positive" and one half is "negative." In this sense, the perceptron can be seen as explicitly finding a good **linear decision boundary**.

4.1 Bio-inspired Learning

Folk biology tells us that our brains are made up of a bunch of little units, called **neurons**, that send electrical signals to one another. The *rate* of firing tells us how "activated" a neuron is. A single neuron, like that shown in Figure 4.1 might have three incoming neurons. These incoming neurons are firing at different rates (i.e., have different **activations**). Based on how much these incoming neurons are firing, and how "strong" the neural connections are, our main neuron will "decide" how strongly it wants to fire. And so on through the whole brain. Learning in the brain happens by neurons becoming connected to other neurons, and the strengths of connections adapting over time.

The real biological world is much more complicated than this. However, our goal isn't to build a brain, but to simply be *inspired* by how they work. We are going to think of our learning algorithm as a *single* neuron. It receives input from D -many other neurons, one for each input feature. The strength of these inputs are the feature values. This is shown schematically in Figure 4.1. Each incoming connection has a weight and the neuron simply sums up all the weighted inputs. Based on this sum, it decides whether to "fire" or

Learning Objectives:

- Describe the biological motivation behind the perceptron.
- Classify learning algorithms based on whether they are error-driven or not.
- Implement the perceptron algorithm for binary classification.
- Draw perceptron weight vectors and the corresponding decision boundaries in two dimensions.
- Contrast the decision boundaries of decision trees, nearest neighbor algorithms and perceptrons.
- Compute the margin of a given weight vector on a given data set.

Dependencies: Chapter 1, Chapter 3

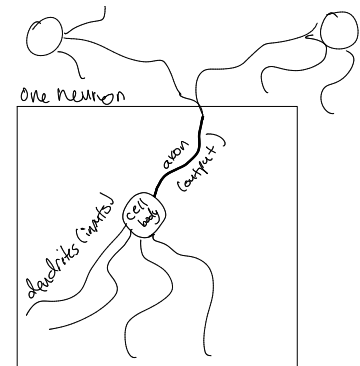


Figure 4.1: a picture of a neuron

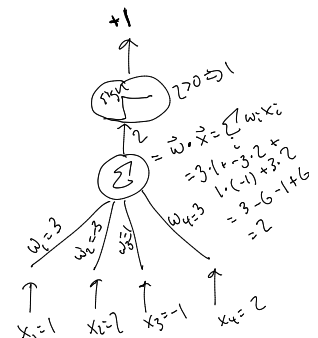


Figure 4.2: figure showing feature vector and weight vector and products and sum

not. Firing is interpreted as being a positive example and not firing is interpreted as being a negative example. In particular, if the weighted sum is positive, it “fires” and otherwise it doesn’t fire. This is shown diagrammatically in Figure 4.2.

Mathematically, an input vector $\mathbf{x} = \langle x_1, x_2, \dots, x_D \rangle$ arrives. The neuron stores D -many weights, w_1, w_2, \dots, w_D . The neuron computes the sum:

$$a = \sum_{d=1}^D w_d x_d \quad (4.1)$$

to determine its amount of “activation.” If this activation is positive (i.e., $a > 0$) it predicts that this example is a positive example. Otherwise it predicts a negative example.

The weights of this neuron are fairly easy to interpret. Suppose that a feature, for instance “is this a System’s class?” gets a zero weight. Then the activation is the same regardless of the value of this feature. So features with zero weight are ignored. Features with positive weights are indicative of positive examples because they cause the activation to increase. Features with negative weights are indicative of negative examples because they cause the activation to decrease.

It is often convenient to have a non-zero **threshold**. In other words, we might want to predict positive if $a > \theta$ for some value θ . The way that is most convenient to achieve this is to introduce a **bias** term into the neuron, so that the activation is always increased by some fixed value b . Thus, we compute:

$$a = \left[\sum_{d=1}^D w_d x_d \right] + b \quad (4.2)$$

This is the complete neural model of learning. The model is parameterized by D -many weights, w_1, w_2, \dots, w_D , and a single scalar bias value b .

?

What would happen if we encoded binary features like “is this a System’s class?” as no=0 and yes=-1 (rather than the standard no=0 and yes=+1)?

?

If you wanted the activation threshold to be $a > \theta$ instead of $a > 0$, what value would b have to be?

4.2 Error-Driven Updating: The Perceptron Algorithm

The **perceptron** is a classic learning algorithm for the neural model of learning. Like K -nearest neighbors, it is one of those frustrating algorithms that is incredibly simple and yet works amazingly well, for some types of problems.

The algorithm is actually quite different than either the decision tree algorithm or the KNN algorithm. First, it is **online**. This means that instead of considering the *entire data set* at the same time, it only ever looks at one example. It processes that example and then goes

Algorithm 5 PERCEPTONTRAIN(\mathbf{D} , MaxIter)

```

1:  $w_d \leftarrow 0$ , for all  $d = 1 \dots D$  // initialize weights
2:  $b \leftarrow 0$  // initialize bias
3: for  $\text{iter} = 1 \dots \text{MaxIter}$  do
4:   for all  $(x, y) \in \mathbf{D}$  do
5:      $a \leftarrow \sum_{d=1}^D w_d x_d + b$  // compute activation for this example
6:     if  $ya \leq 0$  then
7:        $w_d \leftarrow w_d + yx_d$ , for all  $d = 1 \dots D$  // update weights
8:        $b \leftarrow b + y$  // update bias
9:     end if
10:   end for
11: end for
12: return  $w_0, w_1, \dots, w_D, b$ 

```

Algorithm 6 PERCEPTONTEST($w_0, w_1, \dots, w_D, b, \hat{x}$)

```

1:  $a \leftarrow \sum_{d=1}^D w_d \hat{x}_d + b$  // compute activation for the test example
2: return SIGN( $a$ )

```

on to the next one. Second, it is **error driven**. This means that, so long as it is doing well, it doesn't bother updating its parameters.

The algorithm maintains a “guess” at good parameters (weights and bias) as it runs. It processes one example at a time. For a given example, it makes a prediction. It checks to see if this prediction is correct (recall that this is *training data*, so we have access to true labels). If the prediction is correct, it does nothing. Only when the prediction is incorrect does it change its parameters, and it changes them in such a way that it would do better on this example next time around. It then goes on to the next example. Once it hits the last example in the training set, it loops back around for a specified number of iterations.

The training algorithm for the perceptron is shown in Algorithm 4.2 and the corresponding prediction algorithm is shown in Algorithm 4.2. There is one “trick” in the training algorithm, which probably seems silly, but will be useful later. It is in line 6, when we check to see if we want to make an update or not. We want to make an update if the current prediction (just $\text{SIGN}(a)$) is incorrect. The trick is to multiply the true label y by the activation a and compare this against zero. Since the label y is either $+1$ or -1 , you just need to realize that ya is positive whenever a and y have the same sign. In other words, the product ya is positive if the current prediction is correct.

The particular form of update for the perceptron is quite simple. The weight w_d is increased by yx_d and the bias is increased by y . The goal of the update is to adjust the parameters so that they are “better” for the current example. In other words, if we saw this example

? It is very very important to check $ya \leq 0$ rather than $ya < 0$. Why?

twice in a row, we should do a better job the second time around.

To see why this particular update achieves this, consider the following scenario. We have some current set of parameters w_1, \dots, w_D, b . We observe an example (x, y) . For simplicity, suppose this is a positive example, so $y = +1$. We compute an activation a , and make an error. Namely, $a < 0$. We now update our weights and bias. Let's call the new weights w'_1, \dots, w'_D, b' . Suppose we observe the same example again and need to compute a new activation a' . We proceed by a little algebra:

$$a' = \sum_{d=1}^D w'_d x_d + b' \quad (4.3)$$

$$= \sum_{d=1}^D (w_d + x_d) x_d + (b + 1) \quad (4.4)$$

$$= \sum_{d=1}^D w_d x_d + b + \sum_{d=1}^D x_d x_d + 1 \quad (4.5)$$

$$= a + \sum_{d=1}^D x_d^2 + 1 > a \quad (4.6)$$

So the difference between the old activation a and the new activation a' is $\sum_d x_d^2 + 1$. But $x_d^2 \geq 0$, since it's squared. So this value is *always* at least one. Thus, the new activation is always at least the old activation plus one. Since this was a positive example, we have successfully moved the activation in the proper direction. (Though note that there's no guarantee that we will correctly classify this point the second, third or even fourth time around!)

The only **hyperparameter** of the perceptron algorithm is *MaxIter*, the number of passes to make over the training data. If we make many many passes over the training data, then the algorithm is likely to overfit. (This would be like studying *too long* for an exam and just confusing yourself.) On the other hand, going over the data only one time might lead to underfitting. This is shown experimentally in Figure 4.3. The x-axis shows the number of passes over the data and the y-axis shows the training error and the test error. As you can see, there is a “sweet spot” at which test performance begins to degrade due to overfitting.

One aspect of the perceptron algorithm that is left underspecified is line 4, which says: loop over all the training examples. The natural implementation of this would be to loop over them in a constant order. This is actually a bad idea.

Consider what the perceptron algorithm would do on a data set that consisted of 500 positive examples followed by 500 negative examples. After seeing the first few positive examples (maybe five), it would likely decide that *every* example is positive, and would stop

? This analysis holds for the case positive examples ($y = +1$). It should also hold for negative examples. Work it out.

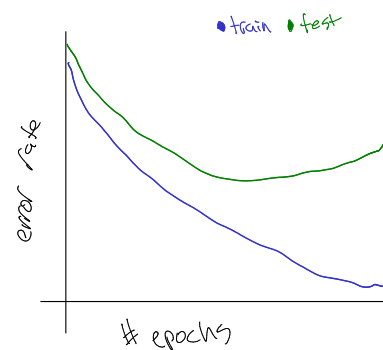


Figure 4.3: training and test error via early stopping

learning anything. It would do well for a while (next 495 examples), until it hit the batch of negative examples. Then it would take a while (maybe ten examples) before it would start predicting everything as negative. By the end of one pass through the data, it would really only have learned from a handful of examples (fifteen in this case).

So one thing you need to avoid is presenting the examples in some fixed order. This can easily be accomplished by permuting the order of examples once in the beginning and then cycling over the data set in the same (permuted) order each iteration. However, it turns out that you can actually do *better* if you re-permute the examples in each iteration. Figure 4.4 shows the effect of re-permuting on convergence speed. In practice, permuting each iteration tends to yield about 20% savings in number of iterations. In theory, you can actually prove that it's expected to be about twice as fast.

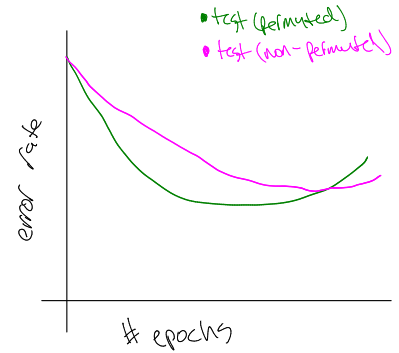


Figure 4.4: training and test error for permuting versus not-permuting

? If permuting the data each iteration saves somewhere between 20% and 50% of your time, are there any cases in which you might *not* want to permute the data every iteration?

4.3 Geometric Interpretation

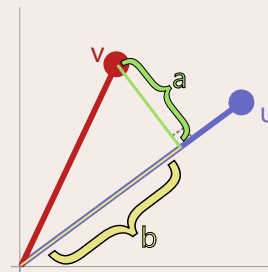
A question you should be asking yourself by now is: what does the decision boundary of a perceptron look like? You can actually answer that question mathematically. For a perceptron, the decision boundary is precisely where the sign of the activation, a , changes from -1 to $+1$. In other words, it is the set of points x that achieve *zero* activation. The points that are not clearly positive nor negative. For simplicity, we'll first consider the case where there is no "bias" term (or, equivalently, the bias is zero). Formally, the decision boundary \mathcal{B} is:

$$\mathcal{B} = \left\{ x : \sum_d w_d x_d = 0 \right\} \quad (4.7)$$

We can now apply some linear algebra. Recall that $\sum_d w_d x_d$ is just the **dot product** between the vector $w = \langle w_1, w_2, \dots, w_D \rangle$ and the vector x . We will write this as $w \cdot x$. Two vectors have a zero dot product if and only if they are **perpendicular**. Thus, if we think of the weights as a vector w , then the decision boundary is simply the plane perpendicular to w .

MATH REVIEW | DOT PRODUCTS

Given two vectors u and v their dot product $u \cdot v$ is $\sum_d u_d v_d$. The dot product grows large and positive when u and v point in same direction, grows large and negative when u and v point in opposite directions, and is zero when they are perpendicular. A useful geometric interpretation of dot products is **projection**. Suppose $\|u\| = 1$, so that u is a **unit vector**. We can think of any other vector v as consisting of two components: (a) a component in the direction of u and (b) a component that's perpendicular to u . This is depicted geometrically to the right: Here, $u = \langle 0.8, 0.6 \rangle$ and $v = \langle 0.37, 0.73 \rangle$. We can think of v as the sum of two vectors, a and b , where a is parallel to u and b is perpendicular. The length of b is exactly $u \cdot v = 0.734$, which is why you can think of dot products as projections: the dot product between u and v is the “projection of v onto u .”



This is shown pictorially in Figure 4.6. Here, the weight vector is shown, together with its perpendicular plane. This plane forms the decision boundary between positive points and negative points. The vector points in the direction of the positive examples and away from the negative examples.

One thing to notice is that the *scale* of the weight vector is irrelevant from the perspective of classification. Suppose you take a weight vector w and replace it with $2w$. All activations are now doubled. But their sign does not change. This makes complete sense geometrically, since all that matters is which side of the plane a test point falls on, now how far it is from that plane. For this reason, it is common to work with **normalized** weight vectors, w , that have length one; i.e., $\|w\| = 1$.

The geometric intuition can help us even more when we realize that dot products compute projections. That is, the value $w \cdot x$ is just the distance of x from the origin when projected *onto* the vector w . This is shown in Figure 4.7. In that figure, all the data points are projected onto w . Below, we can think of this as a one-dimensional version of the data, where each data point is placed according to its projection along w . This distance along w is exactly the *activation* of that example, with no bias.

From here, you can start thinking about the role of the bias term. Previously, the threshold would be at zero. Any example with a negative projection onto w would be classified negative; any example with a positive projection, positive. The bias simply moves this threshold. Now, after the projection is computed, b is added to get the overall activation. The projection *plus* b is then compared against

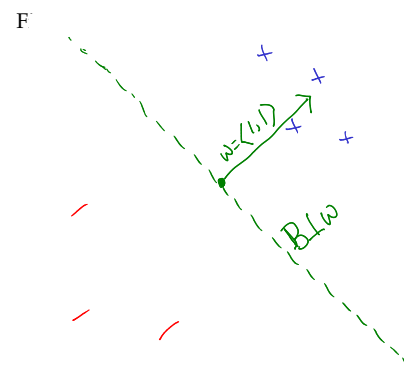


Figure 4.6: picture of data points with hyperplane and weight vector

? If I give you a non-zero weight vector w , how do I compute a weight vector w' that points in the same direction but has a norm of one?

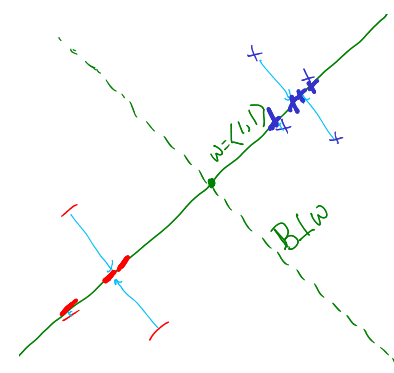


Figure 4.7: The same picture as before, but with projections onto weight vector; then, below, those points along a one-dimensional axis with zero marked.

zero.

Thus, from a geometric perspective, the role of the bias is to *shift* the decision boundary away from the origin, in the direction of w . It is shifted exactly $-b$ units. So if b is positive, the boundary is shifted away from w and if b is negative, the boundary is shifted toward w . This is shown in Figure 4.8. This makes intuitive sense: a positive bias means that more examples should be classified positive. By moving the decision boundary in the negative direction, more space yields a positive classification.

The decision boundary for a perceptron is a very magical thing. In D dimensional space, it is always a $D - 1$ -dimensional hyperplane. (In two dimensions, a 1-d hyperplane is simply a line. In three dimensions, a 2-d hyperplane is like a sheet of paper.) This hyperplane divides space in half. In the rest of this book, we'll refer to the weight vector, and to hyperplane it defines, interchangeably.

The perceptron update can also be considered geometrically. (For simplicity, we will consider the **unbiased** case.) Consider the situation in Figure 4.9. Here, we have a current guess as to the hyperplane, and positive training example comes in that is currently mis-classified. The weights are updated: $w \leftarrow w + yx$. This yields the new weight vector, also shown in the Figure. In this case, the weight vector changed enough that this training example is now correctly classified.

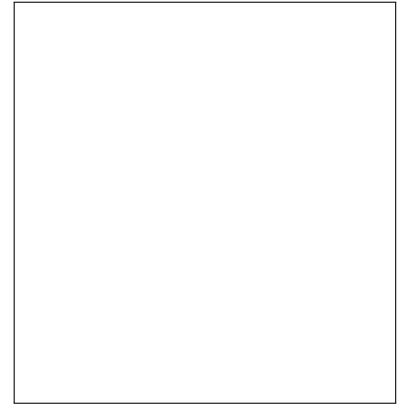


Figure 4.8: perc:bias: perceptron picture with bias

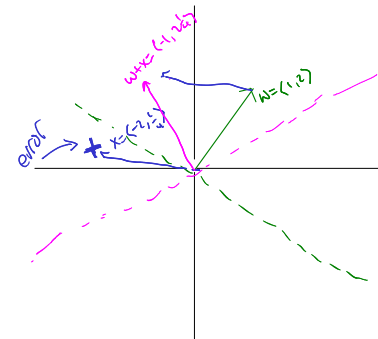


Figure 4.9: perceptron picture with update, no bias

4.4 Interpreting Perceptron Weights

You may find yourself having run the perceptron, learned a really awesome classifier, and then wondering “what the heck is this classifier doing?” You might ask this question because you’re curious to learn something about the underlying data. You might ask this question because you want to make sure that the perceptron is learning “the right thing.” You might ask this question because you want to remove a bunch of features that aren’t very useful because they’re expensive to compute or take a lot of storage.

The perceptron learns a classifier of the form $x \mapsto \text{sign}(\sum_d w_d x_d + b)$. A reasonable question to ask is: how sensitive is the final classification to *small changes* in some particular feature. We can answer this question by taking a derivative. If we arbitrarily take the 7th feature we can compute $\frac{\partial}{\partial x_7} (\sum_d w_d x_d + b) = w_7$. This says: the rate at which the activation changes as a function of the 7th feature is exactly w_7 . This gives rise to a useful heuristic for interpreting perceptron weights: **sort all the weights from largest (positive) to largest (negative), and take the top ten and bottom ten**. The top ten are the features that the perceptron is most sensitive to for making positive

predictions. The bottom ten are the features that the perceptron is most sensitive to for making negative predictions.

This heuristic is useful, especially when the inputs x consist entirely of binary values (like a bag of words representation). The heuristic is less useful when the range of the individual features varies significantly. The issue is that if you have one feat x_5 that's either 0 or 1, and another feature x_7 that's either 0 or 100, but $w_5 = w_7$, it's reasonable to say that w_7 is more important because it is likely to have a much larger influence on the final prediction. The easiest way to compensate for this is simply to scale your features ahead of time: this is another reason why feature scaling is a useful preprocessing step.

4.5 Perceptron Convergence and Linear Separability

You already have an intuitive feeling for why the perceptron works: it moves the decision boundary in the direction of the training examples. A question you should be asking yourself is: does the perceptron converge? If so, what does it converge to? And how long does it take?

It is easy to construct data sets on which the perceptron algorithm will never converge. In fact, consider the (very uninteresting) learning problem with *no features*. You have a data set consisting of one positive example and one negative example. Since there are no features, the only thing the perceptron algorithm will ever do is adjust the bias. Given this data, you can run the perceptron for a bajillion iterations and it will never settle down. As long as the bias is non-negative, the negative example will cause it to decrease. As long as it is non-positive, the positive example will cause it to increase. Ad infinitum. (Yes, this is a very contrived example.)

What does it mean for the perceptron to converge? It means that it can make an entire pass through the training data without making *any* more updates. In other words, it has correctly classified *every* training example. Geometrically, this means that it was found some hyperplane that correctly segregates the data into positive and negative examples, like that shown in Figure 4.10.

In this case, this data is **linearly separable**. This means that there exists *some* hyperplane that puts all the positive examples on one side and all the negative examples on the other side. If the training is *not* linearly separable, like that shown in Figure 4.11, then the perceptron has no hope of converging. It could never possibly classify each point correctly.

The somewhat surprising thing about the perceptron algorithm is that *if* the data is linearly separable, *then* it will converge to a weight

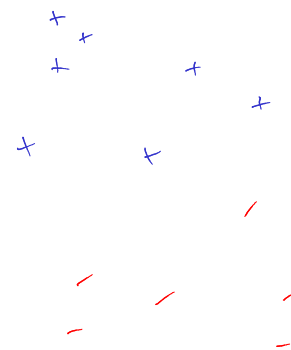


Figure 4.10: separable data

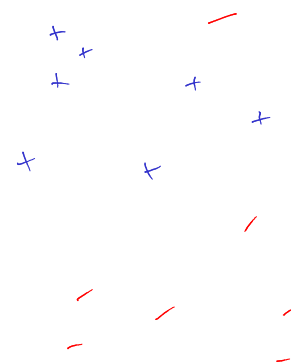


Figure 4.11: inseparable data

vector that separates the data. (And if the data is inseparable, then it will never converge.) This is great news. It means that the perceptron converges whenever it is even remotely possible to converge.

The second question is: how long does it take to converge? By “how long,” what we really mean is “how many updates?” As is the case for much learning theory, you will not be able to get an answer of the form “it will converge after 5293 updates.” This is asking too much. The sort of answer we can hope to get is of the form “it will converge after *at most* 5293 updates.”

What you might expect to see is that the perceptron will converge more quickly for easy learning problems than for hard learning problems. This certainly fits intuition. The question is how to *define* “easy” and “hard” in a meaningful way. One way to make this definition is through the notion of **margin**. If I give you a data set and hyperplane that separates it then the *margin* is the distance between the hyperplane and the nearest point. Intuitively, problems with large margins should be easy (there’s lots of “wiggle room” to find a separating hyperplane); and problems with small margins should be hard (you really have to get a very specific well tuned weight vector).

Formally, given a data set \mathbf{D} , a weight vector w and bias b , the margin of w, b on \mathbf{D} is defined as:

$$\text{margin}(\mathbf{D}, w, b) = \begin{cases} \min_{(x,y) \in \mathbf{D}} y(w \cdot x + b) & \text{if } w \text{ separates } \mathbf{D} \\ -\infty & \text{otherwise} \end{cases} \quad (4.8)$$

In words, the margin is only defined if w, b actually separate the data (otherwise it is just $-\infty$). In the case that it separates the data, we find the point with the minimum activation, after the activation is multiplied by the label.

For some historical reason (that is unknown to the author), margins are always denoted by the Greek letter γ (gamma). One often talks about the **margin of a data set**. The margin of a data set is the largest attainable margin on this data. Formally:

$$\text{margin}(\mathbf{D}) = \sup_{w,b} \text{margin}(\mathbf{D}, w, b) \quad (4.9)$$

In words, to compute the margin of a data set, you “try” every possible w, b pair. For each pair, you compute its margin. We then take the largest of these as the overall margin of the data.¹ If the data is not linearly separable, then the value of the sup, and therefore the value of the margin, is $-\infty$.

There is a famous theorem due to Rosenblatt² that shows that the number of errors that the perceptron algorithm makes is bounded by γ^{-2} . More formally:

? So long as the margin is not $-\infty$, it is always positive. Geometrically this makes sense, but why does Eq (4.8) yield this?

¹ You can read “sup” as “max” if you like: the only difference is a technical difference in how the $-\infty$ case is handled.

² Rosenblatt 1958

Theorem 2 (Perceptron Convergence Theorem). *Suppose the perceptron algorithm is run on a linearly separable data set \mathbf{D} with margin $\gamma > 0$. Assume that $\|x\| \leq 1$ for all $x \in \mathbf{D}$. Then the algorithm will converge after at most $\frac{1}{\gamma^2}$ updates.*

The proof of this theorem is elementary, in the sense that it does not use any fancy tricks: it's all just algebra. The *idea* behind the proof is as follows. If the data is linearly separable with margin γ , then there exists some weight vector w^* that achieves this margin. Obviously we don't know what w^* is, but we know it exists. The perceptron algorithm is trying to find a weight vector w that points roughly in the same direction as w^* . (For large γ , "roughly" can be very rough. For small γ , "roughly" is quite precise.) Every time the perceptron makes an update, the angle between w and w^* changes. What we prove is that the angle actually *decreases*. We show this in two steps. First, the dot product $w \cdot w^*$ increases a lot. Second, the norm $\|w\|$ does not increase very much. Since the dot product is increasing, but w isn't getting too long, the angle between them has to be shrinking. The rest is algebra.

Proof of Theorem 2. The margin $\gamma > 0$ must be realized by some set of parameters, say x^* . Suppose we train a perceptron on this data. Denote by $w^{(0)}$ the initial weight vector, $w^{(1)}$ the weight vector after the *first update*, and $w^{(k)}$ the weight vector after the k th update. (We are essentially ignoring data points on which the perceptron doesn't update itself.) First, we will show that $w^* \cdot w^{(k)}$ grows quickly as a function of k . Second, we will show that $\|w^{(k)}\|$ does not grow quickly.

First, suppose that the k th update happens on example (x, y) . We are trying to show that $w^{(k)}$ is becoming aligned with w^* . Because we updated, know that this example was misclassified: $yw^{(k-1)} \cdot x < 0$. After the update, we get $w^{(k)} = w^{(k-1)} + yx$. We do a little computation:

$$w^* \cdot w^{(k)} = w^* \cdot (w^{(k-1)} + yx) \quad \text{definition of } w^{(k)} \quad (4.10)$$

$$= w^* \cdot w^{(k-1)} + yw^* \cdot x \quad \text{vector algebra} \quad (4.11)$$

$$\geq w^* \cdot w^{(k-1)} + \gamma \quad w^* \text{ has margin } \gamma \quad (4.12)$$

Thus, every time $w^{(k)}$ is updated, its projection onto w^* increases by at least γ . **Therefore:** $w^* \cdot w^{(k)} \geq k\gamma$.

Next, we need to show that the increase of γ along w^* occurs because $w^{(k)}$ is getting closer to w^* , not just because it's getting exceptionally long. To do this, we compute the norm of $w^{(k)}$:

$$\|w^{(k)}\|^2$$

$$= \left\| \mathbf{w}^{(k-1)} + y\mathbf{x} \right\|^2 \quad \text{def. of } \mathbf{w}^{(k)} \quad (4.13)$$

$$= \left\| \mathbf{w}^{(k-1)} \right\|^2 + y^2 \|\mathbf{x}\|^2 + 2y\mathbf{w}^{(k-1)} \cdot \mathbf{x} \quad \text{quadratic rule} \quad (4.14)$$

$$\leq \left\| \mathbf{w}^{(k-1)} \right\|^2 + 1 + 0 \quad \text{assumption and } a < 0 \quad (4.15)$$

Thus, the squared norm of $\mathbf{w}^{(k)}$ increases by at most one every update. **Therefore:** $\left\| \mathbf{w}^{(k)} \right\|^2 \leq k$.

Now we put together the two things we have learned before. By our first conclusion, we know $\mathbf{w}^* \cdot \mathbf{w}^{(k)} \geq k\gamma$. But our second conclusion, $\sqrt{k} \geq \left\| \mathbf{w}^{(k)} \right\|^2$. Finally, because \mathbf{w}^* is a unit vector, we know that $\left\| \mathbf{w}^{(k)} \right\| \geq \mathbf{w}^* \cdot \mathbf{w}^{(k)}$. Putting this together, we have:

$$\sqrt{k} \geq \left\| \mathbf{w}^{(k)} \right\| \geq \mathbf{w}^* \cdot \mathbf{w}^{(k)} \geq k\gamma \quad (4.16)$$

Taking the left-most and right-most terms, we get that $\sqrt{k} \geq k\gamma$.

Dividing both sides by k , we get $\frac{1}{\sqrt{k}} \geq \gamma$ and therefore $k \leq \frac{1}{\gamma^2}$.

This means that once we've made $\frac{1}{\gamma^2}$ updates, we cannot make any more! \square

It is important to keep in mind what this proof shows and what it does not show. It shows that if I give the perceptron data that is linearly separable with margin $\gamma > 0$, then the perceptron will converge to a solution that separates the data. And it will converge quickly when γ is large. It does not say anything about the solution, *other than* the fact that it separates the data. In particular, the proof makes use of the maximum margin separator. But the perceptron is not guaranteed to *find* this maximum margin separator. The data may be separable with margin 0.9 and the perceptron might still find a separating hyperplane with a margin of only 0.000001. Later (in Chapter 7), we will see algorithms that explicitly try to find the maximum margin solution.

Perhaps we don't want to assume that all \mathbf{x} have norm at most 1. If they have all have norm at most R , you can achieve a very similar bound. Modify the perceptron convergence proof to handle this case.

4.6 Improved Generalization: Voting and Averaging

In the beginning of this chapter, there was a comment that the perceptron works amazingly well. This was a half-truth. The “vanilla” perceptron algorithm does well, but not *amazingly* well. In order to make it more competitive with other learning algorithms, you need to modify it a bit to get better generalization. The key issue with the vanilla perceptron is that *it counts later points more than it counts earlier points*.

To see why, consider a data set with 10,000 examples. Suppose that after the first 100 examples, the perceptron has learned a really

Why does the perceptron convergence bound not contradict the earlier claim that poorly ordered data points (e.g., all positives followed by all negatives) will cause the perceptron to take an astronomically long time to learn?

good classifier. It's so good that it goes over the next 9899 examples without making *any* updates. It reaches the 10,000th example and makes an error. It updates. For all we know, the update on this 10,000th example *completely ruins* the weight vector that has done so well on 99.99% of the data!

What we would like is for weight vectors that “survive” a long time to get more say than weight vectors that are overthrown quickly. One way to achieve this is by **voting**. As the perceptron learns, it remembers how long each hyperplane survives. At test time, each hyperplane encountered during training “votes” on the class of a test example. If a particular hyperplane survived for 20 examples, then it gets a vote of 20. If it only survived for one example, it only gets a vote of 1. In particular, let $(\mathbf{w}, b)^{(1)}, \dots, (\mathbf{w}, b)^{(K)}$ be the $K + 1$ weight vectors encountered during training, and $c^{(1)}, \dots, c^{(K)}$ be the survival times for each of these weight vectors. (A weight vector that gets immediately updated gets $c = 1$; one that survives another round gets $c = 2$ and so on.) Then the prediction on a test point is:

$$\hat{y} = \text{sign} \left(\sum_{k=1}^K c^{(k)} \text{sign} \left(\mathbf{w}^{(k)} \cdot \hat{\mathbf{x}} + b^{(k)} \right) \right) \quad (4.17)$$

This algorithm, known as the **voted perceptron** works quite well in practice, and there is some nice theory showing that it is guaranteed to generalize better than the vanilla perceptron. Unfortunately, it is also completely impractical. If there are 1000 updates made during perceptron learning, the voted perceptron requires that you store 1000 weight vectors, together with their counts. This requires an absurd amount of storage, and makes prediction 1000 times slower than the vanilla perceptron.

A much more practical alternative is the **averaged perceptron**. The idea is similar: you maintain a collection of weight vectors and survival times. However, at test time, you predict according to the *average* weight vector, rather than the voting. In particular, the prediction is:

$$\hat{y} = \text{sign} \left(\sum_{k=1}^K c^{(k)} \left(\mathbf{w}^{(k)} \cdot \hat{\mathbf{x}} + b^{(k)} \right) \right) \quad (4.18)$$

The only difference between the voted prediction, Eq (4.17), and the averaged prediction, Eq (4.18), is the presense of the interior sign operator. With a little bit of algebra, we can rewrite the test-time prediction as:

$$\hat{y} = \text{sign} \left(\left(\sum_{k=1}^K c^{(k)} \mathbf{w}^{(k)} \right) \cdot \hat{\mathbf{x}} + \sum_{k=1}^K c^{(k)} b^{(k)} \right) \quad (4.19)$$

The advantage of the averaged perceptron is that we can simply maintain a *running sum* of the averaged weight vector (the blue term)

? The *training algorithm* for the voted perceptron is the same as the vanilla perceptron. In particular, in line 5 of Algorithm 4.2, the activation on a training example is computed based on the *current weight vector*, not based on the voted prediction. Why?

Algorithm 7 AVERAGEDPERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```

1:  $\mathbf{w} \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $b \leftarrow 0$  // initialize weights and bias
2:  $\mathbf{u} \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $\beta \leftarrow 0$  // initialize cached weights and bias
3:  $c \leftarrow 1$  // initialize example counter to one
4: for  $iter = 1 \dots MaxIter$  do
5:   for all  $(x, y) \in \mathbf{D}$  do
6:     if  $y(\mathbf{w} \cdot \mathbf{x} + b) \leq 0$  then
7:        $\mathbf{w} \leftarrow \mathbf{w} + y \mathbf{x}$  // update weights
8:        $b \leftarrow b + y$  // update bias
9:        $\mathbf{u} \leftarrow \mathbf{u} + y c \mathbf{x}$  // update cached weights
10:       $\beta \leftarrow \beta + y c$  // update cached bias
11:     end if
12:    $c \leftarrow c + 1$  // increment counter regardless of update
13: end for
14: end for
15: return  $\mathbf{w} - \frac{1}{c} \mathbf{u}$ ,  $b - \frac{1}{c} \beta$  // return averaged weights and bias

```

and averaged bias (the red term). Test-time prediction is then just as efficient as it is with the vanilla perceptron.

The full training algorithm for the averaged perceptron is shown in Algorithm 4.6. Some of the notation is changed from the original perceptron: namely, vector operations are written as vector operations, and the activation computation is folded into the error checking.

It is probably not immediately apparent from Algorithm 4.6 that the computation unfolding is precisely the calculation of the averaged weights and bias. The most *natural* implementation would be to keep track of an averaged weight vector \mathbf{u} . At the end of every example, you would increase $\mathbf{u} \leftarrow \mathbf{u} + \mathbf{w}$ (and similarly for the bias). However, such an implementation would require that you updated the averaged vector on *every* example, rather than just on the examples that were incorrectly classified! Since we hope that eventually the perceptron learns to do a good job, we would hope that it will not make updates on every example. So, ideally, you would like to only update the averaged weight vector when the actual weight vector changes. The slightly clever computation in Algorithm 4.6 achieves this.

The averaged perceptron is almost always better than the perceptron, in the sense that it generalizes better to test data. However, that does not free you from having to do **early stopping**. It will, eventually, overfit.

By writing out the computation of the averaged weights from Eq (4.18) as a telescoping sum, derive the computation from Algorithm 4.6.

4.7 Limitations of the Perceptron

Although the perceptron is very useful, it is fundamentally limited in a way that neither decision trees nor KNN are. Its limitation is that

its decision boundaries can *only* be linear. The classic way of showing this limitation is through the XOR problem (XOR = exclusive or). The XOR problem is shown graphically in Figure 4.12. It consists of four data points, each at a corner of the unit square. The labels for these points are the same, along the diagonals. You can try, but you will not be able to find a linear decision boundary that perfectly separates these data points.

One question you might ask is: do XOR-like problems exist in the real world? Unfortunately for the perceptron, the answer is yes. Consider a sentiment classification problem that has three features that simply say whether a given word is contained in a review of a course. These features are: **EXCELLENT**, **TERRIBLE** and **NOT**. The **EXCELLENT** feature is indicative of positive reviews and the **TERRIBLE** feature is indicative of negative reviews. But in the presence of the **NOT** feature, this categorization flips.

One way to address this problem is by adding **feature combinations**. We could add two additional features: **EXCELLENT-AND-NOT** and **TERRIBLE-AND-NOT** that indicate a conjunction of these base features. By assigning weights as follows, you can achieve the desired effect:

$$\begin{aligned} w_{\text{EXCELLENT}} &= +1 & w_{\text{TERRIBLE}} &= -1 & w_{\text{NOT}} &= 0 \\ w_{\text{EXCELLENT-AND-NOT}} &= -2 & w_{\text{TERRIBLE-AND-NOT}} &= +2 \end{aligned}$$

In this particular case, we have addressed the problem. However, if we start with D -many features, if we want to add all pairs, we'll blow up to $\binom{D}{2} = \mathcal{O}(D^2)$ features through this **feature mapping**. And there's no guarantee that pairs of features is enough. We might need triples of features, and now we're up to $\binom{D}{3} = \mathcal{O}(D^3)$ features. These additional features will drastically increase computation and will often result in a stronger propensity to overfitting.

In fact, the "XOR problem" is so significant that it basically killed research in classifiers with linear decision boundaries for a decade or two. Later in this book, we will see two alternative approaches to taking key ideas from the perceptron and generating classifiers with non-linear decision boundaries. One approach is to combine multiple perceptrons in a single framework: this is the **neural networks** approach (see Chapter 10). The second approach is to find computationally efficient ways of doing feature mapping in a computationally and statistically efficient way: this is the **kernels** approach (see Chapter 11).

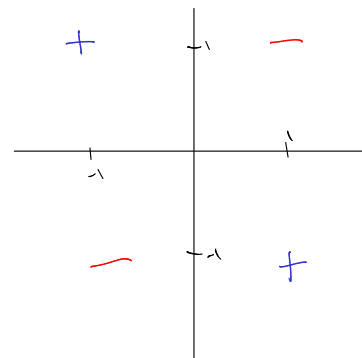


Figure 4.12: picture of xor problem

? Suppose that you took the XOR problem and added one new feature: $x_3 = x_1 \wedge x_2$ (the logical and of the two existing features). Write out feature weights and a bias that would achieve perfect classification on this data.

4.8 Further Reading

TODO further reading

5 | PRACTICAL ISSUES

A ship in port is safe, but that is not what ships are for.

Grace Hopper

AT THIS POINT, you have seen three qualitatively different models for learning: decision trees, nearest neighbors, and perceptrons. You have also learned about clustering with the K-means algorithm. You will shortly learn about more complex models, most of which are variants on things you already know. However, before attempting to understand more complex models of learning, it is important to have a firm grasp on how to use machine learning in practice. This chapter is all about how to go from an abstract learning problem to a concrete implementation. You will see some examples of “best practices” along with justifications of these practices.

In many ways, going from an abstract problem to a concrete learning task is more of an art than a science. However, this art can have a huge impact on the practical performance of learning systems. In many cases, moving to a more complicated learning algorithm will gain you a few percent improvement. Going to a better representation will gain you an order of magnitude improvement. To this end, we will discuss several high level ideas to help you develop a better artistic sensibility.

5.1 The Importance of Good Features

Machine learning is magical. You give it data and it manages to classify that data. For many, it can actually outperform a human! But, like so many problems in the world, there is a significant “garbage in, garbage out” aspect to machine learning. If the data you give it is trash, the learning algorithm is unlikely to be able to overcome it.

Consider a problem of object recognition from images. If you start with a 100×100 pixel image, a very easy feature representation of this image is as a 30,000 dimensional vector, where each dimension corresponds to the red, green or blue component of some pixel in the image. So perhaps feature 1 is the amount of red in pixel (1,1); feature 2 is the amount of green in that pixel; and so on. This is the **pixel representation** of images.

Learning Objectives:

- Translate between a problem description and a concrete learning problem.
- Perform basic feature engineering on image and text data.
- Explain how to use cross-validation to tune hyperparameters and estimate future performance.
- Compare and contrast the differences between several evaluation metrics.
- Explain why feature combinations are important for learning with some models but not others.
- Explain the relationship between the three learning techniques you have seen so far.
- Apply several debugging techniques to learning algorithms.

Dependencies: Chapter 1, Chapter 3, Chapter 4

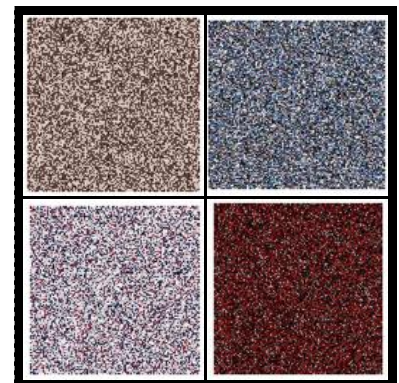


Figure 5.1: object recognition in pixels

One thing to keep in mind is that the pixel representation *throws away* all locality information in the image. Learning algorithms don’t care about features: they only care about feature values. So I can *permute* all of the features, with no effect on the learning algorithm (so long as I apply the same permutation to all training and test examples). Figure 5.1 shows some images whose pixels have been randomly permuted (in this case only the pixels are permuted, not the colors). All of these objects are things that you’ve seen plenty of examples of; can you identify them? Should you expect a machine to be able to?

An alternative representation of images is the **patch representation**, where the unit of interest is a small rectangular block of an image, rather than a single pixel. Again, permuting the patches has no effect on the classifier. Figure 5.2 shows the same images in patch representation. Can you identify them? A final representation is a **shape representation**. Here, we throw out all color and pixel information and simply provide a bounding polygon. Figure 5.3 shows the same images in this representation. Is this now enough to identify them? (If not, you can find the answers in Figure 5.15 at the end of the chapter.)

In the context of **text categorization** (for instance, the sentiment recognition task), one standard representation is the **bag of words** representation. Here, we have one feature for each unique word that appears in a document. For the feature **HAPPY**, the feature value is the number of times that the word “happy” appears in the document. The bag of words (BOW) representation throws away all position information. Table 5.1 shows a BOW representation for two chapters of this book. Can you tell which is which?

5.2 Irrelevant and Redundant Features

One big difference between learning models is how robust they are to the addition of noisy or irrelevant features. Intuitively, an irrelevant feature is one that is completely uncorrelated with the prediction task. A feature f whose expectation does not depend on the label $\mathbb{E}[f \mid Y] = \mathbb{E}[f]$ might be irrelevant. For instance, the presence of the word “the” might be largely irrelevant for predicting whether a course review is positive or negative.

A secondary issue is how well these algorithms deal with **redundant features**. Two features are redundant if they are highly correlated, regardless of whether they are correlated with the task or not. For example, having a bright red pixel in an image at position (20,93) is probably highly redundant with having a bright red pixel at position (21,93). Both might be useful (e.g., for identifying fire hy-

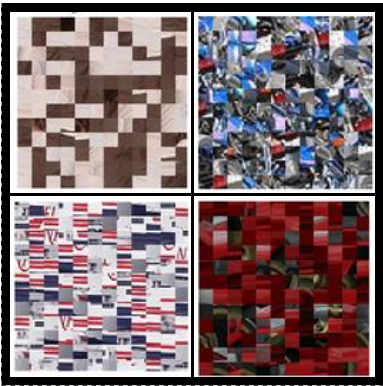


Figure 5.2: object recognition in patches

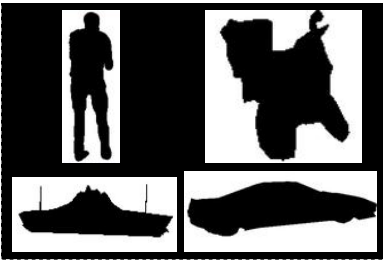


Figure 5.3: object recognition in shapes

data	learning	data	knn
training	set	dimensions	
predict	fea-	points	fea-
ture	function	ture	figure
test	machine	decision	fea-
loss	alice	treetures	point
guess	features	fig	training
algorithm	set	space	
		examples	

Table 5.1: Bag of (most frequent) words representation for the Decision Tree and KNN chapters of this book, after dropping high frequency words like “the”.

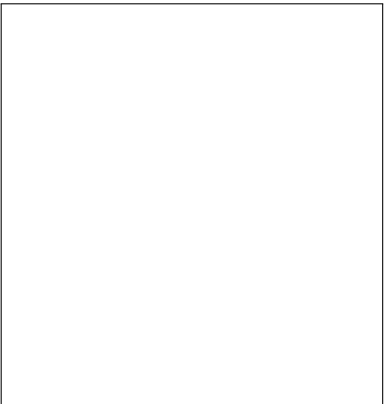


Figure 5.4: prac:baw: BOW repr of one positive and one negative review

drants), but because of how images are structured, these two features are likely to co-occur frequently.

When thinking about robustness to irrelevant or redundant features, it is usually not worthwhile thinking of the case where one has 999 great features and 1 bad feature. The interesting case is when the bad features outnumber the good features, and often outnumber by a large degree. The question is how robust are algorithms in this case.¹

For shallow **decision trees**, the model explicitly selects features that are highly correlated with the label. In particular, by limiting the depth of the decision tree, one can at least *hope* that the model will be able to throw away irrelevant features. Redundant features are almost certainly thrown out: once you select one feature, the second feature now looks mostly useless. The only possible issue with irrelevant features is that even though they're irrelevant, they *happen to* correlate with the class label on the training data, but chance.

As a thought experiment, suppose that we have N training examples, and exactly half are positive examples and half are negative examples. Suppose there's some binary feature, f , that is completely uncorrelated with the label. This feature has a 50/50 chance of appearing in any example, regardless of the label. In principle, the decision tree should *not* select this feature. But, by chance, especially if N is small, the feature might *look* correlated with the label. This is analogous to flipping two coins simultaneously N times. Even though the coins are independent, it's entirely possible that you will observe a sequence like $(H, H), (T, T), (H, H), (H, H)$, which makes them look entirely correlated! The hope is that as N grows, this becomes less and less likely. In fact, we can explicitly compute how likely this is to happen.

To do this, let's fix the sequence of N labels. We now flip a coin N times and consider how likely it is that it exactly matches the label. This is easy: the probability is 0.5^N . Now, we would also be confused if it exactly matched *not* the label, which has the same probability. So the chance that it looks perfectly correlated is $0.5^N + 0.5^N = 0.5^{N-1}$. Thankfully, this shrinks down very small (e.g., 10^{-6}) after only 21 data points, meaning that even with a very small training set, the chance that a random feature *happens* to correlate exactly with the label is tiny.

This makes us happy. The problem is that we don't have one irrelevant feature: we have many! If we randomly pick two irrelevant features, each has the same probability of perfectly correlating: 0.5^{N-1} . But since there are two and they're independent coins, the chance that *either* correlates perfectly is $2 \times 0.5^{N-1} = 0.5^{N-2}$. In general, if we have K irrelevant features, all of which are random independent coins, the chance that at least one of them perfectly correlates is

? Is it possible to have a feature f whose expectation does not depend on the label, but is nevertheless still useful for prediction?

¹ You might think it's absurd to have so many irrelevant features, but the cases you've seen so far (bag of words, bag of pixels) are both reasonable examples of this! How many words, out of the entire English vocabulary (roughly 10,000 – 100,000 words), are actually useful for predicting positive and negative course reviews?

0.5^{N-K} . This suggests that if we have a sizeable number K of irrelevant features, we'd better have at *least* $K + 21$ training examples.

Unfortunately, the situation is actually worse than this. In the above analysis we only considered the case of *perfect* correlation. We could also consider the case of *partial* correlation, which would yield even higher probabilities. Suffice it to say that even decision trees can become confused.

In the case of **K-nearest neighbors**, the situation is perhaps more dire. Since KNN weighs each feature just as much as another feature, the introduction of irrelevant features can completely mess up KNN prediction. In fact, as you saw, in high dimensional space, randomly distributed points all look approximately the same distance apart. If we add lots and lots of randomly distributed features to a data set, then all distances still converge.

In the case of the **perceptron**, one can *hope* that it might learn to assign zero weight to irrelevant features. For instance, consider a binary feature is randomly one or zero independent of the label. If the perceptron makes just as many updates for positive examples as for negative examples, there is a reasonable chance this feature weight will be zero. At the very least, it should be small.

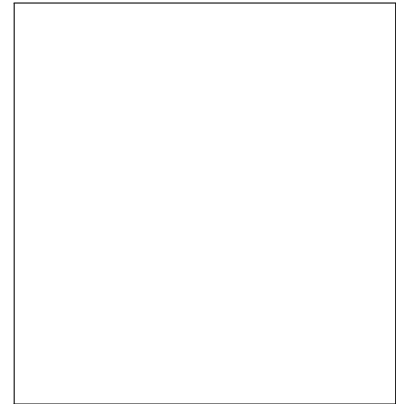


Figure 5.5: `prac:addirrel`: data from high dimensional warning, interpolated

? What happens with the perceptron with truly redundant features (i.e., one is literally a copy of the other)?

5.3 Feature Pruning and Normalization

In text categorization problems, some words simply do not appear very often. Perhaps the word “groovy”² appears in exactly one training document, which is positive. Is it really worth keeping this word around as a feature? It’s a dangerous endeavor because it’s hard to tell with just one training example if it is really correlated with the positive class, or is it just noise. You could hope that your learning algorithm is smart enough to figure it out. Or you could just remove it. That means that (a) the learning algorithm won’t have to figure it out, and (b) you’ve reduced the number of dimensions you have, so things should be more efficient, and less “scary.”

This idea of feature pruning is very useful and applied in many applications. It is easiest in the case of binary features. If a binary feature only appears some small number K times (in the training data: no fair looking at the test data!), you simply remove it from consideration. (You might also want to remove features that appear in all-but- K many documents, for instance the word “the” appears in pretty much every English document ever written.) Typical choices for K are 1, 2, 5, 10, 20, 50, mostly depending on the size of the data. On a text data set with 1000 documents, a cutoff of 5 is probably reasonable. On a text data set the size of the web, a cut of 50 or even 100 or 200 is probably reasonable³. Figure 5.6 shows the effect of

² This is typically positive indicator, or at least it was back in the US in the 1970s.

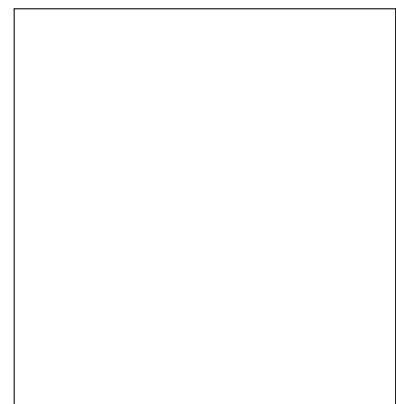


Figure 5.6: `prac:pruning`: effect of pruning on text data

³ According to Google, the following words (among many others) appear 200 times on the web: moudlings, agag-gagctg, setgravity, rogov, prosomeric, spunlaid, piyushtwok, telelesson, nes-mysl, brighnasa. For comparison, the word “the” appears 19,401,194,714 (19 billion) times.

MATH REVIEW | DATA STATISTICS: MEANS AND VARIANCES

We often need to discuss various statistics of a data set. Most often, it is enough to consider univariate (one-dimensional) data. Suppose we have N real valued numbers z_1, z_2, \dots, z_N . The **sample mean** (or just **mean**) of these numbers is just their average value, or expected value: $\mu = \frac{1}{N} \sum_n z_n$. The **sample variance** (or just **variance**) measures how much they vary around their mean: $\sigma^2 = \frac{1}{N-1} \sum_n (z_n - \mu)^2$, where μ is the sample mean.

The mean and variance have convenient interpretations in terms of prediction. Suppose we wanted to choose a single constant value to “predict” the next z , and were minimizing squared error. Call this constant value a . Then $a = \operatorname{argmin}_{a \in \mathbb{R}} \frac{1}{2} \sum_n (a - z_n)^2$. (Here, the $\frac{1}{2}$ is for convenience and does not change the answer.) To solve for a , we can take derivatives and set to zero: $\frac{\partial}{\partial a} \frac{1}{2} \sum_n (a - z_n)^2 = \sum_n (a - z_n) = Na - \sum_n z_n$; therefore $Na = \sum_n z_n$ and $a = \mu$. This means that the sample mean is the number that minimizes squared error to the sample. Moreover, the variance is proportional to the squared error of that “predictor.”

pruning on a sentiment analysis task. In the beginning, pruning does not hurt (and sometimes helps!) but eventually we prune away all the interesting words and performance suffers.

In the case of real-valued features, the question is how to extend the idea of “does not occur much” to real values. A reasonable definition is to look for features with *low variance*. In fact, for binary features, ones that almost never appear or almost always appear will also have low variance. Figure 5.8 shows the result of pruning low-variance features on the digit recognition task. Again, at first pruning does not hurt (and sometimes helps!) but eventually we have thrown out all the useful features.

It is often useful to **normalize** the data so that it is consistent in some way. There are two basic types of normalization: **feature normalization** and **example normalization**. In feature normalization, you go through each feature and adjust it the same way across all examples. In example normalization, each example is adjusted individually.

The goal of both types of normalization is to make it *easier* for your learning algorithm to learn. In feature normalization, there are two standard things to do:

1. Centering: moving the entire data set so that it is centered around the origin.
2. Scaling: rescaling each feature so that one of the following holds:
 - (a) Each feature has variance 1 across the training data.

Figure 5.7:



Figure 5.8: prac:variance: effect of pruning on vision

Earlier we discussed the problem of *scale* of features (e.g., millimeters versus centimeters). Does this have an impact on variance-based feature pruning?

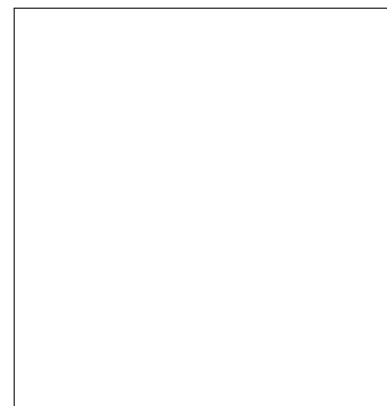


Figure 5.9: prac:transform: picture

- (b) Each feature has maximum absolute value 1 across the training data.

These transformations are shown geometrically in Figure 5.9. The goal of centering is to make sure that no features are arbitrarily large. The goal of scaling is to make sure that all features have roughly the same scale (to avoid the issue of centimeters versus millimeters).

These computations are fairly straightforward. Here, $x_{n,d}$ refers to the d th feature of example n . Since it is very rare to apply scaling without previously applying centering, the expressions below for scaling assume that the data is already centered.

$$\text{Centering:} \quad x_{n,d} \leftarrow x_{n,d} - \mu_d \quad (5.1)$$

$$\text{Variance Scaling:} \quad x_{n,d} \leftarrow x_{n,d} / \sigma_d \quad (5.2)$$

$$\text{Absolute Scaling:} \quad x_{n,d} \leftarrow x_{n,d} / r_d \quad (5.3)$$

$$\text{where:} \quad \mu_d = \frac{1}{N} \sum_n x_{n,d} \quad (5.4)$$

$$\sigma_d = \sqrt{\frac{1}{N-1} \sum_n (x_{n,d} - \mu_d)^2} \quad (5.5)$$

$$r_d = \max_n |x_{n,d}| \quad (5.6)$$

In practice, if the dynamic range of your features is already some subset of $[-2, 2]$ or $[-3, 3]$, then it is probably not worth the effort of centering and scaling. (It's an effort because you have to keep around your centering and scaling calculations so that you can apply them to the test data as well!) However, if some of your features are orders of magnitude larger than others, it might be helpful. Remember that you might know best: if the difference in scale is actually significant for your problem, then rescaling might throw away useful information.

One thing to be wary of is centering binary data. In many cases, binary data is very *sparse*: for a given example, only a few of the features are “on.” For instance, out of a vocabulary of 10,000 or 100,000 words, a given document probably only contains about 100. From a storage and computation perspective, this is very useful. However, after centering, the data will no longer be sparse and you will pay dearly with outrageously slow implementations.

In **example normalization**, you view examples one at a time. The most standard normalization is to ensure that the length of each example vector is one: namely, each example lies somewhere on the unit hypersphere. This is a simple transformation:

$$\text{Example Normalization:} \quad \mathbf{x}_n \leftarrow \mathbf{x}_n / \|\mathbf{x}_n\| \quad (5.7)$$

This transformation is depicted in Figure 5.10.

For the three models you know about (KNN, DT, Perceptron), which are most sensitive to centering? Which are most sensitive to scaling?

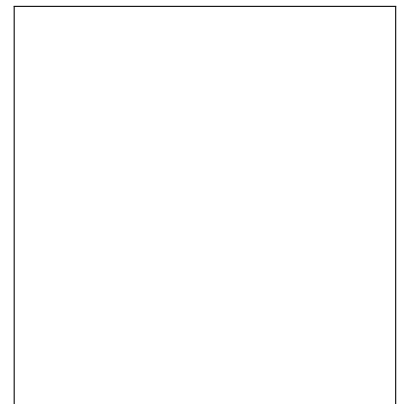


Figure 5.10: `prac:exnorm`: example of example normalization

The main advantage to example normalization is that it makes comparisons more straightforward across data sets. If I hand you two data sets that differ only in the norm of the feature vectors (i.e., one is just a scaled version of the other), it is difficult to compare the learned models. Example normalization makes this more straightforward. Moreover, as you saw in the perceptron convergence proof, it is often just mathematically easier to assume normalized data.

5.4 Combinatorial Feature Explosion

You learned in Chapter 4 that linear models (like the perceptron) cannot solve the XOR problem. You also learned that by performing a combinatorial feature explosion, they could. But that came at the computational expense of gigantic feature vectors.

Of the algorithms that you've seen so far, the perceptron is the one that has the most to gain by feature combination. And the decision tree is the one that has the least to gain. In fact, the decision tree construction is essentially building meta features for you. (Or, at least, it is building meta features constructed purely through “logical ands.”)

This observation leads to a heuristic for constructing meta features *for perceptrons from decision trees*. The idea is to train a decision tree on the training data. From that decision tree, you can extract meta features by looking at feature combinations along branches. You can then add *only* those feature combinations as meta features to the feature set for the perceptron. Figure 5.11 shows a small decision tree and a set of meta features that you might extract from it. There is a hyperparameter here of what length paths to extract from the tree: in this case, only paths of length two are extracted. For bigger trees, or if you have more data, you might benefit from longer paths.

In addition to combinatorial transformations, the **logarithmic transformation** can be quite useful in practice. It seems like a strange thing to be useful, since it doesn't seem to fundamentally change the data. However, since many learning algorithms operate by linear operations on the features (both perceptron and KNN do this), the log-transform is a way to get product-like operations. The question is which of the following feels more applicable to your data: (1) every time this feature increases by one, I'm equally more likely to predict a positive label; (2) every time this feature doubles, I'm equally more likely to predict a positive label. In the first case, you should stick with linear features and in the second case you should switch to a log-transform. This is an important transformation in text data, where the presence of the word “excellent” once is a good indicator of a positive review; seeing “excellent” twice is a better indicator;

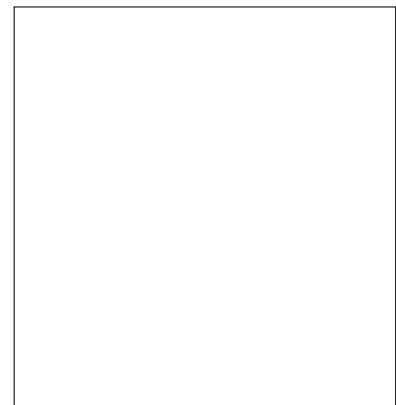


Figure 5.11: `prac:dttoperc`: turning a DT into a set of meta features

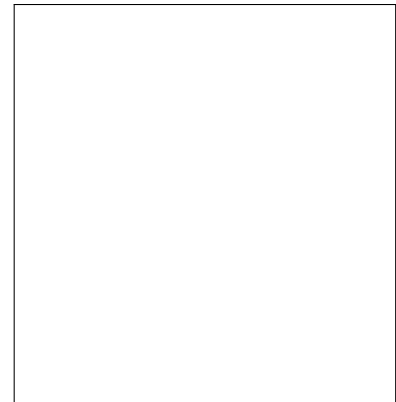


Figure 5.12: `prac:log`: performance on text categorization with word counts versus log word counts

but the difference between seeing “excellent” 10 times and seeing it 11 times really isn’t a big deal any more. A log-transform achieves this. Experimentally, you can see the difference in test performance between word count data and log-word count data in Figure 5.12. Here, the transformation is actually $x_d \mapsto \log_2(x_d + 1)$ to ensure that zeros remain zero and sparsity is retained. In the case that feature values can also be negative, the slightly more complex mapping $x_d \mapsto \log_2(|x_d| + 1)\text{sign}(x_d)$, where $\text{sign}(x_d)$ denotes the sign of x_d .

5.5 Evaluating Model Performance

So far, our focus has been on classifiers that achieve *high accuracy*. In some cases, this is not what you might want. For instance, if you are trying to predict whether a patient has cancer or not, it might be better to err on one side (saying they have cancer when they don’t) than the other (because then they die). Similarly, letting a little spam slip through might be better than accidentally blocking one email from your boss.

There are two major types of binary classification problems. One is “X versus Y.” For instance, positive versus negative sentiment. Another is “X versus not-X.” For instance, spam versus non-spam. (The argument being that there are lots of types of non-spam.) Or in the context of web search, relevant document versus irrelevant document. This is a subtle and subjective decision. But “X versus not-X” problems often have more of the feel of “X spotting” rather than a true distinction between X and Y. (Can you spot the spam? can you spot the relevant documents?)

For spotting problems (X versus not-X), there are often more appropriate success metrics than accuracy. A very popular one from information retrieval is the **precision/recall** metric. Precision asks the question: of all the X’s that you found, how many of them were actually X’s? Recall asks: of all the X’s that were out there, how many of them did you find?⁴ Formally, precision and recall are defined as:

$$P = \frac{I}{S} \quad (5.8)$$

$$R = \frac{I}{T} \quad (5.9)$$

$$S = \text{number of Xs that your system found} \quad (5.10)$$

$$T = \text{number of Xs in the data} \quad (5.11)$$

$$I = \text{number of correct Xs that your system found} \quad (5.12)$$

Here, *S* is mnemonic for “System,” *T* is mnemonic for “Truth” and *I* is mnemonic for “Intersection.” It is generally accepted that $0/0 = 1$ in these definitions. Thus, if you system found nothing, your preci-

⁴ A colleague make the analogy to the US court system’s saying “Do you promise to tell the whole truth and nothing but the truth?” In this case, the “whole truth” means high recall and “nothing but the truth” means high precision.”

sion is always perfect; and if there is nothing to find, your recall is always perfect.

Once you can compute precision and recall, you are often able to produce **precision/recall curves**. Suppose that you are attempting to identify spam. You run a learning algorithm to make predictions on a test set. But instead of just taking a “yes/no” answer, you allow your algorithm to produce its confidence. For instance, in perceptron, you might use the distance from the hyperplane as a confidence measure. You can then sort all of your test emails according to this ranking. You may put the most spam-like emails at the top and the least spam-like emails at the bottom, like in Figure 5.13.

Once you have this sorted list, you can choose how aggressively you want your spam filter to be by setting a threshold *anywhere* on this list. One would hope that if you set the threshold very high, you are likely to have high precision (but low recall). If you set the threshold very low, you’ll have high recall (but low precision). By considering *every possible* place you could put this threshold, you can trace out a curve of precision/recall values, like the one in Figure 5.14. This allows us to ask the question: for some fixed precision, what sort of recall can I get. Obviously, the closer your curve is to the upper-right corner, the better. And when comparing learning algorithms A and B you can say that A **dominates** B if A’s precision/recall curve is always higher than B’s.

Precision/recall curves are nice because they allow us to visualize many ways in which we could use the system. However, sometimes we like to have a *single number* that informs us of the quality of the solution. A popular way of combining precision and recall into a single number is by taking their harmonic mean. This is known as the balanced f-measure (or f-score):

$$F = \frac{2 \times P \times R}{P + R} \quad (5.13)$$

The reason that you want to use a harmonic mean rather than an arithmetic mean (the one you’re more used to) is that it favors systems that achieve roughly equal precision and recall. In the extreme case where $P = R$, then $F = P = R$. But in the imbalanced case, for instance $P = 0.1$ and $R = 0.9$, the overall f-measure is a modest 0.18. Table 5.2 shows f-measures as a function of precision and recall, so that you can see how important it is to get balanced values.

In some cases, you might believe that precision is more important than recall. This idea leads to the *weighted* f-measure, which is parameterized by a weight $\beta \in [0, \infty)$ (beta):

$$F_\beta = \frac{(1 + \beta^2) \times P \times R}{\beta^2 \times P + R} \quad (5.14)$$

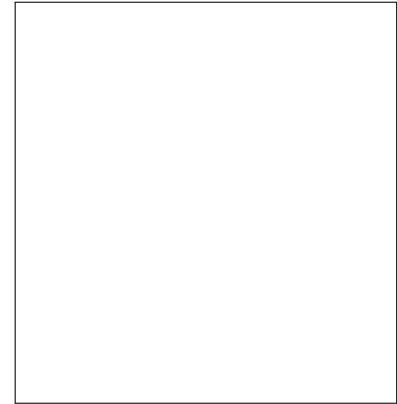


Figure 5.13: prac:spam: show a bunch of emails spam/nospam sorted by model prediction, not perfect

? How would you get a confidence out of a decision tree or KNN?

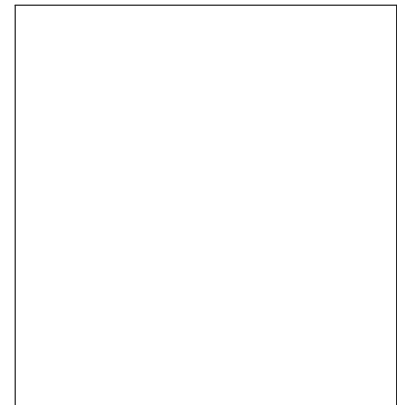


Figure 5.14: prac:prcurve: precision recall curve

	0.0	0.2	0.4	0.6	0.8	1.0
0.0	0.00	0.00	0.00	0.00	0.00	0.00
0.2	0.00	0.20	0.26	0.30	0.32	0.33
0.4	0.00	0.26	0.40	0.48	0.53	0.57
0.6	0.00	0.30	0.48	0.60	0.68	0.74
0.8	0.00	0.32	0.53	0.68	0.80	0.88
1.0	0.00	0.33	0.57	0.74	0.88	1.00

Table 5.2: Table of f-measures when varying precision and recall values.

For $\beta = 1$, this reduces to the standard f-measure. For $\beta = 0$, it focuses entirely on recall and for $\beta \rightarrow \infty$ it focuses entirely on precision. The interpretation of the weight is that F_β measures the performance for a user who cares β times as much about precision as about recall.

One thing to keep in mind is that precision and recall (and hence f-measure) depend crucially on which class is considered the thing you wish to find. In particular, if you take a binary data set and flip what it means to be a positive or negative example, you will end up with completely different precision and recall values. It is *not* the case that precision on the flipped task is equal to recall on the original task (nor vice versa). Consequently, f-measure is also not the same. For some tasks where people are less sure about what they want, they will occasionally report two sets of precision/recall/f-measure numbers, which vary based on which class is considered the thing to spot.

There are other standard metrics that are used in different communities. For instance, the medical community is fond of the **sensitivity/specificity** metric. A sensitive classifier is one which almost always finds everything it is looking for: it has high recall. In fact, sensitivity is exactly the same as recall. A specific classifier is one which does a good job *not* finding the things that it doesn't want to find. Specificity is precision on the negation of the task at hand.

You can compute curves for sensitivity and specificity much like those for precision and recall. The typical plot, referred to as the **receiver operating characteristic** (or **ROC curve**) plots the sensitivity against $1 - \text{specificity}$. Given an ROC curve, you can compute the **area under the curve** (or **AUC**) metric, which also provides a meaningful single number for a system's performance. Unlike f-measures, which tend to be low because they require agreement, AUC scores tend to be very high, even for not great systems. This is because random chance will give you an AUC of 0.5 and the best possible AUC is 1.0.

The main message for evaluation metrics is that you should choose whichever one makes the most sense. In many cases, several might make sense. In that case, you should do whatever is more commonly done in your field. There is no reason to be an outlier without cause.

5.6 Cross Validation

In Chapter 1, you learned about using development data (or held-out data) to set hyperparameters. The main disadvantage to the development data approach is that you throw out some of your training data, just for estimating one or two hyperparameters.

Algorithm 8 `CROSSVALIDATE`(*LearningAlgorithm*, *Data*, *K*)

```

1:  $\hat{\epsilon} \leftarrow \infty$  // store lowest error encountered so far
2:  $\hat{\alpha} \leftarrow \text{unknown}$  // store the hyperparameter setting that yielded it
3: for all hyperparameter settings  $\alpha$  do
4:    $err \leftarrow []$  // keep track of the  $K$ -many error estimates
5:   for  $k = 1$  to  $K$  do
6:      $train \leftarrow \{(x_n, y_n) \in Data : n \bmod K \neq k - 1\}$ 
7:      $test \leftarrow \{(x_n, y_n) \in Data : n \bmod K = k - 1\}$  // test every  $K$ th example
8:      $model \leftarrow \text{Run } LearningAlgorithm \text{ on } train$ 
9:      $err \leftarrow err \oplus \text{error of } model \text{ on } test$  // add current error to list of errors
10:   end for
11:    $avgErr \leftarrow \text{mean of set } err$ 
12:   if  $avgErr < \hat{\epsilon}$  then
13:      $\hat{\epsilon} \leftarrow avgErr$  // remember these settings
14:      $\hat{\alpha} \leftarrow \alpha$  // because they're the best so far
15:   end if
16: end for

```

An alternative is the idea of **cross validation**. In cross validation, you break your training data up into 10 equally-sized partitions. You train a learning algorithm on 9 of them and test it on the remaining 1. You do this 10 times, each time holding out a different partition as the “development” part. You can then average your performance over all ten parts to get an estimate of how well your model will perform in the future. You can repeat this process for every possible choice of hyperparameters to get an estimate of which one performs best. The general K -fold cross validation technique is shown in Algorithm 5.6, where $K = 10$ in the preceeding discussion.

In fact, the development data approach can be seen as an approximation to cross validation, wherein only one of the K loops (line 5 in Algorithm 5.6) is executed.

Typical choices for K are 2, 5, 10 and $N - 1$. By far the most common is $K = 10$: 10-fold cross validation. Sometimes 5 is used for efficiency reasons. And sometimes 2 is used for subtle statistical reasons, but that is quite rare. In the case that $K = N - 1$, this is known as **leave-one-out cross validation** (or abbreviated as **LOO cross validation**). After running cross validation, you have two choices. You can either select one of the K trained models as your final model to make predictions with, or you can train a *new* model on all of the data, using the hyperparameters selected by cross-validation. If you have the time, the latter is probably a better options.

It may seem that LOO cross validation is prohibitively expensive to run. This is true for most learning algorithms *except for* K -nearest neighbors. For KNN, leave-one-out is actually very natural. We loop through each training point and ask ourselves whether this example would be correctly classified for all different possible values of K .

Algorithm 9 KNN-TRAIN-LOO(D)

```

1:  $err_k \leftarrow 0, \forall 1 \leq k \leq N-1$            //  $err_k$  stores how well you do with  $k$ NN
2: for  $n = 1$  to  $N$  do
3:    $S_m \leftarrow \langle ||x_n - x_m||, m \rangle, \forall m \neq n$    // compute distances to other points
4:    $S \leftarrow \text{SORT}(S)$                                // put lowest-distance objects first
5:    $\hat{y} \leftarrow 0$                                      // current label prediction
6:   for  $k = 1$  to  $N-1$  do
7:      $\langle dist, m \rangle \leftarrow S_k$ 
8:      $\hat{y} \leftarrow \hat{y} + y_m$                            // let  $k$ th closest point vote
9:     if  $\hat{y} \neq y_m$  then
10:       $err_k \leftarrow err_k + 1$                        // one more error for  $k$ NN
11:     end if
12:   end for
13: end for
14: return  $\text{argmin}_k err_k$                                // return the  $K$  that achieved lowest error

```

This requires only as much computation as computing the K nearest neighbors for the highest value of K . This is such a popular and effective approach for KNN classification that it is spelled out in Algorithm 5.6.

Overall, the main advantage to cross validation over development data is robustness. The main advantage of development data is speed.

One warning to keep in mind is that the goal of both cross validation and development data is to estimate how well you will do in the future. This is a question of statistics, and holds *only if* your test data really looks like your training data. That is, it is drawn from the same distribution. In many practical cases, this is not entirely true.

For example, in person identification, we might try to classify every pixel in an image based on whether it contains a person or not. If we have 100 training images, each with 10,000 pixels, then we have a total of $1m$ training examples. The classification for a pixel in image 5 is highly dependent on the classification for a neighboring pixel in the same image. So if one of those pixels happens to fall in training data, and the other in development (or cross validation) data, your model will do unreasonably well. In this case, it is important that when you cross validate (or use development data), you do so over *images*, not over *pixels*. The same goes for text problems where you sometimes want to classify things at a word level, but are handed a collection of documents. The important thing to keep in mind is that it is the *images* (or documents) that are drawn independently from your data distribution and *not* the pixels (or words), which are drawn dependently.

5.7 Hypothesis Testing and Statistical Significance

Suppose that you've presented a machine learning solution to your boss that achieves 7% error on cross validation. Your nemesis, Gabe, gives a solution to your boss that achieves 6.9% error on cross validation. How impressed should your boss be? It depends. If this 0.1% improvement was measured over 1000 examples, perhaps not too impressed. It would mean that Gabe got exactly one more example right than you did. (In fact, they probably got 15 more right and 14 more wrong.) If this 0.1% improvement was measured over 1,000,000 examples, perhaps this is more impressive.

This is one of the most fundamental questions in statistics. You have a scientific hypothesis of the form "Gabe's algorithm is better than mine." You wish to test whether this hypothesis is true. You are testing it against the **null hypothesis**, which is that Gabe's algorithm is no better than yours. You've collected data (either 1000 or 1m data points) to measure the strength of this hypothesis. You want to ensure that the difference in performance of these two algorithms is **statistically significant**: i.e., is probably not just due to random luck. (A more common question statisticians ask is whether one drug treatment is better than another, where "another" is either a placebo or the competitor's drug.)

There are about ∞ -many ways of doing **hypothesis testing**. Like evaluation metrics and the number of folds of cross validation, this is something that is very discipline specific. Here, we will discuss two popular tests: the **paired t-test** and **bootstrapping**. These tests, and other statistical tests, have underlying assumptions (for instance, assumptions about the distribution of observations) and strengths (for instance, small or large samples). In most cases, the goal of hypothesis testing is to compute a **p-value**: namely, the probability that the observed difference in performance was by chance. The standard way of reporting results is to say something like "there is a 95% chance that this difference was not by chance." The value 95% is arbitrary, and occasionally people use weaker (90%) test or stronger (99.5%) tests.

The **t-test** is an example of a **parametric test**. It is applicable when the null hypothesis states that the difference between two responses has mean zero and unknown variance. The t-test actually assumes that data is distributed according to a Gaussian distribution, which is probably *not* true of binary responses. Fortunately, for large samples (at least a few hundred), binary samples are well approximated by a Gaussian distribution. So long as your sample is sufficiently large, the t-test is reasonable either for regression or classification problems.

Suppose that you evaluate two algorithm on N -many examples.

t	significance
≥ 1.28	90.0%
≥ 1.64	95.0%
≥ 1.96	97.5%
≥ 2.58	99.5%

Table 5.3: Table of significance values for the t-test.

On each example, you can compute whether the algorithm made the correct prediction. Let a_1, \dots, a_N denote the error of the first algorithm on each example. Let b_1, \dots, b_N denote the error of the second algorithm. You can compute μ_a and μ_b as the means of \mathbf{a} and \mathbf{b} , respectively. Finally, center the data as $\hat{\mathbf{a}} = \mathbf{a} - \mu_a$ and $\hat{\mathbf{b}} = \mathbf{b} - \mu_b$. The t-statistic is defined as:

$$t = (\mu_a - \mu_b) \sqrt{\frac{N(N-1)}{\sum_n (\hat{a}_n - \hat{b}_n)^2}} \quad (5.15)$$

After computing the t -value, you can compare it to a list of values for computing **confidence intervals**. Assuming you have a lot of data (N is a few hundred or more), then you can compare your t -value to Table 5.3 to determine the significance level of the difference.

One disadvantage to the t-test is that it cannot easily be applied to evaluation metrics like f-score. This is because f-score is a computed over an entire test set and does not decompose into a set of individual errors. This means that the t-test cannot be applied.

Fortunately, **cross validation** gives you a way around this problem. When you do K -fold cross validation, you are able to compute K error metrics over the same data. For example, you might run 5-fold cross validation and compute f-score for every fold. Perhaps the f-scores are 92.4, 93.9, 96.1, 92.2 and 94.4. This gives you an average f-score of 93.8 over the 5 folds. The standard deviation of this set of f-scores is:

$$\sigma = \sqrt{\frac{1}{N-1} \sum_n (a_i - \mu)^2} \quad (5.16)$$

$$= \sqrt{\frac{1}{4} (1.96 + 0.01 + 5.29 + 2.56 + 0.36)} \quad (5.17)$$

$$= 1.595 \quad (5.18)$$

You can now assume that the distribution of scores is approximately Gaussian. If this is true, then approximately 70% of the probability mass lies in the range $[\mu - \sigma, \mu + \sigma]$; 95% lies in the range $[\mu - 2\sigma, \mu + 2\sigma]$; and 99.5% lies in the range $[\mu - 3\sigma, \mu + 3\sigma]$. So, if we were comparing our algorithm against one whose average f-score was 90.6%, we could be 95% certain that our superior performance was not due to chance.⁵

WARNING: A confidence of 95% does not mean “There is a 95% chance that I am better.” All it means is that if I reran the same experiment 100 times, then in 95 of those experiments I would still win. These are *very* different statements. If you say the first one, people who know about statistics will get very mad at you!

One disadvantage to cross validation is that it is computationally expensive. More folds typically leads to better estimates, but every

? What does it mean for the means μ_a and μ_b to become further apart? How does this affect the t -value? What happens if the variance of \mathbf{a} increases?

⁵ Had we run 10-fold cross validation we might be been able to get tighter confidence intervals.

Algorithm 10 BOOTSTRAP-EVALUATE(\mathbf{y} , $\hat{\mathbf{y}}$, NumFolds)

```

1:  $\text{scores} \leftarrow []$ 
2: for  $k = 1$  to NumFolds do
3:    $\text{truth} \leftarrow []$                                 // list of values we want to predict
4:    $\text{pred} \leftarrow []$                                 // list of values we actually predicted
5:   for  $n = 1$  to  $N$  do
6:      $m \leftarrow$  uniform random value from 1 to  $N$     // sample a test point
7:      $\text{truth} \leftarrow \text{truth} \oplus y_m$                 // add on the truth
8:      $\text{pred} \leftarrow \text{pred} \oplus \hat{y}_m$                 // add on our prediction
9:   end for
10:   $\text{scores} \leftarrow \text{scores} \oplus \text{F-SCORE}(\text{truth}, \text{pred})$  // evaluate
11: end for
12: return ( $\text{MEAN}(\text{scores})$ ,  $\text{STDDEV}(\text{scores})$ )

```

new fold requires training a new classifier. This can get very time consuming. The technique of **bootstrapping** (and closely related idea of **jack-knifing**) can address this problem.

Suppose that you didn't want to run cross validation. All you have is a single held-out test set with 1000 data points in it. You can run your classifier and get predictions on these 1000 data points. You would like to be able to compute a metric like f-score on this test set, but also get confidence intervals. The idea behind bootstrapping is that this set of 1000 is a random draw from some distribution. We would like to get multiple random draws from this distribution on which to evaluate. We can *simulate* multiple draws by repeatedly subsampling from these 1000 examples, with replacement.

To perform a *single* bootstrap, you will sample 1000 random points from your test set of 1000 random points. This sampling must be done with replacement (so that the same example can be sampled more than once), otherwise you'll just end up with your original test set. This gives you a bootstrapped sample. On this sample, you can compute f-score (or whatever metric you want). You then do this 99 more times, to get a 100-fold bootstrap. For each bootstrapped sample, you will be a different f-score. The mean and standard deviation of this set of f-scores can be used to estimate a confidence interval for your algorithm.

The bootstrap resampling procedure is sketched in Algorithm 5.7. This takes three arguments: the true labels \mathbf{y} , the predicted labels $\hat{\mathbf{y}}$ and the number of folds to run. It returns the mean and standard deviation from which you can compute a confidence interval.

5.8 Debugging Learning Algorithms

Learning algorithms are notoriously hard to debug, as you may have already experienced if you have implemented any of the models

presented so far. The main issue is that when a learning algorithm doesn't learn, it's unclear if this is because there's a bug or because the learning problem is too hard (or there's too much noise, or ...). Moreover, sometimes bugs lead to learning algorithms performing *better* than they should: these are especially hard to catch (and always a bit disappointing when you do catch them).

In order to debug failing learning models, it is useful to revisit the notion of: where can error enter our system? In Chapter 2, we considered a typical design process for machine learning in Figure 2.4. Leaving off the top steps in that are not relevant to machine learning in particular, the basic steps that go into crafting a machine learning system are: collect data, choose features, choose model family, choose training data, train model, evaluate on test data. In each of these steps, things can go wrong. Below are some strategies for isolating the cause of error.

Is the problem with generalization to the test data? We have talked a lot about training error versus test error. In general, it's unrealistic to expect to do *better* on the test data than on the training data. Can your learning system do well on fitting the training data? If so, then the problem is in generalization (perhaps your model family is too complicated, you have too many features or not enough data). If not, then the problem is in representation (you probably need better features or better data).

Do you have train/test mismatch? If you can fit the training data, but it doesn't generalize, it could be because there's something different about your test data. Try shuffling your training data and test data together and then randomly selecting a new test set. If you do well in that condition, then probably the test distribution is strange in some way. If reselecting the test data doesn't help, you have other generalization problems.

Is your learning algorithm implemented correctly? This often means: is it optimizing what you think it's optimizing. Instead of measuring accuracy, try measuring whatever-quantity-your-algorithm-is-supposedly-optimizing (like log loss or hinge loss) and make sure that the optimizer is successfully minimizing this quantity. It is usually useful to hand-craft some datasets on which you know the desired behavior. For instance, you could run KNN on the XOR data. Or you could run perceptron on some easily linearly separable data (for instance positive points along the line $x_2 = x_1 + 1$ and negative points along the line $x_2 = x_1 - 1$). Or a decision tree on nice axis-aligned data. Finally, can you compare against a reference implementation?

Do you have an adequate representation? If you cannot even fit the training data, you might not have a rich enough feature set.

The easiest way to try to get a learning algorithm to overfit is to add a new feature to it. You can call this feature the **CHEATINGISFUN** feature. The feature value associated with this feature is +1 if this is a positive example and −1 (or zero) if this is a negative example. In other words, this feature is a *perfect* indicator of the class of this example. If you add the **CHEATINGISFUN** feature and your algorithm does not get near 0% training error, this could be because there are too many noisy features confusing it. You could either remove a lot of the other features, or make the feature value for **CHEATINGISFUN** either +100 or −100 so that the algorithm *really* looks at it. If you do this and your algorithm still cannot overfit then you likely have a bug. (Remember to remove the **CHEATINGISFUN** feature from your final implementation!) If the **CHEATINGISFUN** technique gets you near 0% error, then you need to work on better feature design or pick another learning model (e.g., decision tree versus linear model). If not, you probably don't have enough data or have too many features; try removing as many features as possible.

Do you have enough data? Try training on 80% of your training data and look at how much this hurts performance. If it hurts a lot, then getting more data is likely to help; if it only hurts a little, you might be data saturated.

5.9 Bias/Variance Trade-off

Because one of the key questions in machine learning is the question of representation, it is common to think about test error in terms of a decomposition into two terms. Let f be the learned classifier, selected from a set \mathcal{F} of “all possible classifiers using a fixed representation,” then:

$$\text{error}(f) = \underbrace{\left[\text{error}(f) - \min_{f^* \in \mathcal{F}} \text{error}(f^*) \right]}_{\text{estimation error}} + \underbrace{\left[\min_{f^* \in \mathcal{F}} \text{error}(f) \right]}_{\text{approximation error}} \quad (5.19)$$

Here, the second term, the **approximation error**, measures the quality of the model family⁶. One way of thinking of approximation error is: suppose someone gave me infinite data to train on—how well could I do with this representation? The first term, the **estimation error**, measures how far the actual learned classifier f is from the optimal classifier f^* . You can think of this term as measuring how much you have to pay for the fact that you *don't* have infinite training data.

Unfortunately, it is nearly impossible to *compute* the estimation error and approximation error, except in constructed cases. This doesn't make the decomposition useless. Decompositions like this

⁶ The “model family” (such as depth 20 decision trees, or linear classifiers) is often referred to as the **hypothesis class**. The hypothesis class \mathcal{F} denotes the set of all possible classifiers we consider, such as all linear classifiers. An classifier $f \in \mathcal{F}$ is sometimes called a **hypothesis**, though we generally avoid this latter terminology here.

are very useful for designing debugging strategies. For instance, the CheatingIsFun strategy is designed explicitly to ensure that the approximation error is zero, and therefore isolating all error into the estimation error.

There is a fundamental trade-off between estimation error and approximation error. As you make your representation more complex, you make \mathcal{F} bigger. This will typically cause a decrease in approximation error, because you can now fit more functions. But you run a risk of increasing the estimation error, because you have added more parameters to fit, and you are likely to suffer from overfitting.

The trade-off between estimation error and approximation error is often called the **bias/variance trade-off**, where “approximation error” is “bias” and “estimation error” is “variance.” To understand this connection, consider a very simple hypothesis class \mathcal{F} that only contains two functions: the always positive classifier (that returns +1 regardless of input) and the always negative classifier. Suppose you have a data generating distribution \mathcal{D} that is 60% positive examples and 40% negative examples. You draw a training set of 41 examples. There’s about a 90% chance that the majority of these training examples will be positive, so on this impoverished hypothesis class \mathcal{F} , there’s a 90% chance that it will learn the “all positive” classifier. That is: 90% of the time, regardless of the training set, the learning algorithm learns the same thing. This is *low variance* as a function of the random draw of the training set. On the other hand, the learned classifier is very insensitive to the input example (in this extreme case, it’s completely insensitive): it is *strongly biased* toward predicting +1 even if everything about the input contradicts this.

5.10 Further Reading

TODO



Figure 5.15: object recognition with full information

6 | BEYOND BINARY CLASSIFICATION

Different general classification methods can give different, but equally plausible, classifications, so you need an application context to choose among them.

– Karen Spärck-Jones

Learning Objectives:

- Represent complex prediction problems in a formal learning setting.
- Be able to artificially “balance” imbalanced data.
- Understand the positive and negative aspects of several reductions from multiclass classification to binary classification.
- Recognize the difference between regression and ordinal regression.

IN THE PRECEEDING CHAPTERS, you have learned all about a very simple form of prediction: predicting bits. In the real world, however, we often need to predict much more complex objects. You may need to categorize a document into one of several categories: sports, entertainment, news, politics, etc. You may need to rank web pages or ads based on relevance to a query. These problems are all commonly encountered, yet fundamentally more complex than binary classification.

In this chapter, you will learn how to use everything you already know about binary classification to solve these more complicated problems. You will see that it’s relatively easy to think of a binary classifier as a black box, which you can reuse for solving these more complex problems. This is a very useful abstraction, since it allows us to reuse knowledge, rather than having to build new learning models and algorithms from scratch.

Dependencies:

6.1 Learning with Imbalanced Data

Your boss tells you to build a classifier that can identify fraudulent transactions in credit card histories. Fortunately, most transactions are legitimate, so perhaps only 0.1% of the data is a positive instance. The **imbalanced data** problem refers to the fact that for a large number of real world problems, the number of positive examples is dwarfed by the number of negative examples (or vice versa). This is actually something of a misnomer: it is not the *data* that is imbalanced, but the *distribution* from which the data is drawn. (And since the distribution is imbalanced, so must the data be.)

Imbalanced data is a problem because machine learning algorithms are too smart for your own good. For most learning algorithms, if you give them data that is 99.9% negative and 0.1% positive, they will simply learn to always predict negative. Why? Because they are trying to minimize error, and they can achieve 0.1% error by doing nothing! If a teacher told you to study for an exam with 1000

true/false questions and only one of them is true, it is unlikely you will study very long.

Really, the problem is not with the data, but rather with the way that you have defined the learning problem. That is to say, what you care about is *not* accuracy: you care about something else. If you want a learning algorithm to do a reasonable job, you have to tell it what you want!

Most likely, what you want is *not* to optimize accuracy, but rather to optimize some other measure, like f-score or AUC. You want your algorithm to make *some* positive predictions, and simply prefer those to be “good.” We will shortly discuss two heuristics for dealing with this problem: subsampling and weighting. In subsampling, you *throw out* some of your negative examples so that you are left with a balanced data set (50% positive, 50% negative). This might scare you a bit since throwing out data seems like a bad idea, but at least it makes learning much more efficient. In weighting, instead of throwing out positive examples, we just give them lower weight. If you assign an **importance weight** of 0.00101 to each of the positive examples, then there will be as much *weight* associated with positive examples as negative examples.

Before formally defining these heuristics, we need to have a mechanism for formally defining supervised learning problems. We will proceed by example, using binary classification as the canonical learning problem.

TASK: BINARY CLASSIFICATION

Given:

1. An input space \mathcal{X}
2. An unknown distribution \mathcal{D} over $\mathcal{X} \times \{-1, +1\}$
3. A training set D sampled from \mathcal{D}

Compute: A function f minimizing: $\mathbb{E}_{(x,y) \sim \mathcal{D}} [f(x) \neq y]$

As in all the binary classification examples you’ve seen, you have some input space (which has always been \mathbb{R}^D). There is some distribution that produces labeled examples over the input space. You do not have access to that distribution, but can obtain samples from it. Your goal is to find a classifier that minimizes error on that distribution.

A small modification on this definition gives a α -weighted classification problem, where you believe that the positive class is α -times as

Algorithm 11 SUBSAMPLEMAP($\mathcal{D}^{\text{weighted}}, \alpha$)

```

1: while true do
2:    $(\mathbf{x}, y) \sim \mathcal{D}^{\text{weighted}}$            // draw an example from the weighted distribution
3:    $u \sim$  uniform random variable in  $[0, 1]$ 
4:   if  $y = +1$  or  $u < \frac{1}{\alpha}$  then
5:     return  $(\mathbf{x}, y)$ 
6:   end if
7: end while

```

Algorithm 12 SUBSAMPLETEST($f^{\text{BINARY}}, \hat{\mathbf{x}}$)

```

1: return  $f^{\text{BINARY}}(\hat{\mathbf{x}})$ 

```

important as the negative class.

TASK: α -WEIGHTED BINARY CLASSIFICATION

Given:

1. An input space \mathcal{X}
2. An unknown distribution \mathcal{D} over $\mathcal{X} \times \{-1, +1\}$
3. A training set D sampled from \mathcal{D}

Compute: A function f minimizing: $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\alpha^{y=1} [f(\mathbf{x}) \neq y]]$

The objects given to you in weighted binary classification are identical to standard binary classification. The only difference is that the *cost* of misprediction for $y = +1$ is α , while the cost of misprediction for $y = -1$ is 1. In what follows, we assume that $\alpha > 1$. If it is not, you can simply swap the labels and use $1/\alpha$.

The question we will ask is: suppose that I have a good algorithm for solving the BINARY CLASSIFICATION problem. Can I turn that into a good algorithm for solving the α -WEIGHTED BINARY CLASSIFICATION problem?

In order to do this, you need to define a *transformation* that maps a concrete weighted problem into a concrete unweighted problem. This transformation needs to happen both at training time and at test time (though it need not be the same transformation!). Algorithm 6.1 sketches a training-time **sub-sampling** transformation and Algorithm 6.1 sketches a test-time transformation (which, in this case, is trivial). All the training algorithm is doing is retaining all positive examples and a $1/\alpha$ fraction of all negative examples. The algorithm is explicitly turning the distribution over weighted examples into a (different) distribution over binary examples. A vanilla binary classifier

is trained on this **induced distribution**.

Aside from the fact that this algorithm throws out a lot of data (especially for large α), it does seem to be doing a reasonable thing. In fact, from a **reductions** perspective, it is an optimal algorithm. You can prove the following result:

Theorem 3 (Subsampling Optimality). *Suppose the binary classifier trained in Algorithm 6.1 achieves a binary error rate of ϵ . Then the error rate of the weighted predictor is equal to $\alpha\epsilon$.*

This theorem states that if your binary classifier does well (on the induced distribution), then the learned predictor will also do well (on the original distribution). Thus, we have successfully converted a weighted learning problem into a plain classification problem! The fact that the error rate of the weighted predictor is exactly α times more than that of the unweighted predictor is unavoidable: the error metric on which it is evaluated is α times bigger!

The proof of this theorem is so straightforward that we will prove it here. It simply involves some algebra on expected values.

Proof of Theorem 3. Let \mathcal{D}^w be the original distribution and let \mathcal{D}^b be the induced distribution. Let f be the binary classifier trained on data from \mathcal{D}^b that achieves a binary error rate of ϵ^b on that distribution. We will compute the expected error ϵ^w of f on the weighted problem:

$$\epsilon^w = \mathbb{E}_{(x,y) \sim \mathcal{D}^w} [\alpha^{y=1} [f(x) \neq y]] \quad (6.1)$$

$$= \sum_{x \in \mathcal{X}} \sum_{y \in \pm 1} \mathcal{D}^w(x, y) \alpha^{y=1} [f(x) \neq y] \quad (6.2)$$

$$= \alpha \sum_{x \in \mathcal{X}} \left(\mathcal{D}^w(x, +1) [f(x) \neq +1] + \mathcal{D}^w(x, -1) \frac{1}{\alpha} [f(x) \neq -1] \right) \quad (6.3)$$

$$= \alpha \sum_{x \in \mathcal{X}} \left(\mathcal{D}^b(x, +1) [f(x) \neq +1] + \mathcal{D}^b(x, -1) [f(x) \neq -1] \right) \quad (6.4)$$

$$= \alpha \mathbb{E}_{(x,y) \sim \mathcal{D}^b} [f(x) \neq y] \quad (6.5)$$

$$= \alpha \epsilon^b \quad (6.6)$$

And we're done! (We implicitly assumed \mathcal{X} is discrete. In the case of continuous data, you need to replace all the sums over x with integrals over x , but the result still holds.) \square

Instead of subsampling the low-cost class, you could alternatively **oversample** the high-cost class. The easiest case is when α is an integer, say 5. Now, whenever you get a positive point, you include 5 copies of it in the induced distribution. Whenever you get a negative point, you include a single copy.

? Why is it unreasonable to expect to be able to achieve, for instance, an error of $\sqrt{\alpha}\epsilon$, or anything that is sublinear in α ?

? How can you handle non-integral α , for instance 5.5?

This oversampling algorithm achieves exactly the same theoretical result as the subsampling algorithm. The main advantage to the oversampling algorithm is that it does not throw out any data. The main advantage to the subsampling algorithm is that it is more computationally efficient.

You might be asking yourself: intuitively, the oversampling algorithm seems like a much better idea than the subsampling algorithm, at least if you don't care about computational efficiency. But the theory tells us that they are the same! What is going on? Of course the theory isn't wrong. It's just that the assumptions are effectively different in the two cases. Both theorems state that if you can get error of ϵ on the binary problem, you automatically get error of $\alpha\epsilon$ on the weighted problem. But they do not say anything about how possible it is to get error ϵ on the binary problem. Since the oversampling algorithm produces more data points than the subsampling algorithm it is very conceivable that you could get lower binary error with oversampling than subsampling.

The primary drawback to oversampling is computational inefficiency. However, for many learning algorithms, it is straightforward to include *weighted* copies of data points at no cost. The idea is to store only the unique data points and maintain a counter saying how many times they are replicated. This is not easy to do for the perceptron (it can be done, but takes work), but it *is* easy for both decision trees and KNN. For example, for decision trees (recall Algorithm 1.3), the only changes are to: (1) ensure that line 1 computes the most frequent *weighted* answer, and (2) change lines 10 and 11 to compute weighted errors.

?

Modify the proof of optimality for the subsampling algorithm so that it applies to the oversampling algorithm.

?

Why is it hard to change the perceptron? (Hint: it has to do with the fact that perceptron is online.)

?

How would you modify KNN to take into account weights?

6.2 Multiclass Classification

Multiclass classification is a natural extension of binary classification. The goal is still to assign a discrete label to examples (for instance, is a document about entertainment, sports, finance or world news?). The difference is that you have $K > 2$ classes to choose from.

Algorithm 13 `ONEVERSUSALLTRAIN`($\mathbf{D}^{\text{multiclass}}$, `BINARYTRAIN`)

```

1: for  $i = 1$  to  $K$  do
2:    $\mathbf{D}^{\text{bin}} \leftarrow$  relabel  $\mathbf{D}^{\text{multiclass}}$  so class  $i$  is positive and  $\neg i$  is negative
3:    $f_i \leftarrow \text{BINARYTRAIN}(\mathbf{D}^{\text{bin}})$ 
4: end for
5: return  $f_1, \dots, f_K$ 

```

Algorithm 14 `ONEVERSUSALLTEST`($f_1, \dots, f_K, \hat{\mathbf{x}}$)

```

1:  $\text{score} \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize  $K$ -many scores to zero
2: for  $i = 1$  to  $K$  do
3:    $y \leftarrow f_i(\hat{\mathbf{x}})$ 
4:    $\text{score}_i \leftarrow \text{score}_i + y$ 
5: end for
6: return  $\text{argmax}_k \text{score}_k$ 

```

TASK: MULTICLASS CLASSIFICATION*Given:*

1. An input space \mathcal{X} and number of classes K
2. An unknown distribution \mathcal{D} over $\mathcal{X} \times [K]$
3. A training set D sampled from \mathcal{D}

Compute: A function f minimizing: $\mathbb{E}_{(x,y) \sim \mathcal{D}} [f(x) \neq y]$

Note that this is *identical* to binary classification, except for the presence of K classes. (In the above, $[K] = \{1, 2, 3, \dots, K\}$.) In fact, if you set $K = 2$ you exactly recover binary classification.

The game we play is the same: someone gives you a binary classifier and you have to use it to solve the multiclass classification problem. A very common approach is the **one versus all** technique (also called **OVA** or **one versus rest**). To perform OVA, you train K -many binary classifiers, f_1, \dots, f_K . Each classifier sees *all* of the training data. Classifier f_i receives all examples labeled class i as positives and all other examples as negatives. At test time, whichever classifier predicts “positive” wins, with ties broken randomly.

The training and test algorithms for OVA are sketched in Algorithms 6.2 and 6.2. In the testing procedure, the prediction of the i th classifier is added to the overall score for class i . Thus, if the prediction is positive, class i gets a vote; if the prediction is negative, everyone else (implicitly) gets a vote. (In fact, if your learning algorithm can output a confidence, as discussed in Section ??, you can often do better by using the confidence as y , rather than a simple ± 1 .)

OVA is quite natural and easy to implement. It also works very

?

Suppose that you have N data points in K classes, evenly divided. How long does it take to train an OVA classifier, if the base binary classifier takes $\mathcal{O}(N)$ time to train? What if the base classifier takes $\mathcal{O}(N^2)$ time?

?

Why would using a confidence help?

well in practice, so long as you do a good job choosing a good binary classification algorithm *tuning* its hyperparameters well. Its weakness is that it can be somewhat brittle. Intuitively, it is not particularly robust to errors in the underlying classifiers. If *one* classifier makes a mistake, it is possible that the entire prediction is erroneous. In fact, it is entirely possible that *none* of the K classifiers predicts positive (which is actually the worst-case scenario from a theoretical perspective)! This is made explicit in the OVA error bound below.

Theorem 4 (OVA Error Bound). *Suppose the average binary error of the K binary classifiers is ϵ . Then the error rate of the OVA multiclass predictor is at most $(K - 1)\epsilon$.*

Proof of Theorem 4. The key question is how erroneous predictions from the binary classifiers lead to multiclass errors. We break it down into false negatives (predicting -1 when the truth is $+1$) and false positives (predicting $+1$ when the truth is -1).

When a false negative occurs, then the testing procedure chooses randomly between available options, which is all labels. This gives a $(K - 1)/K$ probability of multiclass error. Since only *one* binary error is necessary to make this happen, the *efficiency* of this error mode is $[(K - 1)/K]/1 = (K - 1)/K$.

Multiple false positives can occur simultaneously. Suppose there are m false positives. If there is simultaneously a false negative, the error is 1. In order for this to happen, there have to be $m + 1$ errors, so the efficiency is $1/(m + 1)$. In the case that there is not a simultaneous false negative, the error probability is $m/(m + 1)$. This requires m errors, leading to an efficiency of $1/(m + 1)$.

The worse case, therefore, is the false negative case, which gives an efficiency of $(K - 1)/K$. Since we have K -many opportunities to err, we multiply this by K and get a bound of $(K - 1)\epsilon$. \square

The constants in this are relatively unimportant: the aspect that matters is that this scales *linearly* in K . That is, as the number of classes grows, so does your expected error.

To develop alternative approaches, a useful way to think about turning multiclass classification problems into binary classification problems is to think of them like tournaments (football, soccer—aka football, cricket, tennis, or whatever appeals to you). You have K teams entering a tournament, but unfortunately the sport they are playing only allows two to compete at a time. You want to set up a way of pairing the teams and having them compete so that you can figure out which team is best. In learning, the teams are now the classes and you're trying to figure out which class is best.¹

One natural approach is to have every team compete against every other team. The team that wins the majority of its matches is

¹ The sporting analogy breaks down a bit for OVA: K games are played, wherein each team will play simultaneously against all other teams.

Algorithm 15 ALLVERSUSALLTRAIN($\mathbf{D}^{\text{multiclass}}$, BINARYTRAIN)

```

1:  $f_{ij} \leftarrow \emptyset, \forall 1 \leq i < j \leq K$ 
2: for  $i = 1$  to  $K-1$  do
3:    $\mathbf{D}^{\text{pos}} \leftarrow$  all  $\mathbf{x} \in \mathbf{D}^{\text{multiclass}}$  labeled  $i$ 
4:   for  $j = i+1$  to  $K$  do
5:      $\mathbf{D}^{\text{neg}} \leftarrow$  all  $\mathbf{x} \in \mathbf{D}^{\text{multiclass}}$  labeled  $j$ 
6:      $\mathbf{D}^{\text{bin}} \leftarrow \{(x, +1) : x \in \mathbf{D}^{\text{pos}}\} \cup \{(x, -1) : x \in \mathbf{D}^{\text{neg}}\}$ 
7:      $f_{ij} \leftarrow \text{BINARYTRAIN}(\mathbf{D}^{\text{bin}})$ 
8:   end for
9: end for
10: return all  $f_{ij}$ s

```

Algorithm 16 ALLVERSUSALLTEST(all f_{ij} , $\hat{\mathbf{x}}$)

```

1:  $\text{score} \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize  $K$ -many scores to zero
2: for  $i = 1$  to  $K-1$  do
3:   for  $j = i+1$  to  $K$  do
4:      $y \leftarrow f_{ij}(\hat{\mathbf{x}})$ 
5:      $\text{score}_i \leftarrow \text{score}_i + y$ 
6:      $\text{score}_j \leftarrow \text{score}_j - y$ 
7:   end for
8: end for
9: return  $\text{argmax}_k \text{score}_k$ 

```

declared the winner. This is the **all versus all** (or **AVA**) approach (sometimes called **all pairs**). The most natural way to think about it is as training $\binom{K}{2}$ classifiers. Say f_{ij} for $1 \leq i < j \leq k$ is the classifier that pits class i against class j . This classifier receives all of the class i examples as “positive” and all of the class j examples as “negative.” When a test point arrives, it is run through all f_{ij} classifiers. Every time f_{ij} predicts positive, class i gets a point; otherwise, class j gets a point. After running all $\binom{K}{2}$ classifiers, the class with the most votes wins.

The training and test algorithms for AVA are sketched in Algorithms 6.2 and 6.2. In theory, the AVA mapping is more complicated than the weighted binary case. The result is stated below, but the proof is omitted.

Theorem 5 (AVA Error Bound). *Suppose the average binary error of the $\binom{K}{2}$ binary classifiers is ϵ . Then the error rate of the AVA multiclass predictor is at most $2(K-1)\epsilon$.*

At this point, you might be wondering if it’s possible to do better than something linear in K . Fortunately, the answer is yes! The solution, like so much in computer science, is divide and conquer. The idea is to construct a *binary tree* of classifiers. The leaves of this tree correspond to the K labels. Since there are only $\log_2 K$ decisions made to get from the root to a leaf, then there are only $\log_2 K$ chances

Suppose that you have N data points in K classes, evenly divided. How long does it take to train an AVA classifier, if the base binary classifier takes $\mathcal{O}(N)$ time to train? What if the base classifier takes $\mathcal{O}(N^2)$ time? How does this compare to OVA?

The bound for AVA is $2(K-1)\epsilon$; the bound for OVA is $(K-1)\epsilon$. Does this mean that OVA is necessarily better than AVA? Why or why not?

to make an error.

An example of a classification tree for $K = 8$ classes is shown in Figure 6.2. At the root, you distinguish between classes $\{1, 2, 3, 4\}$ and classes $\{5, 6, 7, 8\}$. This means that you will train a binary classifier whose positive examples are all data points with multiclass label $\{1, 2, 3, 4\}$ and whose negative examples are all data points with multiclass label $\{5, 6, 7, 8\}$. Based on what decision is made by this classifier, you can walk down the appropriate path in the tree. When K is not a power of 2, the tree will not be full. This classification tree algorithm achieves the following bound.

Theorem 6 (Tree Error Bound). *Suppose the average binary classifiers error is ϵ . Then the error rate of the tree classifier is at most $\lceil \log_2 K \rceil \epsilon$.*

Proof of Theorem 6. A multiclass error is made if any classifier on the path from the root to the correct leaf makes an error. Each has probability ϵ of making an error and the path consists of at most $\lceil \log_2 K \rceil$ binary decisions. \square

One thing to keep in mind with tree classifiers is that you have control over how the tree is defined. In OVA and AVA you have no say in what classification problems are created. In tree classifiers, the only thing that matters is that, at the root, half of the classes are considered positive and half are considered negative. You want to split the classes in such a way that this classification decision is as easy as possible. You can use whatever you happen to know about your classification problem to try to separate the classes out in a reasonable way.

Can you do better than $\lceil \log_2 K \rceil \epsilon$? It turns out the answer is yes, but the algorithms to do so are relatively complicated. You can actually do as well as 2ϵ using the idea of error-correcting tournaments. Moreover, you can prove a *lower bound* that states that the best you could possibly do is $\epsilon/2$. This means that error-correcting tournaments are at most a factor of four worse than optimal.

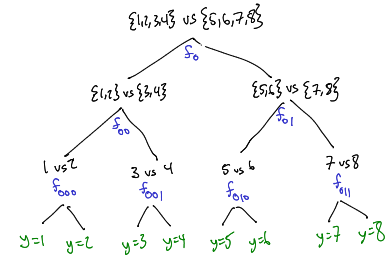


Figure 6.2: example classification tree for $K = 8$

6.3 Ranking

You start a new web search company called Goohooing. Like other search engines, a user inputs a query and a set of documents is retrieved. Your goal is to rank the resulting documents based on relevance to the query. The ranking problem is to take a collection of items and sort them according to some notion of preference. One of the trickiest parts of doing ranking through learning is to properly define the loss function. Toward the end of this section you will see a very general loss function, but before that let's consider a few special cases.

Algorithm 17 NAIVERANKTRAIN(*RankingData*, BINARYTRAIN)

```

1:  $\mathbf{D} \leftarrow []$ 
2: for  $n = 1$  to  $N$  do
3:   for all  $i, j = 1$  to  $M$  and  $i \neq j$  do
4:     if  $i$  is preferred to  $j$  on query  $n$  then
5:        $\mathbf{D} \leftarrow \mathbf{D} \oplus (x_{nij}, +1)$ 
6:     else if  $j$  is preferred to  $i$  on query  $n$  then
7:        $\mathbf{D} \leftarrow \mathbf{D} \oplus (x_{nij}, -1)$ 
8:     end if
9:   end for
10: end for
11: return BINARYTRAIN( $\mathbf{D}$ )

```

Continuing the web search example, you are given a collection of queries. For each query, you are also given a collection of documents, together with a desired ranking over those documents. In the following, we'll assume that you have N -many queries and for each query you have M -many documents. (In practice, M will probably vary by query, but for ease we'll consider the simplified case.) The goal is to train a binary classifier to predict a **preference function**. Given a query q and two documents d_i and d_j , the classifier should predict whether d_i should be preferred to d_j with respect to the query q .

As in all the previous examples, there are two things we have to take care of: (1) how to train the classifier that predicts preferences; (2) how to turn the predicted preferences into a ranking. Unlike the previous examples, the second step is somewhat complicated in the ranking case. This is because we need to predict an entire ranking of a large number of documents, somehow assimilating the preference function into an overall permutation.

For notational simplicity, let x_{nij} denote the features associated with comparing document i to document j on query n . Training is fairly straightforward. For every n and every pair $i \neq j$, we will create a binary classification example based on features x_{nij} . This example is positive if i is preferred to j in the true ranking. It is negative if j is preferred to i . (In some cases the true ranking will not express a preference between two objects, in which case we exclude the i, j and j, i pair from training.)

Now, you might be tempted to evaluate the classification performance of this binary classifier on its own. The problem with this approach is that it's impossible to tell—just by looking at its output on one i, j pair—how good the overall ranking is. This is because there is the intermediate step of turning these pairwise predictions into a coherent ranking. What you need to do is measure how well the ranking based on your predicted preferences compares to the true ordering. Algorithms 6.3 and 6.3 show naive algorithms for training

Algorithm 18 NAIVERANKTEST(f, \hat{x})

```

1:  $score \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize  $M$ -many scores to zero
2: for all  $i, j = 1$  to  $M$  and  $i \neq j$  do
3:    $y \leftarrow f(\hat{x}_{ij})$  // get predicted ranking of  $i$  and  $j$ 
4:    $score_i \leftarrow score_i + y$ 
5:    $score_j \leftarrow score_j - y$ 
6: end for
7: return ARGSORT( $score$ ) // return queries sorted by score

```

and testing a ranking function.

These algorithms actually work quite well in the case of **bipartite ranking problems**. A bipartite ranking problem is one in which you are only ever trying to predict a binary response, for instance “is this document relevant or not?” but are being evaluated according to a metric like **AUC**. This is essentially because the only goal in bipartite problems is to ensure that all the relevant documents are ahead of all the irrelevant documents. There is no notion that one relevant document is *more relevant* than another.

For non-bipartite ranking problems, you can do better. First, when the preferences that you get at training time are more nuanced than “relevant or not,” you can incorporate these preferences at training time. Effectively, you want to give a higher weight to binary problems that are very different in terms of preference than others. Second, rather than producing a list of scores and then calling an arbitrary sorting algorithm, you can actually use the preference function as the sorting function inside your own implementation of quicksort.

We can now formalize the problem. Define a ranking as a function σ that maps the objects we are ranking (documents) to the desired position in the list, $1, 2, \dots, M$. If $\sigma_u < \sigma_v$ then u is preferred to v (i.e., appears earlier on the ranked document list). Given data with observed rankings σ , our goal is to learn to predict rankings for new objects, $\hat{\sigma}$. We define Σ_M as the set of all ranking functions over M objects. We also wish to express the fact that making a mistake on some pairs is worse than making a mistake on others. This will be encoded in a cost function ω (omega), where $\omega(i, j)$ is the cost for accidentally putting something in position j when it should have gone in position i . To be a valid cost function, ω must be (1) symmetric, (2) monotonic and (3) satisfy the triangle inequality. Namely: (1) $\omega(i, j) = \omega(j, i)$; (2) if $i < j < k$ or $i > j > k$ then $\omega(i, j) \leq \omega(i, k)$; (3) $\omega(i, j) + \omega(j, k) \geq \omega(i, k)$. With these definitions, we can properly define the ranking problem.

TASK: ω -RANKING*Given:*

1. An input space \mathcal{X}
2. An unknown distribution \mathcal{D} over $\mathcal{X} \times \Sigma_M$
3. A training set D sampled from \mathcal{D}

Compute: A function $f : \mathcal{X} \rightarrow \Sigma_M$ minimizing:

$$\mathbb{E}_{(x, \sigma) \sim \mathcal{D}} \left[\sum_{u \neq v} [\sigma_u < \sigma_v] [\hat{\sigma}_v < \hat{\sigma}_u] \omega(\sigma_u, \sigma_v) \right] \quad (6.7)$$

where $\hat{\sigma} = f(x)$

In this definition, the only complex aspect is the loss function 6.7. This loss sums over all pairs of objects u and v . If the true ranking (σ) prefers u to v , but the predicted ranking ($\hat{\sigma}$) prefers v to u , then you incur a cost of $\omega(\sigma_u, \sigma_v)$.

Depending on the problem you care about, you can set ω to many “standard” options. If $\omega(i, j) = 1$ whenever $i \neq j$, then you achieve the Kemeny distance measure, which simply counts the number of pairwise misordered items. In many applications, you may only care about getting the top K predictions correct. For instance, your web search algorithm may only display $K = 10$ results to a user. In this case, you can define:

$$\omega(i, j) = \begin{cases} 1 & \text{if } \min\{i, j\} \leq K \text{ and } i \neq j \\ 0 & \text{otherwise} \end{cases} \quad (6.8)$$

In this case, only errors in the top K elements are penalized. Swapping items 55 and 56 is irrelevant (for $K < 55$).

Finally, in the bipartite ranking case, you can express the **area under the curve (AUC)** metric as:

$$\omega(i, j) = \frac{\binom{M}{2}}{M^+(M - M^+)} \times \begin{cases} 1 & \text{if } i \leq M^+ \text{ and } j > M^+ \\ 1 & \text{if } j \leq M^+ \text{ and } i > M^+ \\ 0 & \text{otherwise} \end{cases} \quad (6.9)$$

Here, M is the total number of objects to be ranked and M^+ is the number that are actually “good.” (Hence, $M - M^+$ is the number that are actually “bad,” since this is a bipartite problem.) You are only penalized if you rank a good item in position greater than M^+ or if you rank a bad item in a position less than or equal to M^+ .

In order to *solve* this problem, you can follow a recipe similar to the naive approach sketched earlier. At training time, the biggest

Algorithm 19 RANKTRAIN($\mathbf{D}^{rank}, \omega, \text{BINARYTRAIN}$)

```

1:  $\mathbf{D}^{bin} \leftarrow []$ 
2: for all  $(x, \sigma) \in \mathbf{D}^{rank}$  do
3:   for all  $u \neq v$  do
4:      $y \leftarrow \text{SIGN}(\sigma_v - \sigma_u)$  //  $y$  is +1 if  $u$  is preferred to  $v$ 
5:      $w \leftarrow \omega(\sigma_u, \sigma_v)$  //  $w$  is the cost of misclassification
6:      $\mathbf{D}^{bin} \leftarrow \mathbf{D}^{bin} \oplus (y, w, x_{uv})$ 
7:   end for
8: end for
9: return BINARYTRAIN( $\mathbf{D}^{bin}$ )

```

Algorithm 20 RANKTEST(f, \hat{x}, obj)

```

1: if  $obj$  contains 0 or 1 elements then
2:   return  $obj$ 
3: else
4:    $p \leftarrow$  randomly chosen object in  $obj$  // pick pivot
5:    $left \leftarrow []$  // elements that seem smaller than  $p$ 
6:    $right \leftarrow []$  // elements that seem larger than  $p$ 
7:   for all  $u \in obj \setminus \{p\}$  do
8:      $\hat{y} \leftarrow f(x_{up})$  // what is the probability that  $u$  precedes  $p$ 
9:     if uniform random variable  $< \hat{y}$  then
10:       $left \leftarrow left \oplus u$ 
11:     else
12:       $right \leftarrow right \oplus u$ 
13:     end if
14:   end for
15:    $left \leftarrow \text{RANKTEST}(f, \hat{x}, left)$  // sort earlier elements
16:    $right \leftarrow \text{RANKTEST}(f, \hat{x}, right)$  // sort later elements
17:   return  $left \oplus \langle p \rangle \oplus right$ 
18: end if

```

change is that you can *weight* each training example by how bad it would be to mess it up. This change is depicted in Algorithm 6.3, where the binary classification data has *weights* w provided for saying how important a given example is. These weights are derived from the cost function ω .

At test time, instead of predicting scores and then sorting the list, you essentially run the quicksort algorithm, using f as a comparison function. At each step in Algorithm 6.3, a pivot p is chosen. Every other object u is compared to p using f . If f thinks u is better, then it is sorted on the left; otherwise it is sorted on the right. There is one major difference between this algorithm and quicksort: the comparison function is allowed to be *probabilistic*. If f outputs probabilities, for instance it predicts that u has an 80% probability of being better than p , then it puts it on the left with 80% probability and on the right with 20% probability. (The pseudocode is written in such a way that even if f just predicts $-1, +1$, the algorithm still works.)

This algorithm is better than the naive algorithm in at least two ways. First, it only makes $\mathcal{O}(M \log_2 M)$ calls to f (in expectation), rather than $\mathcal{O}(M^2)$ calls in the naive case. Second, it achieves a better error bound, shown below:

Theorem 7 (Rank Error Bound). *Suppose the average binary error of f is ϵ . Then the ranking algorithm achieves a test error of at most 2ϵ in the general case, and ϵ in the bipartite case.*

6.4 Further Reading

TODO further reading

7 | LINEAR MODELS

The essence of mathematics is not to make simple things complicated, but to make complicated things simple. — Stanley Gudder

IN CHAPTER 4, YOU LEARNED about the perceptron algorithm for linear classification. This was both a *model* (linear classifier) and *algorithm* (the perceptron update rule) in one. In this section, we will separate these two, and consider general ways for optimizing linear models. This will lead us into some aspects of optimization (aka mathematical programming), but not very far. At the end of this chapter, there are pointers to more literature on optimization for those who are interested.

The basic idea of the perceptron is to run a particular algorithm until a linear separator is found. You might ask: are there better algorithms for finding such a linear separator? We will follow this idea and formulate a learning problem as an explicit optimization problem: find me a linear separator that is not too complicated. We will see that finding an “optimal” separator is actually computationally prohibitive, and so will need to “relax” the optimality requirement. This will lead us to a **convex** objective that combines a loss function (how well are we doing on the training data?) and a regularizer (how complicated is our learned model?). This learning framework is known as both **Tikhonov regularization** and **structural risk minimization**.

Learning Objectives:

- Define and plot four surrogate loss functions: squared loss, logistic loss, exponential loss and hinge loss.
- Compare and contrast the optimization of 0/1 loss and surrogate loss functions.
- Solve the optimization problem for squared loss with a quadratic regularizer in closed form.
- Implement and debug gradient descent and subgradient descent.

Dependencies:

7.1 The Optimization Framework for Linear Models

You have already seen the perceptron as a way of finding a weight vector w and bias b that do a good job of separating positive training examples from negative training examples. The perceptron is a **model** and **algorithm** in one. Here, we are interested in *separating* these issues. We will focus on linear models, like the perceptron. But we will think about other, more generic ways of finding good parameters of these models.

The goal of the perceptron was to find a **separating hyperplane** for some training data set. For simplicity, you can ignore the issue of overfitting (but just for now!). Not all data sets are linearly sepa-

table. In the case that your training data *isn't* linearly separable, you might want to find the hyperplane that makes the *fewest errors* on the training data. We can write this down as a formal mathematics **optimization problem** as follows:

$$\min_{w,b} \sum_n \mathbf{1}[y_n(w \cdot x_n + b) > 0] \quad (7.1)$$

In this expression, you are optimizing over two variables, w and b . The **objective function** is the thing you are trying to minimize. In this case, the objective function is simply the **error rate** (or **0/1 loss**) of the linear classifier parameterized by w, b . In this expression, $\mathbf{1}[\cdot]$ is the **indicator function**: it is one when (\cdot) is true and zero otherwise.

We know that the perceptron algorithm is guaranteed to find parameters for this model if the data is linearly separable. In other words, if the optimum of Eq (7.1) is zero, then the perceptron will efficiently find parameters for this model. The notion of “efficiency” depends on the margin of the data for the perceptron.

You might ask: what happens if the data is *not* linearly separable? Is there an efficient algorithm for finding an optimal setting of the parameters? Unfortunately, the answer is *no*. There is no polynomial time algorithm for solving Eq (7.1), unless $P=NP$. In other words, this problem is NP-hard. Sadly, the proof of this is quite complicated and beyond the scope of this book, but it relies on a reduction from a variant of satisfiability. The key idea is to turn a satisfiability problem into an optimization problem where a clause is satisfied exactly when the hyperplane correctly separates the data.

You might then come back and say: okay, well I don't really need an *exact* solution. I'm willing to have a solution that makes one or two more errors than it has to. Unfortunately, the situation is really bad. Zero/one loss is NP-hard to even *approximately minimize*. In other words, there is no efficient algorithm for even finding a solution that's a small constant worse than optimal. (The best known constant at this time is $418/415 \approx 1.007$.)

However, before getting too disillusioned about this whole enterprise (remember: there's an entire chapter about this framework, so it must be going somewhere!), you should remember that optimizing Eq (7.1) perhaps isn't even what you want to do! In particular, all it says is that you will get minimal *training error*. It says nothing about what your *test error* will be like. In order to try to find a solution that will *generalize* well to test data, you need to ensure that you do not overfit the data. To do this, you can introduce a **regularizer** over the parameters of the model. For now, we will be vague about what this regularizer looks like, and simply call it an arbitrary function $R(w, b)$.

? You should remember the $yw \cdot x$ trick from the perceptron discussion. If not, re-convince yourself that this is doing the right thing.

This leads to the following, **regularized objective**:

$$\min_{w,b} \sum_n \mathbf{1}[y_n(w \cdot x_n + b) > 0] + \lambda R(w, b) \quad (7.2)$$

In Eq (7.2), we are now trying to optimize a *trade-off* between a solution that gives low training error (the first term) and a solution that is “simple” (the second term). You can think of the maximum depth hyperparameter of a decision tree as a form of regularization for trees. Here, R is a form of regularization for hyperplanes. In this formulation, λ becomes a **hyperparameter** for the optimization.

The key remaining questions, given this formalism, are:

- How can we adjust the optimization problem so that there *are* efficient algorithms for solving it?
- What are good regularizers $R(w, b)$ for hyperplanes?
- Assuming we can adjust the optimization problem appropriately, what algorithms exist for efficiently solving this regularized optimization problem?

We will address these three questions in the next sections.

Assuming R does the “right thing,” what value(s) of λ will lead to overfitting? What value(s) will lead to underfitting?

7.2 Convex Surrogate Loss Functions

You might ask: why is optimizing zero/one loss so hard? Intuitively, one reason is that small changes to w, b can have a large impact on the value of the objective function. For instance, if there is a positive training example with $w \cdot x + b = -0.0000001$, then adjusting b upwards by 0.00000011 will decrease your error rate by 1. But adjusting it upwards by 0.00000009 will have no effect. This makes it really difficult to figure out good ways to adjust the parameters.

To see this more clearly, it is useful to look at plots that relate *margin* to *loss*. Such a plot for zero/one loss is shown in Figure 7.1. In this plot, the horizontal axis measures the margin of a data point and the vertical axis measures the loss associated with that margin. For zero/one loss, the story is simple. If you get a positive margin (i.e., $y(w \cdot x + b) > 0$) then you get a loss of zero. Otherwise you get a loss of one. By thinking about this plot, you can see how changes to the parameters that change the margin *just a little bit* can have an enormous effect on the overall loss.

You might decide that a reasonable way to address this problem is to replace the non-smooth zero/one loss with a smooth approximation. With a bit of effort, you could probably concoct an “S”-shaped function like that shown in Figure 7.2. The benefit of using such an S-function is that it is smooth, and potentially easier to optimize. The difficulty is that it is not **convex**.

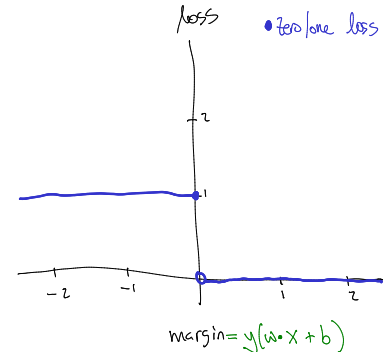


Figure 7.1: plot of zero/one versus margin

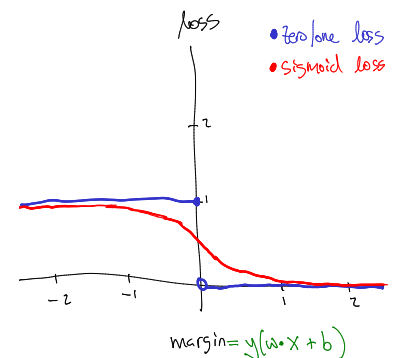


Figure 7.2: plot of zero/one versus margin and an S version of it

If you remember from calculus, a convex function is one that looks like a happy face (☺). (On the other hand, a **concave** function is one that looks like a sad face (☹); an easy mnemonic is that you can hide under a **concave** function.) There are two equivalent definitions of a convex function. The first is that its second derivative is always non-negative. The second, more geometric, definition is that any **chord** of the function lies above it. This is shown in Figure 7.3. There you can see a convex function and a non-convex function, both with two chords drawn in. In the case of the convex function, the chords lie above the function. In the case of the non-convex function, there are parts of the chord that lie below the function.

Convex functions are nice because they are *easy to minimize*. Intuitively, if you drop a ball anywhere in a convex function, it will eventually get to the minimum. This is not true for non-convex functions. For example, if you drop a ball on the very left end of the S-function from Figure 7.2, it will not go anywhere.

This leads to the idea of **convex surrogate loss functions**. Since zero/one loss is hard to optimize, you want to optimize something else, instead. Since convex functions are easy to optimize, we want to approximate zero/one loss with a convex function. This approximating function will be called a **surrogate loss**. The surrogate losses we construct will always be *upper bounds* on the true loss function: this guarantees that if you minimize the surrogate loss, you are also pushing down the real loss.

There are four common surrogate loss functions, each with their own properties: **hinge loss**, **logistic loss**, **exponential loss** and **squared loss**. These are shown in Figure 7.4 and defined below. These are defined in terms of the true label y (which is just $\{-1, +1\}$) and the predicted value $\hat{y} = w \cdot x + b$.

$$\text{Zero/one:} \quad \ell^{(0/1)}(y, \hat{y}) = \mathbf{1}[y\hat{y} \leq 0] \quad (7.3)$$

$$\text{Hinge:} \quad \ell^{(\text{hin})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\} \quad (7.4)$$

$$\text{Logistic:} \quad \ell^{(\text{log})}(y, \hat{y}) = \frac{1}{\log 2} \log(1 + \exp[-y\hat{y}]) \quad (7.5)$$

$$\text{Exponential:} \quad \ell^{(\text{exp})}(y, \hat{y}) = \exp[-y\hat{y}] \quad (7.6)$$

$$\text{Squared:} \quad \ell^{(\text{sq})}(y, \hat{y}) = (y - \hat{y})^2 \quad (7.7)$$

In the definition of logistic loss, the $\frac{1}{\log 2}$ term out front is there simply to ensure that $\ell^{(\text{log})}(y, 0) = 1$. This ensures, like all the other surrogate loss functions, that logistic loss upper bounds the zero/one loss. (In practice, people typically omit this constant since it does not affect the optimization.)

There are two big differences in these loss functions. The first difference is how “upset” they get by erroneous predictions. In the

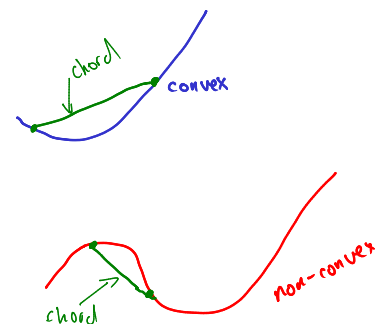


Figure 7.3: plot of convex and non-convex functions with two chords each

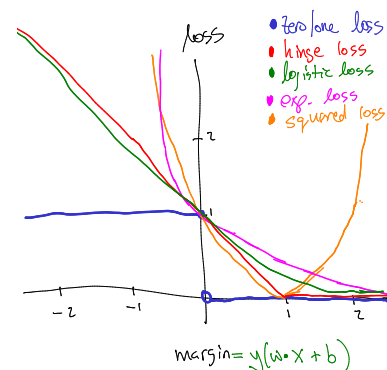


Figure 7.4: surrogate loss fns

case of hinge loss and logistic loss, the growth of the function as \hat{y} goes negative is linear. For squared loss and exponential loss, it is super-linear. This means that exponential loss would rather get a few examples a little wrong than one example really wrong. The other difference is how they deal with very confident correct predictions. Once $y\hat{y} > 1$, hinge loss does not care any more, but logistic and exponential still think you can do better. On the other hand, squared loss thinks it's just as bad to predict $+3$ on a positive example as it is to predict -1 on a positive example.

7.3 Weight Regularization

In our learning objective, Eq (7.2), we had a term correspond to the zero/one loss on the training data, plus a **regularizer** whose goal was to ensure that the learned function didn't get too "crazy." (Or, more formally, to ensure that the function did not overfit.) If you replace to zero/one loss with a surrogate loss, you obtain the following objective:

$$\min_{w,b} \sum_n \ell(y_n, w \cdot x_n + b) + \lambda R(w, b) \quad (7.8)$$

The question is: what should $R(w, b)$ look like?

From the discussion of surrogate loss function, we would like to ensure that R is convex. Otherwise, we will be back to the point where optimization becomes difficult. Beyond that, a common desire is that the components of the weight vector (i.e., the w_d s) should be small (close to zero). This is a form of **inductive bias**.

Why are small values of w_d good? Or, more precisely, why do small values of w_d correspond to *simple functions*? Suppose that we have an example x with label $+1$. We might believe that other examples, x' that are nearby x should also have label $+1$. For example, if I obtain x' by taking x and changing the first component by some small value ϵ and leaving the rest the same, you might think that the classification would be the same. If you do this, the difference between \hat{y} and \hat{y}' will be exactly ϵw_1 . So if w_1 is reasonably small, this is unlikely to have much of an effect on the classification decision. On the other hand, if w_1 is large, this could have a large effect.

Another way of saying the same thing is to look at the derivative of the predictions as a function of w_1 . The derivative of $w \cdot x + b$ with respect to w_1 is:

$$\frac{\partial [w \cdot x + b]}{\partial w_1} = \frac{\partial [\sum_d w_d x_d + b]}{\partial w_1} = x_1 \quad (7.9)$$

Interpreting the derivative as the rate of change, we can see that the rate of change of the prediction function is proportional to the

individual weights. So if you want the function to change slowly, you want to ensure that the weights stay small.

One way to accomplish this is to simply use the norm of the weight vector. Namely $R^{(\text{norm})}(\mathbf{w}, b) = \|\mathbf{w}\| = \sqrt{\sum_d w_d^2}$. This function is convex and smooth, which makes it easy to minimize. In practice, it's often easier to use the squared norm, namely $R^{(\text{sqr})}(\mathbf{w}, b) = \|\mathbf{w}\|^2 = \sum_d w_d^2$ because it removes the ugly square root term and remains convex. An alternative to using the sum of squared weights is to use the sum of absolute weights: $R^{(\text{abs})}(\mathbf{w}, b) = \sum_d |w_d|$. Both of these norms are convex.

In addition to small weights being good, you could argue that *zero* weights are better. If a weight w_d goes to zero, then this means that feature d is not used at all in the classification decision. If there are a large number of irrelevant features, you might want as many weights to go to zero as possible. This suggests an alternative regularizer: $R^{(\text{cnt})}(\mathbf{w}, b) = \sum_d \mathbf{1}[x_d \neq 0]$.

This line of thinking leads to the general concept of ***p*-norms**. (Technically these are called ℓ_p (or “ell *p*”) norms, but this notation clashes with the use of ℓ for “loss.”) This is a family of norms that all have the same general flavor. We write $\|\mathbf{w}\|_p$ to denote the *p*-norm of \mathbf{w} .

$$\|\mathbf{w}\|_p = \left(\sum_d |w_d|^p \right)^{\frac{1}{p}} \quad (7.10)$$

You can check that the 2-norm exactly corresponds to the usual Euclidean norm, and that the 1-norm corresponds to the “absolute” regularizer described above.

When *p*-norms are used to regularize weight vectors, the interesting aspect is how they trade-off multiple features. To see the behavior of *p*-norms in two dimensions, we can plot their **contour** (or **level-set**). Figure 7.5 shows the contours for the same *p* norms in two dimensions. Each line denotes the two-dimensional vectors to which this norm assigns a total value of 1. By changing the value of *p*, you can interpolate between a square (the so-called “max norm”), down to a circle (2-norm), diamond (1-norm) and pointy-star-shaped-thing ($p < 1$ norm).

In general, smaller values of *p* “prefer” sparser vectors. You can see this by noticing that the contours of small *p*-norms “stretch” out along the axes. It is for this reason that small *p*-norms tend to yield weight vectors with many zero entries (aka **sparse** weight vectors). Unfortunately, for $p < 1$ the norm becomes non-convex. As you might guess, this means that the 1-norm is a popular choice for sparsity-seeking applications.

? Why do we not regularize the bias term *b*?

? Why might you not want to use $R^{(\text{cnt})}$ as a regularizer?

? You can actually identify the $R^{(\text{cnt})}$ regularizer with a *p*-norm as well. Which value of *p* gives it to you? (Hint: you may have to take a limit.)

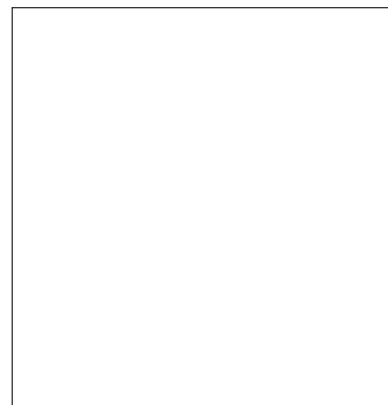


Figure 7.5: loss : norms2d: level sets of the same *p*-norms

? The max norm corresponds to $\lim_{p \rightarrow \infty}$. Why is this called the max norm?

MATH REVIEW | GRADIENTS

A gradient is a multidimensional generalization of a derivative. Suppose you have a function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ that takes a vector $x = \langle x_1, x_2, \dots, x_D \rangle$ as input and produces a scalar value as output. You can differentiate this function according to any one of the inputs; for instance, you can compute $\frac{\partial f}{\partial x_5}$ to get the derivative with respect to the fifth input. The **gradient** of f is just the vector consisting of the derivative f with respect to each of its input coordinates independently, and is denoted ∇f , or, when the input to f is ambiguous, $\nabla_x f$. This is defined as:

$$\nabla_x f = \left\langle \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_D} \right\rangle \quad (7.11)$$

For example, consider the function $f(x_1, x_2, x_3) = x_1^3 + 5x_1x_2 - 3x_2x_3^2$. The gradient is:

$$\nabla_x f = \left\langle 3x_1^2 + 5x_2, 5x_1 - 6x_2x_3, -6x_2x_3 \right\rangle \quad (7.12)$$

Note that if $f : \mathbb{R}^D \rightarrow \mathbb{R}$, then $\nabla f : \mathbb{R}^D \rightarrow \mathbb{R}^D$. If you evaluate $\nabla f(x)$, this will give you the gradient *at* x , a vector in \mathbb{R}^D . This vector can be interpreted as the **direction of steepest ascent**: namely, if you were to travel an infinitesimal amount in the direction of the gradient, you would go uphill (i.e., increase f) the most.

Figure 7.6:

7.4 Optimization with Gradient Descent

Envision the following problem. You're taking up a new hobby: blindfolded mountain climbing. Someone blindfolds you and drops you on the side of a mountain. Your goal is to get to the peak of the mountain as quickly as possible. All you can do is feel the mountain where you are standing, and take steps. How would you get to the top of the mountain? Perhaps you would feel to find out what direction feels the most "upward" and take a step in that direction. If you do this repeatedly, you might hope to get to the top of the mountain. (Actually, if your friend promises always to drop you on purely concave mountains, you *will* eventually get to the peak!)

The idea of gradient-based methods of optimization is exactly the same. Suppose you are trying to find the maximum of a function $f(x)$. The optimizer maintains a current estimate of the parameter of interest, x . At each step, it measures the **gradient** of the function it is trying to optimize. This measurement occurs *at* the current location, x . Call the gradient g . It then takes a step in the direction of the gradient, where the size of the step is controlled by a parameter η (eta). The complete step is $x \leftarrow x + \eta g$. This is the basic idea of **gradient ascent**.

The opposite of gradient ascent is **gradient descent**. All of our

Algorithm 21 GRADIENTDESCENT($\mathcal{F}, K, \eta_1, \dots$)

```

1:  $\mathbf{z}^{(0)} \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize variable we are optimizing
2: for  $k = 1 \dots K$  do
3:    $\mathbf{g}^{(k)} \leftarrow \nabla_{\mathbf{z}} \mathcal{F}|_{\mathbf{z}^{(k-1)}}$  // compute gradient at current location
4:    $\mathbf{z}^{(k)} \leftarrow \mathbf{z}^{(k-1)} - \eta^{(k)} \mathbf{g}^{(k)}$  // take a step down the gradient
5: end for
6: return  $\mathbf{z}^{(K)}$ 

```

learning problems will be framed as *minimization* problems (trying to reach the bottom of a ditch, rather than the top of a hill). Therefore, descent is the primary approach you will use. One of the major conditions for gradient ascent being able to find the true, **global minimum**, of its objective function is convexity. Without convexity, all is lost.

The gradient descent algorithm is sketched in Algorithm 7.4. The function takes as arguments the function \mathcal{F} to be minimized, the number of iterations K to run and a sequence of learning rates η_1, \dots, η_K . (This is to address the case that you might want to start your mountain climbing taking large steps, but only take small steps when you are close to the peak.)

The only real work you need to do to apply a gradient descent method is be able to compute derivatives. For concreteness, suppose that you choose exponential loss as a loss function and the 2-norm as a regularizer. Then, the regularized objective function is:

$$\mathcal{L}(\mathbf{w}, b) = \sum_n \exp[-y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (7.13)$$

The only “strange” thing in this objective is that we have replaced λ with $\frac{\lambda}{2}$. The reason for this change is just to make the gradients cleaner. We can first compute derivatives with respect to b :

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial}{\partial b} \sum_n \exp[-y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] + \frac{\partial}{\partial b} \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (7.14)$$

$$= \sum_n \frac{\partial}{\partial b} \exp[-y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] + 0 \quad (7.15)$$

$$= \sum_n \left(\frac{\partial}{\partial b} - y_n(\mathbf{w} \cdot \mathbf{x}_n + b) \right) \exp[-y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] \quad (7.16)$$

$$= - \sum_n y_n \exp[-y_n(\mathbf{w} \cdot \mathbf{x}_n + b)] \quad (7.17)$$

Before proceeding, it is worth thinking about what this says. From a practical perspective, the optimization will operate by updating $b \leftarrow b - \eta \frac{\partial \mathcal{L}}{\partial b}$. Consider positive examples: examples with $y_n = +1$. We would hope for these examples that the current prediction, $\mathbf{w} \cdot \mathbf{x}_n + b$, is as large as possible. As this value tends toward ∞ , the term in the $\exp[]$ goes to zero. Thus, such points will not contribute to the step.

However, if the current prediction is small, then the $\exp[\cdot]$ term will be positive and non-zero. This means that the bias term b will be *increased*, which is exactly what you would want. Moreover, once all points are very well classified, the derivative goes to zero.

Now that we have done the easy case, let's do the gradient with respect to w .

$$\nabla_w \mathcal{L} = \nabla_w \sum_n \exp[-y_n(w \cdot x_n + b)] + \nabla_w \frac{\lambda}{2} \|w\|^2 \quad (7.18)$$

$$= \sum_n (\nabla_w - y_n(w \cdot x_n + b)) \exp[-y_n(w \cdot x_n + b)] + \lambda w \quad (7.19)$$

$$= - \sum_n y_n x_n \exp[-y_n(w \cdot x_n + b)] + \lambda w \quad (7.20)$$

Now you can repeat the previous exercise. The update is of the form $w \leftarrow w - \eta \nabla_w \mathcal{L}$. For well classified points (ones that tend toward $y_n \infty$), the gradient is near zero. For poorly classified points, the gradient points in the direction $-y_n x_n$, so the update is of the form $w \leftarrow w + c y_n x_n$, where c is some constant. This is just like the perceptron update! Note that c is large for very poorly classified points and small for relatively well classified points.

By looking at the part of the gradient related to the regularizer, the update says: $w \leftarrow w - \lambda w = (1 - \lambda)w$. This has the effect of *shrinking* the weights toward zero. This is exactly what we expect the regularizer to be doing!

The success of gradient descent hinges on appropriate choices for the step size. Figure 7.7 shows what can happen with gradient descent with poorly chosen step sizes. If the step size is too big, you can accidentally step over the optimum and end up oscillating. If the step size is too small, it will take way too long to get to the optimum. For a well-chosen step size, you can show that gradient descent will approach the optimal value at a fast *rate*. The notion of convergence here is that the *objective value* converges to the true minimum.

Theorem 8 (Gradient Descent Convergence). *Under suitable conditions¹, for an appropriately chosen constant step size (i.e., $\eta_1 = \eta_2 = \dots = \eta$), the **convergence rate** of gradient descent is $\mathcal{O}(1/k)$. More specifically, letting z^* be the global minimum of \mathcal{F} , we have: $\mathcal{F}(z^{(k)}) - \mathcal{F}(z^*) \leq \frac{2\|z^{(0)} - z^*\|^2}{\eta k}$.*

The proof of this theorem is a bit complicated because it makes heavy use of some linear algebra. The key is to set the learning rate to $1/L$, where L is the maximum **curvature** of the function that is being optimized. The curvature is simply the “size” of the second derivative. Functions with high curvature have gradients that change

? This considered the case of positive examples. What happens with negative examples?

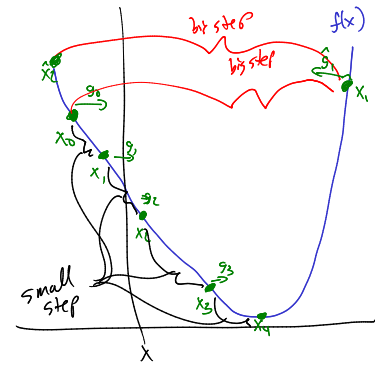


Figure 7.7: good and bad step sizes

¹ Specifically the function to be optimized needs to be **strongly convex**. This is true for all our problems, provided $\lambda > 0$. For $\lambda = 0$ the rate could be as bad as $\mathcal{O}(1/\sqrt{k})$.

? A naive reading of this theorem seems to say that you should choose huge values of η . It should be obvious that this cannot be right. What is missing?

quickly, which means that you need to take small steps to avoid overstepping the optimum.

This convergence result suggests a simple approach to deciding when to stop optimizing: wait until the objective function stops changing by much. An alternative is to wait until the *parameters* stop changing by much. A final example is to do what you did for perceptron: early stopping. Every iteration, you can check the performance of the current model on some held-out data, and stop optimizing when performance plateaus.

7.5 From Gradients to Subgradients

As a good exercise, you should try deriving gradient descent update rules for the different loss functions and different regularizers you've learned about. However, if you do this, you might notice that *hinge loss* and the 1-norm regularizer are not differentiable everywhere! In particular, the 1-norm is not differentiable around $w_d = 0$, and the hinge loss is not differentiable around $y\hat{y} = 1$.

The solution to this is to use **subgradient** optimization. One way to think about subgradients is just to not think about it: you essentially need to just ignore the fact that you forgot that your function wasn't differentiable, and just try to apply gradient descent anyway.

To be more concrete, consider the hinge function $f(z) = \max\{0, 1 - z\}$. This function is differentiable for $z > 1$ and differentiable for $z < 1$, but not differentiable at $z = 1$. You can derive this using differentiation by parts:

$$\frac{\partial}{\partial z} f(z) = \frac{\partial}{\partial z} \begin{cases} 0 & \text{if } z > 1 \\ 1 - z & \text{if } z < 1 \end{cases} \quad (7.21)$$

$$= \begin{cases} \frac{\partial}{\partial z} 0 & \text{if } z > 1 \\ \frac{\partial}{\partial z} (1 - z) & \text{if } z < 1 \end{cases} \quad (7.22)$$

$$= \begin{cases} 0 & \text{if } z \geq 1 \\ -1 & \text{if } z < 1 \end{cases} \quad (7.23)$$

Thus, the derivative is zero for $z < 1$ and -1 for $z > 1$, matching intuition from the Figure. At the non-differentiable point, $z = 1$, we can use a **subderivative**: a generalization of derivatives to non-differentiable functions. Intuitively, you can think of the derivative of f at z as the tangent line. Namely, it is the line that touches f at z that is always below f (for convex functions). The subderivative, denoted ∂f , is the *set* of all such lines. At differentiable positions, this set consists just of the actual derivative. At non-differentiable positions, this contains all slopes that define lines that always lie under the function and make contact at the operating point. This is

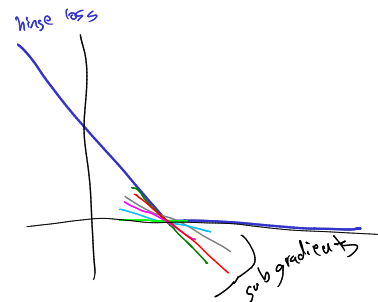


Figure 7.8: hinge loss with sub

Algorithm 22 HINGEREGULARIZEDGD($\mathbf{D}, \lambda, \text{MaxIter}$)

```

1:  $\mathbf{w} \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $b \leftarrow 0$  // initialize weights and bias
2: for  $iter = 1 \dots \text{MaxIter}$  do
3:    $\mathbf{g} \leftarrow \langle 0, 0, \dots, 0 \rangle$  ,  $g \leftarrow 0$  // initialize gradient of weights and bias
4:   for all  $(x, y) \in \mathbf{D}$  do
5:     if  $y(\mathbf{w} \cdot \mathbf{x} + b) \leq 1$  then
6:        $\mathbf{g} \leftarrow \mathbf{g} + y \mathbf{x}$  // update weight gradient
7:        $g \leftarrow g + y$  // update bias derivative
8:     end if
9:   end for
10:   $\mathbf{g} \leftarrow \mathbf{g} - \lambda \mathbf{w}$  // add in regularization term
11:   $\mathbf{w} \leftarrow \mathbf{w} + \eta \mathbf{g}$  // update weights
12:   $b \leftarrow b + \eta g$  // update bias
13: end for
14: return  $\mathbf{w}, b$ 

```

shown pictorially in Figure 7.8, where example subderivatives are shown for the hinge loss function. In the particular case of hinge loss, any value between 0 and -1 is a valid subderivative at $z = 0$. In fact, the subderivative is always a closed set of the form $[a, b]$, where a and b can be derived by looking at limits from the left and right.

This gives you a way of computing derivative-like things for non-differentiable functions. Take hinge loss as an example. For a given example n , the subgradient of hinge loss can be computed as:

$$\partial_{\mathbf{w}} \max\{0, 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b)\} \quad (7.24)$$

$$= \partial_{\mathbf{w}} \begin{cases} 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) > 1 \\ 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases} \quad (7.25)$$

$$= \begin{cases} \partial_{\mathbf{w}} 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) > 1 \\ \partial_{\mathbf{w}} 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases} \quad (7.26)$$

$$= \begin{cases} \mathbf{0} & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) > 1 \\ -y_n \mathbf{x}_n & \text{otherwise} \end{cases} \quad (7.27)$$

If you plug this subgradient form into Algorithm 7.4, you obtain Algorithm 7.5. This is the **subgradient descent** for regularized hinge loss (with a 2-norm regularizer).

7.6 Closed-form Optimization for Squared Loss

Although gradient descent is a good, generic optimization algorithm, there are cases when you can do better. An example is the case of a 2-norm regularizer and squared error loss function. For this, you can actually obtain a *closed form* solution for the optimal weights. However, to obtain this, you need to rewrite the optimization problem in terms of matrix operations. For simplicity, we will only consider the

MATH REVIEW | MATRIX MULTIPLICATION AND INVERSION

If \mathbf{A} and \mathbf{B} are matrices, and \mathbf{A} is $N \times K$ and \mathbf{B} is $K \times M$ (the inner dimensions must match), then the matrix product \mathbf{AB} is a matrix \mathbf{C} that is $N \times M$, with $C_{n,m} = \sum_k A_{n,k} B_{k,m}$. If \mathbf{v} is a vector in \mathbb{R}^D , we will treat it as a *column vector*, or a matrix of size $D \times 1$. Thus, \mathbf{Av} is well defined if \mathbf{A} is $D \times M$, and the resulting product is a vector \mathbf{u} with $u_m = \sum_d A_{d,m} v_d$.

Aside from matrix product, a fundamental matrix operation is inversion. We will often encounter a form like $\mathbf{Ax} = \mathbf{y}$, where \mathbf{A} and \mathbf{y} are known and we want to solve for \mathbf{x} . If \mathbf{A} is square of size $N \times N$, then the inverse of \mathbf{A} , denoted \mathbf{A}^{-1} , is also a square matrix of size $N \times N$, such that $\mathbf{AA}^{-1} = \mathbf{I}_N = \mathbf{A}^{-1}\mathbf{A}$. I.e., multiplying a matrix by its inverse (on either side) gives back the identity matrix. Using this, we can solve $\mathbf{Ax} = \mathbf{y}$ by multiplying both sides by \mathbf{A}^{-1} on the left (recall that order matters in matrix multiplication), yielding $\mathbf{A}^{-1}\mathbf{Ax} = \mathbf{A}^{-1}\mathbf{y}$ from which we can conclude $\mathbf{x} = \mathbf{A}^{-1}\mathbf{y}$. Note that not all square matrices are invertible. For instance, the all zeros matrix does not have an inverse (in the same way that $1/0$ is not defined for scalars). However, there are other matrices that do not have inverses; such matrices are called **singular**.

Figure 7.9:

unbiased version, but the extension is Exercise ?? . This is precisely the **linear regression** setting.

You can think of the training data as a large matrix \mathbf{X} of size $N \times D$, where $X_{n,d}$ is the value of the d th feature on the n th example. You can think of the labels as a column (“tall”) vector \mathbf{Y} of dimension N . Finally, you can think of the weights as a column vector \mathbf{w} of size D . Thus, the matrix-vector product $\mathbf{a} = \mathbf{Xw}$ has dimension N . In particular:

$$a_n = [\mathbf{Xw}]_n = \sum_d X_{n,d} w_d \quad (7.28)$$

This means, in particular, that \mathbf{a} is actually the predictions of the model. Instead of calling this \mathbf{a} , we will call it $\hat{\mathbf{Y}}$. The squared error says that we should minimize $\frac{1}{2} \sum_n (\hat{Y}_n - Y_n)^2$, which can be written in vector form as a minimization of $\frac{1}{2} \|\hat{\mathbf{Y}} - \mathbf{Y}\|^2$.

This can be expanded visually as:

$$\underbrace{\begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,D} \\ x_{2,1} & x_{2,2} & \dots & x_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,D} \end{bmatrix}}_{\mathbf{X}} \underbrace{\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}}_{\mathbf{w}} = \underbrace{\begin{bmatrix} \sum_d x_{1,d} w_d \\ \sum_d x_{2,d} w_d \\ \vdots \\ \sum_d x_{N,d} w_d \end{bmatrix}}_{\hat{\mathbf{Y}}} \approx \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\mathbf{Y}} \quad (7.29)$$



Verify that the squared error can actually be written as this vector norm.

So, compactly, our optimization problem can be written as:

$$\min_w \mathcal{L}(w) = \frac{1}{2} \|Xw - Y\|^2 + \frac{\lambda}{2} \|w\|^2 \quad (7.30)$$

If you recall from calculus, you can minimize a function by setting its derivative to zero. We start with the weights w and take gradients:

$$\nabla_w \mathcal{L}(w) = X^\top (Xw - Y) + \lambda w \quad (7.31)$$

$$= X^\top Xw - X^\top Y + \lambda w \quad (7.32)$$

$$= (X^\top X + \lambda I) w - X^\top Y \quad (7.33)$$

We can equate this to zero and solve, yielding:

$$(X^\top X + \lambda I) w - X^\top Y = 0 \quad (7.34)$$

$$\iff (X^\top X + \lambda I_D) w = X^\top Y \quad (7.35)$$

$$\iff w = (X^\top X + \lambda I_D)^{-1} X^\top Y \quad (7.36)$$

Thus, the *optimal* solution of the weights can be computed by a few matrix multiplications and a matrix inversion. As a sanity check, you can make sure that the dimensions match. The matrix $X^\top X$ has dimension $D \times D$, and therefore so does the inverse term. The inverse is $D \times D$ and X^\top is $D \times N$, so that product is $D \times N$. Multiplying through by the $N \times 1$ vector Y yields a $D \times 1$ vector, which is precisely what we want for the weights.

Note that this gives an *exact solution*, modulo numerical inaccuracies with computing matrix inverses. In contrast, gradient descent will give you progressively better solutions and will “eventually” converge to the optimum at a rate of $1/k$. This means that if you want an answer that’s within an accuracy of $\epsilon = 10^{-4}$, you will need something on the order of one thousand steps.

The question is whether getting this exact solution is always more efficient. To run gradient descent for one step will take $\mathcal{O}(ND)$ time, with a relatively small constant. You will have to run K iterations, yielding an overall runtime of $\mathcal{O}(KND)$. On the other hand, the closed form solution requires constructing $X^\top X$, which takes $\mathcal{O}(D^2N)$ time. The inversion take $\mathcal{O}(D^3)$ time using standard matrix inversion routines. The final multiplications take $\mathcal{O}(ND)$ time. Thus, the overall runtime is on the order $\mathcal{O}(D^3 + D^2N)$. In most standard cases (though this is becoming less true over time), $N > D$, so this is dominated by $\mathcal{O}(D^2N)$.

Thus, the overall question is whether you will need to run more than D -many iterations of gradient descent. If so, then the matrix inversion will be (roughly) faster. Otherwise, gradient descent will be (roughly) faster. For low- and medium-dimensional problems (say,

For those who are keen on linear algebra, you might be worried that the matrix you must invert might not be invertible. Is this actually a problem?

$D \leq 100$), it is probably faster to do the closed form solution via matrix inversion. For high dimensional problems ($D \geq 10,000$), it is probably faster to do gradient descent. For things in the middle, it's hard to say for sure.

7.7 Support Vector Machines

At the beginning of this chapter, you may have looked at the convex surrogate loss functions and asked yourself: where did these come from?! They are all derived from different underlying principles, which essentially correspond to different inductive biases.

Let's start by thinking back to the original goal of linear classifiers: to find a hyperplane that separates the positive training examples from the negative ones. Figure 7.10 shows some data and three potential hyperplanes: red, green and blue. Which one do you like best?

Most likely you chose the green hyperplane. And most likely you chose it because it was furthest away from the closest training points. In other words, it had a large **margin**. The desire for hyperplanes with large margins is a perfect example of an inductive bias. The data does not tell us which of the three hyperplanes is best: we have to choose one using some other source of information.

Following this line of thinking leads us to the **support vector machine** (SVM). This is simply a way of setting up an optimization problem that attempts to find a separating hyperplane with as large a margin as possible. It is written as a **constrained optimization problem**:

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{\gamma(w,b)} \\ \text{subj. to} \quad & y_n (w \cdot x_n + b) \geq 1 \quad (\forall n) \end{aligned} \quad (7.37)$$

In this optimization, you are trying to find parameters that maximize the margin, denoted γ , (i.e., minimize the reciprocal of the margin) subject to the constraint that *all* training examples are correctly classified.

The “odd” thing about this optimization problem is that we require the classification of each point to be greater than *one* rather than simply greater than *zero*. However, the problem doesn't fundamentally change if you replace the “1” with any other positive constant (see Exercise ??). As shown in Figure 7.11, the constant one can be interpreted visually as ensuring that there is a non-trivial margin between the positive points and negative points.

The difficulty with the optimization problem in Eq (7.37) is what happens with data that is not linearly separable. In that case, there is *no* set of parameters w, b that can simultaneously satisfy all the

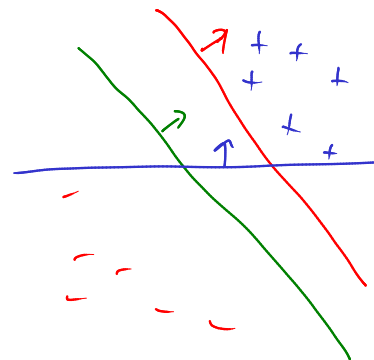


Figure 7.10: picture of data points with three hyperplanes, RGB with G the best

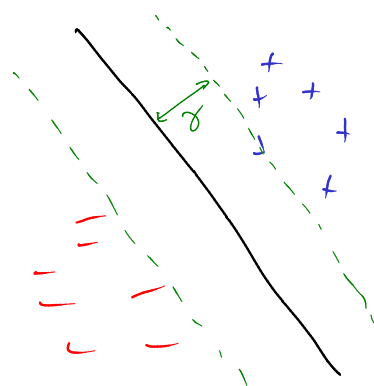


Figure 7.11: hyperplane with margins on sides

constraints. In optimization terms, you would say that the **feasible region** is *empty*. (The feasible region is simply the set of all parameters that satisfy the constraints.) For this reason, this is referred to as the **hard-margin SVM**, because enforcing the margin is a hard constraint. The question is: how to modify this optimization problem so that it can handle inseparable data.

The key idea is the use of **slack parameters**. The intuition behind slack parameters is the following. Suppose we find a set of parameters w, b that do a really good job on 9999 data points. The points are perfectly classified and you achieve a large margin. But there's one pesky data point left that cannot be put on the proper side of the margin: perhaps it is noisy. (See Figure 7.12.) You want to be able to pretend that you can "move" that point across the hyperplane on to the proper side. You will have to pay a little bit to do so, but as long as you aren't moving a *lot* of points around, it should be a good idea to do this. In this picture, the amount that you move the point is denoted ξ (ξ).

By introducing one slack parameter for each training example, and penalizing yourself for having to use slack, you can create an objective function like the following, **soft-margin SVM**:

$$\begin{aligned} \min_{w, b, \xi} \quad & \underbrace{\frac{1}{\gamma(w, b)}}_{\text{large margin}} + C \underbrace{\sum_n \xi_n}_{\text{small slack}} \quad (7.38) \\ \text{subj. to} \quad & y_n (w \cdot x_n + b) \geq 1 - \xi_n \quad (\forall n) \\ & \xi_n \geq 0 \quad (\forall n) \end{aligned}$$

The goal of this objective function is to ensure that all points are correctly classified (the first constraint). But if a point n cannot be correctly classified, then you can set the slack ξ_n to something greater than zero to "move" it in the correct direction. However, for all non-zero slacks, you have to pay in the objective function proportional to the amount of slack. The hyperparameter $C > 0$ controls overfitting versus underfitting. The second constraint simply says that you must not have negative slack.

One major advantage of the soft-margin SVM over the original hard-margin SVM is that the feasible region is *never empty*. That is, there is always going to be some solution, regardless of whether your training data is linearly separable or not.

It's one thing to write down an optimization problem. It's another thing to try to solve it. There are a very large number of ways to optimize SVMs, essentially because they are such a popular learning model. Here, we will talk just about one, very simple way. More complex methods will be discussed later in this book once you have a bit more background.

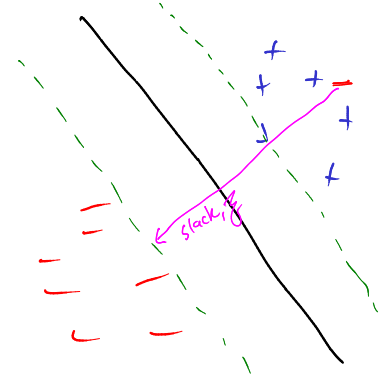


Figure 7.12: one bad point with slack

? What values of C will lead to overfitting? What values will lead to underfitting?

? Suppose I give you a data set. Without even looking at the data, construct for me a feasible solution to the soft-margin SVM. What is the value of the objective for this solution?

To make progress, you need to be able to measure the size of the margin. Suppose someone gives you parameters w, b that optimize the hard-margin SVM. We wish to measure the size of the margin. The first observation is that the hyperplane will lie *exactly* halfway between the nearest positive point and nearest negative point. If not, the margin could be made bigger by simply sliding it one way or the other by adjusting the bias b .

By this observation, there is some positive example that lies exactly 1 unit from the hyperplane. Call it x^+ , so that $w \cdot x^+ + b = 1$. Similarly, there is some negative example, x^- , that lies exactly on the other side of the margin: for which $w \cdot x^- + b = -1$. These two points, x^+ and x^- give us a way to measure the size of the margin. As shown in Figure 7.11, we can measure the size of the margin by looking at the difference between the lengths of projections of x^+ and x^- onto the hyperplane. Since projection requires a normalized vector, we can measure the distances as:

$$d^+ = \frac{1}{\|w\|} w \cdot x^+ + b - 1 \quad (7.39)$$

$$d^- = -\frac{1}{\|w\|} w \cdot x^- - b + 1 \quad (7.40)$$

We can then compute the margin by algebra:

$$\gamma = \frac{1}{2} [d^+ - d^-] \quad (7.41)$$

$$= \frac{1}{2} \left[\frac{1}{\|w\|} w \cdot x^+ + b - 1 - \left(-\frac{1}{\|w\|} w \cdot x^- - b + 1 \right) \right] \quad (7.42)$$

$$= \frac{1}{2} \left[\frac{1}{\|w\|} w \cdot x^+ - \frac{1}{\|w\|} w \cdot x^- \right] \quad (7.43)$$

$$= \frac{1}{2} \left[\frac{1}{\|w\|} (+1) - \frac{1}{\|w\|} (-1) \right] \quad (7.44)$$

$$= \frac{1}{\|w\|} \quad (7.45)$$

This is a remarkable conclusion: the size of the margin is inversely proportional to the norm of the weight vector. Thus, **maximizing the margin is equivalent to minimizing $\|w\|$** ! This serves as an additional justification of the 2-norm regularizer: having small weights means having large margins!

However, our goal wasn't to justify the regularizer: it was to understand hinge loss. So let us go back to the soft-margin SVM and plug in our new knowledge about margins:

$$\min_{w, b, \xi} \underbrace{\frac{1}{2} \|w\|^2}_{\text{large margin}} + C \underbrace{\sum_n \xi_n}_{\text{small slack}} \quad (7.46)$$

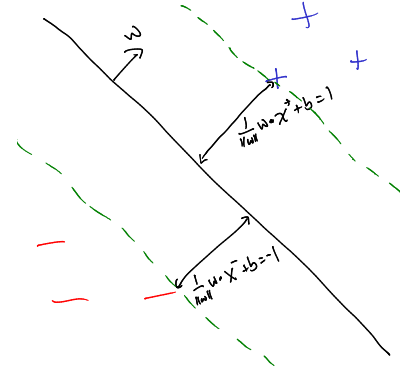


Figure 7.13: copy of figure from p5 of cs544 svm tutorial

$$\begin{aligned} \text{subj. to } y_n (\mathbf{w} \cdot \mathbf{x}_n + b) &\geq 1 - \xi_n & (\forall n) \\ \xi_n &\geq 0 & (\forall n) \end{aligned}$$

Now, let's play a thought experiment. Suppose someone handed you a solution to this optimization problem that consisted of weights (\mathbf{w}) and a bias (b), but they forgot to give you the slacks. Could you recover the slacks from the information you have?

In fact, the answer is yes! For simplicity, let's consider positive examples. Suppose that you look at some positive example \mathbf{x}_n . You need to figure out what the slack, ξ_n , would have been. There are two cases. Either $\mathbf{w} \cdot \mathbf{x}_n + b$ is at least 1 or it is not. If it's large enough, then you want to set $\xi_n = 0$. Why? It cannot be less than zero by the second constraint. Moreover, if you set it greater than zero, you will "pay" unnecessarily in the objective. So in this case, $\xi_n = 0$. Next, suppose that $\mathbf{w} \cdot \mathbf{x}_n + b = 0.2$, so it is not big enough. In order to satisfy the first constraint, you'll need to set $\xi_n \geq 0.8$. But because of the objective, you'll not want to set it any larger than necessary, so you'll set $\xi_n = 0.8$ exactly.

Following this argument through for both positive and negative points, if someone gives you solutions for \mathbf{w}, b , you can automatically compute the optimal ξ variables as:

$$\xi_n = \begin{cases} 0 & \text{if } y_n(\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1 \\ 1 - y_n(\mathbf{w} \cdot \mathbf{x}_n + b) & \text{otherwise} \end{cases} \quad (7.47)$$

In other words, the optimal value for a slack variable is *exactly* the hinge loss on the corresponding example! Thus, we can write the SVM objective as an *unconstrained* optimization problem:

$$\min_{\mathbf{w}, b} \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{large margin}} + C \underbrace{\sum_n \ell^{(\text{hin})}(y_n, \mathbf{w} \cdot \mathbf{x}_n + b)}_{\text{small slack}} \quad (7.48)$$

Multiplying this objective through by λ/C , we obtain exactly the regularized objective from Eq (7.8) with hinge loss as the loss function and the 2-norm as the regularizer!

7.8 Further Reading

TODO further reading

8 | BIAS AND FAIRNESS

Science and everyday life cannot
and should not be separated.

– Rosalind Franklin

Learning Objectives:

-

At the end of Chapter 1, you saw the “Russian Tank” example of a biased data set leading to a classifier that *seemed* like it was doing well, but was really relying on some artifact of the data collection process. As machine learning algorithms have a greater and greater impact on the real world, it is crucially important to ensure that they are making decisions based on the “right” aspects of the input, rather than exploiting arbitrary idiosyncracies of a particular training set.

For the rest of this chapter, we will consider two real world examples of bias issues that have had significant impact: the effect of gender in speech recognition systems and the effect of race in predicting criminal recidivism (i.e., will a convicted criminal commit further crimes if released).¹ The gender issue is that early speech recognition systems in cars failed to recognize the voices of many people who were not men. The race issue is that a specific recidivism predictor based on standard learning algorithms was biased against minorities.

Dependencies: Chapter 1, Chapter 3, Chapter 4, Chapter 5

¹ See [Autoblog](#) and [ProPublica](#) for press coverage of these two issues.

8.1 Train/Test Mismatch

One of the most common issues in bias is a mismatch between the training distribution and the testing distribution. In the running example of speech recognition failing to work on many non-men speakers, a large part of this happened because most of the training data on which the speech recognition system was trained was spoken by men. The learning algorithm learned—very well—how to recognize men’s speech, but its accuracy dropped significantly when faced with a different distribution of test data.

To understand why this happens, recall the Bayes Optimal classifier from Chapter 2. This was the classifier that magically knew the data distribution \mathcal{D} , and then when presented with an example x predicted $\operatorname{argmax}_y \mathcal{D}(x, y)$. This was optimal, but *only* because \mathcal{D} was the correct distribution. Even if one has access to the true distribution for male speakers, say $\mathcal{D}^{(\text{male})}$, the Bayes Optimal classifier under

$\mathcal{D}^{(\text{male})}$ will generally *not* be optimal under the distribution for any other gender.

Another example occurs in sentiment analysis. It is common to train sentiment analysis systems on data collected from reviews: product reviews, restaurant reviews, movie reviews, etc. This data is convenient because it includes both text (the review itself) and also a rating (typically one to five stars). This yields a “free” dataset for training a model to predict rating given text. However, one should be wary when running such a sentiment classifier on text *other than* the types of reviews it was trained on. One can easily imagine that a sentiment analyzer trained on movie reviews may not work so well when trying to detect sentiment in a politician’s speech. Even moving between one type of product and another can fail wildly: “very small” typically expresses positive sentiment for USB drives, but not for hotel rooms.

The issue of **train/test mismatch** has been widely studied under many different names: **covariate shift** (in statistics, the input features are called “covariates”), **sample selection bias** and **domain adaptation** are the most common. We will refer to this problem simply as **adaptation**. The adaptation challenge is: given training data from one “old” distribution, learn a classifier that does a good job on *another* related, but different, “new” distribution.

It’s important to recognize that in general, adaptation is impossible. For example, even if the task remains the same (positive/negative sentiment), if the old distribution is text reviews and the new distribution is images of text reviews, it’s hard to imagine that doing well on the old distribution says anything about the new. As a less extreme example, if the old distribution is movie reviews in English and the new distribution is movie reviews in Mandarin, it’s unlikely that adaptation will be easy.

These examples give rise to the tension in adaptation problems.

1. What does it mean for two distributions to be related? We might believe that “reviews of DVDs” and “reviews of movies” are “highly related” (though somewhat different: DVD reviews often discuss bonus features, quality, etc.); while “reviews of hotels” and “reviews of political policies” are “less related.” But how can we formalize this?
2. When two distributions *are* related, how can we build models that effectively share information between them?

8.2 Unsupervised Adaptation

The first type of adaptation we will cover is **unsupervised adaptation**. The setting is the following. There are two distributions, \mathcal{D}^{old} and \mathcal{D}^{new} . We have *labeled* training data from \mathcal{D}^{old} , say $(x_1, y_1), \dots, (x_N, y_N)$ totalling N examples. We also have M many *unlabeled examples* from \mathcal{D}^{new} : z_1, \dots, z_M . We assume that the examples live in the same space, \mathbb{R}^D . This is called *unsupervised adaptation* because we do not have access to any labels in the new distribution.²

² Sometimes this is called **semi-supervised adaptation** in the literature.

Our goal is to learn a classifier f that achieves low expected loss *under the new distribution*, \mathcal{D}^{new} . The challenge is that we do not have access to any labeled data from \mathcal{D}^{new} . As a warm-up, let's suppose that we have a black box machine learning algorithm \mathcal{A} that takes in *weighted* examples and produces a classifier. At the very least, this can be achieved using either undersampling or oversampling (see Section 6.1). We're going to attempt to *reweigh* the (old distribution) labeled examples based on how similar they are to the new distribution. This is justified using the **importance sampling** trick for switching expectations:

$$\text{test loss} \tag{8.1}$$

$$= \mathbb{E}_{(x,y) \sim \mathcal{D}^{\text{new}}} [\ell(y, f(x))] \quad \text{definition} \tag{8.2}$$

$$= \sum_{(x,y)} \mathcal{D}^{\text{new}}(x, y) \ell(y, f(x)) \quad \text{expand expectation} \tag{8.3}$$

$$= \sum_{(x,y)} \mathcal{D}^{\text{new}}(x, y) \frac{\mathcal{D}^{\text{old}}(x, y)}{\mathcal{D}^{\text{old}}(x, y)} \ell(y, f(x)) \quad \text{times one} \tag{8.4}$$

$$= \sum_{(x,y)} \mathcal{D}^{\text{old}}(x, y) \frac{\mathcal{D}^{\text{new}}(x, y)}{\mathcal{D}^{\text{old}}(x, y)} \ell(y, f(x)) \quad \text{rearrange} \tag{8.5}$$

$$= \mathbb{E}_{(x,y) \sim \mathcal{D}^{\text{old}}} \left[\frac{\mathcal{D}^{\text{new}}(x, y)}{\mathcal{D}^{\text{old}}(x, y)} \ell(y, f(x)) \right] \quad \text{definition} \tag{8.6}$$

What we have achieved here is rewriting the test loss, which is an expectation over \mathcal{D}^{new} , as an expectation over \mathcal{D}^{old} instead.³ This is useful because we have access to labeled examples from \mathcal{D}^{old} but not \mathcal{D}^{new} . The implicit suggested algorithm by this analysis is to train a classifier using our learning algorithm \mathcal{A} , but where each training example (x_n, y_n) is weighted according to the ratio $\mathcal{D}^{\text{new}}(x_n, y_n) / \mathcal{D}^{\text{old}}(x_n, y_n)$. Intuitively, this makes sense: the classifier is being told to pay more attention to training examples that have high probability under the new distribution, and less attention to training that have low probability under the new distribution.

The problem with this approach is that we do not have access to \mathcal{D}^{new} or \mathcal{D}^{old} , so we cannot compute this ratio and therefore cannot run this algorithm. One approach to this problem is to try to explic-

³ In this example, we assumed a discrete distribution; if the distributions are continuous, the sums are simply replaced with integrals.

itly estimate these distributions, a task known as **density estimation**. This is an *incredibly* difficult problem; far harder than the original adaptation problem.

A solution to this problem is to try to estimate the *ratio* directly, rather than separately estimating the two probability distributions⁴. The key idea is to think of the adaptation as follows. All examples are drawn according to some fixed base distribution $\mathcal{D}^{\text{base}}$. Some of these are selected to go into the new distribution, and some of them are selected to go into the old distribution. The mechanism for deciding which ones are kept and which are thrown out is governed by a *selection variable*, which we call s . The choice of selection-or-not, s , is based *only* on the input example x and not on its label. In particular, we define:

$$\mathcal{D}^{\text{old}}(x, y) \propto \mathcal{D}^{\text{base}}(x, y)p(s = 1 | x) \quad (8.7)$$

$$\mathcal{D}^{\text{new}}(x, y) \propto \mathcal{D}^{\text{base}}(x, y)p(s = 0 | x) \quad (8.8)$$

That is, the probability of drawing some pair (x, y) in the old distribution is proportional to the probability of first drawing that example according to the base distribution, and then the probability of selecting that particular example into the old distribution. If we can successfully estimate $p(s = 1 | x)$, then the ratio that we sought, then we can compute the importance ratio as:

$$\frac{\mathcal{D}^{\text{new}}(x, y)}{\mathcal{D}^{\text{old}}(x, y)} = \frac{\frac{1}{Z^{\text{new}}} \mathcal{D}^{\text{base}}(x, y)p(s = 0 | x)}{\frac{1}{Z^{\text{old}}} \mathcal{D}^{\text{base}}(x, y)p(s = 1 | x)} \quad \text{definition} \quad (8.9)$$

$$= \frac{\frac{1}{Z^{\text{new}}} p(s = 0 | x)}{\frac{1}{Z^{\text{old}}} p(s = 1 | x)} \quad \text{cancel base} \quad (8.10)$$

$$= Z \frac{p(s = 0 | x)}{p(s = 1 | x)} \quad \text{consolidate} \quad (8.11)$$

$$= Z \frac{1 - p(s = 1 | x)}{p(s = 1 | x)} \quad \text{binary selection} \quad (8.12)$$

$$= Z \left[\frac{1}{p(s = 1 | x)} - 1 \right] \quad \text{rearrange} \quad (8.13)$$

This means that if we can estimate the selection probability $p(s = 1 | x)$, we're done. We can therefore use $1/p(s = 1 | x_n) - 1$ as an example weight on example (x_n, y_n) when feeding these examples into our learning algorithm \mathcal{A} .

The remaining question is how to estimate $p(s = 1 | x_n)$. Recall that $s = 1$ denotes the case that x is selected into the old distribution and $s = 0$ denotes the case that x is selected into the new distribution. This means that predicting s is *exactly* a binary classification problem, where the “positive” class is the set of N examples from the old distribution and the “negative” class is the set of M examples from the new distribution.

⁴ Bickel et al. 2007

? What could go wrong if s got to look at the label, too?

? As a check: make sure that these weights are always non-negative. Furthermore, why is it okay to ignore the Z factor?

Algorithm 23 SELECTIONADAPTATION($\langle \langle x_n, y_n \rangle \rangle_{n=1}^N, \langle z_m \rangle_{m=1}^M, \mathcal{A}$)

```

1:  $D^{\text{dist}} \leftarrow \langle \langle x_n, +1 \rangle \rangle_{n=1}^N \cup \langle \langle z_m, -1 \rangle \rangle_{m=1}^M$  // assemble data for distinguishing
                                     // between old and new distributions
2:  $\hat{p} \leftarrow$  train logistic regression on  $D^{\text{dist}}$ 
3:  $D^{\text{weighted}} \leftarrow \langle \langle x_n, y_n, \frac{1}{\hat{p}(x_n)} - 1 \rangle \rangle_{n=1}^N$  // assemble weight classification
                                                         // data using selector
4: return  $\mathcal{A}(D^{\text{weighted}})$  // train classifier

```

This analysis gives rise to Algorithm 8.2, which consists of essentially two steps. The first is to train a logistic regression classifier⁵ to distinguish between old and new distributions. The second is to use that classifier to produce weights on the labeled examples from the old distribution and then train whatever learning algorithm you wish on that.

In terms of the questions posed at the beginning of this chapter, this approach to adaptation measures *nearness* of the two distributions by the degree to which the selection probability is constant. In particular, if the selection probability is independent of x , then the two distributions are identical. If the selection probabilities vary significantly as x changes, then the two distributions are considered very different. **More generally, if it is *easy* to train a classifier to distinguish between the old and new distributions, then they are very different.**

In the case of speech recognition failing as a function of gender, a core issue is that speech from men was massively over-represented in the training data but not the test data. When the selection logistic regression is trained, it is likely to say “old” on speech from men and “new” on other speakers, thereby *downweighting* the significance of male speakers and *upweighting* the significance of speakers of other genders on the final learned model. This would (hopefully) address many of the issues confounding that system.

⁵ The use of logistic regression is arbitrary: it need only be a classification algorithm that can produce probabilities.

Make up percentages for fraction of speakers who are male in the old and new distributions; estimate (you’ll have to make some assumptions) what the importance weights would look like in this case.

8.3 Supervised Adaptation

Unsupervised adaptation is very challenging because we never get to see test labels in the new distribution. In many ways, unsupervised adaptation attempts to *guard against* bad things happening. That is, if an old distribution training example looks very unlike the new distribution, it (and consequently its features) are downweighted so much as to be ignored. In supervised adaptation, we can hope for more: we can hope to actually do *better* on the new distribution than the old because we have labeled data there.

The typical setup is similar to the unsupervised case. There are

Algorithm 24 $\text{EASYADAPT}(\langle (x_n^{(old)}, y_n^{(old)}) \rangle_{n=1}^N, \langle (x_m^{(new)}, y_m^{(new)}) \rangle_{m=1}^M, \mathcal{A})$

```

1:  $D \leftarrow \langle \langle (x_n^{(old)}, x_n^{(old)}, 0), y_n^{(old)} \rangle \rangle_{n=1}^N \cup \langle \langle (x_m^{(new)}, 0, x_m^{(new)}), y_m^{(new)} \rangle \rangle_{m=1}^M$  // union
// of transformed data
2: return  $\mathcal{A}(D)$  // train classifier

```

two distributions, \mathcal{D}^{old} and \mathcal{D}^{new} , and our goal is to learn a classifier that does well in expectation on \mathcal{D}^{new} . However, now we have labeled data from both distributions: N labeled examples from \mathcal{D}^{old} and M labeled examples from \mathcal{D}^{new} . Call them $\langle x_n^{(old)}, y_n^{(old)} \rangle_{n=1}^N$ from \mathcal{D}^{old} and $\langle x_m^{(new)}, y_m^{(new)} \rangle_{m=1}^M$ from \mathcal{D}^{new} . Again, suppose that both $x_n^{(old)}$ and $x_m^{(new)}$ both live in \mathbb{R}^D .

One way of answer the question of “how do we share information between the two distributions” is to say: when the distributions agree on the value of a feature, let them share it, but when they disagree, allow them to learn separately. For instance, in a sentiment analysis task, \mathcal{D}^{old} might be reviews of electronics and \mathcal{D}^{new} might be reviews of hotel rooms. In both cases, if the review contains the word “awesome” then it’s probably a positive review, regardless of which distribution we’re in. We would want to *share* this information across distributions. On the other hand, “small” might be positive in electronics and negative in hotels, and we would like the learning algorithm to be able to learn separate information for that feature.

A very straightforward way to accomplish this is the **feature augmentation** approach⁶. This is a simple preprocessing step after which one can apply any learning algorithm. The idea is to create three versions of every feature: one that’s shared (for words like “awesome”), one that’s old-distribution-specific and one that’s new-distribution-specific. The mapping is:

⁶ Daumé III 2007

$$x_n^{(old)} \mapsto \left\langle \begin{array}{ccc} \text{shared} & \text{old-only} & \text{new-only} \\ x_n^{(old)} & , x_n^{(old)} & , \underbrace{0, 0, \dots, 0}_{D\text{-many}} \end{array} \right\rangle \quad (8.14)$$

$$x_m^{(new)} \mapsto \left\langle \begin{array}{ccc} x_m^{(new)} & , \underbrace{0, 0, \dots, 0}_{D\text{-many}} & , x_m^{(new)} \end{array} \right\rangle \quad (8.15)$$

Once you’ve applied this transformation, you can take the union of the (transformed) old and new labeled examples and feed the entire set into your favorite classification algorithm. That classification algorithm can then choose to share strength between the two distributions by using the “shared” features, if possible; or, if not, it can learn distribution-specific properties on the old-only or new-only parts. This is summarized in Algorithm 8.3.

Note that this algorithm can be *combined* with the instance weight-

ing (unsupervised) learning algorithm. In this case, the logistic regression separator should be trained on the *untransformed* data, and then weights can be used exactly as in Algorithm 8.2. This is particularly useful if you have access to *way* more old distribution data than new distribution data, and you don't want the old distribution data to "wash out" the new distribution data.

Although this approach is general, it is most effective when the two distributions are "not too close but not too far":

- If the distributions are too far, and there's little information to share, you're probably better off throwing out the old distribution data and training just on the (untransformed) new distribution data.
- If the distributions are too close, then you might as well just take the union of the (untransformed) old and new distribution data, and training on that.

In general, the interplay between how far the distributions are and how much new distribution data you have is complex, and you should always try "old only" and "new only" and "simple union" as baselines.

? Why is it crucial that the separator be trained on the untransformed data?

8.4 Fairness and Data Bias

Almost any data set in existence is biased in some way, in the sense that it captures an imperfect view of the world. The degree to which this bias is obvious can vary greatly:

- There might be obvious bias in the labeling. For instance, in criminology, learning to predict sentence lengths by predicting the sentences assigned by judges will simply learn to reproduce whatever bias already exists in judicial sentencing.
- There might be sample selection bias, as discussed early. In the same criminology example, the only people for which we have training data are those that have already been arrested, charged with and convicted of a crime; these processes are inherently biased, and so any model learned on this data may exhibit similar biases.
- The task itself might be biased because the designers had blindspots. An intelligent billboard that predicts the gender of the person walking toward it so as to show "better" advertisements may be trained as a binary classifier between male/female and thereby excludes anyone who does not fall in the gender binary. A similar example holds for a classifier that predicts political affiliation in

the US as Democrat/Republican, when in fact there are far more possibilities.

- There might be bias in the features or model structure. Machine translation systems often exhibit incorrect gender stereotypes⁷ because the relevant context, which would tell them the correct gender, is not encoded in the feature representation. Alone, or coupled with biased data, this leads to errors.
- The loss function may favor certain types of errors over others. You've already seen such an example: using zero/one loss on a highly imbalanced class problem will often lead to a learned model that ignores the minority class entirely.
- A deployed system creates feedback loops when it begins to consume its own output as new input. For instance, once a spam filter is in place, spammers will adjust their strategies, leading to distribution shift in the inputs. Or a car guidance system for predicting which roads are likely to be unoccupied may find those roads now occupied by other cars that it directed there.

⁷ For instance, many languages have verbal marking that agrees with the gender of the subject. In such cases, "the doctor treats ..." puts masculine markers on the translation of "treats" but "the nurse treats ..." uses feminine markers.

These are all difficult questions, none of which has easy answers. Nonetheless, it's important to be aware of how our systems may fail, especially as they take over more and more of our lives. Most computing, engineering and mathematical societies (like the ACM, IEEE, BCS, etc.) have codes of ethics, all of which include a statement about avoiding harm; the following is taken from the ACM code⁸:

To minimize the possibility of indirectly harming others, computing professionals must minimize malfunctions by following generally accepted standards for system design and testing. Furthermore, it is often necessary to assess the social consequences of systems to project the likelihood of any serious harm to others.

⁸ ACM Code of Ethics

In addition to ethical questions, there are often related *legal* questions. For example, US law prohibits discrimination by race and gender (among other "protected attributes") in processes like hiring and housing. The current legal mechanism for measuring discrimination is **disparate impact** and the **80% rule**. Informally, the 80% rule says that your rate of hiring women (for instance) must be at least 80% of your rate of hiring men. Formally, the rule states:

$$\Pr(y = +1 \mid G \neq \text{male}) \geq 0.8 \times \Pr(y = +1 \mid G = \text{male}) \quad (8.16)$$

Of course, gender/male can be replaced with any other protected attribute.

One non-solution to the disparate impact problem is to simply throw out protected attributes from your dataset. Importantly, yet unfortunately, this is **not** sufficient. Many other features often correlate

strongly with gender, race, and other demographic information, and just because you've removed *explicit* encodings of these factors does not imply that a model trained as such would satisfy the 80% rule. A natural question to ask is: can we build machine learning algorithms that have high accuracy but simultaneously avoid disparate impact?

Disparate impact is an imperfect measure of (un)fairness, and there are alternatives, each with its own pros and cons⁹. All of these rely on predefined categories of protected attributes, and a natural question is where these come from if not governmental regulation. Regardless of the measure you choose, the most important thing to keep in mind is that just because something comes from data, or is algorithmically implemented, does not mean it's fair.

⁹ Friedler et al. 2016, Hardt et al. 2016

8.5 How Badly can it Go?

To help better understand how badly things can go when the distribution over inputs changes, it's worth thinking about how to analyze this situation formally. Suppose we have two distributions \mathcal{D}^{old} and \mathcal{D}^{new} , and let's assume that the only thing that's different about these is the distribution they put on the inputs \mathcal{X} and not the outputs \mathcal{Y} . (We will return later to the usefulness of this assumption.) That is: $\mathcal{D}^{\text{old}}(x, y) = \mathcal{D}^{\text{old}}(x)\mathcal{D}(y | x)$ and $\mathcal{D}^{\text{new}}(x, y) = \mathcal{D}^{\text{new}}(x)\mathcal{D}(y | x)$, where $\mathcal{D}(y | x)$ is shared between them.

Let's say that you've learned a classifier f that does really well on \mathcal{D}^{old} and achieves some *test error* of $\epsilon^{(\text{old})}$. That is:

$$\epsilon^{(\text{old})} = \mathbb{E}_{\substack{x \sim \mathcal{D}^{\text{old}} \\ y \sim \mathcal{D}(\cdot | x)}} [\mathbf{1}[f(x) \neq y]] \quad (8.17)$$

The question is: how badly can f do on the new distribution?

We can calculate this directly.

$$\epsilon^{(\text{new})} = \mathbb{E}_{\substack{x \sim \mathcal{D}^{\text{new}} \\ y \sim \mathcal{D}(\cdot | x)}} [\mathbf{1}[f(x) \neq y]] \quad (8.18)$$

$$= \int_{\mathcal{X}} \int_{\mathcal{Y}} \mathcal{D}^{\text{new}}(x) \mathcal{D}(y | x) \mathbf{1}[f(x) \neq y] dy dx \quad \text{def. of } \mathbb{E} \quad (8.19)$$

$$= \int_{\mathcal{X}} \left(\mathcal{D}^{\text{new}}(x) - \mathcal{D}^{\text{old}}(x) + \mathcal{D}^{\text{old}}(x) \right) \times \int_{\mathcal{Y}} \mathcal{D}(y | x) \mathbf{1}[f(x) \neq y] dy dx \quad \text{add zero} \quad (8.20)$$

$$= \epsilon^{(\text{old})} + \int_{\mathcal{X}} \left(\mathcal{D}^{\text{new}}(x) - \mathcal{D}^{\text{old}}(x) \right) \times \int_{\mathcal{Y}} \mathcal{D}(y | x) \mathbf{1}[f(x) \neq y] dy dx \quad \text{def. } \epsilon^{\text{old}} \quad (8.21)$$

$$\leq \epsilon^{(\text{old})} + \int_{\mathcal{X}} \left| \mathcal{D}^{\text{new}}(x) - \mathcal{D}^{\text{old}}(x) \right| (1) dx \quad \text{worst case } \int_{\mathcal{Y}} \quad (8.22)$$

$$= \epsilon^{(\text{old})} + 2 \left| \mathcal{D}^{\text{new}} - \mathcal{D}^{\text{old}} \right|_{\text{Var}} \quad \text{def. } |\cdot|_{\text{Var}} \quad (8.23)$$

Here, $|\cdot|_{\text{Var}}$ is the **total variation distance** (or **variational distance**) between two probability distributions, defined as:

$$|P - Q|_{\text{Var}} = \sup_e |P(e) - Q(e)| = \frac{1}{2} \int_{\mathcal{X}} |P(x) - Q(x)| dx \quad (8.24)$$

Is a standard measure of dissimilarity between probability distributions. In particular, the variational distance is the largest difference in probability that P and Q assign to the same event (in this case, the event is an example x).

The bound tells us that we might not be able to hope for very good error on the new distribution, even if we have very good error on the old distribution, when \mathcal{D}^{old} and \mathcal{D}^{new} are very different (assign very different probabilities to the same event). Of course, this is an *upper bound*, and possibly a very loose one.

The second observation is that we barely used the assumption that \mathcal{D}^{new} and \mathcal{D}^{old} share the same label distribution $\mathcal{D}(y | x)$ in the above analysis. (In fact, as an exercise, repeat the analysis *without* this assumption. It will still go through.) In general, this assumption buys us very little. In an extreme case, \mathcal{D}^{new} and \mathcal{D}^{old} can essentially “encode” which distribution a given x came from in one of the features. Once the origin distribution is completely encoded in a feature, then $\mathcal{D}(y | x)$ could look at the encoding, and completely flip the label based on which distribution it’s from. How could \mathcal{D}^{new} and \mathcal{D}^{old} encode the distribution? One way would be to set the 29th decimal digit of feature 1 to an odd value in the old distribution and an even value in the new distribution. This tiny change will be essentially imperceptible if one looks at the data, but would give $\mathcal{D}(y | x)$ enough power to make our lives miserable.

If we want a way out of this dilemma, we need more technology. The core idea is that if we’re learning a function f from some hypothesis class \mathcal{F} , and this hypothesis class isn’t rich enough to peek at the 29th decimal digit of feature 1, then perhaps things are not as bad as they could be. This motivates the idea of looking at a measure of distance between probability distributions that *depends on the hypothesis class*. A popular measure is the **$d_{\mathcal{A}}$ -distance** or the **discrepancy**. The discrepancy measure distances between probability distributions based on how much two function f and f' in the hypothesis class can disagree on their labels. Let:

$$\epsilon_P(f, f') = \mathbb{E}_{x \sim P} [\mathbf{1}[f(x) \neq f'(x)]] \quad (8.25)$$

You can think of $\epsilon_P(f, f')$ as the *error* of f' when the ground truth is given by f , where the error is taken with respect to examples drawn from P . Given a hypothesis class \mathcal{F} , the discrepancy between P and Q is defined as:

$$d_{\mathcal{A}}(P, Q) = \max_{f, f' \in \mathcal{F}} |\epsilon_P(f, f') - \epsilon_Q(f, f')| \quad (8.26)$$

The discrepancy very much has the flavor of a classifier: if you think of f as providing the ground truth labeling, then the discrepancy is large whenever there exists a function f' that agrees strongly with f on distribution P but *disagrees* strongly with f on Q . This feels natural: if all functions behave similarly on P and Q , then the discrepancy will be small, and we also expect to be able to generalize across these two distributions easily.

One very attractive property of the discrepancy is that you can estimate it from finite *unlabeled* samples from \mathcal{D}^{old} and \mathcal{D}^{new} . Although not obvious at first, the discrepancy is very closely related to a quantity we saw earlier in unsupervised adaptation: a classifier that distinguishes between \mathcal{D}^{old} and \mathcal{D}^{new} . In fact, the discrepancy is precisely twice the *accuracy* of the best classifier from \mathcal{H} at separating \mathcal{D}^{old} from \mathcal{D}^{new} .

How does this work in practice? Exactly as in the section on unsupervised adaptation, we train a classifier to distinguish between \mathcal{D}^{old} and \mathcal{D}^{new} . It needn't be a probabilistic classifier; any binary classifier is fine. This classifier, the “domain separator,” will have some (heldout) accuracy, call it acc . The discrepancy is then exactly $d_{\mathcal{A}} = 2(\text{acc} - 0.5)$.

Intuitively the accuracy of the domain separator is a natural measure of how different the two distributions are. If the two distributions are identical, you shouldn't expect to get very good accuracy at separating them. In particular, you expect the accuracy to be around 0.5, which puts the discrepancy at zero. On the other hand, if the two distributions are really far apart, separation is easy, and you expect an accuracy of about 1, yielding a discrepancy also of about 1.

One can, in fact, prove a generalization bound—generalizing from finite samples from \mathcal{D}^{old} to expected loss on \mathcal{D}^{new} —based on the discrepancy¹⁰.

¹⁰ Ben-David et al. 2007

Theorem 9 (Unsupervised Adaptation Bound). *Given a fixed representation and a fixed hypothesis space \mathcal{F} , let $f \in \mathcal{F}$ and let $\epsilon^{(\text{best})} = \min_{f^* \in \mathcal{F}} \frac{1}{2} [\epsilon^{(\text{old})}(f^*) + \epsilon^{(\text{new})}(f^*)]$, then, for all $f \in \mathcal{F}$:*

$$\underbrace{\epsilon^{(\text{new})}(f)}_{\text{error on } \mathcal{D}^{\text{new}}} \leq \underbrace{\epsilon^{(\text{old})}(f)}_{\text{error on } \mathcal{D}^{\text{old}}} + \underbrace{\epsilon^{(\text{best})}}_{\text{minimal avg error}} + \underbrace{d_{\mathcal{A}}(\mathcal{D}^{\text{old}}, \mathcal{D}^{\text{new}})}_{\text{distance}} \quad (8.27)$$

In this bound, $\epsilon^{(\text{best})}$ denotes the error rate of the best possible classifier from \mathcal{F} , where the error rate is measured as the *average*

error rate on \mathcal{D}^{new} and \mathcal{D}^{old} ; this term ensures that at least some good classifier exists that does well on both distributions.

The main practical use of this result is that it suggests a way to look for representations that are good for adaptation: on the one hand, we should try to get our training error (on \mathcal{D}^{old}) as low as possible; on the other hand, we want to make sure that it is *hard* to distinguish between \mathcal{D}^{old} and \mathcal{D}^{new} in this representation.

8.6 *Further Reading*

TODO further reading

9 | PROBABILISTIC MODELING

The world is noisy and messy. You need to deal with the noise and uncertainty.

– Daphne Koller

MANY OF THE MODELS AND ALGORITHMS you have learned about thus far are relatively *disconnected*. There is an alternative view of machine learning that unites and generalizes much of what you have already learned. This is the **probabilistic modeling** framework, in which you will explicitly think of learning as a problem of **statistical inference**.

In this chapter, you will learn about two flavors of probabilistic models: generative and conditional. You will see that many of the approaches (both supervised and unsupervised) we have seen already can be cast as probabilistic models. Through this new view, you will be able to develop learning algorithms that have inductive biases closer to what you, as a designer, believe. Moreover, the two chapters that follow will make heavy use of the probabilistic modeling approach to open doors to other learning problems.

Learning Objectives:

- Define the generative story for a naive Bayes classifier.
- Derive relative frequency as the solution to a constrained optimization problem.
- Compare and contrast generative, conditional and discriminative learning.
- Explain when generative models are likely to fail.
- Derive logistic loss with an ℓ_2 regularizer from a probabilistic perspective.

Dependencies:

9.1 Classification by Density Estimation

Recall from Chapter 2 that if we had access to the underlying probability distribution \mathcal{D} , then we could form a Bayes optimal classifier as:

$$f^{(\text{BO})}(\hat{x}) = \arg \max_{\hat{y} \in \mathcal{Y}} \mathcal{D}(\hat{x}, \hat{y}) \quad (9.1)$$

Unfortunately, no one gave you this distribution, but the optimality of this approach suggests that good way to build a classifier is to try to *estimate* \mathcal{D} . In other words, you try to learn a distribution $\hat{\mathcal{D}}$, which you hope to very similar to \mathcal{D} , and then use this distribution for classification. Just as in the preceding chapters, you can try to form your estimate of \mathcal{D} based on a finite training set.

The most direct way that you can attempt to construct such a probability distribution is to select a *family* of parametric distributions. For instance, a Gaussian (or Normal) distribution is parametric: its parameters are its mean and covariance. The job of learning is

MATH REVIEW | RULES OF PROBABILITY

A probability distribution p specifies the likelihood of an event e , where $p(e) \in [0, 1]$. It's often convenient to think of events as "configurations of the world", so $p(e)$ says "how likely is it that the world is in configuration e ." Often world configurations are built up of smaller pieces, for instance you might say " e = the configuration in which it is rainy, windy and cold." Formally, we might write this as " $e = \{\text{Weather} = \text{rainy}, \text{Wind} = \text{windy}, \text{Temperature} = \text{cold}\}$ ", where we've used a convention that **random variables** (like Temperature) are capitalized and their instantiations (like cold) are lower case. Considering this event, we want to evaluate $p(\text{Weather} = \text{rainy}, \text{Wind} = \text{windy}, \text{Temperature} = \text{cold})$, or more generally $p(A = a, B = b, C = c)$ for some random variables A, B and C , and some instantiations of those random variables a, b and c respectively.

There are a few standard rules of probability that we will use regularly:

sum-to-one: if you sum over all possible configurations of the world, p sums to one: $\sum_e p(E = e) = 1$.

marginalization: you can sum out one random variable to remove it from the world: $\sum_a p(A = a, B = b) = p(B = b)$.

chain rule: if a world configuration consists of two or more random variables, you can evaluate the likelihood of the world one step at a time: $p(A = a, B = b) = p(A = a)p(B = b | A = a)$. Events are unordered, so you can also get $p(A = a, B = b) = p(B = b)p(A = a | B = b)$.

Bayes rule: combining the two chain rule equalities and dividing, we can relate a conditional probability in one direction with that in the other direction: $p(A = a | B = b) = p(A = a)p(B = b | A = a) / p(B = b)$.

Figure 9.1:

then to infer which parameters are "best" as far as the observed training data is concerned, as well as whatever inductive bias you bring. A key assumption that you will need to make is that the training data you have access to is drawn **independently** from \mathcal{D} . In particular, as you draw examples $(x_1, y_1) \sim \mathcal{D}$ then $(x_2, y_2) \sim \mathcal{D}$ and so on, the n th draw (x_n, y_n) is drawn from \mathcal{D} and *does not* otherwise depend on the previous $n - 1$ samples. This assumption is usually false, but is also usually sufficiently close to being true to be useful. Together with the assumption that all the training data is drawn from the same distribution \mathcal{D} leads to the **i.i.d. assumption** or **independently and identically distributed** assumption. This is a key assumption in almost all of machine learning.

9.2 Statistical Estimation

Suppose you need to model a coin that is possibly biased (you can think of this as modeling the *label* in a binary classification problem), and that you observe data HHTH (where H means a flip came up heads

and T means it came up tails). You can assume that all the flips came from the same coin, and that each flip was independent (hence, the data was i.i.d.). Further, you may choose to believe that the coin has a fixed probability β of coming up heads (and hence $1 - \beta$ of coming up tails). Thus, the parameter of your model is simply the scalar β .

The most basic computation you might perform is **maximum likelihood estimation**: namely, select the parameter β that maximizes the probability of the data under that parameter. In order to do so, you need to compute the probability of the data:

$$p_{\beta}(D) = p_{\beta}(\text{HHTH}) \quad \text{definition of } D \quad (9.2)$$

$$= p_{\beta}(\text{H})p_{\beta}(\text{H})p_{\beta}(\text{T})p_{\beta}(\text{H}) \quad \text{data is independent} \quad (9.3)$$

$$= \beta\beta(1 - \beta)\beta \quad (9.4)$$

$$= \beta^3(1 - \beta) \quad (9.5)$$

$$= \beta^3 - \beta^4 \quad (9.6)$$

Thus, if you want the parameter β that maximizes the probability of the data, you can take the derivative of $\beta^3 - \beta^4$ with respect to β , set it equal to zero and solve for β :

$$\frac{\partial}{\partial \beta} [\beta^3 - \beta^4] = 3\beta^2 - 4\beta^3 \quad (9.7)$$

$$4\beta^3 = 3\beta^2 \quad (9.8)$$

$$\iff 4\beta = 3 \quad (9.9)$$

$$\iff \beta = \frac{3}{4} \quad (9.10)$$

Thus, the maximum likelihood β is 0.75, which is probably what you would have selected by intuition. You can solve this problem more generally as follows. If you have H -many heads and T -many tails, the probability of your data sequence is $\beta^H(1 - \beta)^T$. You can try to take the derivative of this with respect to β and follow the same recipe, but all of the products make things difficult. A more friendly solution is to work with the **log likelihood** or **log probability** instead. The log likelihood of this data sequence is $H \log \beta + T \log(1 - \beta)$. Differentiating with respect to β , you get $H/\beta - T/(1 - \beta)$. To solve, you obtain $H/\beta = T/(1 - \beta)$ so $H(1 - \beta) = T\beta$. Thus $H - H\beta = T\beta$ and so $H = (H + T)\beta$, finally yielding that $\beta = H/(H + T)$ or, simply, the fraction of observed data that came up heads. In this case, the maximum likelihood estimate is nothing but the relative frequency of observing heads!

Now, suppose that instead of flipping a coin, you're rolling a K -sided die (for instance, to pick the label for a multiclass classification problem). You might model this by saying that there are parameters $\theta_1, \theta_2, \dots, \theta_K$ specifying, respectively, the probabilities that any given

?

Describe a case in which at least one of the assumptions we are making about the coin flip is false.

?

How do you know that the solution of $\beta = H/(H + T)$ is actually a maximum?

side comes up on a role. Since these are themselves probabilities, each θ_k should be at least zero, and the sum of the θ_k s should be one. Given a data set that consists of x_1 rolls of 1, x_2 rolls of 2 and so on, the probability of this data is $\prod_k \theta_k^{x_k}$, yielding a log probability of $\sum_k x_k \log \theta_k$. If you pick some particular parameter, say θ_3 , the derivative of this with respect to θ_3 is x_3/θ_3 , which you want to equate to zero. This leads to $\dots \theta_3 \rightarrow \infty$.

This is obviously “wrong.” From the mathematical formulation, it’s correct: in fact, setting all of the θ_k s to ∞ *does* maximize $\prod_k \theta_k^{x_k}$ for any (non-negative) x_k s. The problem is that you need to constrain the θ s to sum to one. In particular, you have a constraint that $\sum_k \theta_k = 1$ that you forgot to enforce. A convenient way to enforce such constraints is through the technique of **Lagrange multipliers**. To make this problem consistent with standard minimization problems, it is convenient to minimize negative log probabilities, instead of maximizing log probabilities. Thus, the *constrained* optimization problem is:

$$\begin{aligned} \min_{\theta} \quad & -\sum_k x_k \log \theta_k \\ \text{subj. to} \quad & \sum_k \theta_k - 1 = 0 \end{aligned} \tag{9.11}$$

The Lagrange multiplier approach involves adding a new variable λ to the problem (called the **Lagrange variable**) corresponding to the constraint, and to use that to move the constraint into the objective. The result, in this case, is:

$$\max_{\lambda} \min_{\theta} \quad -\sum_k x_k \log \theta_k - \lambda \left(\sum_k \theta_k - 1 \right) \tag{9.12}$$

Turning a constrained optimization problem into its corresponding **Lagrangian** is straightforward. The mystical aspect is why it works. In this case, the idea is as follows. Think of λ as an adversary: λ is trying to maximize this function (you’re trying to minimize it). If you pick some parameters θ that actually satisfy the constraint, then the green term in Eq (9.12) goes to zero, and therefore λ does not matter: the adversary cannot do anything. On the other hand, if the constraint is even *slightly* unsatisfied, then λ can tend toward $+\infty$ or $-\infty$ to blow up the objective. So, in order to have a non-infinite objective value, the optimizer *must* find values of θ that satisfy the constraint.

If we solve the *inner* optimization of Eq (9.12) by differentiating with respect to θ_1 , we get $x_1/\theta_1 = \lambda$, yielding $\theta_1 = x_1/\lambda$. In general, the solution is $\theta_k = x_k/\lambda$. Remembering that the goal of λ is to enforce the sums-to-one constraint, we can set $\lambda = \sum_k x_k$ and verify

that this is a solution. Thus, our optimal $\theta_k = x_k / \sum_k x_k$, which again completely corresponds to intuition.

9.3 Naive Bayes Models

Now, consider the binary classification problem. You are looking for a *parameterized* probability distribution that can describe the training data you have. To be concrete, your task might be to predict whether a movie review is positive or negative (label) based on what words (features) appear in that review. Thus, the probability for a *single* data point can be written as:

$$p_{\theta}((y, \mathbf{x})) = p_{\theta}(y, x_1, x_2, \dots, x_D) \quad (9.13)$$

The challenge in working with a probability distribution like Eq (9.13) is that it's a distribution over a *lot* of variables. You can try to simplify it by applying the **chain rule** of probabilities:

$$p_{\theta}(x_1, x_2, \dots, x_D, y) = p_{\theta}(y) p_{\theta}(x_1 | y) p_{\theta}(x_2 | y, x_1) p_{\theta}(x_3 | y, x_1, x_2) \dots p_{\theta}(x_D | y, x_1, x_2, \dots, x_{D-1}) \quad (9.14)$$

$$= p_{\theta}(y) \prod_d p_{\theta}(x_d | y, x_1, \dots, x_{d-1}) \quad (9.15)$$

At this point, this equality is *exact* for any probability distribution. However, it might be difficult to craft a probability distribution for the 10000th feature, given the previous 9999. Even if you could, it might be difficult to accurately estimate it. At this point, you can make assumptions. A classic assumption, called the **naive Bayes assumption**, is that *the features are independent, conditioned on the label*. In the movie review example, this is saying that *once you know that it's a positive review*, the probability that the word “excellent” appears is independent of whether “amazing” also appeared. (Note that this does *not* imply that these words are independent when you don't know the label—they most certainly are not.) Formally this assumption states that:

$$\textbf{Assumption: } p(x_d | y, x_{d'}) = p(x_d | y) \quad , \quad \forall d \neq d' \quad (9.16)$$

Under this assumption, you can simplify Eq (9.15) to:

$$p_{\theta}((y, \mathbf{x})) = p_{\theta}(y) \prod_d p_{\theta}(x_d | y) \quad \textbf{naive Bayes assumption} \quad (9.17)$$

At this point, you can start parameterizing p . Suppose, for now, that your labels are binary *and* your features are also binary. In this case, you could model the label as a biased coin, with probability of heads (e.g., positive review) given by θ_0 . Then, for each label, you

can imagine having one (biased) coin for each feature. So if there are D -many features, you'll have $1 + 2D$ total coins: one for the label (call it θ_0) and one for each label/feature combination (call these θ_{+1} and as θ_{-1}). In the movie review example, we might expect $\theta_0 \approx 0.4$ (forty percent of movie reviews are positive) and also that θ_{+1} might give high probability to words like “excellent” and “amazing” and “good” and θ_{-1} might give high probability to words like “terrible” and “boring” and “hate”. You can rewrite the probability of a single example as follows, eventually leading to the log probability of the entire data set:

$$p_{\theta}((y, x)) = p_{\theta}(y) \prod_d p_{\theta}(x_d | y) \quad \text{naive Bayes assumption} \quad (9.18)$$

$$= \left(\theta_0^{[y=+1]} (1 - \theta_0)^{[y=-1]} \right) \prod_d \theta_{(y),d}^{[x_d=1]} (1 - \theta_{(y),d})^{[x_d=0]} \quad \text{model assumptions} \quad (9.19)$$

Solving for θ_0 is identical to solving for the biased coin case from before: it is just the relative frequency of positive labels in your data (because θ_0 doesn't depend on x at all). For the other parameters, you can repeat the same exercise as before for each of the $2D$ coins independently. This yields:

$$\hat{\theta}_0 = \frac{1}{N} \sum_n [y_n = +1] \quad (9.20)$$

$$\hat{\theta}_{(+1),d} = \frac{\sum_n [y_n = +1 \wedge x_{n,d} = 1]}{\sum_n [y_n = +1]} \quad (9.21)$$

$$\hat{\theta}_{(-1),d} = \frac{\sum_n [y_n = -1 \wedge x_{n,d} = 1]}{\sum_n [y_n = -1]} \quad (9.22)$$

In the case that the features are *not* binary, you need to choose a different model for $p(x_d | y)$. The model we chose here is the **Bernoulli distribution**, which is effectively a distribution over independent coin flips. For other types of data, other distributions become more appropriate. The die example from before corresponds to a **discrete distribution**. If the data is continuous, you might choose to use a **Gaussian distribution** (aka **Normal distribution**). The choice of distribution is a form of **inductive bias** by which you can inject your knowledge of the problem into the learning algorithm.

9.4 Prediction

Consider the *predictions* made by the naive Bayes model with Bernoulli features in Eq (9.18). You can better understand this model by considering its decision boundary. In the case of probabilistic models,

MATH REVIEW | COMMON PROBABILITY DISTRIBUTIONS

There are a few common probability distributions that we use in this book. The first is the Bernoulli distribution, which models binary outcomes (like coin flips). A Bernoulli distribution, $\mathcal{Ber}(\theta)$ is parameterized by a single scalar value $\theta \in [0, 1]$ that represents the probability of heads. The likelihood function is $\mathcal{Ber}(x | \theta) = \theta^x(1 - \theta)^{1-x}$. The generalization of the Bernoulli to more than two possible outcomes (like rolls of a die) is the Discrete distribution, $\mathcal{Disc}(\theta)$. If the die has K sides, then $\theta \in \mathbb{R}^K$ with all entries non-negative and $\sum_k \theta_k = 1$. θ_k is the probability that the die comes up on side k . The likelihood function is $\mathcal{Disc}(x | \theta) = \prod_k \theta_k^{1[x=k]}$. The Binomial distribution is just like the Bernoulli distribution but for multiple flips of the rather than a single flip; its likelihood is $\binom{n}{k} \theta^k (1-\theta)^{n-k}$, where n is the number of flips and k is the number of heads. The Multinomial distribution extends the Discrete distribution also to multiple rolls; its likelihood is $\frac{n!}{\prod_k x_k!} \prod_k \theta_k^{x_k}$, where n is the total number of rolls and x_k is the number of times the die came up on side k (so $\sum_k x_k = n$). The preceding distributions are all discrete.

There are two common continuous distributions we need. The first is the Uniform distribution, $\mathcal{Uni}(a, b)$ which is uniform over the closed range $[a, b]$. Its density function is $\mathcal{Uni}(x | a, b) = \frac{1}{b-a} \mathbf{1}[x \in [a, b]]$. Finally, the Gaussian distribution is parameterized by a mean μ and variance σ^2 and has density $\mathcal{Nor}(x | \mu, \sigma^2) = (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right]$.

Figure 9.2:

the decision boundary is the set of inputs for which the likelihood of $y = +1$ is precisely 50%. Or, in other words, the set of inputs x for which $p(y = +1 | x) / p(y = -1 | x) = 1$. In order to do this, the first thing to notice is that $p(y | x) = p(y, x) / p(x)$. In the ratio, the $p(x)$ terms cancel, leaving $p(y = +1, x) / p(y = -1, x)$. Instead of computing this ratio, it is easier to compute the **log-likelihood ratio** (or LLR), $\log p(y = +1, x) - \log p(y = -1, x)$, computed below:

$$\begin{aligned} \text{LLR} &= \log \left[\theta_0 \prod_d \theta_{(+1),d}^{[x_d=1]} (1 - \theta_{(+1),d})^{[x_d=0]} \right] \\ &\quad - \log \left[(1 - \theta_0) \prod_d \theta_{(-1),d}^{[x_d=1]} (1 - \theta_{(-1),d})^{[x_d=0]} \right] \end{aligned} \tag{9.23}$$

model assumptions

$$\begin{aligned} &= \log \theta_0 - \log(1 - \theta_0) + \sum_d [x_d = 1] \left(\log \theta_{(+1),d} - \log \theta_{(-1),d} \right) \\ &\quad + \sum_d [x_d = 0] \left(\log(1 - \theta_{(+1),d}) - \log(1 - \theta_{(-1),d}) \right) \end{aligned} \tag{9.24}$$

take logs and rearrange

$$\begin{aligned} &= \sum_d x_d \log \frac{\theta_{(+1),d}}{\theta_{(-1),d}} + \sum_d (1 - x_d) \log \frac{1 - \theta_{(+1),d}}{1 - \theta_{(-1),d}} + \log \frac{\theta_0}{1 - \theta_0} \end{aligned} \tag{9.25}$$

simplify log terms

$$= \sum_d x_d \left[\log \frac{\theta_{(+1),d}}{\theta_{(-1),d}} - \log \frac{1 - \theta_{(+1),d}}{1 - \theta_{(-1),d}} \right] + \sum_d \log \frac{1 - \theta_{(+1),d}}{1 - \theta_{(-1),d}} + \log \frac{\theta_0}{1 - \theta_0} \quad \text{group } x\text{-terms} \quad (9.26)$$

$$= \mathbf{x} \cdot \mathbf{w} + b \quad (9.27)$$

$$w_d = \log \frac{\theta_{(+1),d}(1 - \theta_{(-1),d})}{\theta_{(-1),d}(1 - \theta_{(+1),d})} \quad , \quad b = \sum_d \log \frac{1 - \theta_{(+1),d}}{1 - \theta_{(-1),d}} + \log \frac{\theta_0}{1 - \theta_0} \quad (9.28)$$

The result of the algebra is that the naive Bayes model has precisely the form of a linear model! Thus, like perceptron and many of the other models you've previous studied, the decision boundary is linear.

9.5 Generative Stories

A useful way to develop probabilistic models is to tell a **generative story**. This is a *fictional* story that explains how you believe your training data came into existence. To make things interesting, consider a multiclass classification problem, with continuous features modeled by independent Gaussians. Since the label can take values $1 \dots K$, you can use a discrete distribution (die roll) to model it (as opposed to the Bernoulli distribution from before):

1. **For each example** $n = 1 \dots N$:
 - (a) **Choose a label** $y_n \sim \text{Disc}(\boldsymbol{\theta})$
 - (b) **For each feature** $d = 1 \dots D$:
 - i. **Choose feature value** $x_{n,d} \sim \text{Nor}(\mu_{y_n,d}, \sigma_{y_n,d}^2)$

This generative story can be directly translated into a likelihood function by replacing the “for each”s with products:

$$p(D) = \overbrace{\prod_n \underbrace{\theta_{y_n}}_{\text{choose label}} \prod_d \underbrace{\frac{1}{\sqrt{2\pi\sigma_{y_n,d}^2}} \exp \left[-\frac{1}{2\sigma_{y_n,d}^2} (x_{n,d} - \mu_{y_n,d})^2 \right]}_{\substack{\text{choose feature value} \\ \text{for each feature}}}}^{\text{for each example}} \quad (9.29)$$

You can take logs to arrive at the log-likelihood:

$$\log p(D) = \sum_n \left[\log \theta_{y_n} + \sum_d -\frac{1}{2} \log(\sigma_{y_n,d}^2) - \frac{1}{2\sigma_{y_n,d}^2} (x_{n,d} - \mu_{y_n,d})^2 \right] + \text{const} \quad (9.30)$$

To optimize for θ , you need to add a “sums to one” constraint as before. This leads to the previous solution where the θ_k s are proportional to the number of examples with label k . In the case of the μ s you can take a derivative with respect to, say $\mu_{k,i}$ and obtain:

$$\frac{\partial \log p(D)}{\partial \mu_{k,i}} = \frac{\partial}{\partial \mu_{k,i}} - \sum_n \sum_d \frac{1}{2\sigma_{y_n,d}^2} (x_{n,d} - \mu_{y_n,d})^2 \quad \text{ignore irrelevant terms} \quad (9.31)$$

$$= \frac{\partial}{\partial \mu_{k,i}} - \sum_{n:y_n=k} \frac{1}{2\sigma_{k,d}^2} (x_{n,i} - \mu_{k,i})^2 \quad \text{ignore irrelevant terms} \quad (9.32)$$

$$= \sum_{n:y_n=k} \frac{1}{\sigma_{k,d}^2} (x_{n,i} - \mu_{k,i}) \quad \text{take derivative} \quad (9.33)$$

Setting this equal to zero and solving yields:

$$\mu_{k,i} = \frac{\sum_{n:y_n=k} x_{n,i}}{\sum_{n:y_n=k} 1} \quad (9.34)$$

Namely, the sample mean of the i th feature of the data points that fall in class k . A similar analysis for $\sigma_{k,i}^2$ yields:

$$\frac{\partial \log p(D)}{\partial \sigma_{k,i}^2} = \frac{\partial}{\partial \sigma_{k,i}^2} - \sum_{y:y_n=k} \left[\frac{1}{2} \log(\sigma_{k,i}^2) + \frac{1}{2\sigma_{k,i}^2} (x_{n,i} - \mu_{k,i})^2 \right] \quad \text{ignore irrelevant terms} \quad (9.35)$$

$$= - \sum_{y:y_n=k} \left[\frac{1}{2\sigma_{k,i}^2} - \frac{1}{2(\sigma_{k,i}^2)^2} (x_{n,i} - \mu_{k,i})^2 \right] \quad \text{take derivative} \quad (9.36)$$

$$= \frac{1}{2\sigma_{k,i}^4} \sum_{y:y_n=k} \left[(x_{n,i} - \mu_{k,i})^2 - \sigma_{k,i}^2 \right] \quad \text{simplify} \quad (9.37)$$

You can now set this equal to zero and solve, yielding:

$$\sigma_{k,i}^2 = \frac{\sum_{n:y_n=k} (x_{n,i} - \mu_{k,i})^2}{\sum_{n:y_n=k} 1} \quad (9.38)$$

Which is just the sample variance of feature i for class k .

9.6 Conditional Models

In the foregoing examples, the task was formulated as attempting to model the **joint** distribution of (x, y) pairs. This may seem wasteful: at prediction time, all you care about is $p(y | x)$, so why not model it directly?

?

What would the estimate be if you decided that, for a given class k , all features had equal variance? What if you assumed feature i had equal variance for each class? Under what circumstances might it be a good idea to make such assumptions?

Starting with the case of regression is actually somewhat simpler than starting with classification in this case. Suppose you “believe” that the relationship between the real value y and the vector x should be linear. That is, you expect that $y = w \cdot x + b$ *should hold* for some parameters (w, b) . Of course, the data that you get does not exactly obey this: that’s fine, you can think of deviations from $y = w \cdot x + b$ as *noise*. To form a probabilistic model, you must assume some distribution over noise; a convenient choice is zero-mean Gaussian noise. This leads to the following generative story:

1. For each example $n = 1 \dots N$:

- (a) Compute $t_n = w \cdot x_n + b$
- (b) Choose noise $e_n \sim \text{Nor}(0, \sigma^2)$
- (c) Return $y_n = t_n + e_n$

In this story, the variable t_n stands for “target.” It is the noiseless variable that you do not get to observe. Similarly e_n is the error (noise) on example n . The value that you actually get to observe is $y_n = t_n + e_n$. See Figure 9.3.

A basic property of the Gaussian distribution is additivity. Namely, that if $a \sim \text{Nor}(\mu, \sigma^2)$ and $b = a + c$, then $b \sim \text{Nor}(\mu + c, \sigma^2)$. Given this, from the generative story above, you can derive a shorter generative story:

1. For each example $n = 1 \dots N$:

- (a) Choose $y_n \sim \text{Nor}(w \cdot x_n + b, \sigma^2)$

Reading off the log likelihood of a dataset from this generative story, you obtain:

$$\log p(D) = \sum_n \left[-\frac{1}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (w \cdot x_n + b - y_n)^2 \right] \quad \text{model assumptions} \quad (9.39)$$

$$= -\frac{1}{2\sigma^2} \sum_n (w \cdot x_n + b - y_n)^2 + \text{const} \quad \text{remove constants} \quad (9.40)$$

This is precisely the linear regression model you encountered in Section 7.6! To minimizing the *negative* log probability, you need only solve for the regression coefficients w, b as before.

In the case of binary classification, using a Gaussian noise model does not make sense. Switching to a Bernoulli model, which describes binary outcomes, makes more sense. The only remaining difficulty is that the parameter of a Bernoulli is a value between zero and one (the probability of “heads”) so your model must produce

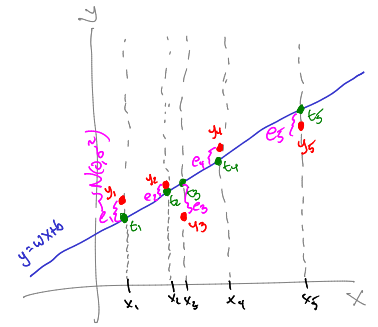


Figure 9.3: pictorial view of targets versus labels

such values. A classic approach is to produce a real-valued target, as before, and then transform this target into a value between zero and one, so that $-\infty$ maps to 0 and $+\infty$ maps to 1. A function that does this is the logistic function¹, defined below and plotted in Figure 9.4:

$$\text{Logistic function: } \sigma(z) = \frac{1}{1 + \exp[-z]} = \frac{\exp z}{1 + \exp z} \quad (9.41)$$

The logistic function has several nice properties that you can verify for yourself: $\sigma(-z) = 1 - \sigma(z)$ and $\partial\sigma/\partial z = z\sigma^2(z)$.

Using the logistic function, you can write down a generative story for binary classification:

1. For each example $n = 1 \dots N$:
 - (a) Compute $t_n = \sigma(\mathbf{w} \cdot \mathbf{x}_n + b)$
 - (b) Compute $z_n \sim \text{Ber}(t_n)$
 - (c) Return $y_n = 2z_n - 1$ (to make it ± 1)

The log-likelihood for this model is:

$$\begin{aligned} \log p(D) &= \sum_n \left[[y_n = +1] \log \sigma(\mathbf{w} \cdot \mathbf{x}_n + b) \right. \\ &\quad \left. + [y_n = -1] \log \sigma(-\mathbf{w} \cdot \mathbf{x}_n + b) \right] \end{aligned} \quad \begin{array}{l} \text{model and properties of } \sigma \\ (9.42) \end{array}$$

$$= \sum_n \log \sigma(y_n (\mathbf{w} \cdot \mathbf{x}_n + b)) \quad \begin{array}{l} \text{join terms} \\ (9.43) \end{array}$$

$$= - \sum_n \log [1 + \exp(-y_n (\mathbf{w} \cdot \mathbf{x}_n + b))] \quad \begin{array}{l} \text{definition of } \sigma \\ (9.44) \end{array}$$

$$= - \sum_n \ell^{(\log)}(y_n, \mathbf{w} \cdot \mathbf{x}_n + b) \quad \begin{array}{l} \text{definition of } \ell^{(\log)} \\ (9.45) \end{array}$$

As you can see, the log-likelihood is *precisely* the negative of (a scaled version of) the logistic loss from Chapter 7. This model is the **logistic regression** model, and this is where logisitic loss originally derived from.

TODO: conditional versus joint

¹ Also called the **sigmoid function** because of its “S”-shape.

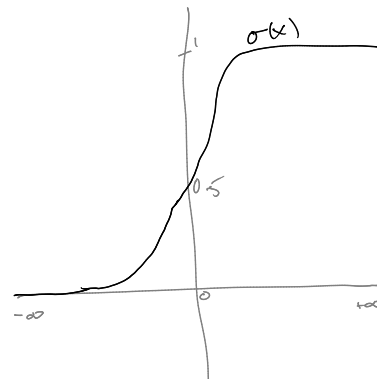


Figure 9.4: sketch of logistic function

9.7 Regularization via Priors

In the foregoing discussion, parameters of the model were selected according to the maximum likelihood criteria: find the parameters θ that maximize $p_\theta(D)$. The trouble with this approach is easy to

see even in a simple coin flipping example. If you flip a coin twice and it comes up heads both times, the maximum likelihood estimate for the bias of the coin is 100%: it will *always* come up heads. This is true even if you had only flipped it once! Of course if you had flipped it one million times and it had come up heads every time, *then* you might find this to be a reasonable solution.

This is clearly undesirable behavior, especially since data is expensive in a machine learning setting. One solution (there are others!) is to seek parameters that balance a tradeoff between the likelihood of the data and some prior belief you have about what values of those parameters are likely. Taking the case of the logistic regression, you might a priori believe that small values of w are more likely than large values, and choose to represent this as a Gaussian prior on each component of w .

The **maximum a posteriori** principle is a method for incorporating both data and prior beliefs to obtain a more balanced parameter estimate. In abstract terms, consider a probabilistic model over data D that is parameterized by parameters θ . If you think of the parameters as just another random variable, then you can write this model as $p(D | \theta)$, and maximum likelihood amounts to choosing θ to maximize $p(D | \theta)$. However, you might instead wish to maximize the probability of the *parameters*, given the data. Namely, maximize $p(\theta | D)$. This term is known as the **posterior** distribution on θ , and can be computed by Bayes' rule:

$$\underbrace{p(\theta | D)}_{\text{posterior}} = \frac{\overbrace{p(\theta)}^{\text{prior}} \overbrace{p(D | \theta)}^{\text{likelihood}}}{\underbrace{p(D)}_{\text{evidence}}}, \text{ where } p(D) = \int d\theta p(\theta) p(D | \theta) \quad (9.46)$$

This reads: the **posterior** is equal to the **prior** times the **likelihood** divided by the **evidence**.² The evidence is a scary-looking term (it has an integral!) but note that from the perspective of seeking parameters θ than maximize the posterior, the evidence is just a constant (it does not depend on θ) and therefore can be ignored.

² The evidence is sometimes called the **marginal likelihood**.

Returning to the logistic regression example with Gaussian priors on the weights, the **log posterior** looks like:

$$\log p(\theta | D) = - \sum_n \ell^{(\log)}(y_n, w \cdot x_n + b) - \sum_d \frac{1}{2\sigma^2} w_d^2 + \text{const} \quad \text{model definition} \quad (9.47)$$

$$= - \sum_n \ell^{(\log)}(y_n, w \cdot x_n + b) - \frac{1}{2\sigma^2} \|w\|^2 \quad (9.48)$$

and therefore reduces to a regularized logistic function, with a

squared 2-norm regularizer on the weights. (A 1-norm regularizer can be obtained by using a Laplace prior on w rather than a Gaussian prior on w .)

9.8 *Further Reading*

TODO

10 | NEURAL NETWORKS

TODO

Learning Objectives:

- Explain the biological inspiration for multi-layer neural networks.
- Construct a two-layer network that can solve the XOR problem.
- Implement the back-propagation algorithm for training multi-layer networks.
- Explain the trade-off between depth and breadth in network structure.
- Contrast neural networks with radial basis functions with k -nearest neighbor learning.

THE FIRST LEARNING MODELS you learned about (decision trees and nearest neighbor models) created complex, **non-linear** decision boundaries. We moved from there to the perceptron, perhaps the most classic linear model. At this point, we will move *back* to non-linear learning models, but using all that we have learned about linear learning thus far.

This chapter presents an extension of perceptron learning to non-linear decision boundaries, taking the biological inspiration of neurons even further. In the perceptron, we thought of the input data point (e.g., an image) as being directly connected to an output (e.g., label). This is often called a **single-layer network** because there is one layer of weights. Now, instead of directly connecting the inputs to the outputs, we will insert a layer of “hidden” nodes, moving from a single-layer network to a **multi-layer network**. But introducing a non-linearity at inner layers, this will give us non-linear decision boundaries. In fact, such networks are able to express almost any function we want, not just linear functions. The trade-off for this flexibility is increased complexity in parameter tuning and model design.

Dependencies:

10.1 Bio-inspired Multi-Layer Networks

One of the major weaknesses of linear models, like perceptron and the regularized linear models from the previous chapter, is that they are linear! Namely, they are unable to learn arbitrary decision boundaries. In contrast, decision trees and KNN *could* learn arbitrarily complicated decision boundaries.

One approach to doing this is to chain together a collection of perceptrons to build more complex **neural networks**. An example of a **two-layer network** is shown in Figure 10.1. Here, you can see five inputs (features) that are fed into two **hidden units**. These hidden units are then fed in to a single **output unit**. Each edge in this figure corresponds to a different weight. (Even though it looks like there are three layers, this is called a two-layer network because we don't count

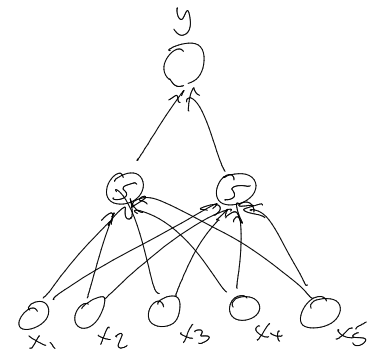


Figure 10.1: picture of a two-layer network with 5 inputs and two hidden units

the inputs as a real layer. That is, it's two layers of *trained* weights.)

Prediction with a neural network is a straightforward generalization of prediction with a perceptron. First you compute activations of the nodes in the hidden unit based on the inputs and the input weights. Then you compute activations of the output unit given the hidden unit activations and the second layer of weights.

The only major difference between this computation and the perceptron computation is that the hidden units compute a *non-linear* function of their inputs. This is usually called the **activation function** or **link function**. More formally, if $w_{i,d}$ is the weights on the edge connecting input d to hidden unit i , then the activation of hidden unit i is computed as:

$$h_i = f(w_i \cdot x) \quad (10.1)$$

Where f is the link function and w_i refers to the vector of weights feeding in to node i .

One example link function is the **sign** function. That is, if the incoming signal is negative, the activation is -1 . Otherwise the activation is $+1$. This is a potentially useful activation function, but you might already have guessed the problem with it: it is non-differentiable.

EXPLAIN BIAS!!!

A more popular link function is the **hyperbolic tangent** function, \tanh . A comparison between the sign function and the \tanh function is in Figure 10.2. As you can see, it is a reasonable approximation to the sign function, but is convenient in that it is differentiable.¹ Because it looks like an "S" and because the Greek character for "S" is "Sigma," such functions are usually called **sigmoid** functions.

Assuming for now that we are using \tanh as the link function, the overall prediction made by a two-layer network can be computed using Algorithm 10.1. This function takes a matrix of weights \mathbf{W} corresponding to the first layer weights and a vector of weights \mathbf{v} corresponding to the second layer. You can write this entire computation out in one line as:

$$\hat{y} = \sum_i v_i \tanh(w_i \cdot \hat{x}) \quad (10.2)$$

$$= \mathbf{v} \cdot \tanh(\mathbf{W}\hat{x}) \quad (10.3)$$

Where the second line is short hand assuming that \tanh can take a vector as input and product a vector as output.

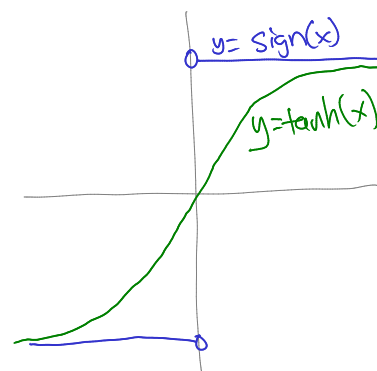


Figure 10.2: picture of sign versus \tanh

¹ Its derivative is just $1 - \tanh^2(x)$.



Is it necessary to use a link function at all? What would happen if you just used the identity function as a link?

Algorithm 25 TWOLAYERNETWORKPREDICT($\mathbf{W}, v, \hat{\mathbf{x}}$)

```

1: for  $i = 1$  to number of hidden units do
2:    $h_i \leftarrow \tanh(\mathbf{w}_i \cdot \hat{\mathbf{x}})$            // compute activation of hidden unit  $i$ 
3: end for
4: return  $v \cdot \mathbf{h}$                        // compute output unit

```

The claim is that two-layer neural networks are more expressive than single layer networks (i.e., perceptrons). To see this, you can construct a very small two-layer network for solving the XOR problem. For simplicity, suppose that the data set consists of four data points, given in Table 10.1. The classification rule is that $y = +1$ if and only if $x_1 = x_2$, where the features are just ± 1 .

You can solve this problem using a two layer network with two hidden units. The key idea is to make the first hidden unit compute an “or” function: $x_1 \vee x_2$. The second hidden unit can compute an “and” function: $x_1 \wedge x_2$. The the output can combine these into a single prediction that mimics XOR. Once you have the first hidden unit activate for “or” and the second for “and,” you need only set the output weights as -2 and $+1$, respectively.

To achieve the “or” behavior, you can start by setting the bias to -0.5 and the weights for the two “real” features as both being 1. You can check for yourself that this will do the “right thing” if the link function were the sign function. Of course it’s not, it’s tanh. To get tanh to mimic sign, you need to make the dot product either really really large or really really small. You can accomplish this by setting the bias to $-500,000$ and both of the two weights to $1,000,000$. Now, the activation of this unit will be just slightly above -1 for $\mathbf{x} = \langle -1, -1 \rangle$ and just slightly below $+1$ for the other three examples.

At this point you’ve seen that one-layer networks (aka perceptrons) can represent any linear function and only linear functions. You’ve also seen that two-layer networks can represent non-linear functions like XOR. A natural question is: do you get additional representational power by moving beyond two layers? The answer is partially provided in the following Theorem, due originally to George Cybenko for one particular type of link function, and extended later by Kurt Hornik to arbitrary link functions.

Theorem 10 (Two-Layer Networks are Universal Function Approximators). *Let F be a continuous function on a bounded subset of D -dimensional space. Then there exists a two-layer neural network \hat{F} with a finite number of hidden units that approximate F arbitrarily well. Namely, for all \mathbf{x} in the domain of F , $|F(\mathbf{x}) - \hat{F}(\mathbf{x})| < \epsilon$.*

Or, in colloquial terms “two-layer networks can approximate any

y	x_0	x_1	x_2
$+1$	$+1$	$+1$	$+1$
$+1$	$+1$	-1	-1
-1	$+1$	$+1$	-1
-1	$+1$	-1	$+1$

Table 10.1: Small XOR data set.

? Verify that these output weights will actually give you XOR.

? This shows how to create an “or” function. How can you create an “and” function?

function.”

This is a remarkable theorem. Practically, it says that if you give me a function F and some error tolerance parameter ϵ , I can construct a two layer network that computes F . In a sense, it says that going from one layer to two layers completely changes the representational capacity of your model.

When working with two-layer networks, the key question is: how many hidden units should I have? If your data is D dimensional and you have K hidden units, then the total number of parameters is $(D + 2)K$. (The first $+1$ is from the bias, the second is from the second layer of weights.) Following on from the heuristic that you should have one to two examples for each parameter you are trying to estimate, this suggests a method for choosing the number of hidden units as roughly $\lfloor \frac{N}{D} \rfloor$. In other words, if you have tons and tons of examples, you can safely have lots of hidden units. If you only have a few examples, you should probably restrict the number of hidden units in your network.

The number of units is both a form of inductive bias and a form of regularization. In both view, the number of hidden units controls how complex your function will be. Lots of hidden units \Rightarrow very complicated function. As the number increases, training performance continues to get better. But at some point, test performance gets worse because the network has overfit the data.

10.2 The Back-propagation Algorithm

The back-propagation algorithm is a classic approach to training neural networks. Although it was not originally seen this way, based on what you know from the last chapter, you can summarize back-propagation as:

$$\text{back-propagation} = \text{gradient descent} + \text{chain rule} \quad (10.4)$$

More specifically, the set up is *exactly* the same as before. You are going to optimize the weights in the network to minimize some objective function. The only difference is that the predictor is no longer linear (i.e., $\hat{y} = w \cdot x + b$) but now non-linear (i.e., $v \cdot \tanh(\mathbf{W}\hat{x})$). The only question is how to do gradient descent on this more complicated objective.

For now, we will ignore the idea of regularization. This is for two reasons. The first is that you already know how to deal with regularization, so everything you’ve learned before applies. The second is that *historically*, neural networks have not been regularized. Instead, people have used **early stopping** as a method for controlling overfitting. Presently, it’s not obvious which is a better solution: both are

valid options.

To be completely explicit, we will focus on optimizing squared error. Again, this is mostly for historic reasons. You could easily replace squared error with your loss function of choice. Our overall objective is:

$$\min_{\mathbf{W}, v} \sum_n \frac{1}{2} \left(y_n - \sum_i v_i f(\mathbf{w}_i \cdot \mathbf{x}_n) \right)^2 \quad (10.5)$$

Here, f is some link function like \tanh .

The easy case is to differentiate this with respect to v : the weights for the output unit. Without even doing any math, you should be able to guess what this looks like. The way to think about it is that from v 's perspective, it is just a linear model, attempting to minimize squared error. The only “funny” thing is that its inputs are the activations \mathbf{h} rather than the examples \mathbf{x} . So the gradient with respect to v is just as for the linear case.

To make things notationally more convenient, let e_n denote the error on the n th example (i.e., the blue term above), and let \mathbf{h}_n denote the vector of hidden unit activations on that example. Then:

$$\nabla_v = - \sum_n e_n \mathbf{h}_n \quad (10.6)$$

This is exactly like the linear case. One way of interpreting this is: how would the output weights have to change to make the prediction better? This is an easy question to answer because they can easily measure how their changes affect the output.

The more complicated aspect to deal with is the weights corresponding to the *first* layer. The reason this is difficult is because the weights in the first layer aren't necessarily trying to produce specific values, say 0 or 5 or -2.1 . They are simply trying to produce activations that get fed to the output layer. So the change they want to make depends crucially on how the output layer interprets them.

Thankfully, the chain rule of calculus saves us. Ignoring the sum over data points, we can compute:

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2} \left(y - \sum_i v_i f(\mathbf{w}_i \cdot \mathbf{x}) \right)^2 \quad (10.7)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_i} = \frac{\partial \mathcal{L}}{\partial f_i} \frac{\partial f_i}{\partial \mathbf{w}_i} \quad (10.8)$$

$$\frac{\partial \mathcal{L}}{\partial f_i} = - \left(y - \sum_i v_i f(\mathbf{w}_i \cdot \mathbf{x}) \right) v_i = -e v_i \quad (10.9)$$

$$\frac{\partial f_i}{\partial \mathbf{w}_i} = f'(\mathbf{w}_i \cdot \mathbf{x}) \mathbf{x} \quad (10.10)$$

Algorithm 26 TWOLAYERNETWORKTRAIN($\mathbf{D}, \eta, K, \text{MaxIter}$)

```

1:  $\mathbf{W} \leftarrow D \times K$  matrix of small random values // initialize input layer weights
2:  $\mathbf{v} \leftarrow K$ -vector of small random values // initialize output layer weights
3: for  $\text{iter} = 1 \dots \text{MaxIter}$  do
4:    $\mathbf{G} \leftarrow D \times K$  matrix of zeros // initialize input layer gradient
5:    $\mathbf{g} \leftarrow K$ -vector of zeros // initialize output layer gradient
6:   for all  $(\mathbf{x}, y) \in \mathbf{D}$  do
7:     for  $i = 1$  to  $K$  do
8:        $a_i \leftarrow \mathbf{w}_i \cdot \hat{\mathbf{x}}$ 
9:        $h_i \leftarrow \tanh(a_i)$  // compute activation of hidden unit  $i$ 
10:    end for
11:     $\hat{y} \leftarrow \mathbf{v} \cdot \mathbf{h}$  // compute output unit
12:     $e \leftarrow y - \hat{y}$  // compute error
13:     $\mathbf{g} \leftarrow \mathbf{g} - e\mathbf{h}$  // update gradient for output layer
14:    for  $i = 1$  to  $K$  do
15:       $\mathbf{G}_i \leftarrow \mathbf{G}_i - ev_i(1 - \tanh^2(a_i))\mathbf{x}$  // update gradient for input layer
16:    end for
17:  end for
18:   $\mathbf{W} \leftarrow \mathbf{W} - \eta \mathbf{G}$  // update input layer weights
19:   $\mathbf{v} \leftarrow \mathbf{v} - \eta \mathbf{g}$  // update output layer weights
20: end for
21: return  $\mathbf{W}, \mathbf{v}$ 

```

Putting this together, we get that the gradient with respect to \mathbf{w}_i is:

$$\nabla_{\mathbf{w}_i} = -ev_i f'(\mathbf{w}_i \cdot \mathbf{x})\mathbf{x} \quad (10.11)$$

Intuitively you can make sense of this. If the overall error of the predictor (e) is small, you want to make small steps. If v_i is small for hidden unit i , then this means that the output is not particularly sensitive to the activation of the i th hidden unit. Thus, its gradient should be small. If v_i flips sign, the gradient at \mathbf{w}_i should also flip signs. The name **back-propagation** comes from the fact that you propagate gradients backward through the network, starting at the end.

The complete instantiation of gradient descent for a two layer network with K hidden units is sketched in Algorithm 10.2. Note that this really is *exactly* a gradient descent algorithm; the only different is that the computation of the gradients of the input layer is moderately complicated.

As a bit of practical advice, implementing the back-propagation algorithm can be a bit tricky. Sign errors often abound. A useful trick is first to keep \mathbf{W} fixed and work on just training \mathbf{v} . Then keep \mathbf{v} fixed and work on training \mathbf{W} . Then put them together.



What would happen to this algorithm if you wanted to optimize exponential loss instead of squared error? What if you wanted to add in weight regularization?



If you like matrix calculus, derive the same algorithm starting from Eq (10.3).

10.3 Initialization and Convergence of Neural Networks

Based on what you know about linear models, you might be tempted to initialize all the weights in a neural network to zero. You might also have noticed that in Algorithm 10.2, this is not what's done: they're initialized to small random values. The question is why?

The answer is because an initialization of $\mathbf{W} = \mathbf{0}$ and $\mathbf{v} = \mathbf{0}$ will lead to “uninteresting” solutions. In other words, if you initialize the model in this way, it will eventually get stuck in a bad local optimum. To see this, first realize that on any example x , the activation h_i of the hidden units will all be zero since $\mathbf{W} = \mathbf{0}$. This means that on the first iteration, the gradient on the output weights (\mathbf{v}) will be zero, so they will stay put. Furthermore, the gradient $w_{1,d}$ for the d th feature on the i th unit will be *exactly* the same as the gradient $w_{2,d}$ for the same feature on the second unit. This means that the weight matrix, after a gradient step, will change in *exactly the same way* for every hidden unit. Thinking through this example for iterations $2 \dots$, the values of the hidden units will *always* be exactly the same, which means that the weights feeding in to any of the hidden units will be exactly the same. Eventually the model will converge, but it will converge to a solution that does not take advantage of having access to the hidden units.

This shows that neural networks are *sensitive* to their initialization. In particular, the function that they optimize is **non-convex**, meaning that it might have plentiful local optima. (One of which is the trivial local optimum described in the preceding paragraph.) In a sense, neural networks *must* have local optima. Suppose you have a two layer network with two hidden units that's been optimized. You have weights w_1 from inputs to the first hidden unit, weights w_2 from inputs to the second hidden unit and weights (v_1, v_2) from the hidden units to the output. If I give you back another network with w_1 and w_2 swapped, and v_1 and v_2 swapped, the network computes *exactly* the same thing, but with a markedly different weight structure. This phenomena is known as **symmetric modes** (“mode” referring to an optima) meaning that there are symmetries in the weight space. It would be one thing if there were lots of modes and they were all symmetric: then finding one of them would be as good as finding any other. Unfortunately there are additional local optima that are *not* global optima.

Random initialization of the weights of a network is a way to address *both* of these problems. By initializing a network with small random weights (say, uniform between -0.1 and 0.1), the network is unlikely to fall into the trivial, symmetric local optimum. Moreover, by training a collection of networks, each with a different random

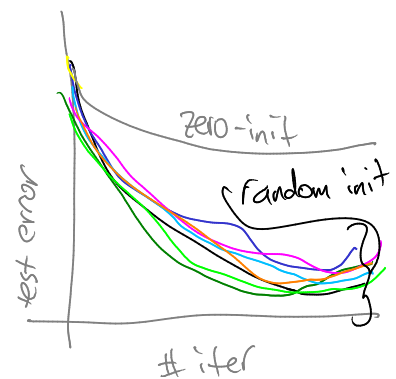


Figure 10.3: convergence of randomly initialized networks

initialization, you can often obtain better solutions than with just one initialization. In other words, you can train ten networks with different random seeds, and then pick the one that does best on held-out data. Figure 10.3 shows prototypical *test-set* performance for ten networks with different random initialization, plus an eleventh plot for the trivial symmetric network initialized with zeros.

One of the typical complaints about neural networks is that they are finicky. In particular, they have a rather large number of knobs to tune:

1. The number of layers
2. The number of hidden units per layer
3. The gradient descent learning rate η
4. The initialization
5. The stopping iteration or weight regularization

The last of these is minor (early stopping is an easy regularization method that does not require much effort to tune), but the others are somewhat significant. Even for two layer networks, having to choose the number of hidden units, and then get the learning rate and initialization “right” can take a bit of work. Clearly it can be automated, but nonetheless it takes time.

Another difficulty of neural networks is that their weights can be difficult to interpret. You’ve seen that, for linear networks, you can often interpret high weights as indicative of positive examples and low weights as indicative of negative examples. In multilayer networks, it becomes very difficult to try to understand what the different hidden units are doing.

10.4 Beyond Two Layers

The definition of neural networks and the back-propagation algorithm can be generalized beyond two layers to any arbitrary directed acyclic graph. In practice, it is most common to use a layered network like that shown in Figure 10.4 unless one has a very strong reason (aka inductive bias) to do something different. However, the view as a directed graph sheds a different sort of insight on the back-propagation algorithm.

Suppose that your network structure is stored in some directed acyclic graph, like that in Figure 10.5. We index nodes in this graph as u, v . The activation *before* applying non-linearity at a node is a_u and after non-linearity is h_u . The graph has a single sink, which is the output node y with activation a_y (no non-linearity is performed

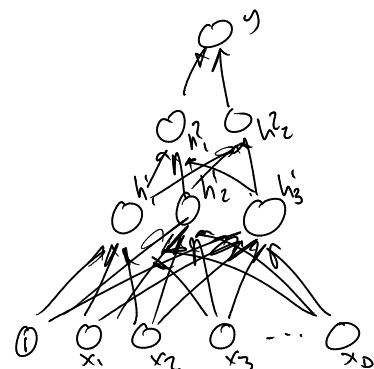


Figure 10.4: multi-layer network

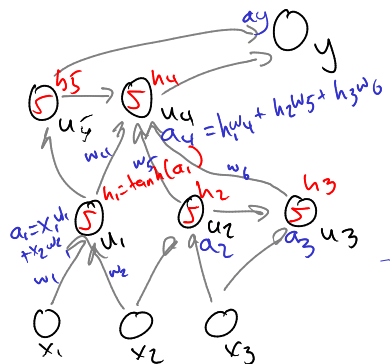


Figure 10.5: DAG network

Algorithm 27 FORWARDPROPAGATION(x)

```

1: for all input nodes  $u$  do
2:    $h_u \leftarrow$  corresponding feature of  $x$ 
3: end for
4: for all nodes  $v$  in the network whose parent's are computed do
5:    $a_v \leftarrow \sum_{u \in \text{par}(v)} w_{u,v} h_u$ 
6:    $h_v \leftarrow \tanh(a_v)$ 
7: end for
8: return  $a_y$ 

```

Algorithm 28 BACKPROPAGATION(x, y)

```

1: run FORWARDPROPAGATION( $x$ ) to compute activations
2:  $e_y \leftarrow y - a_y$  // compute overall network error
3: for all nodes  $v$  in the network whose error  $e_v$  is computed do
4:   for all  $u \in \text{par}(v)$  do
5:      $g_{u,v} \leftarrow -e_v h_u$  // compute gradient of this edge
6:      $e_u \leftarrow e_u + e_v w_{u,v} (1 - \tanh^2(a_u))$  // compute the "error" of the parent node
7:   end for
8: end for
9: return all gradients  $g_e$ 

```

on the output unit). The graph has D -many inputs (i.e., nodes with no parent), whose activations h_u are given by an input example. An edge (u, v) is from a parent to a child (i.e., from an input to a hidden unit, or from a hidden unit to the sink). Each edge has a weight $w_{u,v}$. We say that $\text{par}(u)$ is the set of parents of u .

There are two relevant algorithms: forward-propagation and back-propagation. Forward-propagation tells you how to compute the activation of the sink y given the inputs. Back-propagation computes derivatives of the edge weights for a given input.

The key aspect of the **forward-propagation** algorithm is to iteratively compute activations, going deeper and deeper in the DAG. Once the activations of all the parents of a node u have been computed, you can compute the activation of node u . This is spelled out in Algorithm 10.4. This is also explained pictorially in Figure 10.6.

Back-propagation (see Algorithm 10.4) does the opposite: it computes gradients top-down in the network. The key idea is to compute an *error* for each node in the network. The error at the output unit is the "true error." For any input unit, the error is the amount of gradient that we see coming from our children (i.e., higher in the network). These errors are computed backwards in the network (hence the name **back-propagation**) along with the gradients themselves. This is also explained pictorially in Figure 10.7.

Given the back-propagation algorithm, you can directly run gradient descent, using it as a subroutine for computing the gradients.

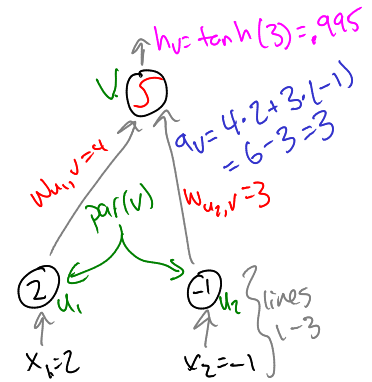


Figure 10.6: picture of forward prop

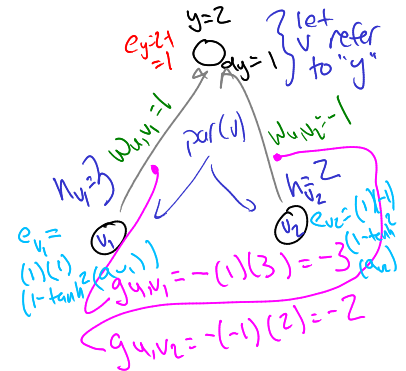


Figure 10.7: picture of back prop

10.5 Breadth versus Depth

At this point, you’ve seen how to train two-layer networks and how to train arbitrary networks. You’ve also seen a theorem that says that two-layer networks are universal function approximators. This begs the question: if two-layer networks are so great, why do we care about deeper networks?

To understand the answer, we can borrow some ideas from CS theory, namely the idea of **circuit complexity**. The goal is to show that there are functions for which it might be a “good idea” to use a deep network. In other words, there are functions that will require a huge number of hidden units if you force the network to be shallow, but can be done in a small number of units if you allow it to be deep. The example that we’ll use is the **parity function** which, ironically enough, is just a generalization of the XOR problem. The function is defined over binary inputs as:

$$\text{parity}(\mathbf{x}) = \sum_d x_d \mod 2 \quad (10.12)$$

$$= \begin{cases} 1 & \text{if the number of 1s in } \mathbf{x} \text{ is odd} \\ 0 & \text{if the number of 1s in } \mathbf{x} \text{ is even} \end{cases} \quad (10.13)$$

It is easy to define a circuit of depth $\mathcal{O}(\log_2 D)$ with $\mathcal{O}(D)$ -many gates for computing the parity function. Each gate is an XOR, arranged in a complete binary tree, as shown in Figure 10.8. (If you want to disallow XOR as a gate, you can fix this by allowing the depth to be doubled and replacing each XOR with an AND, OR and NOT combination, like you did at the beginning of this chapter.)

This shows that if you are allowed to be deep, you can construct a circuit with that computes parity using a number of hidden units that is linear in the dimensionality. So can you do the same with shallow circuits? The answer is no. It’s a famous result of circuit complexity that parity requires exponentially many gates to compute in constant depth. The formal theorem is below:

Theorem 11 (Parity Function Complexity). *Any circuit of depth $K < \log_2 D$ that computes the parity function of D input bits must contain $\mathcal{O}e^D$ gates.*

This is a very famous result because it shows that constant-depth circuits are less powerful than deep circuits. Although a neural network isn’t exactly the same as a circuit, it is generally believed that the same result holds for neural networks. At the very least, this gives a strong indication that depth might be an important consideration in neural networks.

One way of thinking about the issue of breadth versus depth has to do with the number of *parameters* that need to be estimated. By

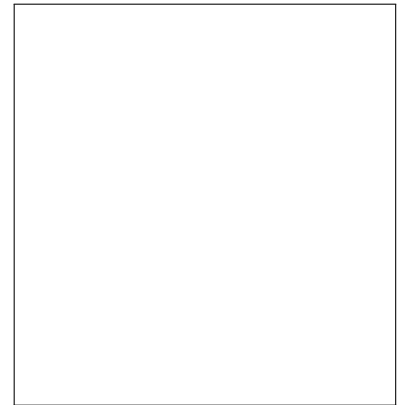


Figure 10.8: nnet:paritydeep: deep function for computing parity

? What is it about neural networks that makes it so that the theorem about circuits does not apply directly?

the heuristic that you need roughly one or two examples for every parameter, a deep model could potentially require exponentially fewer examples to train than a shallow model!

This now flips the question: if deep is potentially so much better, why doesn't everyone use deep networks? There are at least two answers. First, it makes the **architecture selection** problem more significant. Namely, when you use a two-layer network, the only hyperparameter to choose is how many hidden units should go in the middle layer. When you choose a deep network, you need to choose how many layers, and what is the width of all those layers. This can be somewhat daunting.

A second issue has to do with training deep models with back-propagation. In general, as back-propagation works its way down through the model, the sizes of the gradients shrink. You can work this out mathematically, but the intuition is simpler. If you are the beginning of a very deep network, changing one single weight is unlikely to have a significant effect on the output, since it has to go through so many other units before getting there. This directly implies that the derivatives are small. This, in turn, means that back-propagation essentially never moves far from its initialization when run on very deep networks.

Finding good ways to train deep networks is an active research area. There are two general strategies. The first is to attempt to initialize the weights better, often by a **layer-wise** initialization strategy. This can be often done using unlabeled data. After this initialization, back-propagation can be run to tweak the weights for whatever classification problem you care about. A second approach is to use a more complex optimization procedure, rather than gradient descent. You will learn about some such procedures later in this book.



While these small derivatives might make training difficult, they might be *good* for other reasons: what reasons?

10.6 Basis Functions

At this point, we've seen that: (a) neural networks can mimic linear functions and (b) they can learn more complex functions. A reasonable question is whether they can mimic a KNN classifier, and whether they can do it efficiently (i.e., with not-too-many hidden units).

A natural way to train a neural network to mimic a KNN classifier is to replace the sigmoid link function with a **radial basis function** (RBF). In a **sigmoid network** (i.e., a network with sigmoid links), the hidden units were computed as $h_i = \tanh(w_i \cdot x)$. In an **RBF network**, the hidden units are computed as:

$$h_i = \exp \left[-\gamma_i \|w_i - x\|^2 \right] \quad (10.14)$$

In other words, the hidden units behave like little Gaussian “bumps” centered around locations specified by the vectors w_i . A one-dimensional example is shown in Figure 10.9. The parameter γ_i specifies the *width* of the Gaussian bump. If γ_i is large, then only data points that are really close to w_i have non-zero activations. To distinguish sigmoid networks from RBF networks, the hidden units are typically drawn with sigmoids or with Gaussian bumps, as in Figure 10.10.

Training RBF networks involves finding good values for the Gaussian widths, γ_i , the centers of the Gaussian bumps, w_i and the connections between the Gaussian bumps and the output unit, v . This can all be done using back-propagation. The gradient terms for v remain unchanged from before, the the derivates for the other variables differ (see Exercise ??).

One of the big questions with RBF networks is: where should the Gaussian bumps be centered? One can, of course, apply back-propagation to attempt to find the centers. Another option is to specify them ahead of time. For instance, one potential approach is to have one RBF unit per data point, centered on that data point. If you carefully choose the γ s and v s, you can obtain something that looks nearly identical to distance-weighted KNN by doing so. This has the added advantage that you can go futher, and use back-propagation to *learn* good Gaussian widths (γ) and “voting” factors (v) for the nearest neighbor algorithm.

10.7 Further Reading

TODO further reading

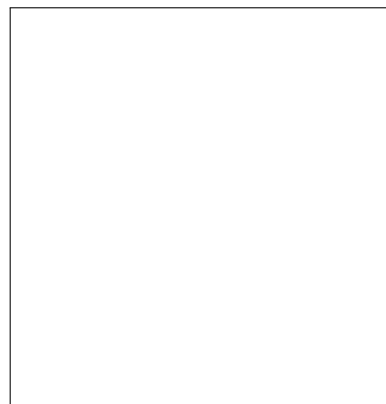


Figure 10.9: `nnet:rbbfpicture`: a one-D picture of RBF bumps

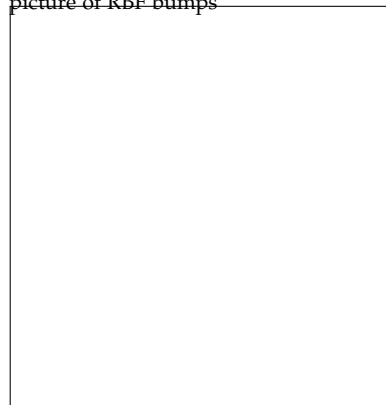


Figure 10.10: `nnet:unitsymbols`: picture of nnet with sigmoid/rbf units

?

Consider an RBF network with one hidden unit per training point, centered at that point. What bad thing might happen if you use back-propagation to estimate the γ s and v on this data if you're not careful? How could you be careful?

11 | KERNEL METHODS

Many who have had an opportunity of knowing any more about mathematics confuse it with arithmetic, and consider it an arid science. In reality, however, it is a science which requires a great amount of imagination.

– Sofia Kovalevskaya

LINEAR MODELS ARE GREAT because they are easy to understand and easy to optimize. They suffer because they can only learn very simple decision boundaries. Neural networks can learn more complex decision boundaries, but lose the nice convexity properties of many linear models.

One way of getting a linear model to behave non-linearly is to transform the input. For instance, by adding feature pairs as additional inputs. Learning a linear model on such a representation is convex, but is computationally prohibitive in all but very low dimensional spaces. You might ask: instead of *explicitly* expanding the feature space, is it possible to stay with our original data representation and do all the feature blow up *implicitly*? Surprisingly, the answer is often “yes” and the family of techniques that makes this possible are known as **kernel** approaches.

Learning Objectives:

- Explain how kernels generalize both feature combinations and basis functions.
- Contrast dot products with kernel products.
- Implement kernelized perceptron.
- Derive a kernelized version of regularized least squares regression.
- Implement a kernelized version of the perceptron.
- Derive the dual formulation of the support vector machine.

Dependencies:

11.1 From Feature Combinations to Kernels

In Section 5.4, you learned one method for increasing the expressive power of linear models: explode the feature space. For instance, a “quadratic” feature explosion might map a feature vector $\mathbf{x} = \langle x_1, x_2, x_3, \dots, x_D \rangle$ to an expanded version denoted $\phi(\mathbf{x})$:

$$\begin{aligned}\phi(\mathbf{x}) = \langle &1, 2x_1, 2x_2, 2x_3, \dots, 2x_D, \\ &x_1^2, x_1x_2, x_1x_3, \dots, x_1x_D, \\ &x_2x_1, x_2^2, x_2x_3, \dots, x_2x_D, \\ &x_3x_1, x_3x_2, x_3^2, \dots, x_3x_D, \\ &\dots, \\ &x_Dx_1, x_Dx_2, x_Dx_3, \dots, x_D^2 \rangle\end{aligned}\tag{11.1}$$

(Note that there are repetitions here, but hopefully most learning algorithms can deal well with redundant features; in particular, the $2x_1$ terms are due to collapsing some repetitions.)

You could then train a classifier on this expanded feature space. There are two primary concerns in doing so. The first is computational: if your learning algorithm scales linearly in the number of features, then you’ve just squared the amount of computation you need to perform; you’ve also squared the amount of memory you’ll need. The second is statistical: if you go by the heuristic that you should have about two examples for every feature, then you will now need quadratically many training examples in order to avoid overfitting.

This chapter is all about dealing with the *computational* issue. It will turn out in Chapter 12 that you can also deal with the statistical issue: for now, you can just hope that regularization will be sufficient to attenuate overfitting.

The key insight in kernel-based learning is that you can *rewrite* many linear models in a way that doesn’t require you to ever *explicitly* compute $\phi(x)$. To start with, you can think of this purely as a computational “trick” that enables you to use the power of a quadratic feature mapping without actually having to compute and store the mapped vectors. Later, you will see that it’s actually quite a bit deeper. Most algorithms we discuss involve a product of the form $w \cdot \phi(x)$, after performing the feature mapping. The goal is to rewrite these algorithms so that they only ever depend on dot products between two examples, say x and z ; namely, they depend on $\phi(x) \cdot \phi(z)$. To understand why this is helpful, consider the quadratic expansion from above, and the dot-product between two vectors. You get:

$$\begin{aligned} \phi(x) \cdot \phi(z) &= 1 + x_1 z_1 + x_2 z_2 + \cdots + x_D z_D + x_1^2 z_1^2 + \cdots + x_1 x_D z_1 z_D + \\ &\quad \cdots + x_D x_1 z_D z_1 + x_D x_2 z_D z_2 + \cdots + x_D^2 z_D^2 \end{aligned} \quad (11.2)$$

$$= 1 + 2 \sum_d x_d z_d + \sum_d \sum_e x_d x_e z_d z_e \quad (11.3)$$

$$= 1 + 2x \cdot z + (x \cdot z)^2 \quad (11.4)$$

$$= (1 + x \cdot z)^2 \quad (11.5)$$

Thus, you can compute $\phi(x) \cdot \phi(z)$ in exactly the same amount of time as you can compute $x \cdot z$ (plus the time it takes to perform an addition and a multiply, about 0.02 nanoseconds on a circa 2011 processor).

The rest of the practical challenge is to rewrite your algorithms so that they only depend on dot products between examples and not on any explicit weight vectors.

11.2 Kernelized Perceptron

Consider the original perceptron algorithm from Chapter 4, repeated in Algorithm 11.2 using linear algebra notation and using feature expansion notation $\phi(x)$. In this algorithm, there are two places

Algorithm 29 PERCEPTONTRAIN(\mathbf{D} , MaxIter)

```

1:  $\mathbf{w} \leftarrow \mathbf{0}, b \leftarrow 0$  // initialize weights and bias
2: for  $\text{iter} = 1 \dots \text{MaxIter}$  do
3:   for all  $(x, y) \in \mathbf{D}$  do
4:      $a \leftarrow \mathbf{w} \cdot \phi(x) + b$  // compute activation for this example
5:     if  $ya \leq 0$  then
6:        $\mathbf{w} \leftarrow \mathbf{w} + y \phi(x)$  // update weights
7:        $b \leftarrow b + y$  // update bias
8:     end if
9:   end for
10: end for
11: return  $\mathbf{w}, b$ 

```

MATH REVIEW | SPANS

If $\mathcal{U} = \{\mathbf{u}_i\}_{i=1}^I$ is a set of vectors in \mathbb{R}^D , then the **span** of \mathcal{U} is the set of vectors that can be written as linear combinations of \mathbf{u}_i s; namely: $\text{span}(\mathcal{U}) = \{\sum_i a_i \mathbf{u}_i : a_1 \in \mathbb{R}, \dots, a_I \in \mathbb{R}\}$. If all of the \mathbf{u}_i s are linearly independent, then the dimension of $\text{span}(\mathcal{U})$ is I ; in particular, if there are D -many linearly independent vectors then they span \mathbb{R}^D .

Figure 11.1:

where $\phi(x)$ is used explicitly. The first is in computing the activation (line 4) and the second is in updating the weights (line 6). The goal is to remove the explicit dependence of this algorithm on ϕ and on the weight vector.

To do so, you can observe that at any point in the algorithm, the weight vector \mathbf{w} can be written as a linear combination of expanded training data. In particular, at any point, $\mathbf{w} = \sum_n \alpha_n \phi(x_n)$ for some parameters α . Initially, $\mathbf{w} = \mathbf{0}$ so choosing $\alpha = 0$ yields this. If the first update occurs on the n th training example, then the resulting weight vector is simply $y_n \phi(x_n)$, which is equivalent to setting $\alpha_n = y_n$. If the second update occurs on the m th training example, then all you need to do is update $\alpha_m \leftarrow \alpha_m + y_m$. This is true, even if you make multiple passes over the data. This observation leads to the following **representer theorem**, which states that the weight vector of the perceptron lies in the **span** of the training data.

Theorem 12 (Perceptron Representer Theorem). *During a run of the perceptron algorithm, the weight vector \mathbf{w} is always in the span of the (assumed non-empty) training data, $\phi(x_1), \dots, \phi(x_N)$.*

Proof of Theorem 12. By induction. Base case: the span of any non-empty set contains the zero vector, which is the initial weight vector. Inductive case: suppose that the theorem is true before the k th update, and suppose that the k th update happens on example n . By the inductive hypothesis, you can write $\mathbf{w} = \sum_i \alpha_i \phi(x_i)$ before

Algorithm 30 KERNELIZEDPERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```

1:  $\alpha \leftarrow 0, b \leftarrow 0$  // initialize coefficients and bias
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(x_n, y_n) \in \mathbf{D}$  do
4:      $a \leftarrow \sum_m \alpha_m \phi(x_m) \cdot \phi(x_n) + b$  // compute activation for this example
5:     if  $y_n a \leq 0$  then
6:        $\alpha_n \leftarrow \alpha_n + y_n$  // update coefficients
7:        $b \leftarrow b + y$  // update bias
8:     end if
9:   end for
10: end for
11: return  $\alpha, b$ 

```

the update. The new weight vector is $[\sum_i \alpha_i \phi(x_i)] + y_n \phi(x_n) = \sum_i (\alpha_i + y_n [i = n]) \phi(x_i)$, which is still in the span of the training data. \square

Now that you know that you can always write $w = \sum_n \alpha_n \phi(x_n)$ for some α_i s, you can additionally compute the activations (line 4) as:

$$w \cdot \phi(x) + b = \left(\sum_n \alpha_n \phi(x_n) \right) \cdot \phi(x) + b \quad \text{definition of } w \quad (11.6)$$

$$= \sum_n \alpha_n [\phi(x_n) \cdot \phi(x)] + b \quad \text{dot products are linear} \quad (11.7)$$

This now depends only on dot-products between data points, and never explicitly requires a weight vector. You can now rewrite the entire perceptron algorithm so that it never refers explicitly to the weights and only ever depends on pairwise dot products between examples. This is shown in Algorithm 11.2.

The advantage to this “kernelized” algorithm is that you can perform feature expansions like the quadratic feature expansion from the introduction for “free.” For example, for exactly the same cost as the quadratic features, you can use a **cubic feature map**, computed as $\phi(\vec{x})\phi(\vec{z}) = (1 + x \cdot z)^3$, which corresponds to three-way interactions between variables. (And, in general, you can do so for any polynomial degree p at the same computational complexity.)

11.3 Kernelized K-means

For a complete change of pace, consider the K -means algorithm from Section 3. This algorithm is for *clustering* where there is no notion of “training labels.” Instead, you want to partition the data into coherent clusters. For data in \mathbb{R}^D , it involves randomly initializing K -many

cluster means $\mu^{(1)}, \dots, \mu^{(K)}$. The algorithm then alternates between the following two steps until convergence, with x replaced by $\phi(x)$ since that is the eventual goal:

1. For each example n , set cluster label $z_n = \arg \min_k \|\phi(x_n) - \mu^{(k)}\|^2$.
2. For each cluster k , update $\mu^{(k)} = \frac{1}{N_k} \sum_{n:z_n=k} \phi(x_n)$, where N_k is the number of n with $z_n = k$.

The question is whether you can perform these steps without explicitly computing $\phi(x_n)$. The **representer theorem** is more straightforward here than in the perceptron. The mean of a set of data is, almost by definition, in the span of that data (choose the a_i s all to be equal to $1/N$). Thus, so long as you *initialize* the means in the span of the data, you are guaranteed always to have the means in the span of the data. Given this, you know that you can write each mean as an expansion of the data; say that $\mu^{(k)} = \sum_n \alpha_n^{(k)} \phi(x_n)$ for some parameters $\alpha_n^{(k)}$ (there are $N \times K$ -many such parameters).

Given this expansion, in order to execute step (1), you need to compute norms. This can be done as follows:

$$z_n = \arg \min_k \|\phi(x_n) - \mu^{(k)}\|^2 \quad \text{definition of } z_n \quad (11.8)$$

$$= \arg \min_k \left\| \phi(x_n) - \sum_m \alpha_m^{(k)} \phi(x_m) \right\|^2 \quad \text{definition of } \mu^{(k)} \quad (11.9)$$

$$= \arg \min_k \|\phi(x_n)\|^2 + \left\| \sum_m \alpha_m^{(k)} \phi(x_m) \right\|^2 + \phi(x_n) \cdot \left[\sum_m \alpha_m^{(k)} \phi(x_m) \right] \quad \text{expand quadratic term} \quad (11.10)$$

$$= \arg \min_k \sum_m \sum_{m'} \alpha_m^{(k)} \alpha_{m'}^{(k)} \phi(x_m) \cdot \phi(x_{m'}) + \sum_m \alpha_m^{(k)} \phi(x_m) \cdot \phi(x_n) + \text{const} \quad \text{linearity and constant} \quad (11.11)$$

This computation can replace the assignments in step (1) of K -means. The mean updates are more direct in step (2):

$$\mu^{(k)} = \frac{1}{N_k} \sum_{n:z_n=k} \phi(x_n) \iff \alpha_n^{(k)} = \begin{cases} \frac{1}{N_k} & \text{if } z_n = k \\ 0 & \text{otherwise} \end{cases} \quad (11.12)$$

11.4 What Makes a Kernel

A **kernel** is just a form of generalized dot product. You can also think of it as simply shorthand for $\phi(x) \cdot \phi(z)$, which is commonly written $K^\phi(x, z)$. Or, when ϕ is clear from context, simply $K(x, z)$.

This is often referred to as the kernel product between x and z (under the mapping ϕ).

In this view, what you've seen in the preceding two sections is that you can rewrite both the perceptron algorithm and the K -means algorithm so that *they only ever depend on kernel products between data points, and never on the actual datapoints themselves*. This is a very powerful notion, as it has enabled the development of a large number of non-linear algorithms essentially “for free” (by applying the so-called **kernel trick**, that you've just seen twice).

This raises an interesting question. If you have rewritten these algorithms so that they only depend on the data through a function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, can you stick *any* function K in these algorithms, or are there some K that are “forbidden?” In one sense, you “could” use any K , but the real question is: for what types of functions K do these algorithms retain the properties that we expect them to have (like convergence, optimality, etc.)?

One way to answer this question is to say that $K(\cdot, \cdot)$ is a valid kernel if it corresponds to the inner product between two vectors. That is, K is valid if there exists a function ϕ such that $K(x, z) = \phi(x) \cdot \phi(z)$. This is a direct definition and it should be clear that if K satisfies this, then the algorithms go through as expected (because this is how we derived them).

You've already seen the general class of **polynomial kernels**, which have the form:

$$K_d^{(\text{poly})}(x, z) = (1 + x \cdot z)^d \quad (11.13)$$

where d is a hyperparameter of the kernel. These kernels correspond to polynomial feature expansions.

There is an alternative characterization of a valid kernel function that is more mathematical. It states that $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a kernel if K is **positive semi-definite** (or, in shorthand, **psd**). This property is also sometimes called **Mercer's condition**. In this context, this means the *for all* functions f that are square integrable (i.e., $\int f(x)^2 dx < \infty$), other than the zero function, the following property holds:

$$\iint f(x)K(x, z)f(z)dx dz > 0 \quad (11.14)$$

This likely seems like it came out of nowhere. Unfortunately, the connection is well beyond the scope of this book, but is covered well in external sources. For now, simply take it as a given that this is an equivalent requirement. (For those so inclined, the appendix of this book gives a proof, but it requires a bit of knowledge of function spaces to understand.)

The question is: why is this alternative characterization useful? It is useful because it gives you an alternative way to construct kernel

functions. For instance, using it you can easily prove the following, which would be difficult from the definition of kernels as inner products after feature mappings.

Theorem 13 (Kernel Addition). *If K_1 and K_2 are kernels, the K defined by $K(x, z) = K_1(x, z) + K_2(x, z)$ is also a kernel.*

Proof of Theorem 13. You need to verify the positive semi-definite property on K . You can do this as follows:

$$\iint f(x)K(x, z)f(z)dx dz = \iint f(x)[K_1(x, z) + K_2(x, z)]f(z)dx dz \quad \text{definition of } K \quad (11.15)$$

$$\begin{aligned} &= \iint f(x)K_1(x, z)f(z)dx dz \\ &\quad + \iint f(x)K_2(x, z)f(z)dx dz \quad \text{distributive rule} \end{aligned} \quad (11.16)$$

$$> 0 + 0 \quad K_1 \text{ and } K_2 \text{ are psd} \quad (11.17)$$

□

More generally, any positive linear combination of kernels is still a kernel. Specifically, if K_1, \dots, K_M are all kernels, and $\alpha_1, \dots, \alpha_M \geq 0$, then $K(x, z) = \sum_m \alpha_m K_m(x, z)$ is also a kernel.

You can also use this property to show that the following **Gaussian kernel** (also called the **RBF kernel**) is also psd:

$$K_\gamma^{(\text{RBF})}(x, z) = \exp \left[-\gamma \|x - z\|^2 \right] \quad (11.18)$$

Here γ is a hyperparameter that controls the width of this Gaussian-like bumps. To gain an intuition for what the RBF kernel is doing, consider what prediction looks like in the perceptron:

$$f(x) = \sum_n \alpha_n K(x_n, x) + b \quad (11.19)$$

$$= \sum_n \alpha_n \exp \left[-\gamma \|x_n - x\|^2 \right] \quad (11.20)$$

In this computation, each training example is getting to “vote” on the label of the test point x . The amount of “vote” that the n th training example gets is proportional to the negative exponential of the distance between the test point and itself. This is very much like an RBF neural network, in which there is a Gaussian “bump” at each training example, with variance $1/(2\gamma)$, and where the α_n s act as the weights connecting these RBF bumps to the output.

Showing that this kernel is positive definite is a bit of an exercise in analysis (particularly, integration by parts), but otherwise not difficult. Again, the proof is provided in the appendix.

So far, you have seen two basic classes of kernels: polynomial kernels ($K(x, z) = (1 + x \cdot z)^d$), which includes the linear kernel ($K(x, z) = x \cdot z$) and RBF kernels ($K(x, z) = \exp[-\gamma \|x - z\|^2]$). The former have a direct connection to feature expansion; the latter to RBF networks. You also know how to combine kernels to get new kernels by addition. In fact, you can do more than that: the product of two kernels is also a kernel.

As far as a “library of kernels” goes, there are many. Polynomial and RBF are by far the most popular. A commonly used, but technically *invalid* kernel, is the hyperbolic-tangent kernel, which mimics the behavior of a two-layer neural network. It is defined as:

$$K^{(\tanh)} = \tanh(1 + x \cdot z) \quad \text{Warning: not psd} \quad (11.21)$$

A final example, which is not very common, but is nonetheless interesting, is the all-subsets kernel. Suppose that your D features are all *binary*: all take values 0 or 1. Let $A \subseteq \{1, 2, \dots, D\}$ be a subset of features, and let $f_A(x) = \bigwedge_{d \in A} x_d$ be the conjunction of all the features in A . Let $\phi(x)$ be a feature vector over *all* such A s, so that there are 2^D features in the vector ϕ . You can compute the kernel associated with this feature mapping as:

$$K^{(\text{subs})}(x, z) = \prod_d (1 + x_d z_d) \quad (11.22)$$

Verifying the relationship between this kernel and the all-subsets feature mapping is left as an exercise (but closely resembles the expansion for the quadratic kernel).

11.5 Support Vector Machines

Kernelization predated support vector machines, but SVMs are definitely the model that popularized the idea. Recall the definition of the soft-margin SVM from Chapter 7.7 and in particular the optimization problem (7.38), which attempts to balance a **large margin** (small $\|w\|^2$) with a **small loss** (small ξ_n s, where ξ_n is the **slack** on the n th training example). This problem is repeated below:

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_n \xi_n & (11.23) \\ \text{subj. to} \quad & y_n (w \cdot x_n + b) \geq 1 - \xi_n & (\forall n) \\ & \xi_n \geq 0 & (\forall n) \end{aligned}$$

Previously, you optimized this by explicitly computing the slack variables ξ_n , given a solution to the decision boundary, w and b . However, you are now an expert with using Lagrange multipliers

to optimize constrained problems! The overall *goal* is going to be to rewrite the SVM optimization problem in a way that it no longer explicitly depends on the weights w and only depends on the examples x_n through kernel products.

There are $2N$ constraints in this optimization, one for each slack constraint and one for the requirement that the slacks are non-negative. Unlike the last time, these constraints are now *inequalities*, which require a slightly different solution. First, you rewrite all the inequalities so that they read as *something* ≥ 0 and then add corresponding Lagrange multipliers. The main difference is that the Lagrange multipliers are now constrained to be non-negative, and their sign in the augmented objective function matters.

The second set of constraints is already in the proper form; the first set can be rewritten as $y_n (w \cdot x_n + b) - 1 + \xi_n \geq 0$. You're now ready to construct the Lagrangian, using multipliers α_n for the first set of constraints and β_n for the second set.

$$\mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \quad (11.24)$$

$$- \sum_n \alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] \quad (11.25)$$

The *new* optimization problem is:

$$\min_{w, b, \xi} \max_{\alpha \geq 0} \max_{\beta \geq 0} \mathcal{L}(w, b, \xi, \alpha, \beta) \quad (11.26)$$

The intuition is exactly the same as before. If you are able to find a solution that satisfies the constraints (e.g., the purple term is properly non-negative), then the β_n s cannot do anything to “hurt” the solution. On the other hand, if the purple term *is* negative, then the corresponding β_n can go to $+\infty$, breaking the solution.

You can solve this problem by taking gradients. This is a bit tedious, but an important step to realize how everything fits together. Since your goal is to remove the dependence on w , the first step is to take a gradient with respect to w , set it equal to zero, and solve for w in terms of the other variables.

$$\nabla_w \mathcal{L} = w - \sum_n \alpha_n y_n x_n = 0 \iff w = \sum_n \alpha_n y_n x_n \quad (11.27)$$

At this point, you should immediately recognize a similarity to the kernelized perceptron: the optimal weight vector takes *exactly* the same form in both algorithms.

You can now take this new expression for w and plug it back in to the expression for \mathcal{L} , thus removing w from consideration. To avoid subscript overloading, you should replace the n in the expression for

w with, say, m . This yields:

$$\mathcal{L}(b, \xi, \alpha, \beta) = \frac{1}{2} \left\| \sum_m \alpha_m y_m \mathbf{x}_m \right\|^2 + C \sum_n \xi_n - \sum_n \beta_n \xi_n \quad (11.28)$$

$$- \sum_n \alpha_n \left[y_n \left(\left[\sum_m \alpha_m y_m \mathbf{x}_m \right] \cdot \mathbf{x}_n + b \right) - 1 + \xi_n \right] \quad (11.29)$$

At this point, it's convenient to rewrite these terms; be sure you understand where the following comes from:

$$\mathcal{L}(b, \xi, \alpha, \beta) = \frac{1}{2} \sum_n \sum_m \alpha_n \alpha_m y_n y_m \mathbf{x}_n \cdot \mathbf{x}_m + \sum_n (C - \beta_n) \xi_n \quad (11.30)$$

$$- \sum_n \sum_m \alpha_n \alpha_m y_n y_m \mathbf{x}_n \cdot \mathbf{x}_m - \sum_n \alpha_n (y_n b - 1 + \xi_n) \quad (11.31)$$

$$= -\frac{1}{2} \sum_n \sum_m \alpha_n \alpha_m y_n y_m \mathbf{x}_n \cdot \mathbf{x}_m + \sum_n (C - \beta_n) \xi_n \quad (11.32)$$

$$- b \sum_n \alpha_n y_n - \sum_n \alpha_n (\xi_n - 1) \quad (11.33)$$

Things are starting to look good: you've successfully removed the dependence on w , and everything is now written in terms of dot products between input vectors! This might still be a difficult problem to solve, so you need to continue and attempt to remove the remaining variables b and ξ .

The derivative with respect to b is:

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_n \alpha_n y_n = 0 \quad (11.34)$$

This doesn't allow you to *substitute* b with something (as you did with w), but it does mean that the fourth term ($b \sum_n \alpha_n y_n$) goes to zero at the optimum.

The last of the original variables is ξ_n ; the derivatives in this case look like:

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = C - \beta_n - \alpha_n \iff C - \beta_n = \alpha_n \quad (11.35)$$

Again, this doesn't allow you to substitute, but it does mean that you can rewrite the second term, which as $\sum_n (C - \beta_n) \xi_n$ as $\sum_n \alpha_n \xi_n$. This then cancels with (most of) the final term. However, you need to be careful to remember something. When we optimize, both α_n and β_n are constrained to be non-negative. What this means is that since we are dropping β from the optimization, we need to ensure that $\alpha_n \leq C$, otherwise the corresponding β will need to be negative, which is not

allowed. You finally wind up with the following, where $\mathbf{x}_n \cdot \mathbf{x}_m$ has been replaced by $K(\mathbf{x}_n, \mathbf{x}_m)$:

$$\mathcal{L}(\boldsymbol{\alpha}) = \sum_n \alpha_n - \frac{1}{2} \sum_n \sum_m \alpha_n \alpha_m y_n y_m K(\mathbf{x}_n, \mathbf{x}_m) \quad (11.36)$$

If you are comfortable with matrix notation, this has a very compact form. Let $\mathbf{1}$ denote the N -dimensional vector of all 1s, let \mathbf{y} denote the vector of labels and let \mathbf{G} be the $N \times N$ matrix, where $G_{n,m} = y_n y_m K(\mathbf{x}_n, \mathbf{x}_m)$, then this has the following form:

$$\mathcal{L}(\boldsymbol{\alpha}) = \boldsymbol{\alpha}^\top \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^\top \mathbf{G} \boldsymbol{\alpha} \quad (11.37)$$

The resulting optimization problem is to *maximize* $\mathcal{L}(\boldsymbol{\alpha})$ as a function of $\boldsymbol{\alpha}$, subject to the constraint that the α_n s are all non-negative and less than C (because of the constraint added when removing the β variables). Thus, your problem is:

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & -\mathcal{L}(\boldsymbol{\alpha}) = \frac{1}{2} \sum_n \sum_m \alpha_n \alpha_m y_n y_m K(\mathbf{x}_n, \mathbf{x}_m) - \sum_n \alpha_n \quad (11.38) \\ \text{subj. to} \quad & 0 \leq \alpha_n \leq C \quad (\forall n) \end{aligned}$$

One way to solve this problem is gradient descent on α . The only complication is making sure that the α s satisfy the constraints. In this case, you can use a **projected gradient** algorithm: after each gradient update, you adjust your parameters to satisfy the constraints by *projecting* them into the feasible region. In this case, the projection is trivial: if, after a gradient step, any $\alpha_n < 0$, simply set it to 0; if any $\alpha_n > C$, set it to C .

11.6 Understanding Support Vector Machines

The prior discussion involved quite a bit of math to derive a representation of the support vector machine in terms of the Lagrange variables. This mapping is actually sufficiently standard that everything in it has a name. The original problem variables ($\mathbf{w}, b, \boldsymbol{\xi}$) are called the **primal variables**; the Lagrange variables are called the **dual variables**. The optimization problem that results after removing all of the primal variables is called the **dual problem**.

A succinct way of saying what you've done is: you found that after converting the SVM into its dual, it is possible to kernelize.

To understand SVMs, a first step is to peek into the dual formulation, Eq (11.38). The objective has two terms: the first depends on the data, and the second depends only on the dual variables. The first thing to notice is that, because of the second term, the α s “want” to

get as large as possible. The constraint ensures that they cannot exceed C , which means that the general tendency is for the α s to grow as close to C as possible.

To further understand the dual optimization problem, it is useful to think of the kernel as being a measure of *similarity* between two data points. This analogy is most clear in the case of RBF kernels, but even in the case of linear kernels, if your examples all have unit norm, then their dot product is still a measure of similarity. Since you can write the prediction function as $f(\hat{x}) = \text{sign}(\sum_n \alpha_n y_n K(x_n, \hat{x}))$, it is natural to think of α_n as the “importance” of training example n , where $\alpha_n = 0$ means that it is not used at all at test time.

Consider two data points that have the same label; namely, $y_n = y_m$. This means that $y_n y_m = +1$ and the objective function has a term that looks like $\alpha_n \alpha_m K(x_n, x_m)$. Since the goal is to make this term small, then one of two things has to happen: either K has to be small, or $\alpha_n \alpha_m$ has to be small. If K is already small, then this doesn't affect the setting of the corresponding α s. But if K is large, then this *strongly* encourages at least one of α_n or α_m to go to zero. So if you have two data points that are very similar *and* have the same label, at least one of the corresponding α s will be small. This makes intuitive sense: if you have two data points that are basically the same (both in the x and y sense) then you only need to “keep” one of them around.

Suppose that you have two data points with different labels: $y_n y_m = -1$. Again, if $K(x_n, x_m)$ is small, nothing happens. But if it is large, then the corresponding α s are encouraged to be as large as possible. In other words, if you have two similar examples with different labels, you are strongly encouraged to keep the corresponding α s as large as C .

An alternative way of understanding the SVM dual problem is geometrically. Remember that the whole point of introducing the variable α_n was to ensure that the n th training example was correctly classified, modulo slack. More formally, the goal of α_n is to ensure that $y_n(w \cdot x_n + b) - 1 + \xi_n \geq 0$. Suppose that this constraint is *not* satisfied. There is an important result in optimization theory, called the **Karush-Kuhn-Tucker conditions** (or **KKT conditions**, for short) that states that at the optimum, the product of the Lagrange multiplier for a constraint, and the value of that constraint, will equal zero. In this case, this says that at the optimum, you have:

$$\alpha_n [y_n (w \cdot x_n + b) - 1 + \xi_n] = 0 \quad (11.39)$$

In order for this to be true, it means that (at least) one of the following must be true:

$$\alpha_n = 0 \quad \text{or} \quad y_n (w \cdot x_n + b) - 1 + \xi_n = 0 \quad (11.40)$$

A reasonable question to ask is: under what circumstances will α_n be *non-zero*? From the KKT conditions, you can discern that α_n can be non-zero *only* when the constraint holds *exactly*; namely, that $y_n(w \cdot x_n + b) - 1 + \xi_n = 0$. When does that constraint hold *exactly*? It holds exactly only for those points *precisely* on the margin of the hyperplane.

In other words, the *only* training examples for which $\alpha_n \neq 0$ are those that lie precisely 1 unit away from the maximum margin decision boundary! (Or those that are “moved” there by the corresponding slack.) These points are called the **support vectors** because they “support” the decision boundary. In general, the number of support vectors is far smaller than the number of training examples, and therefore you naturally end up with a solution that only uses a subset of the training data.

From the first discussion, you know that the points that wind up being support vectors are exactly those that are “confusable” in the sense that you have to examples that are nearby, but have different labels. This is a completely in line with the previous discussion. If you have a decision boundary, it will pass between these “confusable” points, and therefore they will end up being part of the set of support vectors.

11.7 Further Reading

TODO further reading

12 | LEARNING THEORY

The Universe is under no obligation to make sense to you. —

Neil deGrasse Tyson

BY NOW, YOU ARE AN EXPERT at building learning algorithms. You probably understand how they work, intuitively. And you understand why they should generalize. However, there are several basic questions you might want to know the answer to. Is learning always possible? How many training examples will I need to do a good job learning? Is my test performance going to be much worse than my training performance? The key idea that underlies all these answers is that *simple functions generalize well*.

The amazing thing is that you can actually prove strong results that address the above questions. In this chapter, you will learn some of the most important results in learning theory that attempt to answer these questions. The goal of this chapter is not theory for theory's sake, but rather as a way to better understand why learning models work, and how to use this theory to build better algorithms. As a concrete example, we will see how 2-norm regularization provably leads to better generalization performance, thus justifying our common practice!

Learning Objectives:

- Explain why inductive bias is necessary.
- Define the PAC model and explain why both the “P” and “A” are necessary.
- Explain the relationship between complexity measures and regularizers.
- Identify the role of complexity in generalization.
- Formalize the relationship between margins and complexity.

Dependencies:

12.1 The Role of Theory

In contrast to the quote at the start of this chapter, a practitioner friend once said “I would happily give up a few percent performance for an algorithm that I can understand.” Both perspectives are completely valid, and are actually not contradictory. The second statement is presupposing that theory helps you understand, which hopefully you’ll find to be the case in this chapter.

Theory can serve two roles. It can justify and help understand why common practice works. This is the “theory after” view. It can also serve to suggest new algorithms and approaches that turn out to work well in practice. This is the “theory before” view. Often, it turns out to be a mix. Practitioners discover something that works surprisingly well. Theorists figure out why it works and prove something about it. And in the process, they make it better or find new algo-

rithms that more directly exploit whatever property it is that made the theory go through.

Theory can also help you understand what's possible and *what's not possible*. One of the first things we'll see is that, in general, machine learning can not work. Of course it *does* work, so this means that we need to think harder about what it means for learning algorithms to work. By understanding what's not possible, you can focus our energy on things that are.

Probably the biggest *practical* success story for theoretical machine learning is the theory of **boosting**, which you won't actually see in this chapter. (You'll have to wait for Chapter 13.) Boosting is a very simple style of algorithm that came out of theoretical machine learning, and has proven to be incredibly successful in practice. So much so that it is one of the de facto algorithms to run when someone gives you a new data set. In fact, in 2004, Yoav Freund and Rob Schapire won the ACM's Paris Kanellakis Award for their boosting algorithm AdaBoost. This award is given for theoretical accomplishments that have had a significant and demonstrable effect on the practice of computing.¹

¹ In 2008, Corinna Cortes and Vladimir Vapnik won it for support vector machines.

12.2 Induction is Impossible

One nice thing about theory is that it forces you to be *precise* about what you are trying to do. You've already seen a formal definition of binary classification in Chapter 6. But let's take a step back and re-analyze what it means to learn to do binary classification.

From an *algorithmic* perspective, a natural question is whether there is an "ultimate" learning algorithm, $\mathcal{A}^{\text{awesome}}$, that solves the Binary Classification problem above. In other words, have you been wasting your time learning about KNN and Perceptron and decision trees, when $\mathcal{A}^{\text{awesome}}$ is out there.

What would such an ultimate learning algorithm do? You would like it to take in a data set D and produce a function f . No matter what D looks like, this function f should get perfect classification on all future examples drawn from the same distribution that produced D .

A little bit of introspection should demonstrate that this is impossible. For instance, there might be label noise in our distribution. As a very simple example, let $\mathcal{X} = \{-1, +1\}$ (i.e., a one-dimensional, binary distribution. Define the data distribution as:

$$\mathcal{D}(\langle +1 \rangle, +1) = 0.4 \quad \mathcal{D}(\langle -1 \rangle, -1) = 0.4 \quad (12.1)$$

$$\mathcal{D}(\langle +1 \rangle, -1) = 0.1 \quad \mathcal{D}(\langle -1 \rangle, +1) = 0.1 \quad (12.2)$$

In other words, 80% of data points in this distribution have $x = y$

and 20% don't. No matter what function your learning algorithm produces, there's no way that it can do better than 20% error on this data.

Given this, it seems hopeless to have an algorithm $\mathcal{A}^{\text{awesome}}$ that always achieves an error rate of zero. The best that we can hope is that the error rate is not "too large."

Unfortunately, simply weakening our requirement on the error rate is not enough to make learning possible. The second source of difficulty comes from the fact that the only access we have to the data distribution is through sampling. In particular, when trying to learn about a distribution like that in 12.1, you only get to see data points *drawn* from that distribution. You know that "eventually" you will see enough data points that your sample is representative of the distribution, but it might not happen immediately. For instance, even though a fair coin will come up heads only with probability $1/2$, it's completely plausible that in a sequence of four coin flips you never see a tails, or perhaps only see one tails.

So the second thing that we have to give up is the hope that $\mathcal{A}^{\text{awesome}}$ will *always* work. In particular, if we happen to get a lousy sample of data from \mathcal{D} , we need to allow $\mathcal{A}^{\text{awesome}}$ to do something completely unreasonable.

Thus, we cannot hope that $\mathcal{A}^{\text{awesome}}$ will do perfectly, every time. We cannot even hope that it will do pretty well, all of the time. Nor can we hope that it will do perfectly, most of the time. The best *best* we can reasonably hope of $\mathcal{A}^{\text{awesome}}$ is that it will do pretty well, most of the time.

?

It's clear that if your algorithm produces a *deterministic* function that it cannot do better than 20% error. What if it produces a stochastic (aka randomized) function?

12.3 Probably Approximately Correct Learning

Probably Approximately Correct (PAC) learning is a formalism of inductive learning based on the realization that the best we can hope of an algorithm is that it does a good job (i.e., is approximately correct), most of the time (i.e., it is *probably* approximately correct).²

Consider a hypothetical learning algorithm. You run it on ten different binary classification data sets. For each one, it comes back with functions f_1, f_2, \dots, f_{10} . For some reason, whenever you run f_4 on a test point, it crashes your computer. For the other learned functions, their performance on test data is always at most 5% error. If this situation is guaranteed to happen, then this hypothetical learning algorithm is a PAC learning algorithm. It satisfies "probably" because it only failed in one out of ten cases, and it's "approximate" because it achieved low, but non-zero, error on the remainder of the cases.

This leads to the formal definition of an (ϵ, δ) PAC-learning algorithm. In this definition, ϵ plays the role of measuring accuracy (in

² Leslie Valiant invented the notion of PAC learning in 1984. In 2011, he received the Turing Award, the highest honor in computing for his work in learning theory, computational complexity and parallel systems.

the previous example, $\epsilon = 0.05$) and δ plays the role of measuring failure (in the previous, $\delta = 0.1$).

Definitions 1. An algorithm \mathcal{A} is an (ϵ, δ) -PAC learning algorithm if, for all distributions \mathcal{D} : given samples from \mathcal{D} , the probability that it returns a “bad function” is at most δ ; where a “bad” function is one with test error rate more than ϵ on \mathcal{D} .

There are two notions of *efficiency* that matter in PAC learning. The first is the usual notion of *computational complexity*. You would prefer an algorithm that runs quickly to one that takes forever. The second is the notion of **sample complexity**: the number of examples required for your algorithm to achieve its goals. Note that the goal of both of these measure of complexity is to bound how much of a scarce resource your algorithm uses. In the computational case, the resource is CPU cycles. In the sample case, the resource is labeled examples.

Definition: An algorithm \mathcal{A} is an **efficient** (ϵ, δ) -PAC learning algorithm if it is an (ϵ, δ) -PAC learning algorithm whose runtime is polynomial in $\frac{1}{\epsilon}$ and $\frac{1}{\delta}$.

In other words, suppose that you want your algorithm to achieve 4% error rate rather than 5%. The runtime required to do so should not go up by an exponential factor.

12.4 PAC Learning of Conjunctions

To get a better sense of PAC learning, we will start with a completely irrelevant and uninteresting example. The purpose of this example is *only* to help understand how PAC learning works.

The setting is learning conjunctions. Your data points are binary vectors, for instance $x = \langle 0, 1, 1, 0, 1 \rangle$. Someone guarantees for you that there is some boolean conjunction that defines the true labeling of this data. For instance, $x_1 \wedge \neg x_2 \wedge x_5$ (“or” is not allowed). In formal terms, we often call the true underlying classification function the **concept**. So this is saying that the concept you are trying to learn is a conjunction. In this case, the boolean function would assign a negative label to the example above.

Since you know that the concept you are trying to learn is a conjunction, it makes sense that you would represent your function as a conjunction as well. For historical reasons, the function that you learn is often called a **hypothesis** and is often denoted h . However, in keeping with the other notation in this book, we will continue to denote it f .

Formally, the set up is as follows. There is some distribution \mathcal{D}^X over binary data points (vectors) $x = \langle x_1, x_2, \dots, x_D \rangle$. There is a fixed

concept conjunction c that we are trying to learn. There is no noise, so for any example x , its true label is simply $y = c(x)$.

What is a reasonable algorithm in this case? Suppose that you observe the example in Table 12.1. From the first example, we know that the true formula cannot include the term x_1 . If it did, this example would have to be negative, which it is not. By the same reasoning, it cannot include x_2 . By analogous reasoning, it also can neither include the term $\neg x_3$ nor the term $\neg x_4$.

This suggests the algorithm in Algorithm 12.4, colloquially the “Throw Out Bad Terms” algorithm. In this algorithm, you begin with a function that includes all possible $2D$ terms. Note that this function will initially classify everything as negative. You then process each example in sequence. On a negative example, you do nothing. On a positive example, you throw out terms from f that contradict the given positive example.

If you run this algorithm on the data in Table 12.1, the sequence of f s that you cycle through are:

$$f^0(x) = x_1 \wedge \neg x_1 \wedge x_2 \wedge \neg x_2 \wedge x_3 \wedge \neg x_3 \wedge x_4 \wedge \neg x_4 \quad (12.3)$$

$$f^1(x) = \neg x_1 \wedge \neg x_2 \wedge x_3 \wedge x_4 \quad (12.4)$$

$$f^2(x) = \neg x_1 \wedge x_3 \wedge x_4 \quad (12.5)$$

$$f^3(x) = \neg x_1 \wedge x_3 \wedge x_4 \quad (12.6)$$

The first thing to notice about this algorithm is that after processing an example, it is guaranteed to classify that example correctly. This observation *requires* that there is no noise in the data.

The second thing to notice is that it’s very computationally efficient. Given a data set of N examples in D dimensions, it takes $\mathcal{O}(ND)$ time to process the data. This is linear in the size of the data set.

However, in order to be an efficient (ϵ, δ) -PAC learning algorithm, you need to be able to get a bound on the **sample complexity** of this algorithm. Sure, you know that its run time is linear in the number of example N . But *how many* examples N do you need to see in order to guarantee that it achieves an error rate of at most ϵ (in all but δ -many cases)? Perhaps N has to be gigantic (like $2^{D/\epsilon}$) to (probably) guarantee a small error.

The goal is to prove that the number of samples N required to (probably) achieve a small error is not-too-big. The general proof technique for this has essentially the same flavor as almost every PAC learning proof around. First, you define a “bad thing.” In this case, a “bad thing” is that there is some term (say $\neg x_8$) that should have been thrown out, but wasn’t. Then you say: well, bad things happen. Then you notice that if this bad thing happened, you must not have

y	x_1	x_2	x_3	x_4
+1	0	0	1	1
+1	0	1	1	1
-1	1	1	0	1

Table 12.1: Data set for learning conjunctions.

?

Verify that Algorithm 12.4 maintains an *invariant* that it always errs on the side of classifying examples negative and never errs the other way.

Algorithm 31 BINARYCONJUNCTIONTRAIN(**D**)

```

1:  $f \leftarrow x_1 \wedge \neg x_1 \wedge x_2 \wedge \neg x_2 \wedge \dots \wedge x_D \wedge \neg x_D$  // initialize function
2: for all positive examples  $(x, +1)$  in D do
3:   for  $d = 1 \dots D$  do
4:     if  $x_d = 0$  then
5:        $f \leftarrow f$  without term " $x_d$ "
6:     else
7:        $f \leftarrow f$  without term " $\neg x_d$ "
8:     end if
9:   end for
10: end for
11: return  $f$ 

```

seen any positive training examples with $x_8 = 0$. So example with $x_8 = 0$ must have low probability (otherwise you would have seen them). So bad things must not be that common.

Theorem 14. *With probability at least $(1 - \delta)$: Algorithm 12.4 requires at most $N = \dots$ examples to achieve an error rate $\leq \epsilon$.*

Proof of Theorem 14. Let c be the concept you are trying to learn and let \mathcal{D} be the distribution that generates the data.

A learned function f can make a mistake if it contains *any* term t that is not in c . There are initially $2D$ many terms in f , and any (or all!) of them might not be in c . We want to ensure that the probability that f makes an error is at most ϵ . It is sufficient to ensure that

For a term t (e.g., $\neg x_5$), we say that t “negates” an example x if $t(x) = 0$. Call a term t “bad” if (a) it does not appear in c and (b) has probability at least $\epsilon/2D$ of appearing (with respect to the unknown distribution \mathcal{D} over data points).

First, we show that if we have no bad terms left in f , then f has an error rate at most ϵ .

We know that f contains at most $2D$ terms, since it begins with $2D$ terms and throws them out.

The algorithm begins with $2D$ terms (one for each variable and one for each negated variable). Note that f will only make one type of error: it can call positive examples negative, but can never call a negative example positive. Let c be the true concept (true boolean formula) and call a term “bad” if it does not appear in c . A specific bad term (e.g., $\neg x_5$) will cause f to err only on positive examples that contain a corresponding bad value (e.g., $x_5 = 1$). TODO... finish this □

What we’ve shown in this theorem is that: *if* the true underlying concept is a boolean conjunction, *and* there is no noise, *then* the “Throw Out Bad Terms” algorithm needs $N \leq \dots$ examples in order

to learn a boolean conjunction that is $(1 - \delta)$ -likely to achieve an error of at most ϵ . That is to say, that the **sample complexity** of “Throw Out Bad Terms” is Moreover, since the algorithm’s runtime is linear in N , it is an efficient PAC learning algorithm.

12.5 Occam’s Razor: Simple Solutions Generalize

The previous example of boolean conjunctions is mostly just a warm-up exercise to understand PAC-style proofs in a concrete setting.

In this section, you get to generalize the above argument to a much larger range of learning problems. We will still assume that there is no noise, because it makes the analysis much simpler. (Don’t worry: noise will be added eventually.)

William of Occam (c. 1288 – c. 1348) was an English friar and philosopher who is most famous for what later became known as Occam’s razor and popularized by Bertrand Russell. The principle basically states that you should only assume as much as you need. Or, more verbosely, “if one can explain a phenomenon without assuming this or that hypothetical entity, then there is no ground for assuming it i.e. that one should always opt for an explanation in terms of the fewest possible number of causes, factors, or variables.” What Occam actually wrote is the quote that began this chapter.

In a machine learning context, a reasonable paraphrase is “simple solutions generalize well.” In other words, you have 10,000 features you could be looking at. If you’re able to explain your predictions using just 5 of them, or using all 10,000 of them, then you should just use the 5.

The Occam’s razor theorem states that this is a good idea, theoretically. It essentially states that if you are learning some unknown concept, and if you are able to fit your training data perfectly, but you don’t need to resort to a huge class of possible functions to do so, then your learned function will generalize well. It’s an amazing theorem, due partly to the simplicity of its proof. In some ways, the proof is actually *easier* than the proof of the boolean conjunctions, though it follows the same basic argument.

In order to state the theorem explicitly, you need to be able to think about a **hypothesis class**. This is the set of possible hypotheses that your algorithm searches through to find the “best” one. In the case of the boolean conjunctions example, the hypothesis class, \mathcal{H} , is the set of all boolean formulae over D -many variables. In the case of a perceptron, your hypothesis class is the set of all possible linear classifiers. The hypothesis class for boolean conjunctions is *finite*; the hypothesis class for linear classifiers is *infinite*. For Occam’s razor, we can only work with finite hypothesis classes.

Theorem 15 (Occam’s Bound). *Suppose \mathcal{A} is an algorithm that learns a function f from some finite hypothesis class \mathcal{H} . Suppose the learned function always gets zero error on the training data. Then, the sample complexity of f is at most $\log |\mathcal{H}|$.*

TODO COMMENTS

Proof of Theorem 15. TODO □

This theorem applies directly to the “Throw Out Bad Terms” algorithm, since (a) the hypothesis class is finite and (b) the learned function always achieves zero error on the training data. To apply Occam’s Bound, you need only compute the size of the hypothesis class \mathcal{H} of boolean conjunctions. You can compute this by noticing that there are a total of $2D$ possible terms in any formula in \mathcal{H} . Moreover, each term may or may not be in a formula. So there are $2^{2D} = 4^D$ possible formulae; thus, $|\mathcal{H}| = 4^D$. Applying Occam’s Bound, we see that the sample complexity of this algorithm is $N \leq \dots$

Of course, Occam’s Bound is general enough to capture other learning algorithms as well. In particular, it can capture decision trees! In the no-noise setting, a decision tree will always fit the training data perfectly. The only remaining difficulty is to compute the size of the hypothesis class of a decision tree learner.

For simplicity’s sake, suppose that our decision tree algorithm always learns complete trees: i.e., every branch from root to leaf is length D . So the number of split points in the tree (i.e., places where a feature is queried) is 2^{D-1} . (See Figure 12.1.) Each split point needs to be assigned a feature: there D -many choices here. This gives $D2^{D-1}$ trees. The last thing is that there are 2^D leaves of the tree, each of which can take two possible values, depending on whether this leaf is classified as $+1$ or -1 : this is $2 \times 2^D = 2^{D+1}$ possibilities. Putting this all together gives a total number of trees $|\mathcal{H}| = D2^{D-1}2^{D+1} = D2^{2D} = D4^D$. Applying Occam’s Bound, we see that *TODO* examples is enough to learn a decision tree!

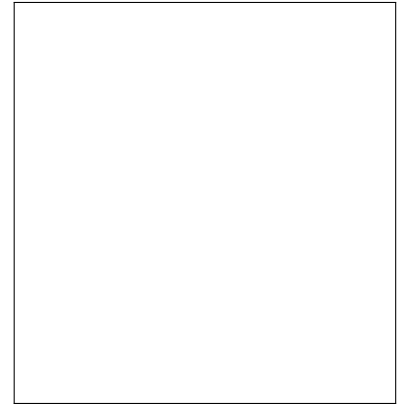


Figure 12.1: thy:dt: picture of full decision tree

12.6 Complexity of Infinite Hypothesis Spaces

Occam’s Bound is a fantastic result for learning over *finite* hypothesis spaces. Unfortunately, it is completely useless when $|\mathcal{H}| = \infty$. This is because the proof works by using each of the N training examples to “throw out” bad hypotheses until only a small number are left. But if $|\mathcal{H}| = \infty$, and you’re throwing out a finite number at each step, there will always be an infinite number remaining.

This means that, if you want to establish sample complexity results for infinite hypothesis spaces, you need some new way of measuring

their “size” or “complexity.” A prototypical way of doing this is to measure the complexity of a hypothesis class as the *number of different things it can do*.

As a silly example, consider boolean conjunctions again. Your input is a vector of binary features. However, instead of representing your hypothesis as a boolean conjunction, you choose to represent it as a conjunction of *inequalities*. That is, instead of writing $x_1 \wedge \neg x_2 \wedge x_5$, you write $[x_1 > 0.2] \wedge [x_2 < 0.77] \wedge [x_5 < \pi/4]$. In this representation, for each feature, you need to choose an inequality ($<$ or $>$) and a threshold. Since the thresholds can be arbitrary real values, there are now infinitely many possibilities: $|\mathcal{H}| = 2^D \times \infty = \infty$. However, you can immediately recognize that on binary features, there really is no difference between $[x_2 < 0.77]$ and $[x_2 < 0.12]$ and any other number of infinitely many possibilities. In other words, *even though there are infinitely many hypotheses, there are only finitely many behaviors*.

The **Vapnik-Chervonenkis dimension** (or **VC dimension**) is a classic measure of complexity of infinite hypothesis classes based on this intuition³. The VC dimension is a very classification-oriented notion of complexity. The idea is to look at a finite set of unlabeled examples, such as those in Figure 12.2. The question is: no matter how these points were labeled, would we be able to find a hypothesis that correctly classifies them. The idea is that as you add more points, being able to represent an arbitrary labeling becomes harder and harder. For instance, regardless of how the three points are labeled, you can find a linear classifier that agrees with that classification. However, for the four points, there exists a labeling for which you cannot find a perfect classifier. The VC dimension is the *maximum* number of points for which you can always find such a classifier.

You can think of VC dimension as a game between you and an adversary. To play this game, *you* choose K unlabeled points however you want. Then your adversary looks at those K points and assigns binary labels to them however they want. You must then find a hypothesis (classifier) that agrees with their labeling. You win if you can find such a hypothesis; they win if you cannot. The VC dimension of your hypothesis class is the *maximum* number of points K so that you can always win this game. This leads to the following formal definition, where you can interpret *there exists* as *your move* and *for all* as *adversary’s move*.

Definitions 2. For data drawn from some space \mathcal{X} , the *VC dimension* of a hypothesis space \mathcal{H} over \mathcal{X} is the maximal K such that: *there exists* a set $X \subseteq \mathcal{X}$ of size $|X| = K$, such that *for all* binary labelings of X , *there exists* a function $f \in \mathcal{H}$ that matches this labeling.



Figure 12.2: thy:vcex: figure with three and four examples

³ Yes, this is the same Vapnik who is credited with the creation of the support vector machine.

? What is that labeling? What is it’s name?

In general, it is much easier to show that the VC dimension is at least some value; it is much harder to show that it is at most some value. For example, following on the example from Figure 12.2, the image of three points (plus a little argumentation) is enough to show that the VC dimension of linear classifiers in two dimension is *at least three*.

To show that the VC dimension is *exactly three* it suffices to show that you *cannot* find a set of four points such that you win this game against the adversary. This is much more difficult. In the proof that the VC dimension is at least three, you simply need to provide an example of three points, and then work through the small number of possible labelings of that data. To show that it is at most three, you need to argue that no matter what set of four point you pick, you cannot win the game.

12.7 Further Reading

TODO

13 | ENSEMBLE METHODS

This is the central illusion in life: that randomness is a risk, that it is a bad thing. . .

– Nassim Nicholas Taleb

Learning Objectives:

- Implement bagging and explain how it reduces variance in a predictor.
- Explain the difference between a weak learner and a strong learner.
- Derive the AdaBoost algorithm.
- Understand the relationship between boosting decision stumps and linear classification.

GROUPS OF PEOPLE CAN OFTEN MAKE BETTER DECISIONS than individuals, especially when group members each come in with their own biases. The same is true in machine learning. Ensemble methods are learning models that achieve performance by combining the opinions of multiple learners. In doing so, you can often get away with using much *simpler* learners and still achieve great performance. Moreover, ensembles are inherently parallel, which can make them much more efficient at training and test time, if you have access to multiple processors.

In this chapter, you will learn about various ways of combining **base learners** into **ensembles**. One of the shocking results we will see is that you can take a learning model that only ever does slightly better than chance, and turn it into an arbitrarily good learning model, though a technique known as **boosting**. You will also learn how ensembles can decrease the variance of predictors as well as perform regularization.

Dependencies:

13.1 Voting Multiple Classifiers

All of the learning algorithms you have seen so far are deterministic. If you train a decision tree multiple times on the same data set, you will always get the same tree back. In order to get an effect out of voting multiple classifiers, they need to differ. There are two primary ways to get variability. You can either change the learning algorithm or change the data set.

Building an ensemble by training different classifiers is the most straightforward approach. As in single-model learning, you are given a data set (say, for classification). Instead of learning a single classifier (e.g., a decision tree) on this data set, you learn multiple different classifiers. For instance, you might train a decision tree, a perceptron, a KNN, and multiple neural networks with different architectures. Call these classifiers f_1, \dots, f_M . At test time, you can make a prediction by *voting*. On a test example \hat{x} , you compute $\hat{y}_1 = f_1(\hat{x}), \dots$,

$\hat{y}_M = f_M(\hat{x})$. If there are more +1s in the list $\langle y_1, \dots, y_M \rangle$ then you predict +1; otherwise you predict -1.

The main advantage of ensembles of different classifiers is that it is unlikely that all classifiers will make the same mistake. In fact, as long as every error is made by a minority of the classifiers, you will achieve optimal classification! Unfortunately, the inductive biases of different learning algorithms are highly correlated. This means that different algorithms are prone to similar types of errors. In particular, ensembles tend to reduce the **variance** of classifiers. So if you have a classification algorithm that tends to be very sensitive to small changes in the training data, ensembles are likely to be useful.

Note that the voting scheme naturally extends to multiclass classification. However, it does not make sense in the contexts of regression, ranking or collective classification. This is because you will rarely see the same exact output predicted twice by two different regression models (or ranking models or collective classification models). For regression, a simple solution is to take the *mean* or *median* prediction from the different models. For ranking and collective classification, different approaches are required.

Instead of training different types of classifiers on the same data set, you can train a single type of classifier (e.g., decision tree) on multiple data sets. The question is: where do these multiple data sets come from, since you're only given one at training time?

One option is to fragment your original data set. For instance, you could break it into 10 pieces and build decision trees on each of these pieces individually. Unfortunately, this means that each decision tree is trained on only a very small part of the entire data set and is likely to perform poorly.

A better solution is to use **bootstrap resampling**. This is a technique from the statistics literature based on the following observation. The data set we are given, D , is a sample drawn i.i.d. from an unknown distribution \mathcal{D} . If we draw a *new* data set \tilde{D} by random sampling from D with replacement¹, then \tilde{D} is *also* a sample from \mathcal{D} . Figure 13.1 shows the process of bootstrap resampling of ten objects.

Applying this idea to ensemble methods yields a technique known as **bagging**. You start with a single data set D that contains N training examples. From this single data set, you create M -many “bootstrapped training sets” $\tilde{D}_1, \dots, \tilde{D}_M$. Each of these bootstrapped sets also contains N training examples, drawn randomly from D with replacement. You can then train a decision tree (or other model) separately on each of these data sets to obtain classifiers f_1, \dots, f_M . As before, you can use these classifiers to vote on new test points.

Note that the bootstrapped data sets will be similar. However, they will not be *too* similar. For example, if N is large then the number of

? Which of the classifiers you've learned about so far have high variance?

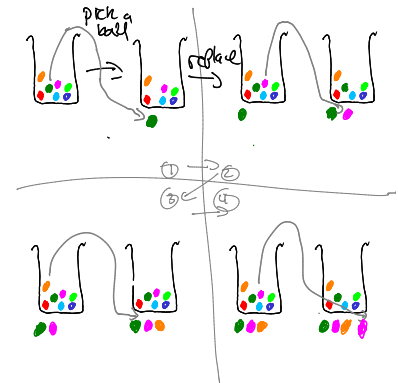


Figure 13.1: picture of sampling with replacement

¹ To sample with replacement, imagine putting all items from D in a hat. To draw a single sample, pick an element at random from that hat, write it down, and then *put it back*.

examples that are not present in any particular bootstrapped sample is relatively large. The probability that the first training example is not selected once is $(1 - 1/N)$. The probability that it is not selected at all is $(1 - 1/N)^N$. As $N \rightarrow \infty$, this tends to $1/e \approx 0.3679$. (Already for $N = 1000$ this is correct to four decimal points.) So only about 63% of the original training examples will be represented in any given bootstrapped set.

Since bagging tends to reduce *variance*, it provides an alternative approach to regularization. That is, even if each of the learned classifiers f_1, \dots, f_M are individually overfit, they are likely to be overfit to different things. Through voting, you are able to overcome a significant portion of this overfitting. Figure 13.2 shows this effect by comparing regularization via hyperparameters to regularization via bagging.

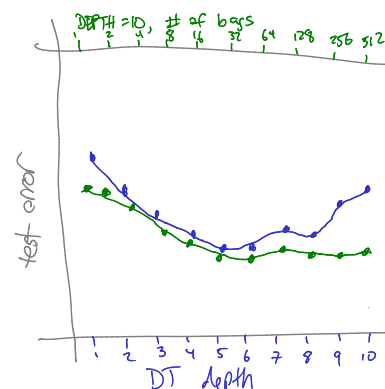


Figure 13.2: graph depicting overfitting using regularization versus bagging

13.2 Boosting Weak Learners

Boosting is the process of taking a crummy learning algorithm (technically called a **weak learner**) and turning it into a great learning algorithm (technically, a **strong learner**). Of all the ideas that originated in the theoretical machine learning community, boosting has had—perhaps—the greatest practical impact. The idea of boosting is reminiscent of what you (like me!) might have thought when you first learned about file compression. If I compress a file, and then re-compress it, and then re-compress it, eventually I’ll end up with a final that’s only one byte in size!

To be more formal, let’s define a **strong learning algorithm** \mathcal{L} as follows. When given a desired error rate ϵ , a failure probability δ and access to “enough” labeled examples from some distribution \mathcal{D} , then, with high probability (at least $1 - \delta$), \mathcal{L} learns a classifier f that has error at most ϵ . This is precisely the definition of **PAC** learning that you learned about in Chapter 12. Building a strong learning algorithm might be difficult. We can as if, instead, it is possible to build a **weak learning algorithm** \mathcal{W} that only has to achieve an error rate of 49%, rather than some arbitrary user-defined parameter ϵ . (49% is arbitrary: anything strictly less than 50% would be fine.)

Boosting is more of a “framework” than an algorithm. It’s a framework for taking a weak learning algorithm \mathcal{W} and turning it into a strong learning algorithm. The particular boosting algorithm discussed here is **AdaBoost**, short for “adaptive boosting algorithm.” AdaBoost is famous because it was one of the first *practical* boosting algorithms: it runs in polynomial time and does not require you to define a large number of hyperparameters. It gets its name from the latter benefit: it automatically *adapts* to the data that you give it.

Algorithm 32 ADABOOST($\mathcal{W}, \mathcal{D}, K$)

```

1:  $\mathbf{d}^{(0)} \leftarrow \langle \frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N} \rangle$  // Initialize uniform importance to each example
2: for  $k = 1 \dots K$  do
3:    $f^{(k)} \leftarrow \mathcal{W}(\mathcal{D}, \mathbf{d}^{(k-1)})$  // Train  $k$ th classifier on weighted data
4:    $\hat{y}_n \leftarrow f^{(k)}(x_n), \forall n$  // Make predictions on training data
5:    $\hat{\epsilon}^{(k)} \leftarrow \sum_n d_n^{(k-1)} [y_n \neq \hat{y}_n]$  // Compute weighted training error
6:    $\alpha^{(k)} \leftarrow \frac{1}{2} \log \left( \frac{1 - \hat{\epsilon}^{(k)}}{\hat{\epsilon}^{(k)}} \right)$  // Compute “adaptive” parameter
7:    $d_n^{(k)} \leftarrow \frac{1}{Z} d_n^{(k-1)} \exp[-\alpha^{(k)} y_n \hat{y}_n], \forall n$  // Re-weight examples and normalize
8: end for
9: return  $f(\hat{x}) = \text{sgn} [\sum_k \alpha^{(k)} f^{(k)}(\hat{x})]$  // Return (weighted) voted classifier

```

The intuition behind AdaBoost is like studying for an exam by using a past exam. You take the past exam and grade yourself. The questions that you got right, you pay less attention to. Those that you got *wrong*, you study more. Then you take the exam again and repeat this process. You continually *down-weight* the importance of questions you routinely answer correctly and *up-weight* the importance of questions you routinely answer incorrectly. After going over the exam multiple times, you hope to have mastered everything.

The precise AdaBoost training algorithm is shown in Algorithm 13.2. The basic functioning of the algorithm is to maintain a *weight distribution* \mathbf{d} , over data points. A weak learner, $f^{(k)}$ is trained on this weighted data. (Note that we implicitly assume that our weak learner can accept weighted training data, a relatively mild assumption that is nearly always true.) The (weighted) error rate of $f^{(k)}$ is used to determine the *adaptive parameter* α , which controls how “important” $f^{(k)}$ is. As long as the weak learner does, indeed, achieve $< 50\%$ error, then α will be greater than zero. As the error drops to zero, α grows without bound.

After the adaptive parameter is computed, the weight distribution is updated for the next iteration. As desired, examples that are correctly classified (for which $y_n \hat{y}_n = +1$) have their weight *decreased* multiplicatively. Examples that are incorrectly classified ($y_n \hat{y}_n = -1$) have their weight *increased* multiplicatively. The Z term is a normalization constant to ensure that the sum of \mathbf{d} is one (i.e., \mathbf{d} can be interpreted as a distribution). The final classifier returned by AdaBoost is a weighted vote of the individual classifiers, with weights given by the adaptive parameters.

To better understand why α is defined as it is, suppose that our weak learner simply returns a *constant* function that returns the (weighted) majority class. So if the total weight of positive examples exceeds that of negative examples, $f(x) = +1$ for all x ; otherwise $f(x) = -1$ for all x . To make the problem moderately interesting, suppose that in the original training set, there are 80 positive ex-

?

What happens if the weak learning assumption is violated and $\hat{\epsilon}$ is equal to 50%? What if it is worse than 50%? What does this mean, in practice?

amples and 20 negative examples. In this case, $f^{(1)}(x) = +1$. It's weighted error rate will be $\hat{\epsilon}^{(1)} = 0.2$ because it gets every negative example wrong. Computing, we get $\alpha^{(1)} = \frac{1}{2} \log 4$. Before normalization, we get the new weight for each positive (correct) example to be $1 \exp[-\frac{1}{2} \log 4] = \frac{1}{2}$. The weight for each negative (incorrect) example becomes $1 \exp[\frac{1}{2} \log 4] = 2$. We can compute $Z = 80 \times \frac{1}{2} + 20 \times 2 = 80$. Therefore, after normalization, the weight distribution on any single positive example is $\frac{1}{160}$ and the weight on any negative example is $\frac{1}{40}$. However, since there are 80 positive examples and 20 negative examples, the *cumulative* weight on all positive examples is $80 \times \frac{1}{160} = \frac{1}{2}$; the cumulative weight on all negative examples is $20 \times \frac{1}{40} = \frac{1}{2}$. Thus, after a single boosting iteration, the data has become precisely evenly weighted. This guarantees that in the next iteration, our weak learner *must* do something more interesting than majority voting if it is to achieve an error rate *less than* 50%, as required.

One of the major attractions of boosting is that it is perhaps easy to design computationally efficient weak learners. A very popular type of weak learner is a **shallow decision tree**: a decision tree with a small depth limit. Figure 13.3 shows test error rates for decision trees of different maximum depths (the different curves) run for differing numbers of boosting iterations (the x-axis). As you can see, if you are willing to boost for many iterations, very shallow trees are quite effective.

In fact, a very popular weak learner is a decision **decision stump**: a decision tree that can only ask *one* question. This may seem like a silly model (and, in fact, it is on its own), but when combined with boosting, it becomes very effective. To understand why, suppose for a moment that our data consists only of binary features, so that any question that a decision tree might ask is of the form “is feature 5 on?” By concentrating on decision stumps, all weak functions must have the form $f(x) = s(2x_d - 1)$, where $s \in \{\pm 1\}$ and d indexes *some* feature.

Now, consider the *final* form of a function learned by AdaBoost. We can expand it as follow, where we let f_k denote the single feature selected by the k th decision stump and let s_k denote its sign:

$$f(x) = \text{sgn} \left[\sum_k \alpha_k f^{(k)}(x) \right] \quad (13.1)$$

$$= \text{sgn} \left[\sum_k \alpha_k s_k (2x_{f_k} - 1) \right] \quad (13.2)$$

$$= \text{sgn} \left[\sum_k 2\alpha_k s_k x_{f_k} - \sum_k \alpha_k s_k \right] \quad (13.3)$$

$$= \text{sgn} [w \cdot x + b] \quad (13.4)$$

This example uses concrete numbers, but the same result holds no matter what the data distribution looks like nor how many examples there are. Write out the general case to see that you will still arrive at an even weighting after one iteration.

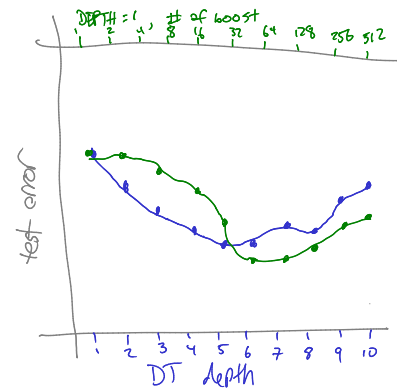


Figure 13.3: perf comparison of depth vs # boost

Why do the functions have this form?

Algorithm 33 `RANDOMFORESTTRAIN`(\mathcal{D} , *depth*, K)

```

1: for  $k = 1 \dots K$  do
2:    $t^{(k)} \leftarrow$  complete binary tree of depth depth with random feature splits
3:    $f^{(k)} \leftarrow$  the function computed by  $t^{(k)}$ , with leaves filled in by  $\mathcal{D}$ 
4: end for
5: return  $f(\hat{x}) = \text{sgn} [\sum_k f^{(k)}(\hat{x})]$            // Return voted classifier

```

$$\text{where } w_d = \sum_{k:f_k=d} 2\alpha_k s_k \quad \text{and} \quad b = - \sum_k \alpha_k s_k \quad (13.5)$$

Thus, when working with decision stumps, AdaBoost actually provides an algorithm for learning **linear classifiers**! In fact, this connection has recently been strengthened: you can show that AdaBoost provides an algorithm for optimizing **exponential loss**. (However, this connection is beyond the scope of this book.)

As a further example, consider the case of boosting a **linear classifier**. In this case, if we let the k th weak classifier be parameterized by $w^{(k)}$ and $b^{(k)}$, the overall predictor will have the form:

$$f(x) = \text{sgn} \left[\sum_k \alpha_k \text{sgn} (w^{(k)} \cdot x + b^{(k)}) \right] \quad (13.6)$$

You can notice that this is nothing but a two-layer **neural network**, with K -many hidden units! Of course it's not a classically trained neural network (once you learn $w^{(k)}$ you never go back and update it), but the structure is identical.

13.3 Random Ensembles

One of the most computationally expensive aspects of ensembles of decision trees is training the decision trees. This is very fast for decision stumps, but for deeper trees it can be prohibitively expensive. The expensive part is choosing the tree structure. Once the tree structure is chosen, it is very cheap to fill in the leaves (i.e., the predictions of the trees) using the training data.

An efficient and surprisingly effective alternative is to use trees with fixed structures and random features. Collections of trees are called forests, and so classifiers built like this are called **random forests**. The random forest training algorithm, shown in Algorithm 13.3 is quite short. It takes three arguments: the data, a desired depth of the decision trees, and a number K of total decision trees to build.

The algorithm generates each of the K trees independently, which makes it very easy to parallelize. For each trees, it constructs a full binary tree of depth *depth*. The features used at the branches of this

tree are selected randomly, typically *with replacement*, meaning that the same feature can appear multiple times, even in one branch. The leaves of this tree, where predictions are made, are filled in based on the training data. This last step is the *only* point at which the training data is used. The resulting classifier is then just a voting of the K -many random trees.

The most amazing thing about this approach is that it actually works remarkably well. It tends to work best when all of the features are at least marginally relevant, since the number of features selected for any given tree is small. An intuitive reason that it works well is the following. Some of the trees will query on useless features. These trees will essentially make random predictions. But some of the trees will happen to query on good features and will make good predictions (because the leaves are estimated based on the training data). If you have enough trees, the random ones will wash out as noise, and only the good trees will have an effect on the final classification.

13.4 *Further Reading*

TODO further reading

14 | EFFICIENT LEARNING

One essential object is to choose that arrangement which shall tend to reduce to a minimum the time necessary for completing the calculation.

– Ada Lovelace

Learning Objectives:

- Understand and be able to implement stochastic gradient descent algorithms.
- Compare and contrast small versus large batch sizes in stochastic optimization.
- Derive subgradients for sparse regularizers.
- Implement feature hashing.

SO FAR, OUR FOCUS HAS BEEN ON *models* of learning and basic algorithms for those models. We have not placed much emphasis on how to learn *quickly*. The basic techniques you learned about so far are enough to get learning algorithms running on tens or hundreds of thousands of examples. But if you want to build an algorithm for web page ranking, you will need to deal with millions or billions of examples, in hundreds of thousands of dimensions. The basic approaches you have seen so far are insufficient to achieve such a massive scale.

In this chapter, you will learn some techniques for scaling learning algorithms. These are useful even when you do not have billions of training examples, because it's always nice to have a program that runs quickly. You will see techniques for speeding up both model training and model prediction. The focus in this chapter is on linear models (for simplicity), but most of what you will learn applies more generally.

Dependencies:

14.1 What Does it Mean to be Fast?

Everyone always wants fast algorithms. In the context of machine learning, this can mean many things. You might want fast training algorithms, or perhaps training algorithms that scale to very large data sets (for instance, ones that will not fit in main memory). You might want training algorithms that can be easily parallelized. Or, you might not care about training efficiency, since it is an offline process, and only care about how quickly your learned functions can make classification decisions.

It is important to separate out these desires. If you care about efficiency at training time, then what you are really asking for are more efficient learning algorithms. On the other hand, if you care about efficiency at test time, then you are asking for *models* that can be quickly evaluated.

One issue that is not covered in this chapter is parallel learning.

This is largely because it is currently not a well-understood area in machine learning. There are many aspects of parallelism that come into play, such as the speed of communication across the network, whether you have shared memory, etc. Right now, this the general, poor-man's approach to parallelization, is to employ ensembles.

14.2 Stochastic Optimization

During training of most learning algorithms, you consider the entire data set simultaneously. This is certainly true of gradient descent algorithms for regularized linear classifiers (recall Algorithm 7.4), in which you first compute a gradient over the entire training data (for simplicity, consider the unbiased case):

$$\mathbf{g} = \sum_n \nabla_{\mathbf{w}} \ell(y_n, \mathbf{w} \cdot \mathbf{x}_n) + \lambda \mathbf{w} \quad (14.1)$$

where $\ell(y, \hat{y})$ is some loss function. Then you update the weights by $\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}$. In this algorithm, in order to make a *single* update, you have to look at *every* training example.

When there are billions of training examples, it is a bit silly to look at every one before doing anything. Perhaps just on the basis of the first few examples, you can already start learning something!

Stochastic optimization involves thinking of your training data as a big distribution over examples. A draw from this distribution corresponds to picking some example (uniformly at random) from your data set. Viewed this way, the optimization problem becomes a **stochastic optimization** problem, because you are trying to optimize some function (say, a regularized linear classifier) over a probability distribution. You can derive this interpretation directly as follows:

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \sum_n \ell(y_n, \mathbf{w} \cdot \mathbf{x}_n) + R(\mathbf{w}) \quad \text{definition} \quad (14.2)$$

$$= \arg \max_{\mathbf{w}} \sum_n \left[\ell(y_n, \mathbf{w} \cdot \mathbf{x}_n) + \frac{1}{N} R(\mathbf{w}) \right] \quad \text{move } R \text{ inside sum} \quad (14.3)$$

$$= \arg \max_{\mathbf{w}} \sum_n \left[\frac{1}{N} \ell(y_n, \mathbf{w} \cdot \mathbf{x}_n) + \frac{1}{N^2} R(\mathbf{w}) \right] \quad \text{divide through by } N \quad (14.4)$$

$$= \arg \max_{\mathbf{w}} \mathbb{E}_{(y, \mathbf{x}) \sim D} \left[\ell(y, \mathbf{w} \cdot \mathbf{x}) + \frac{1}{N} R(\mathbf{w}) \right] \quad \text{write as expectation} \quad (14.5)$$

$$\text{where } D \text{ is the training data distribution} \quad (14.6)$$

Given this framework, you have the following general form of an

Algorithm 34 STOCHASTIC GRADIENT DESCENT($\mathcal{F}, \mathcal{D}, S, K, \eta_1, \dots$)

```

1:  $\mathbf{z}^{(0)} \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize variable we are optimizing
2: for  $k = 1 \dots K$  do
3:    $D^{(k)} \leftarrow S$ -many random data points from  $\mathcal{D}$ 
4:    $\mathbf{g}^{(k)} \leftarrow \nabla_{\mathbf{z}} \mathcal{F}(D^{(k)})|_{\mathbf{z}^{(k-1)}}$  // compute gradient on sample
5:    $\mathbf{z}^{(k)} \leftarrow \mathbf{z}^{(k-1)} - \eta^{(k)} \mathbf{g}^{(k)}$  // take a step down the gradient
6: end for
7: return  $\mathbf{z}^{(K)}$ 

```

optimization problem:

$$\min_{\mathbf{z}} \mathbb{E}_{\zeta}[\mathcal{F}(\mathbf{z}, \zeta)] \quad (14.7)$$

In the example, ζ denotes the random choice of examples over the dataset, \mathbf{z} denotes the weight vector and $\mathcal{F}(\mathbf{w}, \zeta)$ denotes the loss on that example *plus* a fraction of the regularizer.

Stochastic optimization problems are formally *harder* than regular (deterministic) optimization problems because you do not even get access to exact function values and gradients. The only access you have to the function \mathcal{F} that you wish to optimize are noisy measurements, governed by the distribution over ζ . Despite this lack of information, you can still run a gradient-based algorithm, where you simply compute *local* gradients on a current sample of data.

More precisely, you can draw a data point at random from your data set. This is analogous to drawing a single value ζ from its distribution. You can compute the gradient of \mathcal{F} just at that point. In this case of a 2-norm regularized linear model, this is simply $\mathbf{g} = \nabla_{\mathbf{w}} \ell(y, \mathbf{w} \cdot \mathbf{x}) + \frac{1}{N} \mathbf{w}$, where (y, \mathbf{x}) is the random point you selected. Given this *estimate* of the gradient (it's an estimate because it's based on a single random draw), you can take a small gradient step $\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}$.

This is the **stochastic gradient descent** algorithm (**SGD**). In practice, taking gradients with respect to a *single* data point might be too myopic. In such cases, it is useful to use a small **batch** of data. Here, you can draw 10 random examples from the training data and compute a small gradient (estimate) based on those examples: $\mathbf{g} = \sum_{m=1}^{10} \nabla_{\mathbf{w}} \ell(y_m, \mathbf{w} \cdot \mathbf{x}_m) + \frac{10}{N} \mathbf{w}$, where you need to include 10 counts of the regularizer. Popular batch sizes are 1 (single points) and 10. The generic SGD algorithm is depicted in Algorithm 14.2, which takes K -many steps over batches of S -many examples.

In stochastic gradient descent, it is *imperative* to choose good step sizes. It is also very important that the steps get smaller over time at a reasonable slow rate. In particular, convergence can be guaranteed for learning rates of the form: $\eta^{(k)} = \frac{\eta_0}{\sqrt{k}}$, where η_0 is a fixed, initial step size, typically 0.01, 0.1 or 1 depending on how quickly you ex-

pect the algorithm to converge. Unfortunately, in comparison to gradient descent, stochastic gradient is quite sensitive to the selection of a good learning rate.

There is one more practical issue related to the use of SGD as a learning algorithm: do you *really* select a random point (or subset of random points) at each step, or do you stream through the data in order. The answer is akin to the answer of the same question for the perceptron algorithm (Chapter 4). If you do not permute your data at all, very bad things can happen. If you *do* permute your data once and then do multiple passes over that same permutation, it will converge, but more slowly. In theory, you really should permute every iteration. If your data is small enough to fit in memory, this is not a big deal: you will only pay for cache misses. However, if your data is too large for memory and resides on a magnetic disk that has a slow seek time, randomly seeking to new data points for each example is prohibitively slow, and you will likely need to forgo permuting the data. The speed hit in convergence speed will almost certainly be recovered by the speed gain in not having to seek on disk routinely. (Note that the story is very different for solid state disks, on which random accesses really are quite efficient.)

14.3 Sparse Regularization

For many learning algorithms, the test-time efficiency is governed by how many features are used for prediction. This is one reason decision trees tend to be among the fastest predictors: they only use a small number of features. Especially in cases where the actual *computation* of these features is expensive, cutting down on the number that are used at test time can yield huge gains in efficiency. Moreover, the amount of memory used to make predictions is also typically governed by the number of features. (Note: this is *not* true of kernel methods like support vector machines, in which the dominant cost is the number of support vectors.) Furthermore, you may simply *believe* that your learning problem can be solved with a very small number of features: this is a very reasonable form of inductive bias.

This is the idea behind sparse models, and in particular, sparse regularizers. One of the disadvantages of a 2-norm regularizer for linear models is that they tend to never produce weights that are *exactly* zero. They get close to zero, but never hit it. To understand why, as a weight w_d approaches zero, its gradient *also* approaches zero. Thus, even if the weight *should* be zero, it will essentially never get there because of the constantly shrinking gradient.

This suggests that an alternative regularizer is required to yield a sparse inductive bias. An ideal case would be the zero-norm regular-

izer, which simply counts the number of non-zero values in a vector: $\|w\|_0 = \sum_d [w_d \neq 0]$. If you could minimize this regularizer, you would be explicitly minimizing the number of non-zero features. Unfortunately, not only is the zero-norm non-convex, it's also discrete. Optimizing it is NP-hard.

A reasonable middle-ground is the one-norm: $\|w\|_1 = \sum_d |w_d|$. It is indeed convex: in fact, it is the tightest ℓ_p norm that *is* convex. Moreover, its gradients do not go to zero as in the two-norm. Just as hinge-loss is the tightest convex upper bound on zero-one error, the one-norm is the tightest convex upper bound on the zero-norm.

At this point, you should be content. You can take your subgradient optimizer for arbitrary functions and plug in the one-norm as a regularizer. The one-norm is surely non-differentiable at $w_d = 0$, but you can simply choose any value in the range $[-1, +1]$ as a subgradient at that point. (You should choose zero.)

Unfortunately, this does not quite work the way you might expect. The issue is that the gradient might “overstep” zero and you will never end up with a solution that is particularly sparse. For example, at the end of one gradient step, you might have $w_3 = 0.6$. Your gradient might have $g_6 = 0.8$ and your gradient step (assuming $\eta = 1$) will update so that the new $w_3 = -0.2$. In the subsequent iteration, you might have $g_6 = -0.3$ and step to $w_3 = 0.1$.

This observation leads to the idea of **truncated gradients**. The idea is simple: if you have a gradient that would step you over $w_d = 0$, then just set $w_d = 0$. In the easy case when the learning rate is 1, this means that if the *sign* of $w_d - g_d$ is different than the sign of w_d then you truncate the gradient step and simply set $w_d = 0$. In other words, g_d should never be *larger* than w_d . Once you incorporate learning rates, you can express this as:

$$g_d \leftarrow \begin{cases} g_d & \text{if } w_d > 0 \text{ and } g_d \leq \frac{1}{\eta^{(k)}} w_d \\ g_d & \text{if } w_d < 0 \text{ and } g_d \geq \frac{1}{\eta^{(k)}} w_d \\ 0 & \text{otherwise} \end{cases} \quad (14.8)$$

This works quite well in the case of subgradient descent. It works somewhat less well in the case of *stochastic* subgradient descent. The problem that arises in the stochastic case is that wherever you choose to stop optimizing, you will have just touched a single example (or small batch of examples), which will increase the weights for a lot of features, before the regularizer “has time” to shrink them back down to zero. You will still end up with somewhat sparse solutions, but not as sparse as they could be. There are algorithms for dealing with this situation, but they all have a heuristic flavor to them and are beyond the scope of this book.

14.4 Feature Hashing

As much as speed is a bottleneck in prediction, so often is memory usage. If you have a very large number of features, the amount of memory that it takes to store weights for all of them can become prohibitive, especially if you wish to run your algorithm on small devices. Feature hashing is an incredibly simple technique for reducing the memory footprint of linear models, with very small sacrifices in accuracy.

The basic idea is to replace all of your features with hashed versions of those features, thus reducing your space from D -many feature weights to P -many feature weights, where P is the range of the hash function. You can actually think of hashing as a (randomized) feature mapping $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^P$, for some $P \ll D$. The idea is as follows. First, you choose a hash function h whose domain is $[D] = \{1, 2, \dots, D\}$ and whose range is $[P]$. Then, when you receive a feature vector $\mathbf{x} \in \mathbb{R}^D$, you map it to a shorter feature vector $\hat{\mathbf{x}} \in \mathbb{R}^P$. Algorithmically, you can think of this mapping as follows:

1. Initialize $\hat{\mathbf{x}} = \langle 0, 0, \dots, 0 \rangle$
2. For each $d = 1 \dots D$:
 - (a) Hash d to position $p = h(d)$
 - (b) Update the p th position by adding x_d : $\hat{x}_p \leftarrow \hat{x}_p + x_d$
3. Return $\hat{\mathbf{x}}$

Mathematically, the mapping looks like:

$$\phi(\mathbf{x})_p = \sum_d [h(d) = p] x_d = \sum_{d \in h^{-1}(p)} x_d \quad (14.9)$$

where $h^{-1}(p) = \{d : h(d) = p\}$.

In the (unrealistic) case where $P = D$ and h simply encodes a permutation, then this mapping does not change the learning problem at all. All it does is rename all of the features. In practice, $P \ll D$ and there will be collisions. In this context, a collision means that two features, which are really different, end up looking the same to the learning algorithm. For instance, “is it sunny today?” and “did my favorite sports team win last night?” might get mapped to the same location after hashing. The hope is that the learning algorithm is sufficiently robust to noise that it can handle this case well.

Consider the kernel defined by this hash mapping. Namely:

$$K^{(\text{hash})}(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z}) \quad (14.10)$$

$$= \sum_p \left(\sum_d [h(d) = p] x_d \right) \left(\sum_d [h(d) = p] z_d \right) \quad (14.11)$$

$$= \sum_p \sum_{d,e} [h(d) = p][h(e) = p] x_d z_e \quad (14.12)$$

$$= \sum_d \sum_{e \in h^{-1}(h(d))} x_d z_e \quad (14.13)$$

$$= \mathbf{x} \cdot \mathbf{z} + \sum_d \sum_{\substack{e \neq d, \\ e \in h^{-1}(h(d))}} x_d z_e \quad (14.14)$$

This **hash kernel** has the form of a linear kernel plus a small number of quadratic terms. The particular quadratic terms are exactly those given by collisions of the hash function.

There are two things to notice about this. The first is that collisions might not actually be bad things! In a sense, they're giving you a little extra representational power. In particular, if the hash function happens to select out feature pairs that benefit from being paired, then you now have a better representation. The second is that even if this doesn't happen, the quadratic term in the kernel has only a small effect on the overall prediction. In particular, if you assume that your hash function is pairwise independent (a common assumption of hash functions), then the *expected value* of this quadratic term is zero, and its variance decreases at a rate of $\mathcal{O}(P^{-2})$. In other words, if you choose $P \approx 100$, then the variance is on the order of 0.0001.

14.5 Further Reading

TODO further reading

15 | UNSUPERVISED LEARNING

[My father] advised me to sit every few months in my reading chair for an entire evening, close my eyes and try to think of new problems to solve. I took his advice very seriously and have been glad ever since that he did.

– Luis Walter Alvarez

IF YOU HAVE ACCESS TO LABELED TRAINING DATA, you know what to do. This is the “supervised” setting, in which you have a teacher telling you the right answers. Unfortunately, finding such a teacher is often difficult, expensive, or down right impossible. In those cases, you might still want to be able to analyze your data, even though you do not have labels.

Unsupervised learning is learning without a teacher. One basic thing that you might want to do with data is to **visualize** it. Sadly, it is difficult to visualize things in more than two (or three) dimensions, and most data is in hundreds of dimensions (or more). **Dimensionality reduction** is the problem of taking high dimensional data and **embedding** it in a lower dimension space. Another thing you might want to do is automatically derive a partitioning of the data into **clusters**. You’ve already learned a basic approach for doing this: the *k*-means algorithm (Chapter 3). Here you will analyze this algorithm to see why it works. You will also learn more advanced clustering approaches.

Learning Objectives:

- Explain the difference between linear and non-linear dimensionality reduction.
- Relate the view of PCA as maximizing variance with the view of it as minimizing reconstruction error.
- Implement latent semantic analysis for text data.
- Motivate manifold learning from the perspective of reconstruction error.
- Understand *K*-means clustering as distance minimization.
- Explain the importance of initialization in *k*-means and furthest-first heuristic.
- Implement agglomerative clustering.
- Argue whether spectral clustering is a clustering algorithm or a dimensionality reduction algorithm.

Dependencies:

15.1 *K*-Means Clustering, Revisited

The *K*-means clustering algorithm is re-presented in Algorithm 15.1. There are two very basic questions about this algorithm: (1) does it converge (and if so, how quickly); (2) how sensitive it is to initialization? The answers to these questions, detailed below, are: (1) yes it converges, and it converges very quickly in practice (though slowly in theory); (2) yes it is sensitive to initialization, but there are good ways to initialize it.

Consider the question of convergence. The following theorem states that the *K*-Means algorithm converges, though it does not say how quickly it happens. The method of proving the convergence is to specify a **clustering quality** objective function, and then to show that the *K*-Means algorithm converges to a (local) optimum of that objective function. The particular objective function that *K*-Means

Algorithm 35 K-MEANS(\mathbf{D}, K)

```

1: for  $k = 1$  to  $K$  do
2:    $\mu_k \leftarrow$  some random location      // randomly initialize mean for  $k$ th cluster
3: end for
4: repeat
5:   for  $n = 1$  to  $N$  do
6:      $z_n \leftarrow \operatorname{argmin}_k \|\mu_k - x_n\|$       // assign example  $n$  to closest center
7:   end for
8:   for  $k = 1$  to  $K$  do
9:      $\mu_k \leftarrow \operatorname{MEAN}(\{x_n : z_n = k\})$       // re-estimate mean of cluster  $k$ 
10:  end for
11: until converged
12: return  $z$                                 // return cluster assignments

```

is optimizing is the *sum of squared distances from any data point to its assigned center*. This is a natural generalization of the definition of a mean: the mean of a set of points is the single point that minimizes the sum of squared distances from the mean to every point in the data. Formally, the K-Means objective is:

$$\mathcal{L}(z, \mu; \mathbf{D}) = \sum_n \left\| x_n - \mu_{z_n} \right\|^2 = \sum_k \sum_{n: z_n=k} \|x_n - \mu_k\|^2 \quad (15.1)$$

Theorem 16 (K-Means Convergence Theorem). *For any dataset \mathbf{D} and any number of clusters K , the K-means algorithm converges in a finite number of iterations, where convergence is measured by \mathcal{L} ceasing the change.*

Proof of Theorem 16. The proof works as follows. There are only two points in which the K-means algorithm changes the values of μ or z : lines 6 and 9. We will show that both of these operations can never increase the value of \mathcal{L} . Assuming this is true, the rest of the argument is as follows. After the first pass through the data, there are only finitely many possible assignments to z and μ , because z is discrete and because μ can only take on a finite number of values: means of some subset of the data. Furthermore, \mathcal{L} is lower-bounded by zero. Together, this means that \mathcal{L} cannot decrease more than a finite number of times. Thus, it must stop decreasing at some point, and at that point the algorithm has converged.

It remains to show that lines 6 and 9 decrease \mathcal{L} . For line 6, when looking at example n , suppose that the previous value of z_n is a and the new value is b . It must be the case that $\|x_n - \mu_b\| \leq \|x_n - \mu_a\|$. Thus, changing from a to b can only decrease \mathcal{L} . For line 9, consider the second form of \mathcal{L} . Line 9 computes μ_k as the mean of the data points for which $z_n = k$, which is precisely the point that minimizes squared sitances. Thus, this update to μ_k can only decrease \mathcal{L} . \square

There are several aspects of K-means that are unfortunate. First, the convergence is only to a local optimum of \mathcal{L} . In practice, this

means that you should usually run it 10 times with different initializations and pick the one with minimal resulting \mathcal{L} . Second, one can show that there are input datasets and initializations on which it might take an exponential amount of time to converge. Fortunately, these cases almost never happen in practice, and in fact it has recently been shown that (roughly) if you limit the floating point precision of your machine, K -means *will* converge in polynomial time (though still only to a local optimum), using techniques of **smoothed analysis**.

The biggest practical issue in K -means is initialization. If the cluster means are initialized poorly, you often get convergence to uninteresting solutions. A useful heuristic is the **furthest-first heuristic**. This gives a way to perform a semi-random initialization that attempts to pick initial means as far from each other as possible. The heuristic is sketched below:

1. Pick a random example m and set $\mu_1 = x_m$.
2. For $k = 2 \dots K$:
 - (a) Find the example m that is as far as possible from *all* previously selected means; namely: $m = \arg \max_m \min_{k' < k} \|x_m - \mu_{k'}\|^2$ and set $\mu_k = x_m$

In this heuristic, the only bit of randomness is the selection of the first data point. After that, it is completely deterministic (except in the rare case that there are multiple equidistant points in step 2a). It is extremely important that when selecting the 3rd mean, you select that point that *maximizes* the *minimum* distance to the closest other mean. You want the point that's as far away from *all* previous means as possible.

The furthest-first heuristic is just that: a heuristic. It works very well in practice, though can be somewhat sensitive to outliers (which will often get selected as some of the initial means). However, this outlier sensitivity is usually reduced after one iteration through the K -means algorithm. Despite being just a heuristic, it is quite useful in practice.

You can turn the heuristic into an algorithm by adding a bit more randomness. This is the idea of the K -means++ algorithm, which is a simple randomized tweak on the furthest-first heuristic. The idea is that when you select the k th mean, instead of choosing the *absolute furthest* data point, you choose a data point at random, with probability proportional to its distance squared. This is made formal in Algorithm 15.1.

If you use K -means++ as an initialization for K -means, then you are able to achieve an approximation guarantee on the final value

Algorithm 36 K-MEANS++(D, K)

```

1:  $\mu_1 \leftarrow x_m$  for  $m$  chosen uniformly at random // randomly initialize first point
2: for  $k = 2$  to  $K$  do
3:    $d_n \leftarrow \min_{k' < k} \|x_n - \mu_{k'}\|^2, \forall n$  // compute distances
4:    $p \leftarrow \frac{1}{\sum_n n_d} d$  // normalize to probability distribution
5:    $m \leftarrow$  random sample from  $p$  // pick an example at random
6:    $\mu_k \leftarrow x_m$ 
7: end for
8: run K-MEANS using  $\mu$  as initial centers

```

of the objective. This doesn't tell you that you will reach the global optimum, but it *does* tell you that you will get reasonably close. In particular, if $\hat{\mathcal{L}}$ is the value obtained by running K-means++, then this will not be "too far" from $\mathcal{L}^{(\text{opt})}$, the true global minimum.

Theorem 17 (K-means++ Approximation Guarantee). *The expected value of the objective returned by K-means++ is never more than $\mathcal{O}(\log K)$ from optimal and can be as close as $\mathcal{O}(1)$ from optimal. Even in the former case, with 2^K random restarts, one restart will be $\mathcal{O}(1)$ from optimal (with high probability). Formally: $\mathbb{E}[\hat{\mathcal{L}}] \leq 8(\log K + 2)\mathcal{L}^{(\text{opt})}$. Moreover, if the data is "well suited" for clustering, then $\mathbb{E}[\hat{\mathcal{L}}] \leq \mathcal{O}(1)\mathcal{L}^{(\text{opt})}$.*

The notion of "well suited" for clustering informally states that the advantage of going from $K - 1$ clusters to K clusters is "large." Formally, it means that $\mathcal{L}_K^{(\text{opt})} \leq \epsilon^2 \mathcal{L}_{K-1}^{(\text{opt})}$, where $\mathcal{L}_K^{(\text{opt})}$ is the optimal value for clustering with K clusters, and ϵ is the desired degree of approximation. The idea is that if this condition does *not* hold, then you shouldn't bother clustering the data.

One of the biggest practical issues with K-means clustering is "choosing K ." Namely, if someone just hands you a dataset and asks you to cluster it, how many clusters should you produce? This is difficult, because increasing K will always decrease $\mathcal{L}_K^{(\text{opt})}$ (until $K > N$), and so simply using \mathcal{L} as a notion of goodness is insufficient (analogous to overfitting in a supervised setting). A number of "information criteria" have been proposed to try to address this problem. They all effectively boil down to "regularizing" K so that the model cannot grow to be too complicated. The two most popular are the Bayes Information Criteria (BIC) and the Akaike Information Criteria (AIC), defined below in the context of K-means:

$$\text{BIC: } \arg \min_K \hat{\mathcal{L}}_K + K \log D \quad (15.2)$$

$$\text{AIC: } \arg \min_K \hat{\mathcal{L}}_K + 2KD \quad (15.3)$$

The informal intuition behind these criteria is that increasing K is going to make \mathcal{L}_K go down. However, if it doesn't go down "by enough" then it's not worth doing. In the case of BIC, "by enough"

means by an amount proportional to $\log D$; in the case of AIC, it's proportional to $2D$. Thus, AIC provides a much stronger penalty for many clusters than does BIC, especially in high dimensions.

A more formal intuition for BIC is the following. You ask yourself the question “if I wanted to send this data across a network, how many bits would I need to send?” Clearly you could simply send all of the N examples, each of which would take roughly $\log D$ bits to send. This gives $N \log D$ to send all the data. Alternatively, you could first cluster the data and send the cluster centers. This will take $K \log D$ bits. Then, for each data point, you send its center as well as its deviation from that center. It turns out this will cost exactly $\hat{\mathcal{L}}_K$ bits. Therefore, the BIC is precisely measuring how many bits it will take to send your data using K clusters. The K that minimizes this number of bits is the optimal value.

15.2 Linear Dimensionality Reduction

Dimensionality reduction is the task of taking a dataset in high dimensions (say 10000) and reducing it to low dimensions (say 2) while retaining the “important” characteristics of the data. Since this is an unsupervised setting, the notion of important characteristics is difficult to define.

Consider the dataset in Figure 15.1, which lives in high dimensions (two) and you want to reduce to low dimensions (one). In the case of linear dimensionality reduction, the only thing you can do is to project the data onto a vector and use the projected distances as the embeddings. Figure 15.2 shows a projection of this data onto the vector that points in the direction of *maximal variance* of the original dataset. Intuitively, this is a reasonable notion of importance, since this is the direction in which most information is encoded in the data.

For the rest of this section, assume that the data is centered: namely, the mean of all the data is at the origin. (This will simply make the math easier.) Suppose the two dimensional data is x_1, \dots, x_N and you're looking for a vector u that points in the direction of maximal variance. You can compute this by projecting each point onto u and looking at the variance of the result. In order for the projection to make sense, you need to constrain $\|u\|^2 = 1$. In this case, the projections are $x_1 \cdot u, x_2 \cdot u, \dots, x_N \cdot u$. Call these values p_1, \dots, p_N .

The goal is to compute the variance of the $\{p_n\}$ s and then choose u to maximize this variance. To compute the variance, you first need to compute the mean. Because the mean of the x_n s was zero, the

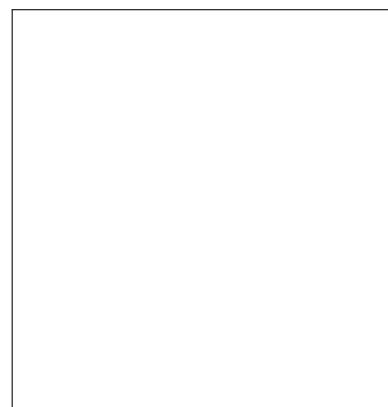


Figure 15.1: `unsup:pcadata`: Data in two dimensions

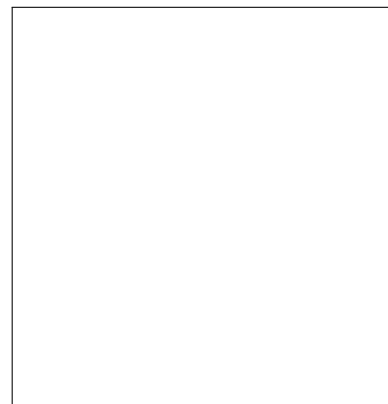


Figure 15.2: `unsup:pcadata2`: Projection of that data down to one dimension by PCA

MATH REVIEW | EIGENVALUES AND EIGENVECTORS

todo the usual...

Figure 15.3:

mean of the p_n is also zero. This can be seen as follows:

$$\sum_n p_n = \sum_n x_n \cdot u = \left(\sum_n x_n \right) \cdot u = \mathbf{0} \cdot u = 0 \quad (15.4)$$

The variance of the $\{p_n\}$ is then just $\sum_n p_n^2$. Finding the optimal u (from the perspective of variance maximization) reduces to the following optimization problem:

$$\max_u \sum_n (x_n \cdot u)^2 \quad \text{subj. to} \quad ||u||^2 = 1 \quad (15.5)$$

In this problem it becomes apparent why keeping u unit length is important: if not, u would simply stretch to have infinite length to maximize the objective.

It is now helpful to write the collection of datapoints x_n as a $N \times D$ matrix \mathbf{X} . If you take this matrix \mathbf{X} and multiply it by u , which has dimensions $D \times 1$, you end up with a $N \times 1$ vector whose values are exactly the values p . The objective in Eq (15.5) is then just the squared norm of p . This simplifies Eq (15.5) to:

$$\max_u ||\mathbf{X}u||^2 \quad \text{subj. to} \quad ||u||^2 - 1 = 0 \quad (15.6)$$

where the constraint has been rewritten to make it amenable to constructing the Lagrangian. Doing so and taking gradients yields:

$$\mathcal{L}(u, \lambda) = ||\mathbf{X}u||^2 - \lambda (||u||^2 - 1) \quad (15.7)$$

$$\nabla_u \mathcal{L} = 2\mathbf{X}^\top \mathbf{X}u - 2\lambda u \quad (15.8)$$

$$\implies \lambda u = (\mathbf{X}^\top \mathbf{X}) u \quad (15.9)$$

You can solve this expression ($\lambda u = \mathbf{X}^\top \mathbf{X}u$) by computing the first eigenvector and eigenvalue of the matrix $\mathbf{X}^\top \mathbf{X}$.

This gives you the solution to a projection into a one-dimensional space. To get a second dimension, you want to find a new vector v on which the data has maximal variance. However, to avoid redundancy, you want v to be orthogonal to u ; namely $u \cdot v = 0$. This gives:

$$\max_v ||\mathbf{X}v||^2 \quad \text{subj. to} \quad ||v||^2 = 1, \text{ and } u \cdot v = 0 \quad (15.10)$$

Following the same procedure as before, you can construct a La-

Algorithm 37 PCA(\mathbf{D}, K)

```

1:  $\mu \leftarrow \text{MEAN}(\mathbf{X})$  // compute data mean for centering
2:  $\mathbf{D} \leftarrow (\mathbf{X} - \mu \mathbf{1}^\top)^\top (\mathbf{X} - \mu \mathbf{1}^\top)$  // compute covariance,  $\mathbf{1}$  is a vector of ones
3:  $\{\lambda_k, \mathbf{u}_k\} \leftarrow$  top  $K$  eigenvalues/eigenvectors of  $\mathbf{D}$ 
4: return  $(\mathbf{X} - \mu \mathbf{1}) \mathbf{U}$  // project data using  $\mathbf{U}$ 

```

grangian and differentiate:

$$\mathcal{L}(\mathbf{v}, \lambda_1, \lambda_2) = \|\mathbf{X}\mathbf{v}\|^2 - \lambda_1 (\|\mathbf{v}\|^2 - 1) - \lambda_2 \mathbf{u} \cdot \mathbf{v} \quad (15.11)$$

$$\nabla_{\mathbf{u}} \mathcal{L} = 2\mathbf{X}^\top \mathbf{X} \mathbf{v} - 2\lambda_1 \mathbf{v} - \lambda_2 \mathbf{u} \quad (15.12)$$

$$\implies \lambda_1 \mathbf{v} = (\mathbf{X}^\top \mathbf{X}) \mathbf{v} - \frac{\lambda_2}{2} \mathbf{u} \quad (15.13)$$

However, you know that \mathbf{u} is the first eigenvector of $\mathbf{X}^\top \mathbf{X}$, so the solution to this problem for λ_1 and \mathbf{v} is given by the *second* eigenvalue/eigenvector pair of $\mathbf{X}^\top \mathbf{X}$.

Repeating this analysis inductively tells you that if you want to project onto K mutually orthogonal dimensions, you simply need to take the first K eigenvectors of the matrix $\mathbf{X}^\top \mathbf{X}$. This matrix is often called the **data covariance matrix** because $[\mathbf{X}^\top \mathbf{X}]_{i,j} = \sum_n \sum_m x_{n,i} x_{m,j}$, which is the sample covariance between features i and j .

This leads to the technique of **principle components analysis**, or **PCA**. For completeness, the is depicted in Algorithm 15.2. The important thing to note is that the eigenanalysis only gives you the projection directions. It does not give you the embedded data. To embed a data point \mathbf{x} you need to compute its embedding as $\langle \mathbf{x} \cdot \mathbf{u}_1, \mathbf{x} \cdot \mathbf{u}_2, \dots, \mathbf{x} \cdot \mathbf{u}_K \rangle$. If you write \mathbf{U} for the $D \times K$ matrix of \mathbf{u} s, then this is just \mathbf{XU} .

There is an alternative derivation of PCA that can be informative, based on *reconstruction error*. Consider the one-dimensional case again, where you are looking for a single projection direction \mathbf{u} . If you were to use this direction, your projected data would be $\mathbf{Z} = \mathbf{Xu}$. Each Z_n gives the position of the n th datapoint along \mathbf{u} . You can project this one-dimensional data back into the original space by multiplying it by \mathbf{u}^\top . This gives you *reconstructed* values \mathbf{Zu}^\top . Instead of maximizing variance, you might instead want to minimize the **reconstruction error**, defined by:

$$\|\mathbf{X} - \mathbf{Zu}^\top\|^2 = \|\mathbf{X} - \mathbf{Xu}\mathbf{u}^\top\|^2 \quad \text{definition of } \mathbf{Z} \quad (15.14)$$

$$= \|\mathbf{X}\|^2 + \|\mathbf{Xu}\mathbf{u}^\top\|^2 - 2\mathbf{X}^\top \mathbf{Xu}\mathbf{u}^\top \quad \text{quadratic rule} \quad (15.15)$$

$$= ||\mathbf{X}||^2 + ||\mathbf{X}\mathbf{u}\mathbf{u}^\top||^2 - 2\mathbf{u}^\top \mathbf{X}^\top \mathbf{X} \mathbf{u} \quad \text{quadratic rule}$$

(15.16)

$$= ||\mathbf{X}||^2 + ||\mathbf{X}||^2 - 2\mathbf{u}^\top \mathbf{X}^\top \mathbf{X} \mathbf{u} \quad \text{\textit{u} is a unit vector}$$

(15.17)

$$= C - 2 ||\mathbf{X}\mathbf{u}||^2 \quad \text{join constants, rewrite last term}$$

(15.18)

Minimizing this final term is equivalent to maximizing $||\mathbf{X}\mathbf{u}||^2$, which is exactly the form of the maximum variance derivation of PCA.

Thus, you can see that maximizing variance is identical to minimizing reconstruction error.

The same question of “what should K be” arises in dimensionality reduction as in clustering. If the purpose of dimensionality reduction is to visualize, then K should be 2 or 3. However, an alternative purpose of dimensionality reduction is to avoid the curse of dimensionality. For instance, even if you have labeled data, it might be worthwhile to reduce the dimensionality before applying supervised learning, essentially as a form of regularization. In this case, the question of an optimal K comes up again. In this case, the same criteria (AIC and BIC) that can be used for clustering can be used for PCA. The only difference is the quality measure changes from a sum of squared distances to means (for clustering) to a sum of squared distances to original data points (for PCA). In particular, for BIC you get the reconstruction error plus $K \log D$; for AIC, you get the reconstruction error plus $2KD$.

15.3 Autoencoders

TODO

15.4 Further Reading

TODO

16 | EXPECTATION MAXIMIZATION

A hen is only an egg's way of making another egg. – Samuel Butler

SUPPOSE YOU WERE BUILDING a naive Bayes model for a text categorization problem. After you were done, your boss told you that it became prohibitively expensive to obtain labeled data. You now have a probabilistic model that assumes access to labels, but you don't have any labels! Can you still do something?

Amazingly, you can. You can treat the labels as **hidden variables**, and attempt to learn them at the same time as you learn the parameters of your model. A very broad family of algorithms for solving problems just like this is the **expectation maximization** family. In this chapter, you will derive expectation maximization (EM) algorithms for clustering and dimensionality reduction, and then see why EM works.

Learning Objectives:

- Explain the relationship between parameters and hidden variables.
- Construct generative stories for clustering and dimensionality reduction.
- Draw a graph explaining how EM works by constructing convex lower bounds.
- Implement EM for clustering with mixtures of Gaussians, and contrasting it with k -means.
- Evaluate the differences between EM and gradient descent for hidden variable models.

Dependencies:

16.1 Grading an Exam without an Answer Key

Alice's machine learning professor Carlos gives out an exam that consists of 50 true/false questions. Alice's class of 100 students takes the exam and Carlos goes to grade their solutions. If Carlos made an answer key, this would be easy: he would just count the fraction of correctly answered questions each student got, and that would be their score. But, like many professors, Carlos was really busy and didn't have time to make an answer key. Can he still grade the exam?

There are two insights that suggest that he might be able to. Suppose he knew ahead of time that Alice was an awesome student, and is basically guaranteed to get 100% on the exam. In that case, Carlos can simply use Alice's answers as the ground truth. More generally, if Carlos assumes that *on average* students are better than random guessing, he can hope that the majority answer for each question is likely to be correct. Combining this with the previous insight, when doing the "voting", he might want to pay more attention to the answers of the better students.

To be a bit more pedantic, suppose there are $N = 100$ students and $M = 50$ questions. Each student n has a score s_n , between 0 and

1 that denotes how well they do on the exam. The score is what we really want to compute. For each question m and each student n , the student has provided an answer $a_{n,m}$, which is either zero or one. There is also an unknown ground truth answer for each question m , which we'll call t_m , which is also either zero or one.

As a starting point, let's consider a simple heuristic and then complexify it. The heuristic is the "majority vote" heuristic and works as follows. First, we estimate t_m as the most common answer for question m : $t_m = \operatorname{argmax}_t \sum_n \mathbf{1}[a_{n,m} = t]$. Once we have a guess for each true answer, we estimate each students' score as how many answers they produced that match this guessed key: $s_n = \frac{1}{M} \sum_m \mathbf{1}[a_{n,m} = t_m]$.

Once we have these scores, however, we might want to trust some of the students more than others. In particular, answers from students with high scores are perhaps more likely to be correct, so we can *recompute* the ground truth, according to weighted votes. The weight of the votes will be precisely the score the corresponding each student:

$$t_m = \operatorname{argmax}_t \sum_n s_n \mathbf{1}[a_{n,m} = t] \quad (16.1)$$

You can recognize this as a *chicken and egg problem*. If you knew the student's scores, you could estimate an answer key. If you had an answer key, you could compute student scores. A very common strategy in computer science for dealing with such chicken and egg problems is to iterate. Take a guess at the first, compute the second, recompute the first, and so on.

In order to develop this idea formally, we have to case the problem in terms of a probabilistic model with a generative story. The generative story we'll use is:

1. For each question m , choose a true answer $t_m \sim \mathcal{Ber}(0.5)$
2. For each student n , choose a score $s_n \sim \mathcal{Uni}(0, 1)$
3. For each question m and each student n , choose an answer $a_{n,m} \sim \mathcal{Ber}(s_n)^{t_m} \mathcal{Ber}(1 - s_n)^{1-t_m}$

In the first step, we generate the true answers independently by flipping a fair coin. In the second step, each students' overall score is determined to be a uniform random number between zero and one. The tricky step is step three, where each students' answer is generated for each question. Consider student n answering question m , and suppose that $s_n = 0.9$. If $t_m = 1$, then $a_{n,m}$ should be 1 (i.e., correct) 90% of the time; this can be accomplished by drawing the answer from $\mathcal{Ber}(0.9)$. On the other hand, if $t_m = 0$, then $a_{n,m}$ should be 1 (i.e., incorrect) 10% of the time; this can be accomplished by drawing

the answer from $\mathcal{Ber}(0.1)$. The exponent in step 3 selects which of two Bernoulli distributions to draw from, and then implements this rule.

This can be translated into the following likelihood:

$$\begin{aligned}
 p(\mathbf{a}, \mathbf{t}, \mathbf{s}) &= \left[\prod_m 0.5^{t_m} 0.5^{1-t_m} \right] \times \left[\prod_n 1 \right] \\
 &\times \left[\prod_n \prod_m s_n^{a_{n,m} t_m} (1-s_n)^{(1-a_{n,m}) t_m} \right. \\
 &\quad \left. s_n^{(1-a_{n,m})(1-t_m)} (1-s_n)^{a_{n,m}(1-t_m)} \right] \tag{16.2}
 \end{aligned}$$

$$= 0.5^M \prod_n \prod_m s_n^{a_{n,m} t_m} (1-s_n)^{(1-a_{n,m}) t_m} s_n^{(1-a_{n,m})(1-t_m)} (1-s_n)^{a_{n,m}(1-t_m)} \tag{16.3}$$

Suppose we knew the true labels \mathbf{t} . We can take the log of this likelihood and differentiate it with respect to the score s_n of some student (note: we can drop the 0.5^M term because it is just a constant):

$$\begin{aligned}
 \log p(\mathbf{a}, \mathbf{t}, \mathbf{s}) &= \sum_n \sum_m \left[a_{n,m} t_m \log s_n + (1-a_{n,m})(1-t_m) \log(s_n) \right. \\
 &\quad \left. + (1-a_{n,m}) t_m \log(1-s_n) + a_{n,m}(1-t_m) \log(1-s_n) \right] \tag{16.4}
 \end{aligned}$$

$$\frac{\partial \log p(\mathbf{a}, \mathbf{t}, \mathbf{s})}{\partial s_n} = \sum_m \left[\frac{a_{n,m} t_m + (1-a_{n,m})(1-t_m)}{s_n} - \frac{(1-a_{n,m}) t_m + a_{n,m}(1-t_m)}{1-s_n} \right] \tag{16.5}$$

The derivative has the form $\frac{A}{s_n} - \frac{B}{1-s_n}$. If we set this equal to zero and solve for s_n , we get an optimum of $s_n = \frac{A}{A+B}$. In this case:

$$A = \sum_m [a_{n,m} t_m + (1-a_{n,m})(1-t_m)] \tag{16.6}$$

$$B = \sum_m [(1-a_{n,m}) t_m + a_{n,m}(1-t_m)] \tag{16.7}$$

$$A + B = \sum_m [1] = M \tag{16.8}$$

Putting this together, we get:

$$s_n = \frac{1}{M} \sum_m [a_{n,m} t_m + (1-a_{n,m})(1-t_m)] \tag{16.9}$$

In the case of known \mathbf{t} s, this matches exactly what we had in the heuristic.

However, we do not know \mathbf{t} , so instead of using the “true” values of \mathbf{t} , we’re going to use their *expectations*. In particular, we will compute s_n by maximizing its likelihood under the *expected* values

of t , hence the name **expectation maximization**. If we are going to compute expectations of t , we have to say: expectations according to which probability distribution? We will use the distribution $p(t_m | \mathbf{a}, \mathbf{s})$. Let \tilde{t}_m denote $\mathbb{E}_{t_m \sim p(t_m | \mathbf{a}, \mathbf{s})}[t_m]$. Because t_m is a binary variable, its expectation is equal to its probability; namely: $\tilde{t}_m = p(t_m | \mathbf{a}, \mathbf{s})$.

How can we compute this? We will compute $C = p(t_m = 1, \mathbf{a}, \mathbf{s})$ and $D = p(t_m = 0, \mathbf{a}, \mathbf{s})$ and then compute $\tilde{t}_m = C/(C + D)$. The computation is straightforward:

$$C = 0.5 \prod_n s_n^{a_{n,m}} (1 - s_n)^{1-a_{n,m}} = 0.5 \prod_{\substack{n: \\ a_{n,m}=1}} s_n \prod_{\substack{n: \\ a_{n,m}=0}} (1 - s_n) \quad (16.10)$$

$$D = 0.5 \prod_n s_n^{1-a_{n,m}} (1 - s_n)^{a_{n,m}} = 0.5 \prod_{\substack{n: \\ a_{n,m}=1}} (1 - s_n) \prod_{\substack{n: \\ a_{n,m}=0}} s_n \quad (16.11)$$

If you inspect the value of C , it is basically “voting” (in a product form, not a sum form) the scores of those students who agree that the answer is 1 with one-minus-the-score of those students who do not. The value of D is doing the reverse. This is a form of multiplicative voting, which has the effect that if a given student has a perfect score of 1.0, their results will carry the vote completely.

We now have a way to:

1. Compute expected ground truth values \tilde{t}_m , given scores.
2. Optimize scores s_n given expected ground truth values.

The full solution is then to alternate between these two. You can start by initializing the ground truth values at the majority vote (this seems like a safe initialization). Given those, compute new scores. Given those new scores, compute new ground truth values. And repeat until tired.

In the next two sections, we will consider a more complex unsupervised learning model for clustering, and then a generic mathematical framework for expectation maximization, which will answer questions like: will this process converge, and, if so, to what?

16.2 Clustering with a Mixture of Gaussians

In Chapter 9, you learned about probabilistic models for classification based on density estimation. Let’s start with a fairly simple classification model that *assumes* we have labeled data. We will shortly remove this assumption. Our model will state that we have K classes, and data from class k is drawn from a Gaussian with mean μ_k and variance σ_k^2 . The choice of classes is parameterized by θ . The generative story for this model is:

1. For each example $n = 1 \dots N$:

- (a) Choose a label $y_n \sim \text{Disc}(\theta)$
- (b) Choose example $x_n \sim \mathcal{N}(\mu_{y_n}, \sigma_{y_n}^2)$

This generative story can be directly translated into a likelihood as before:

$$p(D) = \prod_n \text{Mult}(y_n \mid \theta) \text{Nor}(x_n \mid \mu_{y_n}, \sigma_{y_n}^2) \quad (16.12)$$

$$= \prod_n \underbrace{\theta_{y_n}}_{\text{choose label}} \underbrace{\left[2\pi\sigma_{y_n}^2 \right]^{-\frac{D}{2}} \exp \left[-\frac{1}{2\sigma_{y_n}^2} \|x_n - \mu_{y_n}\|^2 \right]}_{\text{choose feature values}} \quad (16.13)$$

If you had access to labels, this would be all well and good, and you could obtain closed form solutions for the maximum likelihood estimates of all parameters by taking a log and then taking gradients of the log likelihood:

$$\theta_k = \text{fraction of training examples in class } k \quad (16.14)$$

$$= \frac{1}{N} \sum_n [y_n = k]$$

$$\mu_k = \text{mean of training examples in class } k \quad (16.15)$$

$$= \frac{\sum_n [y_n = k] x_n}{\sum_n [y_n = k]}$$

$$\sigma_k^2 = \text{variance of training examples in class } k \quad (16.16)$$

$$= \frac{\sum_n [y_n = k] \|x_n - \mu_k\|^2}{\sum_n [y_n = k]}$$

Suppose that you *don't* have labels. Analogously to the *K*-means algorithm, one potential solution is to iterate. You can start off with guesses for the values of the unknown variables, and then iteratively improve them over time. In *K*-means, the approach was the *assign* examples to labels (or clusters). This time, instead of making hard assignments (“example 10 belongs to cluster 4”), we’ll make **soft assignments** (“example 10 belongs half to cluster 4, a quarter to cluster 2 and a quarter to cluster 5”). So as not to confuse ourselves too much, we’ll introduce a new variable, $z_n = \langle z_{n,1}, \dots, z_{n,K} \rangle$ (that sums to one), to denote a fractional assignment of examples to clusters.

This notion of soft-assignments is visualized in Figure 16.1. Here, we’ve depicted each example as a pie chart, and its coloring denotes the degree to which it’s been assigned to each (of three) clusters. The size of the pie pieces correspond to the z_n values.



You should be able to derive the maximum likelihood solution results formally by now.

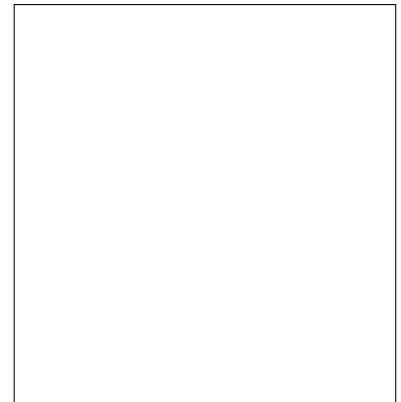


Figure 16.1: em:piecharts: A figure showing pie charts

Formally, $z_{n,k}$ denotes the probability that example n is assigned to cluster k :

$$z_{n,k} = p(y_n = k \mid \mathbf{x}_n) \quad (16.17)$$

$$= \frac{p(y_n = k, \mathbf{x}_n)}{p(\mathbf{x}_n)} \quad (16.18)$$

$$= \frac{1}{Z_n} \text{Mult}(k \mid \boldsymbol{\theta}) \text{Nor}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \sigma_k^2) \quad (16.19)$$

Here, the normalizer Z_n is to ensure that z_n sums to one.

Given a set of parameters (the θ s, μ s and σ^2 s), the **fractional assignments** $z_{n,k}$ are easy to compute. Now, akin to K -means, given fractional assignments, you need to recompute estimates of the model parameters. In analogy to the maximum likelihood solution (Eqs (??)-(??)), you can do this by counting fractional points rather than full points. This gives the following re-estimation updates:

$$\theta_k = \text{fraction of training examples in class } k \quad (16.20)$$

$$= \frac{1}{N} \sum_n z_{n,k}$$

$$\boldsymbol{\mu}_k = \text{mean of fractional examples in class } k \quad (16.21)$$

$$= \frac{\sum_n z_{n,k} \mathbf{x}_n}{\sum_n z_{n,k}}$$

$$\sigma_k^2 = \text{variance of fractional examples in class } k \quad (16.22)$$

$$= \frac{\sum_n z_{n,k} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2}{\sum_n z_{n,k}}$$

All that has happened here is that the hard assignments “[$y_n = k$]” have been replaced with soft assignments “ $z_{n,k}$ ”. As a bit of foreshadowing of what is to come, what we’ve done is essentially replace known labels with *expected labels*, hence the name “expectation maximization.”

Putting this together yields Algorithm 16.2. This is the **GMM** (“**Gaussian Mixture Models**”) algorithm, because the probabilistic model being learned describes a dataset as being drawn from a mixture distribution, where each component of this distribution is a Gaussian.

Just as in the K -means algorithm, this approach is susceptible to local optima and quality of initialization. The heuristics for computing better initializers for K -means are also useful here.

Aside from the fact that GMMs use soft assignments and K -means uses hard assignments, there are other differences between the two approaches. What are they?

16.3 The Expectation Maximization Framework

At this point, you’ve seen a method for learning in a particular probabilistic model with hidden variables. Two questions remain: (1) can

Algorithm 38 GMM(\mathbf{X}, K)

```

1: for  $k = 1$  to  $K$  do
2:    $\mu_k \leftarrow$  some random location           // randomly initialize mean for  $k$ th cluster
3:    $\sigma_k^2 \leftarrow 1$                        // initialize variances
4:    $\theta_k \leftarrow 1/K$                      // each cluster equally likely a priori
5: end for
6: repeat
7:   for  $n = 1$  to  $N$  do
8:     for  $k = 1$  to  $K$  do
9:        $z_{n,k} \leftarrow \theta_k [2\pi\sigma_k^2]^{-\frac{D}{2}} \exp \left[ -\frac{1}{2\sigma_k^2} \|x_n - \mu_k\|^2 \right]$  // compute
// (unnormalized) fractional assignments
10:    end for
11:     $z_n \leftarrow \frac{1}{\sum_k z_{n,k}} z_n$  // normalize fractional assignments
12:  end for
13:  for  $k = 1$  to  $K$  do
14:     $\theta_k \leftarrow \frac{1}{N} \sum_n z_{n,k}$  // re-estimate prior probability of cluster  $k$ 
15:     $\mu_k \leftarrow \frac{\sum_n z_{n,k} x_n}{\sum_n z_{n,k}}$  // re-estimate mean of cluster  $k$ 
16:     $\sigma_k^2 \leftarrow \frac{\sum_n z_{n,k} \|x_n - \mu_k\|^2}{\sum_n z_{n,k}}$  // re-estimate variance of cluster  $k$ 
17:  end for
18: until converged
19: return  $z$  // return cluster assignments

```

you apply this idea more generally and (2) why is it even a reasonable thing to do? Expectation maximization is a *family* of algorithms for performing maximum likelihood estimation in probabilistic models with hidden variables.

The general flavor of how we will proceed is as follows. We want to maximize the log likelihood \mathcal{L} , but this will turn out to be difficult to do directly. Instead, we'll pick a surrogate function $\tilde{\mathcal{L}}$ that's a lower bound on \mathcal{L} (i.e., $\tilde{\mathcal{L}} \leq \mathcal{L}$ everywhere) that's (hopefully) easier to maximize. We'll construct the surrogate in such a way that increasing it will force the true likelihood to also go up. After maximizing $\tilde{\mathcal{L}}$, we'll construct a *new* lower bound and optimize that. This process is shown pictorially in Figure 16.2.

To proceed, consider an arbitrary probabilistic model $p(x, y | \theta)$, where x denotes the observed data, y denotes the hidden data and θ denotes the parameters. In the case of Gaussian Mixture Models, x was the data points, y was the (unknown) labels and θ included the cluster prior probabilities, the cluster means and the cluster variances. Now, given access *only* to a number of examples x_1, \dots, x_N , you would like to estimate the parameters (θ) of the model.

Probabilistically, this means that some of the variables are unknown and therefore you need to marginalize (or sum) over their possible values. Now, your data consists only of $\mathbf{X} = \langle x_1, x_2, \dots, x_N \rangle$,

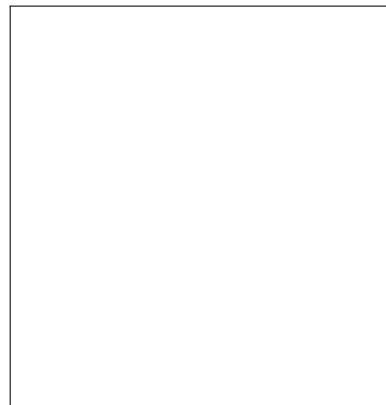


Figure 16.2: em:lowerbound: A figure showing successive lower bounds

not the (x, y) pairs in D . You can then write the likelihood as:

$$p(\mathbf{X} \mid \theta) = \sum_{y_1} \sum_{y_2} \cdots \sum_{y_N} p(\mathbf{X}, y_1, y_2, \dots, y_N \mid \theta) \quad \text{marginalization} \quad (16.23)$$

$$= \sum_{y_1} \sum_{y_2} \cdots \sum_{y_N} \prod_n p(x_n, y_n \mid \theta) \quad \text{examples are independent} \quad (16.24)$$

$$= \prod_n \sum_{y_n} p(x_n, y_n \mid \theta) \quad \text{algebra} \quad (16.25)$$

At this point, the natural thing to do is to take logs and then start taking gradients. However, once you start taking logs, you run into a problem: the log cannot eat the sum!

$$\mathcal{L}(\mathbf{X} \mid \theta) = \sum_n \log \sum_{y_n} p(x_n, y_n \mid \theta) \quad (16.26)$$

Namely, the log gets “stuck” outside the sum and cannot move in to decompose the rest of the likelihood term!

The next step is to apply the somewhat strange, but strangely useful, trick of multiplying by 1. In particular, let $q(\cdot)$ be an arbitrary probability distribution. We will multiply the $p(\dots)$ term above by $q(y_n)/q(y_n)$, a valid step so long as q is never zero. This leads to:

$$\mathcal{L}(\mathbf{X} \mid \theta) = \sum_n \log \sum_{y_n} q(y_n) \frac{p(x_n, y_n \mid \theta)}{q(y_n)} \quad (16.27)$$

We will now construct a lower bound using **Jensen’s inequality**.

This is a very useful (and easy to prove!) result that states that $f(\sum_i \lambda_i x_i) \geq \sum_i \lambda_i f(x_i)$, so long as (a) $\lambda_i \geq 0$ for all i , (b) $\sum_i \lambda_i = 1$, and (c) f is concave. If this looks familiar, that’s just because it’s a direct result of the definition of **concavity**. Recall that f is concave if $f(ax + by) \geq af(x) + bf(y)$ whenever $a + b = 1$.

You can now apply Jensen’s inequality to the log likelihood by identifying the list of $q(y_n)$ s as the λ s, log as f (which is, indeed, concave) and each “ x ” as the p/q term. This yields:

$$\mathcal{L}(\mathbf{X} \mid \theta) \geq \sum_n \sum_{y_n} q(y_n) \log \frac{p(x_n, y_n \mid \theta)}{q(y_n)} \quad (16.28)$$

$$= \sum_n \sum_{y_n} \left[q(y_n) \log p(x_n, y_n \mid \theta) - q(y_n) \log q(y_n) \right] \quad (16.29)$$

$$\triangleq \tilde{\mathcal{L}}(\mathbf{X} \mid \theta) \quad (16.30)$$

Note that this inequality holds for *any* choice of function q , so long as its non-negative and sums to one. In particular, it needn’t even be the

? Prove Jensen’s inequality using the definition of concavity and induction.

same function q for each n . We will need to take advantage of both of these properties.

We have succeeded in our first goal: constructing a lower bound on \mathcal{L} . When you go to optimize this lower bound for θ , the only part that matters is the first term. The second term, $q \log q$, drops out as a function of θ . This means that the the maximization you need to be able to compute, for fixed q_n s, is:

$$\theta^{(\text{new})} \leftarrow \arg \max_{\theta} \sum_n \sum_{y_n} q_n(y_n) \log p(x_n, y_n \mid \theta) \quad (16.31)$$

This is *exactly* the sort of maximization done for Gaussian mixture models when we recomputed new means, variances and cluster prior probabilities.

The second question is: what should $q_n(\cdot)$ actually be? Any reasonable q will lead to a lower bound, so in order to choose one q over another, we need another criterion. Recall that we are hoping to maximize \mathcal{L} by instead maximizing a lower bound. In order to ensure that an increase in the lower bound implies an increase in \mathcal{L} , we need to ensure that $\mathcal{L}(\mathbf{X} \mid \theta) = \tilde{\mathcal{L}}(\mathbf{X} \mid \theta)$. In words: $\tilde{\mathcal{L}}$ should be a lower bound on \mathcal{L} that makes contact at the current point, θ .

16.4 Further Reading

TODO further reading

17 | STRUCTURED PREDICTION

It is often the case that instead of predicting a *single* output, you need to predict multiple, correlated outputs simultaneously. In natural language processing, you might want to assign a syntactic label (like noun, verb, adjective, etc.) to words in a sentence: there is clear correlation among these labels. In computer vision, you might want to label regions in an image with object categories; again, there is correlation among these regions. The branch of machine learning that studies such questions is **structured prediction**.

In this chapter, we will cover two of the most common algorithms for structured prediction: the structured perceptron and the structured support vector machine. We will consider two types of structure. The first is the “sequence labeling” problem, typified by the natural language processing example above, but also common in computational biology (labeling amino acids in DNA) and robotics (labeling actions in a sequence). For this, we will develop specialized prediction algorithms that take advantage of the sequential nature of the task. We will also consider more general structures beyond sequences, and discuss how to cast them in a generic optimization framework: **integer linear programming** (or **ILP**).

The general framework we will explore is that of *jointly scoring input/output configurations*. We will construct algorithms that learn a function $s(x, \hat{y})$ (s for “score”), where x is an input (like an image) and \hat{y} is some predicted output (like a segmentation of that image). For any given image, there are a *lot* of possible segmentations (i.e., a lot of possible \hat{y} s), and the goal of s is to rank them in order of “how good” they are: how compatible they are with the input x . The most important thing is that the scoring function s ranks the *true* segmentation y higher than any other imposter segmentation \hat{y} . That is, we want to ensure that $s(x, y) > s(x, \hat{y})$ for all $\hat{y} \neq y$. The main challenge we will face is how to do this *efficiently*, given that there are so many imposter \hat{y} s.

Learning Objectives:

- Recognize when a problem should be solved using a structured prediction technique.
- Implement the structured perceptron algorithm for sequence labeling.
- Map “argmax” problems to integer linear programs.
- Augment the structured perceptron with losses to derive structured SVMs.

Dependencies:

17.1 Multiclass Perceptron

In order to build up to structured problems, let's begin with a simplified by pedagogically useful stepping stone: multiclass classification with a perceptron. As discussed earlier, in multiclass classification we have inputs $\mathbf{x} \in \mathbb{R}^D$ and output labels $y \in \{1, 2, \dots, K\}$. Our goal is to learn a scoring function s so that $s(\mathbf{x}, y) > s(\mathbf{x}, \hat{y})$ for all $\hat{y} \neq y$, where y is the true label and \hat{y} is an imposter label. The general form of scoring function we consider is a linear function of a joint feature vector $\phi(\mathbf{x}, y)$:

$$s(\mathbf{x}, y) = \mathbf{w} \cdot \phi(\mathbf{x}, y) \quad (17.1)$$

Here, the features $\phi(\mathbf{x}, y)$ should denote how “compatible” the input \mathbf{x} is with the label y . We keep track of a single weight vector \mathbf{w} that learns how to weigh these different “compatibility” features.

A natural way to represent ϕ , if you know nothing else about the problem, is an *outer product* between \mathbf{x} and the label space. This yields the following representation:

$$\phi(\mathbf{x}, k) = \left\langle \underbrace{0, 0, \dots, 0}_{D(k-1) \text{ zeros}}, \underbrace{\mathbf{x}}_{\in \mathbb{R}^D}, \underbrace{0, 0, \dots, 0}_{D(K-k) \text{ zeros}} \right\rangle \in \mathbb{R}^{DK} \quad (17.2)$$

In this representation, \mathbf{w} effectively encodes a separate weight for every feature/label pair.

How are we going to learn \mathbf{w} ? We will start with $\mathbf{w} = \mathbf{0}$ and then process each input one at a time. Suppose we get an input \mathbf{x} with gold standard label y . We will use the current scoring function to predict a label. In particular, we will predict the label \hat{y} that maximizes the score:

$$\hat{y} = \operatorname{argmax}_{\hat{y} \in [1, K]} s(\mathbf{x}, \hat{y}) \quad (17.3)$$

$$= \operatorname{argmax}_{\hat{y} \in [1, K]} \mathbf{w} \cdot \phi(\mathbf{x}, \hat{y}) \quad (17.4)$$

If this predicted output is correct (i.e., $\hat{y} = y$), then, per the normal perceptron, we will do nothing. Suppose that $\hat{y} \neq y$. This means that the score of \hat{y} is greater than the score of y , so we want to update \mathbf{w} so that the score of \hat{y} is decreased and the score of y is increased. We do this by:

$$\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}, y) - \phi(\mathbf{x}, \hat{y}) \quad (17.5)$$

To make sure this is doing what we expect, let's consider what would happen if we computed scores under the *updated* value of \mathbf{w} . To make the notation clear, let's say $\mathbf{w}^{(\text{old})}$ are the weights before update, and

Algorithm 39 MULTICLASSPERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```

1:  $\mathbf{w} \leftarrow \mathbf{0}$  // initialize weights
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(\mathbf{x}, y) \in \mathbf{D}$  do
4:      $\hat{y} \leftarrow \operatorname{argmax}_k \mathbf{w} \cdot \phi(\mathbf{x}, k)$  // compute prediction
5:     if  $\hat{y} \neq y$  then
6:        $\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}, y) - \phi(\mathbf{x}, \hat{y})$  // update weights
7:     end if
8:   end for
9: end for
10: return  $\mathbf{w}$  // return learned weights

```

$\mathbf{w}^{(new)}$ are the weights after update. Then:

$$\mathbf{w}^{(new)} \cdot \phi(\mathbf{x}, y) \quad (17.6)$$

$$= \left(\mathbf{w}^{(old)} + \phi(\mathbf{x}, y) - \phi(\mathbf{x}, \hat{y}) \right) \cdot \phi(\mathbf{x}, y) \quad (17.7)$$

$$= \underbrace{\mathbf{w}^{(old)} \cdot \phi(\mathbf{x}, y)}_{\text{old prediction}} + \underbrace{\phi(\mathbf{x}, y) \cdot \phi(\mathbf{x}, y)}_{\geq 0} - \underbrace{\phi(\mathbf{x}, \hat{y}) \cdot \phi(\mathbf{x}, y)}_{=0} \quad (17.8)$$

Here, the first term is the old prediction. The second term is of the form $\mathbf{a} \cdot \mathbf{a}$ which is non-negative (and, unless $\phi(\mathbf{x}, y)$ is the zero vector, positive). The third term is the dot product between ϕ for two different labels, which by definition of ϕ is zero (see Eq (17.2)).

This gives rise to the updated multiclass perceptron specified in Algorithm 17.1. As with the normal perceptron, the generalization of the multiclass perceptron increases dramatically if you do weight averaging.

An important note is that MulticlassPerceptronTrain is actually more powerful than suggested so far. For instance, suppose that you have three categories, but believe that two of them are tightly related, while the third is very different. For instance, the categories might be {music, movies, oncology}. You can encode this relatedness by defining a feature expansion ϕ that reflects this:

$$\phi(\mathbf{x}, \text{music}) = \langle \mathbf{x}, \mathbf{0}, \mathbf{0}, \mathbf{x} \rangle \quad (17.9)$$

$$\phi(\mathbf{x}, \text{movies}) = \langle \mathbf{0}, \mathbf{x}, \mathbf{0}, \mathbf{x} \rangle \quad (17.10)$$

$$\phi(\mathbf{x}, \text{oncology}) = \langle \mathbf{0}, \mathbf{0}, \mathbf{x}, \mathbf{0} \rangle \quad (17.11)$$

This encoding is identical to the normal encoding in the first three positions, but includes an extra copy of the features at the end, shared between music and movies. By doing so, if the perceptron wants to learn something common to music and movies, it can use this final shared position.



Verify the score of \hat{y} , $\mathbf{w}^{(new)} \cdot \phi(\mathbf{x}, \hat{y})$, decreases after an update, as we would want.



Suppose you have a *hierarchy* of classes arranged in a tree. How could you use that to construct a feature representation. You can think of the music/movies/oncology example as a binary tree: the left branch of the root splits into music and movies; the right branch of the root is just oncology.

17.2 Structured Perceptron

Let us now consider the sequence labeling task. In sequence labeling, the outputs are themselves variable-length vectors. An input/output pair (which must have the same length) might look like:

$$\mathbf{x} = \text{" monsters eat tasty bunnies "} \quad (17.12)$$

$$\mathbf{y} = \quad \text{noun verb adj noun} \quad (17.13)$$

To set terminology, we will refer to the *entire sequence* \mathbf{y} as the “output” and a single label within \mathbf{y} as a “label”. As before, our goal is to learn a scoring function that scores the true output sequence \mathbf{y} higher than any imposter output sequence.

As before, despite the fact that \mathbf{y} is now a vector, we can *still* define feature functions over the *entire* input/output pair. For instance, we might want to count the number of times “monsters” has been tagged as “noun” in a given output. Or the number of times “verb” is followed by “noun” in an output. Both of these are features that are likely indicative of a *correct* output. We might also count the number of times “tasty” has been tagged as a verb (probably a negative feature) and the number of times two verbs are adjacent (again, probably a negative feature).

More generally, a very standard set of features would be:

- the number of times word w has been labeled with tag l , for all words w and all syntactic tags l
- the number of times tag l is adjacent to tag l' in the output, for all tags l and l'

The first set of features are often called **unary features**, because they talk only about the relationship between the input (sentence) and a *single* (unit) label in the output sequence. The second set of features are often called **Markov features**, because they talk about adjacent labels in the output sequence, which is reminiscent of Markov models which only have short term memory.

Note that for a given input x of length L (in the example, $L = 4$), the number of possible outputs is K^L , where K is the number of syntactic tags. This means that the number of possible outputs *grows exponentially* in the length of the input. In general, we write $\mathcal{Y}(x)$ to mean “the set of all possible structured outputs for the input x ”. We have just seen that $|\mathcal{Y}(x)| = K^{\text{len}(x)}$.

Despite the fact that the inputs and outputs have variable length, the size of the *feature representation* is constant. If there are V words in your vocabulary and K labels for a given word, the the number of unary features is VK and the number of Markov features is K^2 , so

Algorithm 40 STRUCTUREDPERCEPTRONTRAIN(\mathbf{D} , $MaxIter$)

```

1:  $\mathbf{w} \leftarrow \mathbf{0}$  // initialize weights
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(\mathbf{x}, \mathbf{y}) \in \mathbf{D}$  do
4:      $\hat{\mathbf{y}} \leftarrow \operatorname{argmax}_{\hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x})} \mathbf{w} \cdot \phi(\mathbf{x}, \hat{\mathbf{y}})$  // compute prediction
5:     if  $\hat{\mathbf{y}} \neq \mathbf{y}$  then
6:        $\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x}, \hat{\mathbf{y}})$  // update weights
7:     end if
8:   end for
9: end for
10: return  $\mathbf{w}$  // return learned weights

```

the total number of features is $K(V + K)$. Of course, more complex feature representations are possible and, in general, are a good idea. For example, it is often useful to have unary features of neighboring words like “the number of times the word immediately preceding a verb was ‘monsters’.”

Now that we have a fixed size feature representation, we can develop a perceptron-style algorithm for sequence labeling. The core idea is the same as before. We will maintain a *single* weight vector \mathbf{w} . We will make predictions by choosing the (entire) output sequence $\hat{\mathbf{y}}$ that maximizes a score given by $\mathbf{w} \cdot \phi(\mathbf{x}, \hat{\mathbf{y}})$. And if this output sequence is incorrect, we will adjust the weights *word* the correct output sequence \mathbf{y} and away from the incorrect output sequence $\hat{\mathbf{y}}$. This is summarized in Algorithm 17.2

You may have noticed that Algorithm 17.2 for the structured perceptron is *identical* to Algorithm 17.1, aside from the fact that in the multiclass perceptron the argmax is over the K possible classes, while in the structured perceptron, the argmax is over the K^L possible output sequences!

The only difficulty in this algorithm is in line 4:

$$\hat{\mathbf{y}} \leftarrow \operatorname{argmax}_{\hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x})} \mathbf{w} \cdot \phi(\mathbf{x}, \hat{\mathbf{y}}) \quad (17.14)$$

In principle, this requires you to search over K^L possible output sequences $\hat{\mathbf{y}}$ to find the one that maximizes the dot product. Except for very small K or very small L , this is computationally infeasible. Because of its difficulty, this is often referred to as the **argmax problem** in structured prediction. Below, we consider how to solve the argmax problem for sequences.

17.3 Argmax for Sequences

We now face an *algorithmic* question, not a machine learning question: how to compute the argmax in Eq 17.14 efficiently. In general,

this is not possible. However, under somewhat restrictive assumptions about the form of our features ϕ , we can solve this problem efficiently, by casting it as the problem of computing a maximum weight path through a specifically constructed lattice. This is a variant of the Viterbi algorithm for hidden Markov models, a classic example of dynamic programming. (Later, in Section 17.6, we will consider argmax for more general problems.)

The key observation for sequences is that—so long as we restrict our attention to unary features and Markov features—the feature function ϕ decomposes over the input. This is easiest to see with an example. Consider the input/output sequence from before: $x =$ “monsters eat tasty bunnies” and $y = [\text{noun verb adj noun}]$. If we want to compute the number of times “bunnies” is tagged as “noun” in this pair, we can do this by:

1. count the number of times “bunnies” is tagged as “noun” in the first three words of the sentence
2. add to that the number of times “bunnies” is tagged as “noun” in the final word

We can do a similar exercise for Markov features, like the number of times “adj” is followed by “noun”.

However, we don’t actually *need* these counts. All we need for computing the argmax sequence is the dot product between the weights w and these counts. In particular, we can compute $w \cdot \phi(x, y)$ as the dot product on all-but-the-last word plus the dot product on the last word: $w \cdot \phi_{1:3}(x, y) + w \cdot \phi_4(x, y)$. Here, $\phi_{1:3}$ means “features for everything up to and including position 3” and ϕ_4 means “features for position 4.”

More generally, we can write $\phi(x, y) = \sum_{l=1}^L \phi_l(x, y)$, where $\phi_l(x, y)$ only includes features about position l .¹ In particular, we’re taking advantage of the associative law for addition:

$$w \cdot \phi(x, y) = w \cdot \sum_{l=1}^L \phi_l(x, y) \quad \text{decomposition of structure} \quad (17.15)$$

$$= \sum_{l=1}^L w \cdot \phi_l(x, y) \quad \text{associative law} \quad (17.16)$$

What this means is that we can build a graph like that in Figure ??, with one verticle slice per time step ($l = 1 \dots L$).² Each *edge* in this graph will receive a weight, constructed in such a way that if you take a complete path through the lattice, and add up all the weights, this will correspond exactly to $w \cdot \phi(x, y)$.

To complete the construction, let $\phi_l(x, \dots \circ y \circ y')$ denote the *unary* features at position l together with the *Markov* features that end at

¹ In the case of Markov features, we think of them as pairs that *end* at position l , so “verb adj” would be the active feature for ϕ_3 .

² A graph of this sort is called a **trellis**, and sometimes a **lattice** in the literature.

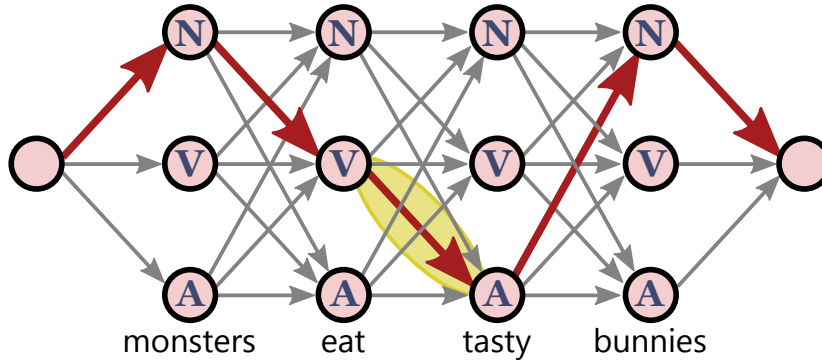


Figure 17.1: A picture of a trellis sequence labeling. At each time step l the corresponding word can have any of the three possible labels. Any path through this trellis corresponds to a unique labeling of this sentence. The gold standard path is drawn with bold red arrows. The highlighted edge corresponds to the edge between $l = 2$ and $l = 3$ for verb/adj as described in the text. That edge has weight $w \cdot \phi_3(x, \dots \circ \text{verb} \circ \text{adj})$.

position l . These features depend *only* on x , y and y' , and *not* any of the previous parts of the output.

For example, in the running example “monsters/noun eat/verb tasty/adj bunnies/noun”, consider the edge between $l = 2$ and $l = 3$ going from “verb” to “adj”. (Note: this is a “correct” edge, in the sense that it belongs to the ground truth output.) The features associated with this edge will be unary features about “tasty/adj” as well as Markov features about “verb/adj”. The *weight* of this edge will be exactly the total score (according to w) of those features.

Formally, consider an edge in the trellis that goes from time $l - 1$ to l , and transitions from y to y' . Set the weight of this edge to exactly $w \cdot \phi_l(x, \dots \circ y \circ y')$. By doing so, we guarantee that the sum of weights along any path through this lattice is exactly equal to the score of that path. Once we have constructed the graph as such, we can run any max-weight path algorithm to compute the highest scoring output. For trellises, this can be computed by the Viterbi algorithm, or by applying any of a number of path finding algorithms for more general graphs. A complete derivation of the dynamic program in this case is given in Section 17.7 for those who want to implement it directly.

The main benefit of this construction is that it is guaranteed to exactly compute the argmax output for sequences required in the structured perceptron algorithm, *efficiently*. In particular, it’s runtime is $O(LK^2)$, which is an exponential improvement on the naive $O(K^L)$ runtime if one were to enumerate every possible output sequence. The algorithm can be naturally extended to handle “higher order” Markov assumptions, where features depend on triples or quadruples of the output. The trellis becomes larger, but the algorithm remains essentially the same. In order to handle a length M Markov features, the resulting algorithm will take $O(LK^M)$ time. In practice, it’s rare that $M > 3$ is necessary or useful.

17.4 Structured Support Vector Machines

In Section 7.7 we saw the support vector machine as a very useful general framework for binary classification. In this section, we will develop a related framework for structured support vector machines. The two main advantages of structured SVMs over the structured perceptron are (1) it is regularized (though averaging in structured perceptron achieves a similar effect) and (2) we can incorporate more complex loss functions.

In particular, one suboptimal thing about the structured perceptron is that all errors are considered equally bad. For structured problems, we often have much more nuanced and elaborate loss functions that we want to optimize. Even for sequence labeling, it is typically far worse to label every word incorrectly than to just label one word incorrectly. It is very common to use **Hamming loss** as a general loss function for structured prediction. Hamming loss simply counts: of all the predictions you made, how many were incorrect? For sequence labeling, it is:

$$\ell^{(\text{Ham})}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{l=1}^L \mathbf{1}[y_l \neq \hat{y}_l] \quad (17.17)$$

In order to build up to structured SVMs, recall that SVMs began with the following optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{large margin}} + C \underbrace{\sum_n \xi_n}_{\text{small slack}} \quad (17.18) \\ \text{subj. to} \quad & y_n (\mathbf{w} \cdot \mathbf{x}_n + b) \geq 1 - \xi_n \quad (\forall n) \\ & \xi_n \geq 0 \quad (\forall n) \end{aligned}$$

After a bit of work, we were able to reformulate this in terms of a standard loss optimization algorithm with hinge loss:

$$\min_{\mathbf{w}} \quad \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{large margin}} + C \underbrace{\sum_n \ell^{(\text{hin})}(y_n, \mathbf{w} \cdot \mathbf{x}_n + b)}_{\text{small slack}} \quad (17.19)$$

We can do a similar derivation in the structured case. The question is: exactly what should the slack be measuring? Our *goal* is for the score of the true output \mathbf{y} to beat the score of any imposter output $\hat{\mathbf{y}}$. To incorporate loss, we will say that we want the score of the true output to beat the score of any imposter output by *at least* the loss that would be suffered if we were to predict that imposter output. An alternative view is the ranking view: we want the true output to be ranked above any imposter by an amount at least equal to the loss.

To keep notation simple, we will write $s_w(\mathbf{x}, \mathbf{y})$ to denote the score of the pair \mathbf{x}, \mathbf{y} , namely $\mathbf{w} \cdot \phi(\mathbf{x}, \mathbf{y})$. This suggests a set of constraints of the form:

$$s_w(\mathbf{x}, \mathbf{y}) - s_w(\mathbf{x}, \hat{\mathbf{y}}) \geq \ell^{(\text{Ham})}(\mathbf{y}, \hat{\mathbf{y}}) - \xi_{\hat{\mathbf{y}}} \quad (\forall n, \forall \hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x})) \quad (17.20)$$

The rest of the optimization problem remains the same, yielding:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_n \sum_{\hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)} \xi_{n, \hat{\mathbf{y}}} \\ \text{subj. to} \quad & s_w(\mathbf{x}, \mathbf{y}) - s_w(\mathbf{x}, \hat{\mathbf{y}}) \\ & \geq \ell^{(\text{Ham})}(\mathbf{y}_n, \hat{\mathbf{y}}) - \xi_{n, \hat{\mathbf{y}}} \quad (\forall n, \forall \hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)) \\ & \xi_{n, \hat{\mathbf{y}}} \geq 0 \quad (\forall n, \forall \hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)) \end{aligned} \quad (17.21)$$

This optimization problem asks for a large margin and small slack, where there is a slack very for every training example and every possible incorrect output associated with that training example. In general, this is *way too many* slack variables and *way too many* constraints!

There is a very useful, general trick we can apply. If you focus on the first constraint, it roughly says (letting $s(\cdot)$ denote score): $s(\mathbf{y}) \geq [s(\hat{\mathbf{y}}) + \ell(\mathbf{y}, \hat{\mathbf{y}})]$ for all $\hat{\mathbf{y}}$, modulo slack. We'll refer to the thing in brackets as the "loss-augmented score." But if we want to guarantee that the score of the true \mathbf{y} beats the loss-augmented score of *all* $\hat{\mathbf{y}}$, it's enough to ensure that it beats the loss-augmented score of the most confusing imposter. Namely, it is sufficient to require that $s(\mathbf{y}) \geq \max_{\hat{\mathbf{y}}} [s(\hat{\mathbf{y}}) + \ell(\mathbf{y}, \hat{\mathbf{y}})]$, modulo slack. Expanding out the definition of $s(\cdot)$ and adding slack back in, we can replace the exponentially large number of constraints in Eq (17.21) with the simpler set of constraints:

$$s_w(\mathbf{x}_n, \mathbf{y}_n) \geq \max_{\hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)} [s_w(\mathbf{x}_n, \hat{\mathbf{y}}) + \ell^{(\text{Ham})}(\mathbf{y}_n, \hat{\mathbf{y}})] - \xi_n \quad (\forall n)$$

We can now apply the same trick as before to remove ξ_n from the analysis. In particular, because ξ_n is constrained to be ≥ 0 and because we are trying to minimize it's sum, we can figure out that out the optimum, it will be the case that:

$$\xi_n = \max \left\{ 0, \max_{\hat{\mathbf{y}} \in \mathcal{Y}(\mathbf{x}_n)} [s_w(\mathbf{x}_n, \hat{\mathbf{y}}) + \ell^{(\text{Ham})}(\mathbf{y}_n, \hat{\mathbf{y}})] - s_w(\mathbf{x}_n, \mathbf{y}_n) \right\} \quad (17.22)$$

$$= \ell^{(\text{s-h})}(\mathbf{y}_n, \mathbf{x}_n, \mathbf{w}) \quad (17.23)$$

This value is referred to as the **structured hinge loss**, which we have denoted as $\ell^{(\text{s-h})}(\mathbf{y}_n, \mathbf{x}_n, \mathbf{w})$. This is because, although it is more complex, it bears a striking resemblance to the **hinge loss** from Chapter 7.

In particular, if the score of the true output beats the score of every the best imposter by at least its loss, then ζ_n will be zero. On the other hand, if some imposter (plus its loss) beats the true output, the loss scales linearly as a function of the difference. At this point, there is nothing special about Hamming loss, so we will replace it with some arbitrary structured loss ℓ .

Plugging this back into the objective function of Eq (17.21), we can write the structured SVM as an *unconstrained* optimization problem, akin to Eq (17.19), as:

$$\min_w \frac{1}{2} \|w\|^2 + C \sum_n \ell^{(s-h)}(y_n, x_n, w) \quad (17.24)$$

This is now in a form that we can optimize using subgradient descent (Chapter 7) or stochastic subgradient descent (Chapter 14).

In order to compute subgradients of Eq (17.24), we need to be able to compute subgradients of the structured hinge loss. Mathematically this is straightforward. If the structured hinge loss on an example (x, y) is zero, then the gradient with respect to w is also zero. If the structured hinge loss is positive, then the gradient is:

$$\nabla_w \ell^{(s-h)}(y, x, w) \quad \text{if the loss is } > 0 \quad (17.25)$$

expand definition using arbitrary structured loss ℓ

$$= \nabla_w \left\{ \max_{\hat{y} \in \mathcal{Y}(x_n)} [w \cdot \phi(x_n, \hat{y}) + \ell(y_n, \hat{y})] - w \cdot \phi(x_n, y_n) \right\} \quad (17.26)$$

define \hat{y}_n to be the output that attains the maximum above, rearrange

$$= \nabla_w \left\{ w \cdot \phi(x_n, \hat{y}) - w \cdot \phi(x_n, y_n) + \ell(y_n, \hat{y}) \right\} \quad (17.27)$$

take gradient

$$= \phi(x_n, \hat{y}) - \phi(x_n, y_n) \quad (17.28)$$

Putting this together, we get the full gradient as:

$$\nabla_w \ell^{(s-h)}(y_n, x_n, w) = \begin{cases} 0 & \text{if } \ell^{(s-h)}(y_n, x_n, w) = 0 \\ \phi(x_n, \hat{y}_n) - \phi(x_n, y_n) & \text{otherwise} \end{cases}$$

$$\text{where } \hat{y}_n = \operatorname{argmax}_{\hat{y}_n \in \mathcal{Y}(x_n)} [w \cdot \phi(x_n, \hat{y}_n) + \ell(y_n, \hat{y}_n)] \quad (17.29)$$

The form of this gradient is very simple: it is equal to the features of the worst imposter minus the features of the truth, unless the truth beats all imposters, in which case it's zero. When plugged into stochastic subgradient descent, you end up with an update that looks very much like the structured perceptron: if the current prediction (\hat{y}_n) is correct, there is no gradient step. But if the current prediction is incorrect, you step w toward the truth and away from the imposter.

Algorithm 41 `STOCHSUBGRADSTRUCTSVM(D, MaxIter, λ , ℓ)`

```

1:  $w \leftarrow 0$  // initialize weights
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(x, y) \in D$  do
4:      $\hat{y} \leftarrow \operatorname{argmax}_{\hat{y} \in \mathcal{Y}(x)} w \cdot \phi(x, \hat{y}) + \ell(y, \hat{y})$  // loss-augmented prediction
5:     if  $\hat{y} \neq y$  then
6:        $w \leftarrow w + \phi(x, y) - \phi(x, \hat{y})$  // update weights
7:     end if
8:      $w \leftarrow w - \frac{\lambda}{N} w$  // shrink weights due to regularizer
9:   end for
10: end for
11: return  $w$  // return learned weights

```

We will consider how to compute the loss-augmented argmax in the next section, but before that we summarize an algorithm for optimizing structured SVMs using stochastic subgradient descent: Algorithm 17.4. Of course there are other possible optimization strategies; we are highlighting this one because it is nearly identical to the structured perceptron. The only differences are: (1) on line 4 you use loss-augmented argmax instead of argmax; and (2) on line 8 the weights are shrunk slightly corresponding to the ℓ_2 regularizer on w . (Note: we have used $\lambda = 1/(2C)$ to make the connection to linear models clearer.)

17.5 Loss-Augmented Argmax

The challenge that arises is that we now have a more complicated argmax problem than before. In structured perceptron, we only needed to compute \hat{y}_n as the output that maximized its score (see Eq 17.14). Here, we need to find the output that maximizes its score *plus* its loss (Eq (17.29)). This optimization problem is referred to as **loss-augmented search** or **loss-augmented inference**.

Before solving the loss-augmented inference problem, it's worth thinking about why it makes sense. What is \hat{y}_n ? It's the output that has the highest score among all outputs, *after* adding the output's corresponding loss to that score. In other words, every incorrect output gets an artificial boost to its score, equal to its loss. The loss is serving to make imposters look *even better* than they really are, so if the truth is to beat an imposter, it has to beat it by a *lot*. In fact, this loss augmentation is essentially playing the role of a margin, where the required margin scales according to the loss.

The algorithmic question, then, is how to compute \hat{y}_n . In the fully general case, this is at least as hard as the normal argmax problem, so we cannot expect a general solution. Moreover, even in cases where the argmax problem is easy (like for sequences), the loss-augmented

argmax problem can still be difficult. In order to make it easier, we need to assume that the loss *decomposes* of the input in a way that's consistent with the features. In particular, if the structured loss function is Hamming loss, this is often straightforward.

As a concrete example, let's consider loss-augmented argmax for sequences under Hamming loss. In comparison to the trellis problem solved in Section 17.7, the only difference is that we want to *reward* paths that go through incorrect nodes in the trellis! In particular, in Figure 17.1, all of the edges that are not part of the gold standard path—those that are thinner and grey—get a free “+1” added to their weights. Since Hamming loss adds one to the score for any word that's predicted incorrectly, this means that every edge in the trellis that leads to an *incorrect* node (i.e., one that does not match the gold truth label) gets a “+1” added to its weight.

Again, consider an edge in the trellis that goes from time $l - 1$ to l , and transitions from y to y' . In the non-loss-augmented, the weight of this edge was exactly $w \cdot \phi_l(x, \dots \circ y \circ y')$. In the loss-augmented cases, the weight of this edge becomes:

$$\underbrace{w \cdot \phi_l(x, \dots \circ y \circ y')}_{\text{edge score, as before}} + \underbrace{\mathbf{1}[y' \neq y_l]}_{+1 \text{ for mispredictions}} \quad (17.30)$$

Once this loss-augmented graph has been constructed, the same max-weight path algorithm can be run to find the loss-augmented argmax sequence.

17.6 Argmax in General

The general argmax problem for structured perceptron is the algorithmic question of whether the following can be efficiently computed:

$$\hat{y} \leftarrow \operatorname{argmax}_{\hat{y} \in \mathcal{Y}(x)} w \cdot \phi(x, \hat{y}) \quad (17.31)$$

We have seen that *if* the output space $\mathcal{Y}(x)$ is sequences *and* the only types of features are unary features and Markov features, then this can be computed efficiently. There are a small number of other structured output spaces and feature restrictions for which efficient problem-specific algorithms exist:

- Binary trees, with context-free features: use the CKY algorithm
- 2d image segmentation, with adjacent-pixel features: use a form of graph cuts
- Spanning trees, with edge-based features: use Kruskal's algorithm (or for directed spanning trees, use Chu-Liu/Edmonds algorithm)

These special cases are often very useful, and many problems can be cast in one of these frameworks. However, it is often the case that you need a more general solution.

One of the most generally useful solutions is to cast the argmax problem as an **integer linear program**, or **ILP**. ILPs are a specific type of mathematical program/optimization problem, in which the objective function being optimized is linear and the constraints are linear. However, unlike “normal” linear programs, in an ILP you are allowed to have integer constraints and disallow fractional values. The general form of an ILP is, for a fixed vector a :

$$\max_z \quad a \cdot z \quad \text{subj. to} \quad \text{linear constraints on } z \quad (17.32)$$

The main point is that the constraints on z are allowed to include constraints like $z_3 \in \{0, 1\}$, which is considered an integer constraint.

Being able to cast your argmax problem as an ILP has the advantage that there are very good, efficiently, well-engineered ILP solvers out there in the world.³ ILPs are not a panacea though: in the worst case, the ILP solver will be horribly inefficient. But for prototyping, or if there are no better options, it’s a very handy technique.

³ I like Gurobi best, and it’s free for academic use. It also has a really nice Python interface.

Figuring out how exactly to cast your argmax problem as an ILP can be a bit challenging. Let’s start with an example of encoding sequence labeling with Markov features as an ILP. We first need to decide what the variables will be. Because we need to encode pairwise features, we will let our variables be of the form:

$$z_{l,k',k} = \mathbf{1}[\text{label } l \text{ is } k \text{ and label } l-1 \text{ is } k'] \quad (17.33)$$

These z s will all be binary indicator variables.

Our next task is to construct the linear objective function. To do so, we need to assign a value to $a_{l,k',k}$ in such a way that $a \cdot z$ will be exactly equal to $w \cdot \phi(x, y(z))$, where $y(z)$ denotes the sequence that we can read off of the variables z . With a little thought, we arrive at:

$$a_{l,k',k} = w \cdot \phi_l(x, \langle \dots, k', k \rangle) \quad (17.34)$$

Finally, we need to construct constraints. There are a few things that these constraints need to enforce:

1. That all the z s are binary. That’s easy: just say $z_{l,k',k} \in \{0, 1\}$, for all l, k', k .
2. That for a given position l , there is exactly one active z . We can do this with an equality constraint: $\sum_k \sum_{k'} z_{l,k',k} = 1$ for all l .
3. That the z s are internally consistent: if the label at position 5 is supposed to be “noun” then both $z_{5,..}$ and $z_{6,..}$ need to agree on

this. We can do this as: $\sum_{k'} z_{l,k',k} = \sum_{k''} z_{l+1,k,k''}$ for all l, k . Effectively what this is saying is that $z_{5,?,\text{verb}} = z_{6,\text{verb},?}$ where the “?” means “sum over all possibilities.”

This fully specifies an ILP that you can relatively easily implement (arguably more easily than the dynamic program in Algorithm 17.7) and which will solve the argmax problem for you. Will it be efficient? In this case, probably yes. Will it be as efficient as the dynamic program? Probably not.

It takes a bit of effort and time to get used to casting optimization problems as ILPs, and certainly not all can be, but most can and it’s a very nice alternative.

In the case of loss-augmented search for structured SVMs (as opposed to structured perceptron), the objective function of the ILP will need to be modified to include terms corresponding to the loss.

17.7 Dynamic Programming for Sequences

Recall the decomposition we derived earlier:

$$w \cdot \phi(x, y) = w \cdot \sum_{l=1}^L \phi_l(x, y) \quad \text{decomposition of structure} \quad (17.35)$$

$$= \sum_{l=1}^L w \cdot \phi_l(x, y) \quad \text{associative law} \quad (17.36)$$

This decomposition allows us to construct the following dynamic program. We will compute $\alpha_{l,k}$ as the score of the *best possible* output prefix up to and including position l that labels the l th word with label k . More formally:

$$\alpha_{l,k} = \max_{\hat{y}_{1:l-1}} w \cdot \phi_{1:l}(x, \hat{y} \circ k) \quad (17.37)$$

Here, \hat{y} is a sequence of length $l - 1$, and $\hat{y} \circ k$ denotes the sequence of length l obtained by adding k onto the end. The max denotes the fact that we are seeking the *best possible* prefix up to position $l - 1$, and the forcing the label for position l to be k .

Before working through the details, let’s consider an example. Suppose that we’ve computing the α s up to $l = 2$, and have: $\alpha_{2,\text{noun}} = 2$, $\alpha_{2,\text{verb}} = 9$, $\alpha_{2,\text{adj}} = -1$ (recall: position $l = 2$ is “eat”). We want to extend this to position 3; for example, we want to compute $\alpha_{3,\text{adj}}$. Let’s assume there’s a single unary feature here, “tasty/adj” and three possible Markov features of the form “?:adj”. Assume these weights are as given to the right.⁴ Now, the question for $\alpha_{3,\text{adj}}$ is: what’s the score of the best prefix that labels “tasty” as “adj”? We can obtain this by taking the best prefix up to “eat” and then appending

⁴ $w_{\text{“tasty/adj”}} = 1.2$
 $w_{\text{“noun:adj”}} = -5$
 $w_{\text{“verb:adj”}} = 2.5$
 $w_{\text{“adj:adj”}} = 2.2$

each possible label. Whichever combination is best is the winner. The relevant computation is:

$$\alpha_{3,\text{adj}} = \max \left\{ \begin{aligned} &\alpha_{2,\text{noun}} + w_{\text{tasty/adj}} + w_{\text{noun:adj}} \\ &\alpha_{2,\text{verb}} + w_{\text{tasty/adj}} + w_{\text{verb:adj}} \\ &\alpha_{2,\text{adj}} + w_{\text{tasty/adj}} + w_{\text{adj:adj}} \end{aligned} \right\} \quad (17.38)$$

$$= \max \left\{ 2 + 1.2 - 5, \quad 9 + 1.2 + 2.5, \quad -1 + 1.2 + 2.2 \right\} \quad (17.39)$$

$$= \max \left\{ -1.8, \quad 12.7, \quad 2.4 \right\} = 12.7 \quad (17.40)$$

This means that (a) the score for the prefix ending at position 3 labeled as adjective is 12.7, and (b) the “winning” previous label was “verb”. We will need to record these winning previous labels so that we can extract the best path at the end. Let’s denote by $\zeta_{l,k}$ the label at position $l-1$ that achieves the max.

From here, we can formally compute the α s recursively. The main observation that will be necessary is that, because we have limited ourselves to Markov features, $\phi_{l+1}(\mathbf{x}, \langle y_1, y_2, \dots, y_l, y_{l+1} \rangle)$ depends only on the last two terms of y , and does *not* depend on y_1, y_2, \dots, y_{l-1} . The full recursion is derived as:

$$\alpha_{0,k} = 0 \quad \forall k \quad (17.41)$$

$$\zeta_{0,k} = \emptyset \quad \forall k \quad (17.42)$$

the score for any empty sequence is zero

$$\alpha_{l+1,k} = \max_{\hat{y}_{1:l}} w \cdot \phi_{l+1}(\mathbf{x}, \hat{y} \circ k) \quad (17.43)$$

separate score of prefix from score of position l+1

$$= \max_{\hat{y}_{1:l}} w \cdot \left(\phi_{1:l}(\mathbf{x}, \hat{y}) + \phi_{l+1}(\mathbf{x}, \hat{y} \circ k) \right) \quad (17.44)$$

distributive law over dot products

$$= \max_{\hat{y}_{1:l}} \left[w \cdot \phi_{1:l}(\mathbf{x}, \hat{y}) + w \cdot \phi_{l+1}(\mathbf{x}, \hat{y} \circ k) \right] \quad (17.45)$$

separate out final label from prefix, call it k'

$$= \max_{\hat{y}_{1:l-1}} \max_{k'} \left[w \cdot \phi_{1:l}(\mathbf{x}, \hat{y} \circ k') + w \cdot \phi_{l+1}(\mathbf{x}, \hat{y} \circ k' \circ k) \right] \quad (17.46)$$

swap order of maxes, and last term doesn't depend on prefix

$$= \max_{k'} \left[\left[\max_{\hat{y}_{1:l-1}} w \cdot \phi_{1:l}(\mathbf{x}, \hat{y} \circ k') \right] + w \cdot \phi_{l+1}(\mathbf{x}, \langle \dots, k', k \rangle) \right] \quad (17.47)$$

apply recursive definition

$$= \max_{k'} \left[\alpha_{l,k'} + w \cdot \phi_{l+1}(\mathbf{x}, \langle \dots, k', k \rangle) \right] \quad (17.48)$$

Algorithm 42 ARGMAXFORSEQUENCES(x, w)

```

1:  $L \leftarrow \text{LEN}(x)$ 
2:  $\alpha_{l,k} \leftarrow 0, \quad \zeta_{k,l} \leftarrow 0, \quad \forall k = 1 \dots K, \quad \forall l = 0 \dots L$  // initialize variables
3: for  $l = 0 \dots L-1$  do
4:   for  $k = 1 \dots K$  do
5:      $\alpha_{l+1,k} \leftarrow \max_{k'} [\alpha_{l,k'} + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle)]$  // recursion:
        // here,  $\phi_{l+1}(\dots k', k \dots)$  is the set of features associated with
        // output position  $l+1$  and two adjacent labels  $k'$  and  $k$  at that position
6:      $\zeta_{l+1,k} \leftarrow$  the  $k'$  that achieves the maximum above // store backpointer
7:   end for
8: end for
9:  $y \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize predicted output to L-many zeros
10:  $y_L \leftarrow \text{argmax}_k \alpha_{L,k}$  // extract highest scoring final label
11: for  $l = L-1 \dots 1$  do
12:    $y_l \leftarrow \zeta_{l, y_{l+1}}$  // traceback  $\zeta$  based on  $y_{l+1}$ 
13: end for
14: return  $y$  // return predicted output

```

and record a backpointer to the k' that achieves the max

$$\zeta_{l+1,k} = \text{argmax}_{k'} [\alpha_{l,k'} + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle)] \quad (17.49)$$

At the end, we can take $\max_k \alpha_{L,k}$ as the score of the best output sequence. To extract the final sequence, we know that the best label for the last word is $\text{argmax}_k \alpha_{L,k}$. Let's call this \hat{y}_L . Once we know that, the best *previous* label is ζ_{L-1, \hat{y}_L} . We can then follow a path through ζ back to the beginning. Putting this all together gives Algorithm 17.7.

The main benefit of Algorithm 17.7 is that it is guaranteed to exactly compute the argmax output for sequences required in the structured perceptron algorithm, *efficiently*. In particular, its runtime is $O(LK^2)$, which is an exponential improvement on the naive $O(K^L)$ runtime if one were to enumerate every possible output sequence. The algorithm can be naturally extended to handle “higher order” Markov assumptions, where features depend on triples or quadruples of the output. The memoization becomes notationally cumbersome, but the algorithm remains essentially the same. In order to handle a length M Markov features, the resulting algorithm will take $O(LK^M)$ time. In practice, it's rare that $M > 3$ is necessary or useful.

In the case of loss-augmented search for structured SVMs (as opposed to structured perceptron), we need to include the scores coming from the loss augmentation in the dynamic program. The only thing that changes between the standard argmax solution (Algorithm 17.7, and derivation in Eq (17.48)) is that the any time an incorrect label is used, the (loss-augmented) score increases by one. Recall that in the non-loss-augmented case, we have the α recursion

as:

$$\alpha_{l+1,k} = \max_{\hat{\mathbf{y}}_{1:l}} \mathbf{w} \cdot \phi_{1:l+1}(\mathbf{x}, \hat{\mathbf{y}} \circ k) \quad (17.50)$$

$$= \max_{k'} \left[\alpha_{l,k'} + \mathbf{w} \cdot \phi_{l+1}(\mathbf{x}, \langle \dots, k', k \rangle) \right] \quad (17.51)$$

If we define $\tilde{\alpha}$ to be the loss-augmented score, the corresponding recursion is (differences highlighted in blue):

$$\tilde{\alpha}_{l+1,k} = \max_{\hat{\mathbf{y}}_{1:l}} \mathbf{w} \cdot \phi_{1:l+1}(\mathbf{x}, \hat{\mathbf{y}} \circ k) + \ell_{1:l+1}^{(\text{Ham})}(\mathbf{y}, \hat{\mathbf{y}} \circ k) \quad (17.52)$$

$$= \max_{k'} \left[\tilde{\alpha}_{l,k'} + \mathbf{w} \cdot \phi_{l+1}(\mathbf{x}, \langle \dots, k', k \rangle) \right] + \mathbf{1}[k \neq \mathbf{y}_{l+1}] \quad (17.53)$$

In other words, when computing $\tilde{\alpha}$ in the loss-augmented case, whenever the output prediction is forced to pass through an incorrect label, the score for that cell in the dynamic program gets increased by one. The resulting algorithm is identical to Algorithm 17.7, except that Eq (17.53) is used for computing α s.

17.8 Further Reading

TODO

18 | IMITATION LEARNING

Programming is a skill best acquired by practice and example rather than from books.

— Alan Turing

SO FAR, WE HAVE LARGELY CONSIDERED MACHINE LEARNING PROBLEMS in which the goal of the learning algorithm is to make a *single* prediction. In many real world problems, however, an algorithm must make a *sequence* of decisions, with the world possibly changing during that sequence. Such problems are often called **sequential decision making** problems. A straightforward example—which will be the running example for this chapter—is that of self-driving cars. We want to train a machine learning algorithm to drive a car. But driving a car is not a single prediction: it's a sequence of predictions over time. And as the machine is driving the car, the world around it is changing, often based on its own behavior. This creates complicated feedback loops, and one of the major challenges we will face is how to deal with these feedback loops.

To make this discussion more concrete, let's consider the case of a self-driving car. And let's consider a very simplistic car, in which the only decision that has to be made is how to steer, and that's between one of three options: {left, right, none}. In the imitation learning setting, we assume that we have access to an **expert** or **oracle** that already knows how to drive. We want to watch the expert driving, and learn to imitate their behavior. Hence: **imitation learning** (sometimes called **learning by demonstration** or **programming by example**, in the sense that programs are learned, and not implemented).

At each point in time $t = 1 \dots T$, the car receives sensor information x_t (for instance, a camera photo ahead of the car, or radar readings). It then has to take an action, a_t ; in the case of the car, this is one of the three available steering actions. The car then suffers some loss ℓ_t ; this might be zero in the case that it's driving well, or large in the case that it crashes. The world then changes, moves to time step $t + 1$, sensor readings x_{t+1} are observed, action a_{t+1} is taken, loss ℓ_{t+1} is suffered, and the process continues.

The goal is to learn a function f that maps from sensor readings x_t to actions. Because of close connections to the field of **reinforcement learning**, this function is typically called a **policy**. The measure of

Learning Objectives:

- Be able to formulate imitation learning problems.
- Understand the failure cases of simple classification approaches to imitation learning.
- Implement solutions to those problems based on either classification or dataset aggregation.
- Relate structured prediction and imitation learning.

Dependencies:

success of a policy is: if we were to run this policy, how much total loss would be suffered. In particular, suppose that the **trajectory** (denoted τ) of observation/action/reward triples encountered by your policy is:

$$\tau = x_1, \underbrace{a_1}_{=f(x_1)}, \ell_1, x_2, \underbrace{a_2}_{=f(x_2)}, \ell_2, \dots, x_T, \underbrace{a_T}_{=f(x_T)}, \ell_T \quad (18.1)$$

The losses ℓ_t depend implicitly on the state of the world and the actions of the policy. The goal of f is to minimize the expected loss $\mathbb{E}_{\tau \sim f} \left[\sum_{t=1}^T \ell_t \right]$, where the expectation is taken over all randomness in the world, and the sequence of actions taken is according to f .¹

18.1 Imitation Learning by Classification

We will begin with a straightforward, but brittle, approach to imitation learning. We assume access to a set of *training trajectories* taken by an expert. For example, consider a self-driving car, like that in Figure 18.1. A single trajectory τ consists of a sequence of observations (what is seen from the car's sensors) and a sequence of actions (what action did the expert take at that point in time). The idea in imitation learning by classification is to learn a classifier that attempts to mimic the expert's action based on the observations at that time.

In particular, we have $\tau_1, \tau_2, \dots, \tau_N$. Each of the N trajectories comprises a sequence of T -many observation/action/loss triples, where the action is the action taken by the expert. T , the length of the trajectory is typically called the **time horizon** (or just **horizon**). For instance, we may ask an expert human driver to drive $N = 20$ different routes, and record the observations and actions that driver saw and took during those routes. These are our training trajectories. We assume for simplicity that each of these trajectories is of fixed length T , though this is mostly for notational convenience.

The most straightforward thing we can do is convert this expert data into a big multiclass classification problem. We take our favorite multiclass classification algorithm, and use it to learn a mapping from x to a . The data on which it is trained is the set of all observation/action pairs visited during any of the N trajectories. In total, this would be NT examples. This approach is summarized in Algorithm 18.1 for training and Algorithm 18.1 for prediction.

How well does this approach work?

The first question is: how good is the expert? If we learn to mimic an expert, but the expert is no good, we lose. In general, it also seems unrealistic to believe this algorithm should be able to *improve* on the expert. Similarly, if our multiclass classification algorithm \mathcal{A} is crummy, we also lose. So part of the question “how well does

¹ It's completely okay for f to look at more than just x_t when making predictions; for instance, it might want to look at x_{t-1} , or a_{t-1} and a_{t-2} . As long as it only references information from the *past*, this is fine. For notational simplicity, we will assume that all of this relevant information is summarized in x_t .

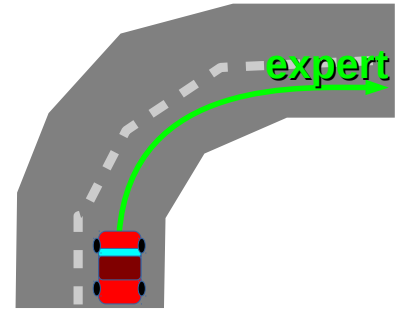


Figure 18.1: A single expert trajectory in a self-driving car.

Algorithm 43 `SUPERVISEDIMITATIONTRAIN`($\mathcal{A}, \tau_1, \tau_2, \dots, \tau_N$)

```

1:  $D \leftarrow \langle (x, a) : \forall n, \forall (x, a, \ell) \in \tau_n \rangle$  // collect all observation/action pairs
2: return  $\mathcal{A}(D)$  // train multiclass classifier on D

```

Algorithm 44 `SUPERVISEDIMITATIONTEST`(f)

```

1: for  $t = 1 \dots T$  do
2:    $x_t \leftarrow$  current observation
3:    $a_t \leftarrow f(x_t)$  // ask policy to choose an action
4:   take action  $a_t$ 
5:    $\ell_t \leftarrow$  observe instantaneous loss
6: end for
7: return  $\sum_{t=1}^T \ell_t$  // return total loss

```

this work” is the more basic question of: what are we even trying to measure?

There is a nice theorem² that gives an upper bound on the loss suffered by the SupervisedIL algorithm (Algorithm 18.1) as a function of (a) the quality of the expert, and (b) the error rate of the learned classifier. To be clear, we need to distinguish between the loss of the policy when run for T steps to form a full trajectory, and the error rate of the learned classifier, which is just its average multiclass classification error. The theorem states, roughly, that the loss of the learned *policy* is at most the loss of the expert plus T^2 times the error rate of the classifier.

² Ross et al. 2011

Theorem 18 (Loss of SupervisedIL). *Suppose that one runs Algorithm 18.1 using a multiclass classifier that optimizes the 0-1 loss (or an upperbound thereof). Let ϵ be the error rate of the underlying classifier (in expectation) and assume that all instantaneous losses are in the range $[0, \ell^{(\max)}]$. Let f be the learned policy; then:*

$$\underbrace{\mathbb{E}_{\tau \sim f} \left[\sum_t \ell_t \right]}_{\text{loss of learned policy}} \leq \underbrace{\mathbb{E}_{\tau \sim \text{expert}} \left[\sum_t \ell_t \right]}_{\text{loss of expert}} + \ell^{(\max)} T^2 \epsilon \quad (18.2)$$

Intuitively, this bound on the loss is about a factor of T away from what we might hope for. In particular, the multiclass classifier makes errors on an ϵ fraction of its actions, measured by zero/one loss. In the worst case, this will lead to a loss of $\ell^{(\max)} \epsilon$ for a single step. Summing all these errors over the entire trajectory would lead to a loss on the order of $\ell^{(\max)} T \epsilon$, which is a factor T better than this theorem provides. A natural question (addressed in the next section) is whether this analysis is tight. A related question (addressed in the section after that) is whether we can do better. Before getting there, though, it’s worth highlighting that an extra factor of T is *really*

bad. It can cause even very small multiclass error rates to blow up; in particular, if $\epsilon \geq 1/T$, we lose, and T can be in the hundreds or more.

18.2 Failure Analysis

The biggest single issue with the supervised learning approach to imitation learning is that it cannot learn to recover from failures. That is: it has only been trained based on expert trajectories. This means that the only training data it has seen is that of an expert driver. If it *ever* veers from that state distribution, it may have no idea how to recover. As a concrete example, perhaps the expert driver never ever gets themselves into a state where they are directly facing a wall. Moreover, the expert driver probably tends to drive forward more than backward. If the imperfect learner manages to make a few errors and get stuck next to a wall, it's likely to resort to the general “drive forward” rule and stay there forever. This is the problem of **compounding error**; and yes, it does happen in practice.

It turns out that it's possible to construct an imitation learning problem on which the T^2 compounding error is unavoidable. Consider the following somewhat artificial problem. At time $t = 1$ you're shown a picture of either a zero or a one. You have two possible actions: press a button marked “zero” or press a button marked “one.” The “correct” thing to do at $t = 1$ is to press the button that corresponds to the image you've been shown. Pressing the correct button leads to $\ell_1 = 0$; the incorrect leads to $\ell_1 = 1$. Now, at time $t = 2$ you are shown another image, again of a zero or one. The correct thing to do in this time step is the xor of (a) the number written on the picture you see right now, and (b) the correct answer from the previous time step. This holds in general for $t > 1$.

There are two important things about this construction. The first is that the expert can easily get zero loss. The second is that once the learned policy makes a single mistake, this can cause it to make *all* future decisions incorrectly. (At least until it “luckily” makes another “mistake” to get it back on track.)

Based on this construction, you can show the following theorem³.

³ Kääriäinen 2006

Theorem 19 (Lower Bound for SupervisedIL). *There exist imitation learning problems on which Algorithm 18.1 is able to achieve small classification error $\epsilon \in [0, 1/T]$ under an optimal expert, but for which the test loss is lower bounded as:*

$$\underbrace{\mathbb{E}_{\tau \sim f} \left[\sum_t \ell_t \right]}_{\text{loss of learned policy}} \geq \frac{T+1}{2} - \frac{1}{4\epsilon} \left[1 - (1-2\epsilon)^{T+1} \right] \quad (18.3)$$

which is bounded by $T^2\epsilon$ and, for small ϵ , grows like $T^2\epsilon$.

Up to constants, this gives matching upper and lower bounds for the loss of a policy learned by supervised imitation learning that is pretty far (a factor of T) from what we might hope for.

18.3 Dataset Aggregation

Supervised imitation learning fails because once it gets “off the expert path,” things can go really badly. Ideally, we might want to train our policy to deal with *any* possible situation it could encounter. Unfortunately, this is unrealistic: we cannot hope to be able to train on every possible configuration of the world; and if we could, we wouldn’t really need to learn anyway, we could just memorize. So we want to train f on a subset of world configurations, but using “configurations visited by the expert” fails because f cannot learn to recover from its own errors. Somehow what we’d like to do is train f to do well on the configurations that it, itself, encounters!

This is a classic chicken-and-egg problem. We want a policy f that does well in a bunch of world configurations. What set of configurations? The configurations that f encounters! A very classic approach to solving chicken-and-egg problems is iteration. Start with some policy f . Run f and see what configurations it visits. Train a new f to do well there. Repeat.

This is exactly what the Dataset Aggregation algorithm (“Dagger”) does. Continuing with the self-driving car analogy, we first let a human expert drive a car for a while, and learn an initial policy f_0 by running standard supervised imitation learning (Algorithm 18.1) on the trajectories visited by the human. We then do something unusual. We put the human expert in the car, and record their actions, but the car behaves not according to the expert’s behavior, but according to f_0 . That is, f_0 is in control of the car, and the expert is trying to steer, but the car is ignoring them⁴ and simply recording their actions as training data. This is shown in Figure 18.2.

Based on trajectories generated by f_0 but actions given by the expert, we generate a new dataset that contains information about how to recover from the errors of f_0 . We now will train a new policy, f_1 . Because we don’t want f_1 to “forget” what f_0 already knows, f_1 is trained on the union of the initial expert-only trajectories together with the new trajectories generated by f_0 . We repeat this process a number of times MaxIter , yielding Algorithm 18.3.

This algorithm returns the list of *all* policies generated during its run. A very practical question is: which one should you use? There are essentially two choices. The first choice is just to use the final policy learned. The problem with this approach is that Dagger can be somewhat unstable in practice, and policies do not monotonically

⁴ This is possibly terrifying for the expert!

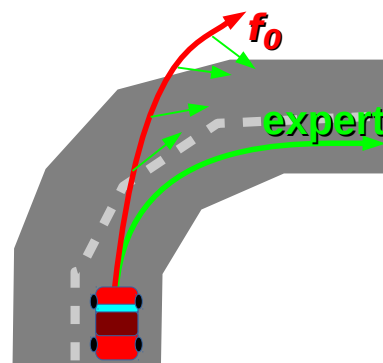


Figure 18.2: In DAgger, the trajectory (red) is generated according to the previously learned policy, f_0 , but the gold standard actions are given by the expert.

Algorithm 45 DAGGERTRAIN(\mathcal{A} , MaxIter, N , expert)

```

1:  $\langle \tau_n^{(0)} \rangle_{n=1}^N \leftarrow$  run the expert  $N$  many times
2:  $D_0 \leftarrow \langle \langle \mathbf{x}, a \rangle : \forall n, \forall (\mathbf{x}, a, \ell) \in \tau_n^{(0)} \rangle$  // collect all pairs (same as supervised)
3:  $f_0 \leftarrow \mathcal{A}(D_0)$  // train initial policy (multiclass classifier) on  $D_0$ 
4: for  $i = 1 \dots \text{MaxIter}$  do
5:    $\langle \tau_n^{(i)} \rangle_{n=1}^N \leftarrow$  run policy  $f_{i-1}$   $N$ -many times // trajectories by  $f_{i-1}$ 
6:    $D_i \leftarrow \langle \langle \mathbf{x}, \text{expert}(\mathbf{x}) \rangle : \forall n, \forall (\mathbf{x}, a, \ell) \in \tau_n^{(i)} \rangle$  // collect data set
   // observations  $\mathbf{x}$  visited by  $f_{i-1}$ 
   // but actions according to the expert
7:    $f_i \leftarrow \mathcal{A}(\cup_{j=0}^i D_j)$  // train policy  $f_i$  on union of all data so far
8: end for
9: return  $\langle f_0, f_1, \dots, f_{\text{MaxIter}} \rangle$  // return collection of all learned policies

```

improve. A safer alternative (as we'll see by theory below) is to test all of them on some held-out “development” tasks, and pick the one that does best there. This requires a bit more computation, but is a much better approach in general.

One major difference in *requirements* between Dagger (Algorithm 18.3) and SupervisedIL (Algorithm 18.1) is the requirement of interaction with the expert. In SupervisedIL, you only need access to a bunch of trajectories taken by the expert, *passively*. In Dagger, you need access to them expert themselves, so you can ask questions like “if you saw configuration \mathbf{x} , what would you do?” This puts *much* more demand on the expert.

Another question that arises is: what should N , the number of trajectories generated in each round, be? In practice, the initial N should probably be reasonably large, so that the initial policy f_0 is pretty good. The number of trajectories generated by iteration subsequently can be much smaller, perhaps even just one.

Intuitively, Dagger should be less sensitive to compounding error than SupervisedIL, precisely because it gets trained on observations that it is likely to see at test time. This is formalized in the following theorem:

Theorem 20 (Loss of Dagger). *Suppose that one runs Algorithm 18.3 using a multiclass classifier that optimizes the 0-1 loss (or an upperbound thereof). Let ϵ be the error rate of the underlying classifier (in expectation) and assume that all instantaneous losses are in the range $[0, \ell^{(\max)}]$. Let f be the learned policy; then:*

$$\underbrace{\mathbb{E}_{\tau \sim f} \left[\sum_t \ell_t \right]}_{\text{loss of learned policy}} \leq \underbrace{\mathbb{E}_{\tau \sim \text{expert}} \left[\sum_t \ell_t \right]}_{\text{loss of expert}} + \ell^{(\max)} T \epsilon + O \left(\frac{\ell^{(\max)} T \log T}{\text{MaxIter}} \right)$$

(18.4)

Furthermore, if the loss function is strongly convex in f , and MaxIter is

$\tilde{O}(T/\epsilon)$, then:

$$\underbrace{\mathbb{E}_{\tau \sim f} \left[\sum_t \ell_t \right]}_{\text{loss of learned policy}} \leq \underbrace{\mathbb{E}_{\tau \sim \text{expert}} \left[\sum_t \ell_t \right]}_{\text{loss of expert}} + \ell^{(\max)} T \epsilon + O(\epsilon) \quad (18.5)$$

Both of these results show that, assuming MaxIter is large enough, the loss of the learned policy f (here, taken to be the best on of all the MaxIter policies learned) grows like $T\epsilon$, which is what we hope for. Note that the final term in the first bound gets small so long as MaxIter is at least $T \log T$.

18.4 Expensive Algorithms as Experts

Because of the strong requirement on the expert in Dagger (i.e., that you need to be able to query it many times during training), one of the most substantial use cases for Dagger is to learn to (quickly) imitate otherwise slow algorithms. Here are two prototypical examples:

1. Game playing. When a game (like chess or minecraft) can be run in simulation, you can often explicitly compute a semi-optimal expert behavior with brute-force search. But this search might be too computationally expensive to play in real time, so you can use it during training time, learning a fast policy that attempts to mimic the expensive search. This learned policy can then be applied at test time.
2. Discrete optimizers. Many discrete optimization problems can be computationally expensive to run in real time; for instance, even shortest path search on a large graph can be too slow for real time use. We can compute shortest paths offline as “training data” and then use imitation learning to try to build shortest path optimizers that will run sufficiently efficiently in real time.

Consider the game playing example, and for concreteness, suppose you are trying to learn to play solitaire (this is an easier example because it’s a single player game). When running `DaggerTrain` (Algorithm 18.3 to learn a chess-playing policy, the algorithm will repeatedly ask for `expert(x)`, where x is the current state of the game. What should this function return? Ideally, it should return the/an action a such that, if a is taken, and then the rest of the game is played optimally, the player wins. Computing this exactly is going to be very difficult for anything except the simplest games, so we need to resort to an approximation.

Algorithm 46 DEPTHLIMITEDDFS($x, h, \text{MaxDepth}$)

```

1: if  $x$  is a terminal state or  $\text{MaxDepth} \leq 0$  then
2:   return  $(\perp, h(x))$  // if we cannot search deeper
                        // return “no action” ( $\perp$ ) and the current heuristic score
3: else
4:    $\text{BestAction}, \text{BestScore} \leftarrow \perp, -\infty$  // keep track of best action & its score
5:   for all actions  $a$  from  $x$  do
6:      $(\_ \text{score}) \leftarrow \text{DEPTHLIMITEDDFS}(x \circ a, h, \text{MaxDepth} - 1)$  // get score
                        // for action  $a$ , depth reduced by one by appending  $a$  to  $x$ 
7:     if  $\text{score} > \text{BestScore}$  then
8:        $\text{BestAction}, \text{BestScore} \leftarrow a, \text{score}$  // update tracked best action & score
9:     end if
10:  end for
11: end if
12: return  $(\text{BestAction}, \text{BestScore})$  // return best found action and its score

```

A common strategy is to run a depth-limited depth first search, starting at state x , and terminating after at most three or four moves (see Figure 18.3). This will generate a search tree. Unless you are very near the end of the game, none of the leaves of this tree will correspond to the end of the game. So you’ll need some heuristic, h , for evaluating states that are non-terminals. You can propagate this heuristic score up to the root, and choose the action that looks best with this depth four search. This is not necessarily going to be the *optimal* action, and there’s a speed/accuracy trade-off for searching deeper, but this is typically effective. This approach summarized in Algorithm 18.4.

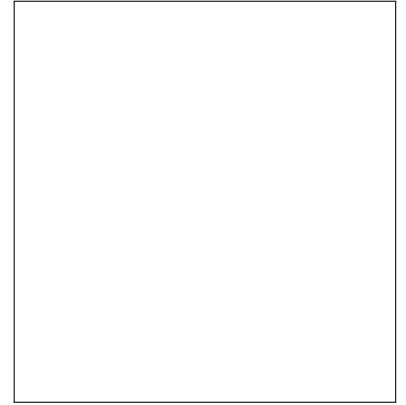


Figure 18.3: `imit:dldfs`: Depth limited depth-first search

18.5 Structured Prediction via Imitation Learning

A final case where an expert can often be computed algorithmically arises when one solves structured prediction (see Chapter 17) via imitation learning. It is clearest how this can work in the case of sequence labeling. Recall there that predicted outputs should be *sequences* of labels. The running example from the earlier chapter was:

$$x = \text{“ monsters eat tasty bunnies “} \quad (18.6)$$

$$y = \quad \text{noun verb adj noun} \quad (18.7)$$

One can easily cast the prediction of y as a sequential decision making problem, by treating the production of y in a left-to-right manner. In this case, we have a time horizon $T = 4$. We want to learn a policy f that first predicts “noun” then “verb” then “adj” then “noun” on this input.

Let's suppose that the input to f consists of features extracted both from the input (x) and the current predicted output prefix \hat{y} , denoted $\phi(x, \hat{y})$. For instance, $\phi(x, \hat{y})$ might represent a similar set of features to those use in Chapter 17. It is perhaps easiest to think of f as just a classifier: given some features of the input sentence x ("monsters eat tasty bunnies"), and some features about previous predictions in the output prefix (so far, produced "noun verb"), the goal of f is to predict the tag for the next word ("tasty") in this context.

An important question is: what is the "expert" in this case? Intuitively, the expert should provide the correct next label, but what does this mean? That depends on the loss function being optimized. Under Hamming loss (sum zero/one loss over each individual prediction), the expert is straightforward. When the expert is asked to produce an action for the third word, the expert's response is *always* "adj" (or whatever happens to be the correct label for the third word in the sentence it is currently training on).

More generally, the expert gets to look at x , y and a prefix \hat{y} of the output. Note, *importantly*, that the prefix *might be wrong*! In particular, after the first iteration of Dagger, the prefix will be predicted by the learned policy, which may make mistakes! The expert also has some structured loss function ℓ that it is trying to minimize. Like in the previous section, the expert's goal is to choose the action that minimizes the long-term loss according to ℓ on this example.

To be more formal, we need a bit of notation. Let $\text{best}(\ell, y, \hat{y})$ denote the loss (according to ℓ and the ground truth y) of the best reachable output starting at \hat{y} . For instance, if y is "noun verb adj noun" and \hat{y} is "noun noun", and the loss is Hamming loss, then the best achievable output (predicting left-to-right) is "noun noun adj noun" which has a loss of 1. Thus, best for this situation is 1.

Given that notion of best, the expert is easy to define:

$$\text{expert}(\ell, y, \hat{y}) = \underset{a}{\operatorname{argmin}} \text{best}(\ell, y, \hat{y} \circ a) \quad (18.8)$$

Namely, it is the action that leads to the best possible completion *after* taking that action. So in the example above, the expert action is "adj". For some problems and some loss functions, computing the expert is easy. In particular, for sequence labeling under Hamming loss, it's trivial. In the case that you can compute the expert exactly, it is often called an **oracle**.⁵ For some other problems, exactly computing an oracle is computationally expensive or intractable. In those cases, one can often resort to depth limited depth-first-search (Algorithm 18.4) to compute an approximate oracle as an expert.

To be very concrete, a typical implementation of Dagger applied to sequence labeling would go as follows. Each structured training example (a pair of sentence and tag-sequence) gives rise to one trajec-

⁵ Some literature calls it a "dynamic oracle", though the extra word is unnecessary.

tory. At training time, a predict tag sequence is generated left-to-right, starting with the empty sequence. At any given time step, you are attempting to predict the label of the t th word in the input. You define a feature vector $\phi(x, \hat{y})$, which will typically consist of: (a) the t th word, (b) left and right neighbors of the t th word, (c) the last few predictions in \hat{y} , and (d) anything else you can think of. *In particular*, the features are *not* limited to Markov style features, because we’re not longer trying to do dynamic programming. The expert label for the t th word is just the corresponding label in the ground truth y . Given all this, one can run Dagger (Algorithm 18.4) exactly as specified.

Moving to structured prediction problems other than sequence labeling problems is beyond the scope of this book. The general framework is to cast your structured prediction problem as a sequential decision making problem. Once you’ve done that, you need to decide on features (this is the easy part) and an expert (this is often the harder part). However, once you’ve done so, there are generic libraries for “compiling” your specification down to code.

18.6 Further Reading

TODO further reading

CODE AND DATASETS

Rating	Easy?	AI?	Sys?	Thy?	Morning?
+2	y	y	n	y	n
+2	y	y	n	y	n
+2	n	y	n	n	n
+2	n	n	n	y	n
+2	n	y	y	n	y
+1	y	y	n	n	n
+1	y	y	n	y	n
+1	n	y	n	y	n
0	n	n	n	n	y
0	y	n	n	y	y
0	n	y	n	y	n
0	y	y	y	y	y
-1	y	y	y	n	y
-1	n	n	y	y	n
-1	n	n	y	n	y
-1	y	n	y	n	y
-2	n	n	y	y	n
-2	n	y	y	n	y
-2	y	n	y	n	n
-2	y	n	y	n	y

Table 1: Course rating data set

BIBLIOGRAPHY

Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. *Advances in neural information processing systems*, 19:137, 2007.

Steffen Bickel, Michael Bruckner, and Tobias Scheffer. Discriminative learning for differing training and test distributions. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2007.

Sergey Brin. Near neighbor search in large metric spaces. In *Conference on Very Large Databases (VLDB)*, 1995.

Hal Daumé III. Frustratingly easy domain adaptation. In *Conference of the Association for Computational Linguistics (ACL)*, Prague, Czech Republic, 2007.

Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. On the (im)possibility of fairness. *arXiv preprint arXiv:1609.07236*, 2016.

Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*, pages 3315–3323, 2016.

Matti Kääriäinen. Lower bounds for reductions. Talk at the Atomic Learning Workshop (TTI-C), March 2006.

Tom M. Mitchell. *Machine Learning*. McGraw Hill, 1997.

J. Ross Quinlan. Induction of decision trees. *Machine learning*, 1(1): 81–106, 1986.

Frank Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65:386–408, 1958. Reprinted in *Neurocomputing* (MIT Press, 1998).

Stéphane Ross, Geoff J. Gordon, and J. Andrew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the Workshop on Artificial Intelligence and Statistics (AISTats)*, 2011.

- K -nearest neighbors, 58
- d_A -distance, 113
- p -norms, 92
- 0/1 loss, 88
- 80% rule, 111

- absolute loss, 14
- activation function, 130
- activations, 41
- AdaBoost, 166
- adaptation, 105
- algorithm, 87
- all pairs, 80
- all versus all, 80
- approximation error, 71
- architecture selection, 139
- area under the curve, 64, 84
- argmax problem, 199
- AUC, 64, 83, 84
- AVA, 80
- averaged perceptron, 52

- back-propagation, 134, 137
- bag of words, 56
- bagging, 165
- base learner, 164
- batch, 173
- Bayes error rate, 20
- Bayes optimal classifier, 19
- Bayes optimal error rate, 20
- Bayes rule, 117
- Bernoulli distribution, 121
- bias, 42
- bias/variance trade-off, 72
- binary features, 30
- bipartite ranking problems, 83
- boosting, 155, 164
- bootstrap resampling, 165
- bootstrapping, 67, 69

- categorical features, 30
- chain rule, 117, 120
- chord, 90
- circuit complexity, 138
- clustering, 35, 178
- clustering quality, 178
- complexity, 34
- compounding error, 215
- concave, 90
- concavity, 193
- concept, 157
- confidence intervals, 68
- constrained optimization problem, 100
- contour, 92
- convergence rate, 95
- convex, 87, 89
- covariate shift, 105
- cross validation, 65, 68
- cubic feature map, 144
- curvature, 95

- data covariance matrix, 184
- data generating distribution, 15
- decision boundary, 34
- decision stump, 168
- decision tree, 8, 10
- decision trees, 57
- density estimation, 107
- development data, 26
- dimensionality reduction, 178
- discrepancy, 113
- discrete distribution, 121
- disparate impact, 111
- distance, 31
- domain adaptation, 105
- dominates, 63
- dot product, 45
- dual problem, 151

- dual variables, 151

- early stopping, 53, 132
- embedding, 178
- ensemble, 164
- error driven, 43
- error rate, 88
- estimation error, 71
- Euclidean distance, 31
- evidence, 127
- example normalization, 59, 60
- examples, 9
- expectation maximization, 186, 189
- expected loss, 16
- expert, 212
- exponential loss, 90, 169

- feasible region, 101
- feature augmentation, 109
- feature combinations, 54
- feature mapping, 54
- feature normalization, 59
- feature scale, 33
- feature space, 31
- feature values, 11, 29
- feature vector, 29, 31
- features, 11, 29
- forward-propagation, 137
- fractional assignments, 191
- furthest-first heuristic, 180

- Gaussian distribution, 121
- Gaussian kernel, 147
- Gaussian Mixture Models, 191
- generalize, 9, 17
- generative story, 123
- geometric view, 29
- global minimum, 94
- GMM, 191

- gradient, 93
- gradient ascent, 93
- gradient descent, 93
- Hamming loss, 202
- hard-margin SVM, 101
- hash kernel, 177
- held-out data, 26
- hidden units, 129
- hidden variables, 186
- hinge loss, 90, 203
- histogram, 12
- horizon, 213
- hyperbolic tangent, 130
- hypercube, 38
- hyperparameter, 26, 44, 89
- hyperplane, 41
- hyperspheres, 38
- hypothesis, 71, 157
- hypothesis class, 71, 160
- hypothesis testing, 67
- i.i.d. assumption, 117
- identically distributed, 24
- ILP, 195, 207
- imbalanced data, 73
- imitation learning, 212
- importance sampling, 106
- importance weight, 74
- independent, 24
- independently, 117
- independently and identically distributed, 117
- indicator function, 88
- induce, 16
- induced distribution, 76
- induction, 9
- inductive bias, 20, 31, 33, 91, 121
- integer linear program, 207
- integer linear programming, 195
- iteration, 36
- jack-knifing, 69
- Jensen's inequality, 193
- joint, 124
- K-nearest neighbors, 32
- Karush-Kuhn-Tucker conditions, 152
- kernel, 141, 145
- kernel trick, 146
- kernels, 54
- KKT conditions, 152
- label, 11
- Lagrange multipliers, 119
- Lagrange variable, 119
- Lagrangian, 119
- lattice, 200
- layer-wise, 139
- learning by demonstration, 212
- leave-one-out cross validation, 65
- level-set, 92
- license, 2
- likelihood, 127
- linear classifier, 169
- linear classifiers, 169
- linear decision boundary, 41
- linear regression, 98
- linearly separable, 48
- link function, 130
- log likelihood, 118
- log posterior, 127
- log probability, 118
- log-likelihood ratio, 122
- logarithmic transformation, 61
- logistic loss, 90
- logistic regression, 126
- LOO cross validation, 65
- loss function, 14
- loss-augmented inference, 205
- loss-augmented search, 205
- margin, 49, 100
- margin of a data set, 49
- marginal likelihood, 127
- marginalization, 117
- Markov features, 198
- maximum a posteriori, 127
- maximum depth, 26
- maximum likelihood estimation, 118
- mean, 59
- Mercer's condition, 146
- model, 87
- modeling, 25
- multi-layer network, 129
- naive Bayes assumption, 120
- nearest neighbor, 29, 31
- neural network, 169
- neural networks, 54, 129
- neurons, 41
- noise, 21
- non-convex, 135
- non-linear, 129
- Normal distribution, 121
- normalize, 46, 59
- null hypothesis, 67
- objective function, 88
- one versus all, 78
- one versus rest, 78
- online, 42
- optimization problem, 88
- oracle, 212, 220
- oracle experiment, 28
- output unit, 129
- OVA, 78
- overfitting, 23
- oversample, 76
- p-value, 67
- PAC, 156, 166
- paired t-test, 67
- parametric test, 67
- parity, 21
- parity function, 138
- patch representation, 56
- PCA, 184
- perceptron, 41, 42, 58
- perpendicular, 45
- pixel representation, 55
- policy, 212
- polynomial kernels, 146
- positive semi-definite, 146
- posterior, 127
- precision, 62
- precision/recall curves, 63
- predict, 9
- preference function, 82
- primal variables, 151
- principle components analysis, 184
- prior, 127
- probabilistic modeling, 116
- Probably Approximately Correct, 156
- programming by example, 212
- projected gradient, 151
- projection, 46
- psd, 146
- radial basis function, 139
- random forests, 169
- random variable, 117
- RBF kernel, 147

- RBF network, 139
- recall, 62
- receiver operating characteristic, 64
- reconstruction error, 184
- reductions, 76
- redundant features, 56
- regularized objective, 89
- regularizer, 88, 91
- reinforcement learning, 212
- representer theorem, 143, 145
- ROC curve, 64

- sample complexity, 157, 158, 160
- sample mean, 59
- sample selection bias, 105
- sample variance, 59
- semi-supervised adaptation, 106
- sensitivity, 64
- separating hyperplane, 87
- sequential decision making, 212
- SGD, 173
- shallow decision tree, 21, 168
- shape representation, 56
- sigmoid, 130
- sigmoid function, 126
- sigmoid network, 139
- sign, 130
- single-layer network, 129
- singular, 98
- slack, 148
- slack parameters, 101
- smoothed analysis, 180
- soft assignments, 190
- soft-margin SVM, 101

- span, 143
- sparse, 92
- specificity, 64
- squared loss, 14, 90
- statistical inference, 116
- statistically significant, 67
- steepest ascent, 93
- stochastic gradient descent, 173
- stochastic optimization, 172
- strong law of large numbers, 24
- strong learner, 166
- strong learning algorithm, 166
- strongly convex, 95
- structural risk minimization, 87
- structured hinge loss, 203
- structured prediction, 195
- sub-sampling, 75
- subderivative, 96
- subgradient, 96
- subgradient descent, 97
- sum-to-one, 117
- support vector machine, 100
- support vectors, 153
- surrogate loss, 90
- symmetric modes, 135

- t-test, 67
- test data, 25
- test error, 25
- test set, 9
- text categorization, 56
- the curse of dimensionality, 37
- threshold, 42
- Tikhonov regularization, 87

- time horizon, 213
- total variation distance, 113
- train/test mismatch, 105
- training data, 9, 16, 24
- training error, 16
- trajectory, 213
- trellis, 200
- truncated gradients, 175
- two-layer network, 129

- unary features, 198
- unbiased, 47
- underfitting, 23
- unit hypercube, 39
- unit vector, 46
- unsupervised adaptation, 106
- unsupervised learning, 35

- validation data, 26
- Vapnik-Chernovenkis dimension, 162
- variance, 59, 165
- variational distance, 113
- VC dimension, 162
- vector, 31
- visualize, 178
- vote, 32
- voted perceptron, 52
- voting, 52

- weak learner, 166
- weak learning algorithm, 166
- weights, 41

- zero/one loss, 14