

# CS229 Lecture Notes

Andrew Ng

Updated by Tengyu Ma

# Contents

<b>I</b>	<b>Supervised learning</b>	<b>3</b>
<b>1</b>	<b>Linear Regression</b>	<b>6</b>
1.1	LMS algorithm . . . . .	7
1.2	The normal equations . . . . .	11
1.2.1	Matrix derivatives . . . . .	11
1.2.2	Least squares revisited . . . . .	12
1.3	Probabilistic interpretation . . . . .	13
1.4	Locally weighted linear regression (optional reading) . . . . .	15
<b>2</b>	<b>Classification and logistic regression</b>	<b>18</b>
2.1	Logistic regression . . . . .	18
2.2	Digression: The perceptron learning algorithm . . . . .	21
2.3	Another algorithm for maximizing $\ell(\theta)$ . . . . .	22
<b>3</b>	<b>Generalized Linear Models</b>	<b>24</b>
3.1	The exponential family . . . . .	24
3.2	Constructing GLMs . . . . .	26
3.2.1	Ordinary Least Squares . . . . .	27
3.2.2	Logistic Regression . . . . .	28
3.2.3	Softmax Regression . . . . .	28
<b>4</b>	<b>Generative Learning algorithms</b>	<b>33</b>
4.1	Gaussian discriminant analysis . . . . .	34
4.1.1	The multivariate normal distribution . . . . .	34
4.1.2	The Gaussian Discriminant Analysis model . . . . .	37
4.1.3	Discussion: GDA and logistic regression . . . . .	39
4.2	Naive Bayes . . . . .	40
4.2.1	Laplace smoothing . . . . .	43
4.2.2	Event models for text classification . . . . .	45

<b>5</b>	<b>Kernel Methods</b>	<b>47</b>
5.1	Feature maps . . . . .	47
5.2	LMS (least mean squares) with features . . . . .	48
5.3	LMS with the kernel trick . . . . .	48
5.4	Properties of kernels . . . . .	52
<b>6</b>	<b>Support Vector Machines</b>	<b>58</b>
6.1	Margins: Intuition . . . . .	58
6.2	Notation (option reading) . . . . .	60
6.3	Functional and geometric margins (option reading) . . . . .	60
6.4	The optimal margin classifier (option reading) . . . . .	62
6.5	Lagrange duality (optional reading) . . . . .	64
6.6	Optimal margin classifiers: the dual form (option reading) . . . . .	67
6.7	Regularization and the non-separable case (optional reading) . . . . .	71
6.8	The SMO algorithm (optional reading) . . . . .	72
6.8.1	Coordinate ascent . . . . .	73
6.8.2	SMO . . . . .	74
<b>II</b>	<b>Deep Learning</b>	<b>78</b>
<b>7</b>	<b>Deep Learning</b>	<b>79</b>
7.1	Supervised Learning with Non-linear Models . . . . .	79
7.2	Neural Networks . . . . .	81
7.3	Backpropagation . . . . .	90
7.3.1	Preliminary: chain rule . . . . .	91
7.3.2	One-neuron neural networks . . . . .	91
7.3.3	Two-layer neural networks: a low-level unpacked computation . . . . .	92
7.3.4	Two-layer neural network with vector notation . . . . .	95
7.3.5	Multi-layer neural networks . . . . .	97
7.4	Vectorization Over Training Examples . . . . .	97

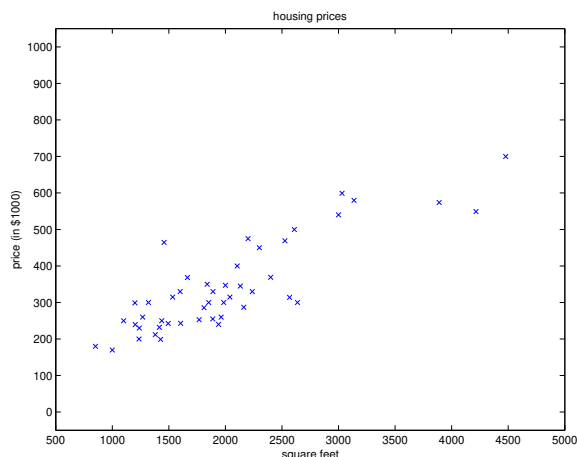
# Part I

## Supervised learning

Let’s start by talking about a few examples of supervised learning problems. Suppose we have a dataset giving the living areas and prices of 47 houses from Portland, Oregon:

Living area (feet <sup>2</sup> )	Price (1000\$)
2104	400
1600	330
2400	369
1416	232
3000	540
$\vdots$	$\vdots$

We can plot this data:

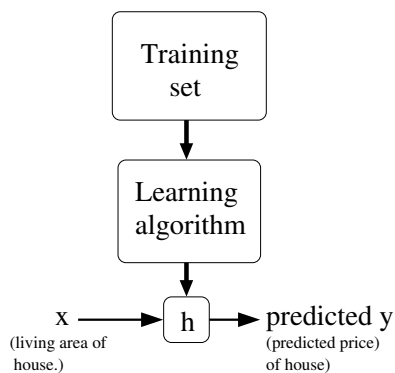


Given data like this, how can we learn to predict the prices of other houses in Portland, as a function of the size of their living areas?

To establish notation for future use, we’ll use  $x^{(i)}$  to denote the “input” variables (living area in this example), also called input **features**, and  $y^{(i)}$  to denote the “output” or **target** variable that we are trying to predict (price). A pair  $(x^{(i)}, y^{(i)})$  is called a **training example**, and the dataset that we’ll be using to learn—a list of  $n$  training examples  $\{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$ —is called a **training set**. Note that the superscript “ $(i)$ ” in the notation is simply an index into the training set, and has nothing to do with exponentiation. We will also use  $\mathcal{X}$  denote the space of input values, and  $\mathcal{Y}$  the space of output values. In this example,  $\mathcal{X} = \mathcal{Y} = \mathbb{R}$ .

To describe the supervised learning problem slightly more formally, our goal is, given a training set, to learn a function  $h : \mathcal{X} \mapsto \mathcal{Y}$  so that  $h(x)$  is a “good” predictor for the corresponding value of  $y$ . For historical reasons, this

function  $h$  is called a **hypothesis**. Seen pictorially, the process is therefore like this:



When the target variable that we're trying to predict is continuous, such as in our housing example, we call the learning problem a **regression** problem. When  $y$  can take on only a small number of discrete values (such as if, given the living area, we wanted to predict if a dwelling is a house or an apartment, say), we call it a **classification** problem.

# Chapter 1

## Linear Regression

To make our housing example more interesting, let's consider a slightly richer dataset in which we also know the number of bedrooms in each house:

Living area (feet <sup>2</sup> )	#bedrooms	Price (1000\$)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
$\vdots$	$\vdots$	$\vdots$

Here, the  $x$ 's are two-dimensional vectors in  $\mathbb{R}^2$ . For instance,  $x_1^{(i)}$  is the living area of the  $i$ -th house in the training set, and  $x_2^{(i)}$  is its number of bedrooms. (In general, when designing a learning problem, it will be up to you to decide what features to choose, so if you are out in Portland gathering housing data, you might also decide to include other features such as whether each house has a fireplace, the number of bathrooms, and so on. We'll say more about feature selection later, but for now let's take the features as given.)

To perform supervised learning, we must decide how we're going to represent functions/hypotheses  $h$  in a computer. As an initial choice, let's say we decide to approximate  $y$  as a linear function of  $x$ :

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

Here, the  $\theta_i$ 's are the **parameters** (also called **weights**) parameterizing the space of linear functions mapping from  $\mathcal{X}$  to  $\mathcal{Y}$ . When there is no risk of

confusion, we will drop the  $\theta$  subscript in  $h_\theta(x)$ , and write it more simply as  $h(x)$ . To simplify our notation, we also introduce the convention of letting  $x_0 = 1$  (this is the **intercept term**), so that

$$h(x) = \sum_{i=0}^d \theta_i x_i = \theta^T x,$$

where on the right-hand side above we are viewing  $\theta$  and  $x$  both as vectors, and here  $d$  is the number of input variables (not counting  $x_0$ ).

Now, given a training set, how do we pick, or learn, the parameters  $\theta$ ? One reasonable method seems to be to make  $h(x)$  close to  $y$ , at least for the training examples we have. To formalize this, we will define a function that measures, for each value of the  $\theta$ 's, how close the  $h(x^{(i)})$ 's are to the corresponding  $y^{(i)}$ 's. We define the **cost function**:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2.$$

If you've seen linear regression before, you may recognize this as the familiar least-squares cost function that gives rise to the **ordinary least squares** regression model. Whether or not you have seen it previously, let's keep going, and we'll eventually show this to be a special case of a much broader family of algorithms.

## 1.1 LMS algorithm

We want to choose  $\theta$  so as to minimize  $J(\theta)$ . To do so, let's use a search algorithm that starts with some "initial guess" for  $\theta$ , and that repeatedly changes  $\theta$  to make  $J(\theta)$  smaller, until hopefully we converge to a value of  $\theta$  that minimizes  $J(\theta)$ . Specifically, let's consider the **gradient descent** algorithm, which starts with some initial  $\theta$ , and repeatedly performs the update:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta).$$

(This update is simultaneously performed for all values of  $j = 0, \dots, d$ .) Here,  $\alpha$  is called the **learning rate**. This is a very natural algorithm that repeatedly takes a step in the direction of steepest decrease of  $J$ .

In order to implement this algorithm, we have to work out what is the partial derivative term on the right hand side. Let's first work it out for the



case of if we have only one training example  $(x, y)$ , so that we can neglect the sum in the definition of  $J$ . We have:

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2 \\
 &= 2 \cdot \frac{1}{2} (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \\
 &= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^d \theta_i x_i - y \right) \\
 &= (h_\theta(x) - y) x_j
 \end{aligned}$$

For a single training example, this gives the update rule:<sup>1</sup>

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}.$$

The rule is called the **LMS** update rule (LMS stands for “least mean squares”), and is also known as the **Widrow-Hoff** learning rule. This rule has several properties that seem natural and intuitive. For instance, the magnitude of the update is proportional to the **error** term  $(y^{(i)} - h_\theta(x^{(i)}))$ ; thus, for instance, if we are encountering a training example on which our prediction nearly matches the actual value of  $y^{(i)}$ , then we find that there is little need to change the parameters; in contrast, a larger change to the parameters will be made if our prediction  $h_\theta(x^{(i)})$  has a large error (i.e., if it is very far from  $y^{(i)}$ ).

We’d derived the LMS rule for when there was only a single training example. There are two ways to modify this method for a training set of more than one example. The first is replace it with the following algorithm:

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^n (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}, \text{ (for every } j \text{)} \quad (1.1)$$

}

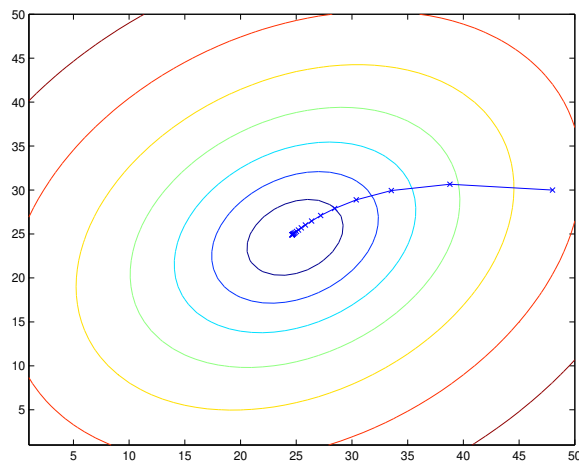
---

<sup>1</sup>We use the notation “ $a := b$ ” to denote an operation (in a computer program) in which we *set* the value of a variable  $a$  to be equal to the value of  $b$ . In other words, this operation overwrites  $a$  with the value of  $b$ . In contrast, we will write “ $a = b$ ” when we are asserting a statement of fact, that the value of  $a$  is equal to the value of  $b$ .

By grouping the updates of the coordinates into an update of the vector  $\theta$ , we can rewrite update (1.1) in a slightly more succinct way:

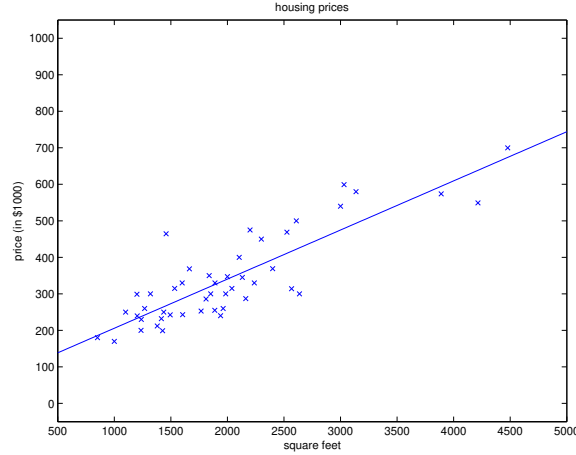
$$\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x^{(i)}$$

The reader can easily verify that the quantity in the summation in the update rule above is just  $\partial J(\theta)/\partial \theta_j$  (for the original definition of  $J$ ). So, this is simply gradient descent on the original cost function  $J$ . This method looks at every example in the entire training set on every step, and is called **batch gradient descent**. Note that, while gradient descent can be susceptible to local minima in general, the optimization problem we have posed here for linear regression has only one global, and no other local, optima; thus gradient descent always converges (assuming the learning rate  $\alpha$  is not too large) to the global minimum. Indeed,  $J$  is a convex quadratic function. Here is an example of gradient descent as it is run to minimize a quadratic function.



The ellipses shown above are the contours of a quadratic function. Also shown is the trajectory taken by gradient descent, which was initialized at (48,30). The  $x$ 's in the figure (joined by straight lines) mark the successive values of  $\theta$  that gradient descent went through.

When we run batch gradient descent to fit  $\theta$  on our previous dataset, to learn to predict housing price as a function of living area, we obtain  $\theta_0 = 71.27$ ,  $\theta_1 = 0.1345$ . If we plot  $h_{\theta}(x)$  as a function of  $x$  (area), along with the training data, we obtain the following figure:



If the number of bedrooms were included as one of the input features as well, we get  $\theta_0 = 89.60$ ,  $\theta_1 = 0.1392$ ,  $\theta_2 = -8.738$ .

The above results were obtained with batch gradient descent. There is an alternative to batch gradient descent that also works very well. Consider the following algorithm:

```

Loop {
    for  $i = 1$  to  $n$ , {


$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}, \quad (\text{for every } j) \quad (1.2)$$


    }
}

```

By grouping the updates of the coordinates into an update of the vector  $\theta$ , we can rewrite update (1.2) in a slightly more succinct way:

$$\theta := \theta + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x^{(i)}$$

In this algorithm, we repeatedly run through the training set, and each time we encounter a training example, we update the parameters according to the gradient of the error with respect to that single training example only. This algorithm is called **stochastic gradient descent** (also **incremental gradient descent**). Whereas batch gradient descent has to scan through the entire training set before taking a single step—a costly operation if  $n$  is large—stochastic gradient descent can start making progress right away, and

continues to make progress with each example it looks at. Often, stochastic gradient descent gets  $\theta$  “close” to the minimum much faster than batch gradient descent. (Note however that it may never “converge” to the minimum, and the parameters  $\theta$  will keep oscillating around the minimum of  $J(\theta)$ ; but in practice most of the values near the minimum will be reasonably good approximations to the true minimum.<sup>2</sup>) For these reasons, particularly when the training set is large, stochastic gradient descent is often preferred over batch gradient descent.

## 1.2 The normal equations

Gradient descent gives one way of minimizing  $J$ . Let’s discuss a second way of doing so, this time performing the minimization explicitly and without resorting to an iterative algorithm. In this method, we will minimize  $J$  by explicitly taking its derivatives with respect to the  $\theta_j$ ’s, and setting them to zero. To enable us to do this without having to write reams of algebra and pages full of matrices of derivatives, let’s introduce some notation for doing calculus with matrices.

### 1.2.1 Matrix derivatives

For a function  $f : \mathbb{R}^{n \times d} \mapsto \mathbb{R}$  mapping from  $n$ -by- $d$  matrices to the real numbers, we define the derivative of  $f$  with respect to  $A$  to be:

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1d}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{n1}} & \cdots & \frac{\partial f}{\partial A_{nd}} \end{bmatrix}$$

Thus, the gradient  $\nabla_A f(A)$  is itself an  $n$ -by- $d$  matrix, whose  $(i, j)$ -element is  $\partial f / \partial A_{ij}$ . For example, suppose  $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$  is a 2-by-2 matrix, and the function  $f : \mathbb{R}^{2 \times 2} \mapsto \mathbb{R}$  is given by

$$f(A) = \frac{3}{2}A_{11} + 5A_{12}^2 + A_{21}A_{22}.$$

---

<sup>2</sup>By slowly letting the learning rate  $\alpha$  decrease to zero as the algorithm runs, it is also possible to ensure that the parameters will converge to the global minimum rather than merely oscillate around the minimum.

Here,  $A_{ij}$  denotes the  $(i, j)$  entry of the matrix  $A$ . We then have

$$\nabla_A f(A) = \begin{bmatrix} \frac{3}{2} & 10A_{12} \\ A_{22} & A_{21} \end{bmatrix}.$$

### 1.2.2 Least squares revisited

Armed with the tools of matrix derivatives, let us now proceed to find in closed-form the value of  $\theta$  that minimizes  $J(\theta)$ . We begin by re-writing  $J$  in matrix-vectorial notation.

Given a training set, define the **design matrix**  $X$  to be the  $n$ -by- $d$  matrix (actually  $n$ -by- $d + 1$ , if we include the intercept term) that contains the training examples' input values in its rows:

$$X = \begin{bmatrix} - & (x^{(1)})^T & - \\ - & (x^{(2)})^T & - \\ & \vdots & \\ - & (x^{(n)})^T & - \end{bmatrix}.$$

Also, let  $\vec{y}$  be the  $n$ -dimensional vector containing all the target values from the training set:

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}.$$

Now, since  $h_\theta(x^{(i)}) = (x^{(i)})^T \theta$ , we can easily verify that

$$\begin{aligned} X\theta - \vec{y} &= \begin{bmatrix} (x^{(1)})^T \theta \\ \vdots \\ (x^{(n)})^T \theta \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \\ &= \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(x^{(n)}) - y^{(n)} \end{bmatrix}. \end{aligned}$$

Thus, using the fact that for a vector  $z$ , we have that  $z^T z = \sum_i z_i^2$ :

$$\begin{aligned} \frac{1}{2}(X\theta - \vec{y})^T(X\theta - \vec{y}) &= \frac{1}{2} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= J(\theta) \end{aligned}$$

Finally, to minimize  $J$ , let's find its derivatives with respect to  $\theta$ . Hence,

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) \\
 &= \frac{1}{2} \nabla_{\theta} ((X\theta)^T X\theta - (X\theta)^T \vec{y} - \vec{y}^T (X\theta) + \vec{y}^T \vec{y}) \\
 &= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X) \theta - \vec{y}^T (X\theta) - \vec{y}^T (X\theta)) \\
 &= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X) \theta - 2(X^T \vec{y})^T \theta) \\
 &= \frac{1}{2} (2X^T X\theta - 2X^T \vec{y}) \\
 &= X^T X\theta - X^T \vec{y}
 \end{aligned}$$

In the third step, we used the fact that  $a^T b = b^T a$ , and in the fifth step used the facts  $\nabla_x b^T x = b$  and  $\nabla_x x^T A x = 2Ax$  for symmetric matrix  $A$  (for more details, see Section 4.3 of “Linear Algebra Review and Reference”). To minimize  $J$ , we set its derivatives to zero, and obtain the **normal equations**:

$$X^T X\theta = X^T \vec{y}$$

Thus, the value of  $\theta$  that minimizes  $J(\theta)$  is given in closed form by the equation

$$\theta = (X^T X)^{-1} X^T \vec{y}.^3$$

## 1.3 Probabilistic interpretation

When faced with a regression problem, why might linear regression, and specifically why might the least-squares cost function  $J$ , be a reasonable choice? In this section, we will give a set of probabilistic assumptions, under which least-squares regression is derived as a very natural algorithm.

Let us assume that the target variables and the inputs are related via the equation

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)},$$

---

<sup>3</sup>Note that in the above step, we are implicitly assuming that  $X^T X$  is an invertible matrix. This can be checked before calculating the inverse. If either the number of linearly independent examples is fewer than the number of features, or if the features are not linearly independent, then  $X^T X$  will not be invertible. Even in such cases, it is possible to “fix” the situation with additional techniques, which we skip here for the sake of simplicity.

where  $\epsilon^{(i)}$  is an error term that captures either unmodeled effects (such as if there are some features very pertinent to predicting housing price, but that we'd left out of the regression), or random noise. Let us further assume that the  $\epsilon^{(i)}$  are distributed IID (independently and identically distributed) according to a Gaussian distribution (also called a Normal distribution) with mean zero and some variance  $\sigma^2$ . We can write this assumption as " $\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$ ." I.e., the density of  $\epsilon^{(i)}$  is given by

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

This implies that

$$p(y^{(i)}|x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right).$$

The notation " $p(y^{(i)}|x^{(i)}; \theta)$ " indicates that this is the distribution of  $y^{(i)}$  given  $x^{(i)}$  and parameterized by  $\theta$ . Note that we should not condition on  $\theta$  (" $p(y^{(i)}|x^{(i)}, \theta)$ "), since  $\theta$  is not a random variable. We can also write the distribution of  $y^{(i)}$  as  $y^{(i)} | x^{(i)}; \theta \sim \mathcal{N}(\theta^T x^{(i)}, \sigma^2)$ .

Given  $X$  (the design matrix, which contains all the  $x^{(i)}$ 's) and  $\theta$ , what is the distribution of the  $y^{(i)}$ 's? The probability of the data is given by  $p(\vec{y}|X; \theta)$ . This quantity is typically viewed a function of  $\vec{y}$  (and perhaps  $X$ ), for a fixed value of  $\theta$ . When we wish to explicitly view this as a function of  $\theta$ , we will instead call it the **likelihood** function:

$$L(\theta) = L(\theta; X, \vec{y}) = p(\vec{y}|X; \theta).$$

Note that by the independence assumption on the  $\epsilon^{(i)}$ 's (and hence also the  $y^{(i)}$ 's given the  $x^{(i)}$ 's), this can also be written

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right). \end{aligned}$$

Now, given this probabilistic model relating the  $y^{(i)}$ 's and the  $x^{(i)}$ 's, what is a reasonable way of choosing our best guess of the parameters  $\theta$ ? The principal of **maximum likelihood** says that we should choose  $\theta$  so as to make the data as high probability as possible. I.e., we should choose  $\theta$  to maximize  $L(\theta)$ .

Instead of maximizing  $L(\theta)$ , we can also maximize any strictly increasing function of  $L(\theta)$ . In particular, the derivations will be a bit simpler if we instead maximize the **log likelihood**  $\ell(\theta)$ :

$$\begin{aligned}\ell(\theta) &= \log L(\theta) \\ &= \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2.\end{aligned}$$

Hence, maximizing  $\ell(\theta)$  gives the same answer as minimizing

$$\frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2,$$

which we recognize to be  $J(\theta)$ , our original least-squares cost function.

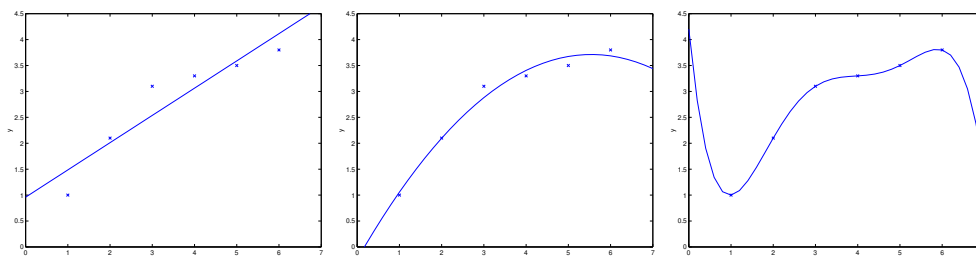
To summarize: Under the previous probabilistic assumptions on the data, least-squares regression corresponds to finding the maximum likelihood estimate of  $\theta$ . This is thus one set of assumptions under which least-squares regression can be justified as a very natural method that's just doing maximum likelihood estimation. (Note however that the probabilistic assumptions are by no means *necessary* for least-squares to be a perfectly good and rational procedure, and there may—and indeed there are—other natural assumptions that can also be used to justify it.)

Note also that, in our previous discussion, our final choice of  $\theta$  did not depend on what was  $\sigma^2$ , and indeed we'd have arrived at the same result even if  $\sigma^2$  were unknown. We will use this fact again later, when we talk about the exponential family and generalized linear models.

## 1.4 Locally weighted linear regression (optional reading)

Consider the problem of predicting  $y$  from  $x \in \mathbb{R}$ . The leftmost figure below shows the result of fitting a  $y = \theta_0 + \theta_1 x$  to a dataset. We see that the data doesn't really lie on straight line, and so the fit is not very good.





Instead, if we had added an extra feature  $x^2$ , and fit  $y = \theta_0 + \theta_1 x + \theta_2 x^2$ , then we obtain a slightly better fit to the data. (See middle figure) Naively, it might seem that the more features we add, the better. However, there is also a danger in adding too many features: The rightmost figure is the result of fitting a 5-th order polynomial  $y = \sum_{j=0}^5 \theta_j x^j$ . We see that even though the fitted curve passes through the data perfectly, we would not expect this to be a very good predictor of, say, housing prices ( $y$ ) for different living areas ( $x$ ). Without formally defining what these terms mean, we'll say the figure on the left shows an instance of **underfitting**—in which the data clearly shows structure not captured by the model—and the figure on the right is an example of **overfitting**. (Later in this class, when we talk about learning theory we'll formalize some of these notions, and also define more carefully just what it means for a hypothesis to be good or bad.)

As discussed previously, and as shown in the example above, the choice of features is important to ensuring good performance of a learning algorithm. (When we talk about model selection, we'll also see algorithms for automatically choosing a good set of features.) In this section, let us briefly talk about the locally weighted linear regression (LWR) algorithm which, assuming there is sufficient training data, makes the choice of features less critical. This treatment will be brief, since you'll get a chance to explore some of the properties of the LWR algorithm yourself in the homework.

In the original linear regression algorithm, to make a prediction at a query point  $x$  (i.e., to evaluate  $h(x)$ ), we would:

1. Fit  $\theta$  to minimize  $\sum_i (y^{(i)} - \theta^T x^{(i)})^2$ .
2. Output  $\theta^T x$ .

In contrast, the locally weighted linear regression algorithm does the following:

1. Fit  $\theta$  to minimize  $\sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$ .
2. Output  $\theta^T x$ .

Here, the  $w^{(i)}$ 's are non-negative valued **weights**. Intuitively, if  $w^{(i)}$  is large for a particular value of  $i$ , then in picking  $\theta$ , we'll try hard to make  $(y^{(i)} - \theta^T x^{(i)})^2$  small. If  $w^{(i)}$  is small, then the  $(y^{(i)} - \theta^T x^{(i)})^2$  error term will be pretty much ignored in the fit.

A fairly standard choice for the weights is<sup>4</sup>

$$w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$

Note that the weights depend on the particular point  $x$  at which we're trying to evaluate  $x$ . Moreover, if  $|x^{(i)} - x|$  is small, then  $w^{(i)}$  is close to 1; and if  $|x^{(i)} - x|$  is large, then  $w^{(i)}$  is small. Hence,  $\theta$  is chosen giving a much higher “weight” to the (errors on) training examples close to the query point  $x$ . (Note also that while the formula for the weights takes a form that is cosmetically similar to the density of a Gaussian distribution, the  $w^{(i)}$ 's do not directly have anything to do with Gaussians, and in particular the  $w^{(i)}$  are not random variables, normally distributed or otherwise.) The parameter  $\tau$  controls how quickly the weight of a training example falls off with distance of its  $x^{(i)}$  from the query point  $x$ ;  $\tau$  is called the **bandwidth** parameter, and is also something that you'll get to experiment with in your homework.

Locally weighted linear regression is the first example we're seeing of a **non-parametric** algorithm. The (unweighted) linear regression algorithm that we saw earlier is known as a **parametric** learning algorithm, because it has a fixed, finite number of parameters (the  $\theta_i$ 's), which are fit to the data. Once we've fit the  $\theta_i$ 's and stored them away, we no longer need to keep the training data around to make future predictions. In contrast, to make predictions using locally weighted linear regression, we need to keep the entire training set around. The term “non-parametric” (roughly) refers to the fact that the amount of stuff we need to keep in order to represent the hypothesis  $h$  grows linearly with the size of the training set.

---

<sup>4</sup>If  $x$  is vector-valued, this is generalized to be  $w^{(i)} = \exp(-(x^{(i)} - x)^T(x^{(i)} - x)/(2\tau^2))$ , or  $w^{(i)} = \exp(-(x^{(i)} - x)^T \Sigma^{-1}(x^{(i)} - x)/(2\tau^2))$ , for an appropriate choice of  $\tau$  or  $\Sigma$ .

## Chapter 2

# Classification and logistic regression

Let's now talk about the classification problem. This is just like the regression problem, except that the values  $y$  we now want to predict take on only a small number of discrete values. For now, we will focus on the **binary classification** problem in which  $y$  can take on only two values, 0 and 1. (Most of what we say here will also generalize to the multiple-class case.) For instance, if we are trying to build a spam classifier for email, then  $x^{(i)}$  may be some features of a piece of email, and  $y$  may be 1 if it is a piece of spam mail, and 0 otherwise. 0 is also called the **negative class**, and 1 the **positive class**, and they are sometimes also denoted by the symbols “-” and “+.” Given  $x^{(i)}$ , the corresponding  $y^{(i)}$  is also called the **label** for the training example.

### 2.1 Logistic regression

We could approach the classification problem ignoring the fact that  $y$  is discrete-valued, and use our old linear regression algorithm to try to predict  $y$  given  $x$ . However, it is easy to construct examples where this method performs very poorly. Intuitively, it also doesn't make sense for  $h_{\theta}(x)$  to take values larger than 1 or smaller than 0 when we know that  $y \in \{0, 1\}$ .

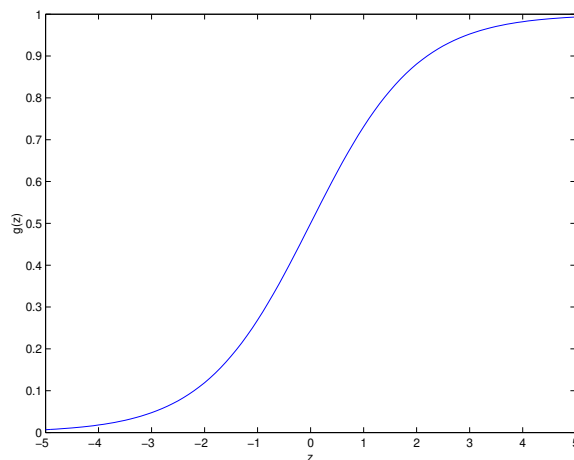
To fix this, let's change the form for our hypotheses  $h_{\theta}(x)$ . We will choose

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}},$$

where

$$g(z) = \frac{1}{1 + e^{-z}}$$

is called the **logistic function** or the **sigmoid function**. Here is a plot showing  $g(z)$ :



Notice that  $g(z)$  tends towards 1 as  $z \rightarrow \infty$ , and  $g(z)$  tends towards 0 as  $z \rightarrow -\infty$ . Moreover,  $g(z)$ , and hence also  $h(x)$ , is always bounded between 0 and 1. As before, we are keeping the convention of letting  $x_0 = 1$ , so that  $\theta^T x = \theta_0 + \sum_{j=1}^d \theta_j x_j$ .

For now, let's take the choice of  $g$  as given. Other functions that smoothly increase from 0 to 1 can also be used, but for a couple of reasons that we'll see later (when we talk about GLMs, and when we talk about generative learning algorithms), the choice of the logistic function is a fairly natural one. Before moving on, here's a useful property of the derivative of the sigmoid function, which we write as  $g'$ :

$$\begin{aligned} g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\ &= \frac{1}{(1 + e^{-z})^2} (e^{-z}) \\ &= \frac{1}{(1 + e^{-z})} \cdot \left( 1 - \frac{1}{(1 + e^{-z})} \right) \\ &= g(z)(1 - g(z)). \end{aligned}$$

So, given the logistic regression model, how do we fit  $\theta$  for it? Following how we saw least squares regression could be derived as the maximum likelihood estimator under a set of assumptions, let's endow our classification model with a set of probabilistic assumptions, and then fit the parameters via maximum likelihood.

Let us assume that

$$\begin{aligned} P(y = 1 \mid x; \theta) &= h_\theta(x) \\ P(y = 0 \mid x; \theta) &= 1 - h_\theta(x) \end{aligned}$$

Note that this can be written more compactly as

$$p(y \mid x; \theta) = (h_\theta(x))^y (1 - h_\theta(x))^{1-y}$$

Assuming that the  $n$  training examples were generated independently, we can then write down the likelihood of the parameters as

$$\begin{aligned} L(\theta) &= p(\vec{y} \mid X; \theta) \\ &= \prod_{i=1}^n p(y^{(i)} \mid x^{(i)}; \theta) \\ &= \prod_{i=1}^n (h_\theta(x^{(i)}))^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{1-y^{(i)}} \end{aligned}$$

As before, it will be easier to maximize the log likelihood:

$$\begin{aligned} \ell(\theta) &= \log L(\theta) \\ &= \sum_{i=1}^n y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \end{aligned}$$

How do we maximize the likelihood? Similar to our derivation in the case of linear regression, we can use gradient ascent. Written in vectorial notation, our updates will therefore be given by  $\theta := \theta + \alpha \nabla_\theta \ell(\theta)$ . (Note the positive rather than negative sign in the update formula, since we're maximizing, rather than minimizing, a function now.) Let's start by working with just one training example  $(x, y)$ , and take derivatives to derive the stochastic gradient ascent rule:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \ell(\theta) &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x) (1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= (y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x)) x_j \\ &= (y - h_\theta(x)) x_j \end{aligned}$$

Above, we used the fact that  $g'(z) = g(z)(1 - g(z))$ . This therefore gives us the stochastic gradient ascent rule

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

If we compare this to the LMS update rule, we see that it looks identical; but this is *not* the same algorithm, because  $h_\theta(x^{(i)})$  is now defined as a non-linear function of  $\theta^T x^{(i)}$ . Nonetheless, it's a little surprising that we end up with the same update rule for a rather different algorithm and learning problem. Is this coincidence, or is there a deeper reason behind this? We'll answer this when we get to GLM models.

## 2.2 Digression: The perceptron learning algorithm

We now digress to talk briefly about an algorithm that's of some historical interest, and that we will also return to later when we talk about learning theory. Consider modifying the logistic regression method to “force” it to output values that are either 0 or 1 or exactly. To do so, it seems natural to change the definition of  $g$  to be the threshold function:

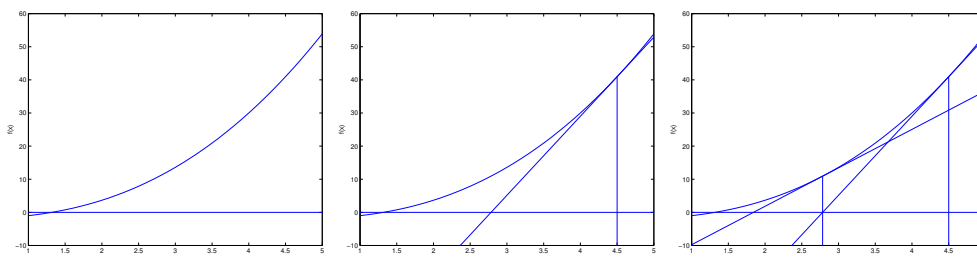
$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

If we then let  $h_\theta(x) = g(\theta^T x)$  as before but using this modified definition of  $g$ , and if we use the update rule

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}.$$

then we have the **perceptron learning algorithm**.

In the 1960s, this “perceptron” was argued to be a rough model for how individual neurons in the brain work. Given how simple the algorithm is, it will also provide a starting point for our analysis when we talk about learning theory later in this class. Note however that even though the perceptron may be cosmetically similar to the other algorithms we talked about, it is actually a very different type of algorithm than logistic regression and least squares linear regression; in particular, it is difficult to endow the perceptron's predictions with meaningful probabilistic interpretations, or derive the perceptron as a maximum likelihood estimation algorithm.



## 2.3 Another algorithm for maximizing $\ell(\theta)$

Returning to logistic regression with  $g(z)$  being the sigmoid function, let's now talk about a different algorithm for maximizing  $\ell(\theta)$ .

To get us started, let's consider Newton's method for finding a zero of a function. Specifically, suppose we have some function  $f : \mathbb{R} \mapsto \mathbb{R}$ , and we wish to find a value of  $\theta$  so that  $f(\theta) = 0$ . Here,  $\theta \in \mathbb{R}$  is a real number. Newton's method performs the following update:

$$\theta := \theta - \frac{f(\theta)}{f'(\theta)}.$$

This method has a natural interpretation in which we can think of it as approximating the function  $f$  via a linear function that is tangent to  $f$  at the current guess  $\theta$ , solving for where that linear function equals to zero, and letting the next guess for  $\theta$  be where that linear function is zero.

Here's a picture of the Newton's method in action:

In the leftmost figure, we see the function  $f$  plotted along with the line  $y = 0$ . We're trying to find  $\theta$  so that  $f(\theta) = 0$ ; the value of  $\theta$  that achieves this is about 1.3. Suppose we initialized the algorithm with  $\theta = 4.5$ . Newton's method then fits a straight line tangent to  $f$  at  $\theta = 4.5$ , and solves for the where that line evaluates to 0. (Middle figure.) This give us the next guess for  $\theta$ , which is about 2.8. The rightmost figure shows the result of running one more iteration, which the updates  $\theta$  to about 1.8. After a few more iterations, we rapidly approach  $\theta = 1.3$ .

Newton's method gives a way of getting to  $f(\theta) = 0$ . What if we want to use it to maximize some function  $\ell$ ? The maxima of  $\ell$  correspond to points where its first derivative  $\ell'(\theta)$  is zero. So, by letting  $f(\theta) = \ell'(\theta)$ , we can use the same algorithm to maximize  $\ell$ , and we obtain update rule:

$$\theta := \theta - \frac{\ell'(\theta)}{\ell''(\theta)}.$$

(Something to think about: How would this change if we wanted to use Newton's method to minimize rather than maximize a function?)

Lastly, in our logistic regression setting,  $\theta$  is vector-valued, so we need to generalize Newton's method to this setting. The generalization of Newton's method to this multidimensional setting (also called the Newton-Raphson method) is given by

$$\theta := \theta - H^{-1} \nabla_{\theta} \ell(\theta).$$

Here,  $\nabla_{\theta} \ell(\theta)$  is, as usual, the vector of partial derivatives of  $\ell(\theta)$  with respect to the  $\theta_i$ 's; and  $H$  is an  $d$ -by- $d$  matrix (actually,  $d+1$ -by- $d+1$ , assuming that we include the intercept term) called the **Hessian**, whose entries are given by

$$H_{ij} = \frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j}.$$

Newton's method typically enjoys faster convergence than (batch) gradient descent, and requires many fewer iterations to get very close to the minimum. One iteration of Newton's can, however, be more expensive than one iteration of gradient descent, since it requires finding and inverting an  $d$ -by- $d$  Hessian; but so long as  $d$  is not too large, it is usually much faster overall. When Newton's method is applied to maximize the logistic regression log likelihood function  $\ell(\theta)$ , the resulting method is also called **Fisher scoring**.



# Chapter 3

## Generalized Linear Models

So far, we've seen a regression example, and a classification example. In the regression example, we had  $y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$ , and in the classification one,  $y|x; \theta \sim \text{Bernoulli}(\phi)$ , for some appropriate definitions of  $\mu$  and  $\phi$  as functions of  $x$  and  $\theta$ . In this section, we will show that both of these methods are special cases of a broader family of models, called Generalized Linear Models (GLMs).<sup>1</sup> We will also show how other models in the GLM family can be derived and applied to other classification and regression problems.

### 3.1 The exponential family

To work our way up to GLMs, we will begin by defining exponential family distributions. We say that a class of distributions is in the exponential family if it can be written in the form

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta)) \quad (3.1)$$

Here,  $\eta$  is called the **natural parameter** (also called the **canonical parameter**) of the distribution;  $T(y)$  is the **sufficient statistic** (for the distributions we consider, it will often be the case that  $T(y) = y$ ); and  $a(\eta)$  is the **log partition function**. The quantity  $e^{-a(\eta)}$  essentially plays the role of a normalization constant, that makes sure the distribution  $p(y; \eta)$  sums/integrates over  $y$  to 1.

A fixed choice of  $T$ ,  $a$  and  $b$  defines a *family* (or set) of distributions that is parameterized by  $\eta$ ; as we vary  $\eta$ , we then get different distributions within this family.

---

<sup>1</sup>The presentation of the material in this section takes inspiration from Michael I. Jordan, *Learning in graphical models* (unpublished book draft), and also McCullagh and Nelder, *Generalized Linear Models* (2nd ed.).

We now show that the Bernoulli and the Gaussian distributions are examples of exponential family distributions. The Bernoulli distribution with mean  $\phi$ , written  $\text{Bernoulli}(\phi)$ , specifies a distribution over  $y \in \{0, 1\}$ , so that  $p(y = 1; \phi) = \phi$ ;  $p(y = 0; \phi) = 1 - \phi$ . As we vary  $\phi$ , we obtain Bernoulli distributions with different means. We now show that this class of Bernoulli distributions, ones obtained by varying  $\phi$ , is in the exponential family; i.e., that there is a choice of  $T$ ,  $a$  and  $b$  so that Equation (3.1) becomes exactly the class of Bernoulli distributions.

We write the Bernoulli distribution as:

$$\begin{aligned} p(y; \phi) &= \phi^y (1 - \phi)^{1-y} \\ &= \exp(y \log \phi + (1 - y) \log(1 - \phi)) \\ &= \exp \left( \left( \log \left( \frac{\phi}{1 - \phi} \right) \right) y + \log(1 - \phi) \right). \end{aligned}$$

Thus, the natural parameter is given by  $\eta = \log(\phi/(1 - \phi))$ . Interestingly, if we invert this definition for  $\eta$  by solving for  $\phi$  in terms of  $\eta$ , we obtain  $\phi = 1/(1 + e^{-\eta})$ . This is the familiar sigmoid function! This will come up again when we derive logistic regression as a GLM. To complete the formulation of the Bernoulli distribution as an exponential family distribution, we also have

$$\begin{aligned} T(y) &= y \\ a(\eta) &= -\log(1 - \phi) \\ &= \log(1 + e^\eta) \\ b(y) &= 1 \end{aligned}$$

This shows that the Bernoulli distribution can be written in the form of Equation (3.1), using an appropriate choice of  $T$ ,  $a$  and  $b$ .

Let's now move on to consider the Gaussian distribution. Recall that, when deriving linear regression, the value of  $\sigma^2$  had no effect on our final choice of  $\theta$  and  $h_\theta(x)$ . Thus, we can choose an arbitrary value for  $\sigma^2$  without changing anything. To simplify the derivation below, let's set  $\sigma^2 = 1$ .<sup>2</sup> We

---

<sup>2</sup>If we leave  $\sigma^2$  as a variable, the Gaussian distribution can also be shown to be in the exponential family, where  $\eta \in \mathbb{R}^2$  is now a 2-dimension vector that depends on both  $\mu$  and  $\sigma$ . For the purposes of GLMs, however, the  $\sigma^2$  parameter can also be treated by considering a more general definition of the exponential family:  $p(y; \eta, \tau) = b(a, \tau) \exp((\eta^T T(y) - a(\eta))/c(\tau))$ . Here,  $\tau$  is called the **dispersion parameter**, and for the Gaussian,  $c(\tau) = \sigma^2$ ; but given our simplification above, we won't need the more general definition for the examples we will consider here.

then have:

$$\begin{aligned} p(y; \mu) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - \mu)^2\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \cdot \exp\left(\mu y - \frac{1}{2}\mu^2\right) \end{aligned}$$

Thus, we see that the Gaussian is in the exponential family, with

$$\begin{aligned} \eta &= \mu \\ T(y) &= y \\ a(\eta) &= \mu^2/2 \\ &= \eta^2/2 \\ b(y) &= (1/\sqrt{2\pi}) \exp(-y^2/2). \end{aligned}$$

There're many other distributions that are members of the exponential family: The multinomial (which we'll see later), the Poisson (for modelling count-data; also see the problem set); the gamma and the exponential (for modelling continuous, non-negative random variables, such as time-intervals); the beta and the Dirichlet (for distributions over probabilities); and many more. In the next section, we will describe a general "recipe" for constructing models in which  $y$  (given  $x$  and  $\theta$ ) comes from any of these distributions.

## 3.2 Constructing GLMs

Suppose you would like to build a model to estimate the number  $y$  of customers arriving in your store (or number of page-views on your website) in any given hour, based on certain features  $x$  such as store promotions, recent advertising, weather, day-of-week, etc. We know that the Poisson distribution usually gives a good model for numbers of visitors. Knowing this, how can we come up with a model for our problem? Fortunately, the Poisson is an exponential family distribution, so we can apply a Generalized Linear Model (GLM). In this section, we will describe a method for constructing GLM models for problems such as these.

More generally, consider a classification or regression problem where we would like to predict the value of some random variable  $y$  as a function of  $x$ . To derive a GLM for this problem, we will make the following three assumptions about the conditional distribution of  $y$  given  $x$  and about our model:

1.  $y \mid x; \theta \sim \text{ExponentialFamily}(\eta)$ . I.e., given  $x$  and  $\theta$ , the distribution of  $y$  follows some exponential family distribution, with parameter  $\eta$ .
2. Given  $x$ , our goal is to predict the expected value of  $T(y)$  given  $x$ . In most of our examples, we will have  $T(y) = y$ , so this means we would like the prediction  $h(x)$  output by our learned hypothesis  $h$  to satisfy  $h(x) = \mathbb{E}[y|x]$ . (Note that this assumption is satisfied in the choices for  $h_\theta(x)$  for both logistic regression and linear regression. For instance, in logistic regression, we had  $h_\theta(x) = p(y = 1|x; \theta) = 0 \cdot p(y = 0|x; \theta) + 1 \cdot p(y = 1|x; \theta) = \mathbb{E}[y|x; \theta]$ .)
3. The natural parameter  $\eta$  and the inputs  $x$  are related linearly:  $\eta = \theta^T x$ . (Or, if  $\eta$  is vector-valued, then  $\eta_i = \theta_i^T x$ .)

The third of these assumptions might seem the least well justified of the above, and it might be better thought of as a “design choice” in our recipe for designing GLMs, rather than as an assumption per se. These three assumptions/design choices will allow us to derive a very elegant class of learning algorithms, namely GLMs, that have many desirable properties such as ease of learning. Furthermore, the resulting models are often very effective for modelling different types of distributions over  $y$ ; for example, we will shortly show that both logistic regression and ordinary least squares can both be derived as GLMs.

### 3.2.1 Ordinary Least Squares

To show that ordinary least squares is a special case of the GLM family of models, consider the setting where the target variable  $y$  (also called the **response variable** in GLM terminology) is continuous, and we model the conditional distribution of  $y$  given  $x$  as a Gaussian  $\mathcal{N}(\mu, \sigma^2)$ . (Here,  $\mu$  may depend  $x$ .) So, we let the *ExponentialFamily*( $\eta$ ) distribution above be the Gaussian distribution. As we saw previously, in the formulation of the Gaussian as an exponential family distribution, we had  $\mu = \eta$ . So, we have

$$\begin{aligned}
 h_\theta(x) &= \mathbb{E}[y|x; \theta] \\
 &= \mu \\
 &= \eta \\
 &= \theta^T x.
 \end{aligned}$$

The first equality follows from Assumption 2, above; the second equality follows from the fact that  $y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$ , and so its expected value is given

by  $\mu$ ; the third equality follows from Assumption 1 (and our earlier derivation showing that  $\mu = \eta$  in the formulation of the Gaussian as an exponential family distribution); and the last equality follows from Assumption 3.

### 3.2.2 Logistic Regression

We now consider logistic regression. Here we are interested in binary classification, so  $y \in \{0, 1\}$ . Given that  $y$  is binary-valued, it therefore seems natural to choose the Bernoulli family of distributions to model the conditional distribution of  $y$  given  $x$ . In our formulation of the Bernoulli distribution as an exponential family distribution, we had  $\phi = 1/(1 + e^{-\eta})$ . Furthermore, note that if  $y|x; \theta \sim \text{Bernoulli}(\phi)$ , then  $E[y|x; \theta] = \phi$ . So, following a similar derivation as the one for ordinary least squares, we get:

$$\begin{aligned} h_{\theta}(x) &= E[y|x; \theta] \\ &= \phi \\ &= 1/(1 + e^{-\eta}) \\ &= 1/(1 + e^{-\theta^T x}) \end{aligned}$$

So, this gives us hypothesis functions of the form  $h_{\theta}(x) = 1/(1 + e^{-\theta^T x})$ . If you are previously wondering how we came up with the form of the logistic function  $1/(1 + e^{-z})$ , this gives one answer: Once we assume that  $y$  conditioned on  $x$  is Bernoulli, it arises as a consequence of the definition of GLMs and exponential family distributions.

To introduce a little more terminology, the function  $g$  giving the distribution's mean as a function of the natural parameter ( $g(\eta) = E[T(y); \eta]$ ) is called the **canonical response function**. Its inverse,  $g^{-1}$ , is called the **canonical link function**. Thus, the canonical response function for the Gaussian family is just the identity function; and the canonical response function for the Bernoulli is the logistic function.<sup>3</sup>

### 3.2.3 Softmax Regression

Let's look at one more example of a GLM. Consider a classification problem in which the response variable  $y$  can take on any one of  $k$  values, so  $y \in \{1, 2, \dots, k\}$ . For example, rather than classifying email into the two classes

---

<sup>3</sup>Many texts use  $g$  to denote the link function, and  $g^{-1}$  to denote the response function; but the notation we're using here, inherited from the early machine learning literature, will be more consistent with the notation used in the rest of the class.

spam or not-spam—which would have been a binary classification problem—we might want to classify it into three classes, such as spam, personal mail, and work-related mail. The response variable is still discrete, but can now take on more than two values. We will thus model it as distributed according to a multinomial distribution.

Let's derive a GLM for modelling this type of multinomial data. To do so, we will begin by expressing the multinomial as an exponential family distribution.

To parameterize a multinomial over  $k$  possible outcomes, one could use  $k$  parameters  $\phi_1, \dots, \phi_k$  specifying the probability of each of the outcomes. However, these parameters would be redundant, or more formally, they would not be independent (since knowing any  $k - 1$  of the  $\phi_i$ 's uniquely determines the last one, as they must satisfy  $\sum_{i=1}^k \phi_i = 1$ ). So, we will instead parameterize the multinomial with only  $k - 1$  parameters,  $\phi_1, \dots, \phi_{k-1}$ , where  $\phi_i = p(y = i; \phi)$ , and  $p(y = k; \phi) = 1 - \sum_{i=1}^{k-1} \phi_i$ . For notational convenience, we will also let  $\phi_k = 1 - \sum_{i=1}^{k-1} \phi_i$ , but we should keep in mind that this is not a parameter, and that it is fully specified by  $\phi_1, \dots, \phi_{k-1}$ .

To express the multinomial as an exponential family distribution, we will define  $T(y) \in \mathbb{R}^{k-1}$  as follows:

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, T(2) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, T(3) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, T(k-1) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}, T(k) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$

Unlike our previous examples, here we do *not* have  $T(y) = y$ ; also,  $T(y)$  is now a  $k - 1$  dimensional vector, rather than a real number. We will write  $(T(y))_i$  to denote the  $i$ -th element of the vector  $T(y)$ .

We introduce one more very useful piece of notation. An indicator function  $1\{\cdot\}$  takes on a value of 1 if its argument is true, and 0 otherwise ( $1\{\text{True}\} = 1$ ,  $1\{\text{False}\} = 0$ ). For example,  $1\{2 = 3\} = 0$ , and  $1\{3 = 5 - 2\} = 1$ . So, we can also write the relationship between  $T(y)$  and  $y$  as  $(T(y))_i = 1\{y = i\}$ . (Before you continue reading, please make sure you understand why this is true!) Further, we have that  $E[(T(y))_i] = P(y = i) = \phi_i$ .

We are now ready to show that the multinomial is a member of the

exponential family. We have:

$$\begin{aligned}
p(y; \phi) &= \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_k^{1\{y=k\}} \\
&= \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_k^{1 - \sum_{i=1}^{k-1} 1\{y=i\}} \\
&= \phi_1^{(T(y))_1} \phi_2^{(T(y))_2} \dots \phi_k^{1 - \sum_{i=1}^{k-1} (T(y))_i} \\
&= \exp((T(y))_1 \log(\phi_1) + (T(y))_2 \log(\phi_2) + \\
&\quad \dots + \left(1 - \sum_{i=1}^{k-1} (T(y))_i\right) \log(\phi_k)) \\
&= \exp((T(y))_1 \log(\phi_1/\phi_k) + (T(y))_2 \log(\phi_2/\phi_k) + \\
&\quad \dots + (T(y))_{k-1} \log(\phi_{k-1}/\phi_k) + \log(\phi_k)) \\
&= b(y) \exp(\eta^T T(y) - a(\eta))
\end{aligned}$$

where

$$\begin{aligned}
\eta &= \begin{bmatrix} \log(\phi_1/\phi_k) \\ \log(\phi_2/\phi_k) \\ \vdots \\ \log(\phi_{k-1}/\phi_k) \end{bmatrix}, \\
a(\eta) &= -\log(\phi_k) \\
b(y) &= 1.
\end{aligned}$$

This completes our formulation of the multinomial as an exponential family distribution.

The link function is given (for  $i = 1, \dots, k$ ) by

$$\eta_i = \log \frac{\phi_i}{\phi_k}.$$

For convenience, we have also defined  $\eta_k = \log(\phi_k/\phi_k) = 0$ . To invert the link function and derive the response function, we therefore have that

$$\begin{aligned}
e^{\eta_i} &= \frac{\phi_i}{\phi_k} \\
\phi_k e^{\eta_i} &= \phi_i \\
\phi_k \sum_{i=1}^k e^{\eta_i} &= \sum_{i=1}^k \phi_i = 1
\end{aligned} \tag{3.2}$$

This implies that  $\phi_k = 1/\sum_{i=1}^k e^{\eta_i}$ , which can be substituted back into Equation (3.2) to give the response function

$$\phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}$$

This function mapping from the  $\eta$ 's to the  $\phi$ 's is called the **softmax** function.

To complete our model, we use Assumption 3, given earlier, that the  $\eta_i$ 's are linearly related to the  $x$ 's. So, have  $\eta_i = \theta_i^T x$  (for  $i = 1, \dots, k-1$ ), where  $\theta_1, \dots, \theta_{k-1} \in \mathbb{R}^{d+1}$  are the parameters of our model. For notational convenience, we can also define  $\theta_k = 0$ , so that  $\eta_k = \theta_k^T x = 0$ , as given previously. Hence, our model assumes that the conditional distribution of  $y$  given  $x$  is given by

$$\begin{aligned} p(y = i|x; \theta) &= \phi_i \\ &= \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}} \\ &= \frac{e^{\theta_i^T x}}{\sum_{j=1}^k e^{\theta_j^T x}} \end{aligned} \tag{3.3}$$

This model, which applies to classification problems where  $y \in \{1, \dots, k\}$ , is called **softmax regression**. It is a generalization of logistic regression.

Our hypothesis will output

$$\begin{aligned} h_\theta(x) &= E[T(y)|x; \theta] \\ &= E \left[ \begin{array}{c|c} \begin{matrix} 1\{y = 1\} \\ 1\{y = 2\} \\ \vdots \\ 1\{y = k-1\} \end{matrix} & x; \theta \end{array} \right] \\ &= \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\exp(\theta_1^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)} \\ \frac{\exp(\theta_2^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)} \\ \vdots \\ \frac{\exp(\theta_{k-1}^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)} \end{bmatrix}. \end{aligned}$$

In other words, our hypothesis will output the estimated probability that  $p(y = i|x; \theta)$ , for every value of  $i = 1, \dots, k$ . (Even though  $h_\theta(x)$  as defined above is only  $k-1$  dimensional, clearly  $p(y = k|x; \theta)$  can be obtained as  $1 - \sum_{i=1}^{k-1} \phi_i$ .)



Lastly, let's discuss parameter fitting. Similar to our original derivation of ordinary least squares and logistic regression, if we have a training set of  $n$  examples  $\{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$  and would like to learn the parameters  $\theta_i$  of this model, we would begin by writing down the log-likelihood

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^n \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^n \log \prod_{l=1}^k \left( \frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \right)^{1_{\{y^{(i)}=l\}}}\end{aligned}$$

To obtain the second line above, we used the definition for  $p(y|x; \theta)$  given in Equation (3.3). We can now obtain the maximum likelihood estimate of the parameters by maximizing  $\ell(\theta)$  in terms of  $\theta$ , using a method such as gradient ascent or Newton's method.

## Chapter 4

# Generative Learning algorithms

So far, we've mainly been talking about learning algorithms that model  $p(y|x; \theta)$ , the conditional distribution of  $y$  given  $x$ . For instance, logistic regression modeled  $p(y|x; \theta)$  as  $h_\theta(x) = g(\theta^T x)$  where  $g$  is the sigmoid function. In these notes, we'll talk about a different type of learning algorithm.

Consider a classification problem in which we want to learn to distinguish between elephants ( $y = 1$ ) and dogs ( $y = 0$ ), based on some features of an animal. Given a training set, an algorithm like logistic regression or the perceptron algorithm (basically) tries to find a straight line—that is, a decision boundary—that separates the elephants and dogs. Then, to classify a new animal as either an elephant or a dog, it checks on which side of the decision boundary it falls, and makes its prediction accordingly.

Here's a different approach. First, looking at elephants, we can build a model of what elephants look like. Then, looking at dogs, we can build a separate model of what dogs look like. Finally, to classify a new animal, we can match the new animal against the elephant model, and match it against the dog model, to see whether the new animal looks more like the elephants or more like the dogs we had seen in the training set.

Algorithms that try to learn  $p(y|x)$  directly (such as logistic regression), or algorithms that try to learn mappings directly from the space of inputs  $\mathcal{X}$  to the labels  $\{0, 1\}$ , (such as the perceptron algorithm) are called **discriminative** learning algorithms. Here, we'll talk about algorithms that instead try to model  $p(x|y)$  (and  $p(y)$ ). These algorithms are called **generative** learning algorithms. For instance, if  $y$  indicates whether an example is a dog (0) or an elephant (1), then  $p(x|y = 0)$  models the distribution of dogs' features, and  $p(x|y = 1)$  models the distribution of elephants' features.

After modeling  $p(y)$  (called the **class priors**) and  $p(x|y)$ , our algorithm

can then use Bayes rule to derive the posterior distribution on  $y$  given  $x$ :

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}.$$

Here, the denominator is given by  $p(x) = p(x|y = 1)p(y = 1) + p(x|y = 0)p(y = 0)$  (you should be able to verify that this is true from the standard properties of probabilities), and thus can also be expressed in terms of the quantities  $p(x|y)$  and  $p(y)$  that we've learned. Actually, if we were calculating  $p(y|x)$  in order to make a prediction, then we don't actually need to calculate the denominator, since

$$\begin{aligned} \arg \max_y p(y|x) &= \arg \max_y \frac{p(x|y)p(y)}{p(x)} \\ &= \arg \max_y p(x|y)p(y). \end{aligned}$$

## 4.1 Gaussian discriminant analysis

The first generative learning algorithm that we'll look at is Gaussian discriminant analysis (GDA). In this model, we'll assume that  $p(x|y)$  is distributed according to a multivariate normal distribution. Let's talk briefly about the properties of multivariate normal distributions before moving on to the GDA model itself.

### 4.1.1 The multivariate normal distribution

The multivariate normal distribution in  $d$ -dimensions, also called the multivariate Gaussian distribution, is parameterized by a **mean vector**  $\mu \in \mathbb{R}^d$  and a **covariance matrix**  $\Sigma \in \mathbb{R}^{d \times d}$ , where  $\Sigma \geq 0$  is symmetric and positive semi-definite. Also written " $\mathcal{N}(\mu, \Sigma)$ ", its density is given by:

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right).$$

In the equation above, " $|\Sigma|$ " denotes the determinant of the matrix  $\Sigma$ .

For a random variable  $X$  distributed  $\mathcal{N}(\mu, \Sigma)$ , the mean is (unsurprisingly) given by  $\mu$ :

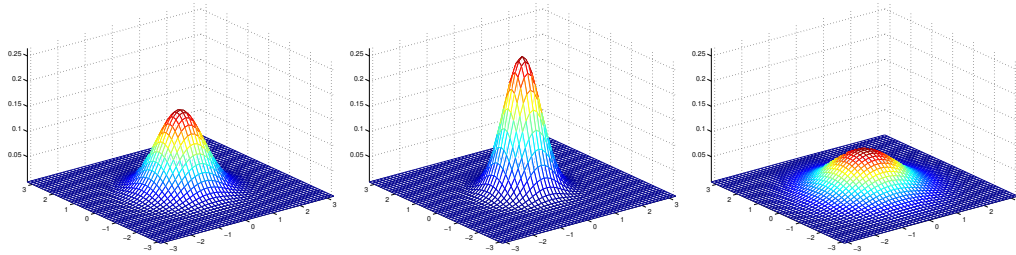
$$\mathbb{E}[X] = \int_x x p(x; \mu, \Sigma) dx = \mu$$

The **covariance** of a vector-valued random variable  $Z$  is defined as  $\text{Cov}(Z) = \mathbb{E}[(Z - \mathbb{E}[Z])(Z - \mathbb{E}[Z])^T]$ . This generalizes the notion of the variance of a

real-valued random variable. The covariance can also be defined as  $\text{Cov}(Z) = E[ZZ^T] - (E[Z])(E[Z])^T$ . (You should be able to prove to yourself that these two definitions are equivalent.) If  $X \sim \mathcal{N}(\mu, \Sigma)$ , then

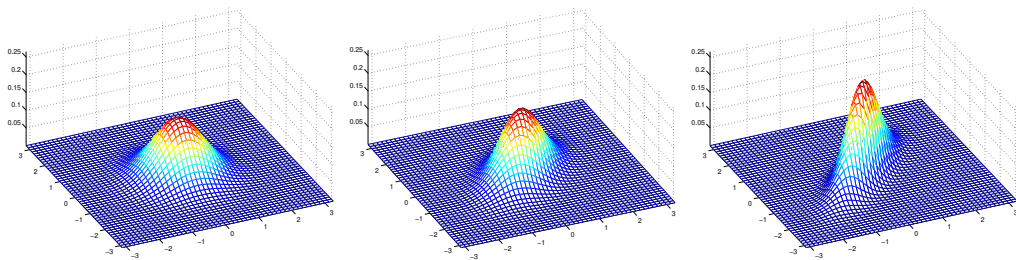
$$\text{Cov}(X) = \Sigma.$$

Here are some examples of what the density of a Gaussian distribution looks like:



The left-most figure shows a Gaussian with mean zero (that is, the 2x1 zero-vector) and covariance matrix  $\Sigma = I$  (the 2x2 identity matrix). A Gaussian with zero mean and identity covariance is also called the **standard normal distribution**. The middle figure shows the density of a Gaussian with zero mean and  $\Sigma = 0.6I$ ; and in the rightmost figure shows one with  $\Sigma = 2I$ . We see that as  $\Sigma$  becomes larger, the Gaussian becomes more “spread-out,” and as it becomes smaller, the distribution becomes more “compressed.”

Let’s look at some more examples.

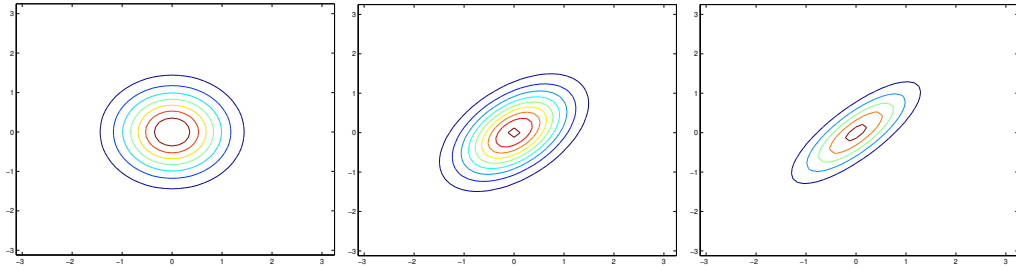


The figures above show Gaussians with mean 0, and with covariance matrices respectively

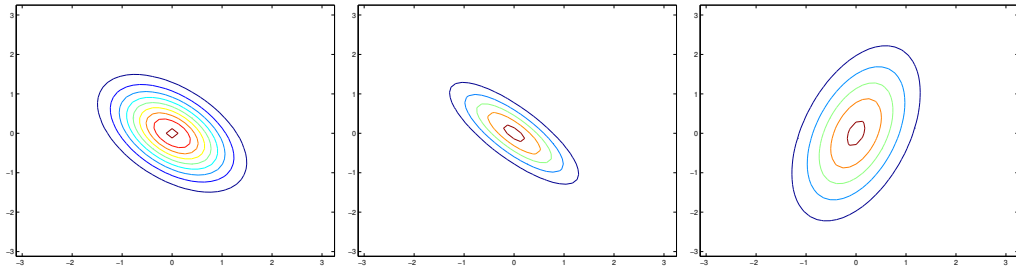
$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}; \quad \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}; \quad \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}.$$

The leftmost figure shows the familiar standard normal distribution, and we see that as we increase the off-diagonal entry in  $\Sigma$ , the density becomes more

“compressed” towards the  $45^\circ$  line (given by  $x_1 = x_2$ ). We can see this more clearly when we look at the contours of the same three densities:



Here’s one last set of examples generated by varying  $\Sigma$ :

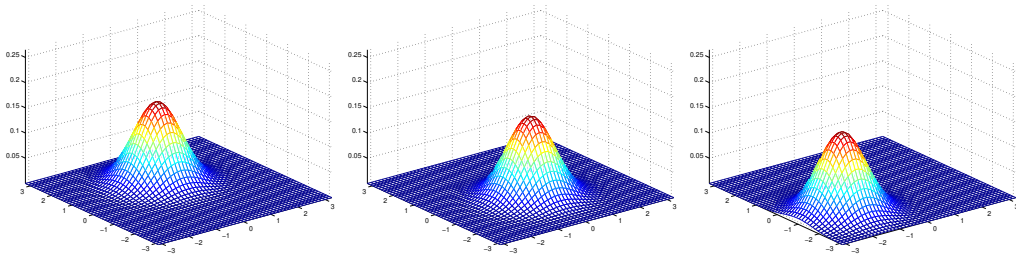


The plots above used, respectively,

$$\Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}; \quad \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}; \quad \Sigma = \begin{bmatrix} 3 & 0.8 \\ 0.8 & 1 \end{bmatrix}.$$

From the leftmost and middle figures, we see that by decreasing the off-diagonal elements of the covariance matrix, the density now becomes “compressed” again, but in the opposite direction. Lastly, as we vary the parameters, more generally the contours will form ellipses (the rightmost figure showing an example).

As our last set of examples, fixing  $\Sigma = I$ , by varying  $\mu$ , we can also move the mean of the density around.



The figures above were generated using  $\Sigma = I$ , and respectively

$$\mu = \begin{bmatrix} 1 \\ 0 \end{bmatrix}; \quad \mu = \begin{bmatrix} -0.5 \\ 0 \end{bmatrix}; \quad \mu = \begin{bmatrix} -1 \\ -1.5 \end{bmatrix}.$$

### 4.1.2 The Gaussian Discriminant Analysis model

When we have a classification problem in which the input features  $x$  are continuous-valued random variables, we can then use the Gaussian Discriminant Analysis (GDA) model, which models  $p(x|y)$  using a multivariate normal distribution. The model is:

$$\begin{aligned} y &\sim \text{Bernoulli}(\phi) \\ x|y=0 &\sim \mathcal{N}(\mu_0, \Sigma) \\ x|y=1 &\sim \mathcal{N}(\mu_1, \Sigma) \end{aligned}$$

Writing out the distributions, this is:

$$\begin{aligned} p(y) &= \phi^y(1-\phi)^{1-y} \\ p(x|y=0) &= \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu_0)^T\Sigma^{-1}(x-\mu_0)\right) \\ p(x|y=1) &= \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu_1)^T\Sigma^{-1}(x-\mu_1)\right) \end{aligned}$$

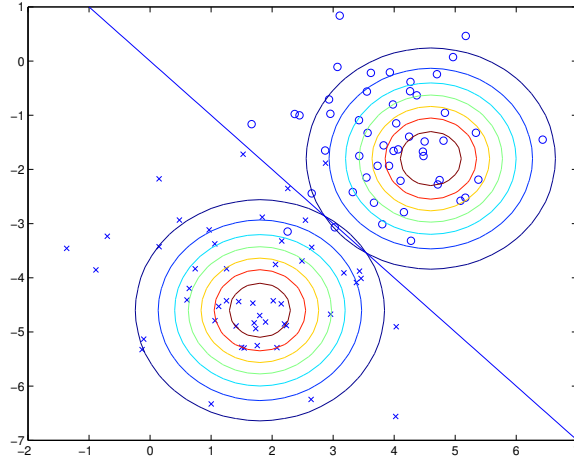
Here, the parameters of our model are  $\phi$ ,  $\Sigma$ ,  $\mu_0$  and  $\mu_1$ . (Note that while there're two different mean vectors  $\mu_0$  and  $\mu_1$ , this model is usually applied using only one covariance matrix  $\Sigma$ .) The log-likelihood of the data is given by

$$\begin{aligned} \ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^n p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi). \end{aligned}$$

By maximizing  $\ell$  with respect to the parameters, we find the maximum likelihood estimate of the parameters (see problem set 1) to be:

$$\begin{aligned}\phi &= \frac{1}{n} \sum_{i=1}^n 1\{y^{(i)} = 1\} \\ \mu_0 &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}} \\ \mu_1 &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}} \\ \Sigma &= \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T.\end{aligned}$$

Pictorially, what the algorithm is doing can be seen in as follows:



Shown in the figure are the training set, as well as the contours of the two Gaussian distributions that have been fit to the data in each of the two classes. Note that the two Gaussians have contours that are the same shape and orientation, since they share a covariance matrix  $\Sigma$ , but they have different means  $\mu_0$  and  $\mu_1$ . Also shown in the figure is the straight line giving the decision boundary at which  $p(y = 1|x) = 0.5$ . On one side of the boundary, we'll predict  $y = 1$  to be the most likely outcome, and on the other side, we'll predict  $y = 0$ .

### 4.1.3 Discussion: GDA and logistic regression

The GDA model has an interesting relationship to logistic regression. If we view the quantity  $p(y = 1|x; \phi, \mu_0, \mu_1, \Sigma)$  as a function of  $x$ , we'll find that it can be expressed in the form

$$p(y = 1|x; \phi, \Sigma, \mu_0, \mu_1) = \frac{1}{1 + \exp(-\theta^T x)},$$

where  $\theta$  is some appropriate function of  $\phi, \Sigma, \mu_0, \mu_1$ .<sup>1</sup> This is exactly the form that logistic regression—a discriminative algorithm—used to model  $p(y = 1|x)$ .

When would we prefer one model over another? GDA and logistic regression will, in general, give different decision boundaries when trained on the same dataset. Which is better?

We just argued that if  $p(x|y)$  is multivariate gaussian (with shared  $\Sigma$ ), then  $p(y|x)$  necessarily follows a logistic function. The converse, however, is not true; i.e.,  $p(y|x)$  being a logistic function does not imply  $p(x|y)$  is multivariate gaussian. This shows that GDA makes *stronger* modeling assumptions about the data than does logistic regression. It turns out that when these modeling assumptions are correct, then GDA will find better fits to the data, and is a better model. Specifically, when  $p(x|y)$  is indeed gaussian (with shared  $\Sigma$ ), then GDA is **asymptotically efficient**. Informally, this means that in the limit of very large training sets (large  $n$ ), there is no algorithm that is strictly better than GDA (in terms of, say, how accurately they estimate  $p(y|x)$ ). In particular, it can be shown that in this setting, GDA will be a better algorithm than logistic regression; and more generally, even for small training set sizes, we would generally expect GDA to better.

In contrast, by making significantly weaker assumptions, logistic regression is also more *robust* and less sensitive to incorrect modeling assumptions. There are many different sets of assumptions that would lead to  $p(y|x)$  taking the form of a logistic function. For example, if  $x|y = 0 \sim \text{Poisson}(\lambda_0)$ , and  $x|y = 1 \sim \text{Poisson}(\lambda_1)$ , then  $p(y|x)$  will be logistic. Logistic regression will also work well on Poisson data like this. But if we were to use GDA on such data—and fit Gaussian distributions to such non-Gaussian data—then the results will be less predictable, and GDA may (or may not) do well.

To summarize: GDA makes stronger modeling assumptions, and is more data efficient (i.e., requires less training data to learn “well”) when the modeling assumptions are correct or at least approximately correct. Logistic

---

<sup>1</sup>This uses the convention of redefining the  $x^{(i)}$ 's on the right-hand-side to be  $(d + 1)$ -dimensional vectors by adding the extra coordinate  $x_0^{(i)} = 1$ ; see problem set 1.



regression makes weaker assumptions, and is significantly more robust to deviations from modeling assumptions. Specifically, when the data is indeed non-Gaussian, then in the limit of large datasets, logistic regression will almost always do better than GDA. For this reason, in practice logistic regression is used more often than GDA. (Some related considerations about discriminative vs. generative models also apply for the Naive Bayes algorithm that we discuss next, but the Naive Bayes algorithm is still considered a very good, and is certainly also a very popular, classification algorithm.)

## 4.2 Naive Bayes

In GDA, the feature vectors  $x$  were continuous, real-valued vectors. Let's now talk about a different learning algorithm in which the  $x_j$ 's are discrete-valued.

For our motivating example, consider building an email spam filter using machine learning. Here, we wish to classify messages according to whether they are unsolicited commercial (spam) email, or non-spam email. After learning to do this, we can then have our mail reader automatically filter out the spam messages and perhaps place them in a separate mail folder. Classifying emails is one example of a broader set of problems called **text classification**.

Let's say we have a training set (a set of emails labeled as spam or non-spam). We'll begin our construction of our spam filter by specifying the features  $x_j$  used to represent an email.

We will represent an email via a feature vector whose length is equal to the number of words in the dictionary. Specifically, if an email contains the  $j$ -th word of the dictionary, then we will set  $x_j = 1$ ; otherwise, we let  $x_j = 0$ . For instance, the vector

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{array}{l} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

is used to represent an email that contains the words “a” and “buy,” but not

“aardvark,” “aardwolf” or “zygmurgy.”<sup>2</sup> The set of words encoded into the feature vector is called the **vocabulary**, so the dimension of  $x$  is equal to the size of the vocabulary.

Having chosen our feature vector, we now want to build a generative model. So, we have to model  $p(x|y)$ . But if we have, say, a vocabulary of 50000 words, then  $x \in \{0, 1\}^{50000}$  ( $x$  is a 50000-dimensional vector of 0’s and 1’s), and if we were to model  $x$  explicitly with a multinomial distribution over the  $2^{50000}$  possible outcomes, then we’d end up with a  $(2^{50000} - 1)$ -dimensional parameter vector. This is clearly too many parameters.

To model  $p(x|y)$ , we will therefore make a very strong assumption. We will assume that the  $x_i$ ’s are conditionally independent given  $y$ . This assumption is called the **Naive Bayes (NB) assumption**, and the resulting algorithm is called the **Naive Bayes classifier**. For instance, if  $y = 1$  means spam email; “buy” is word 2087 and “price” is word 39831; then we are assuming that if I tell you  $y = 1$  (that a particular piece of email is spam), then knowledge of  $x_{2087}$  (knowledge of whether “buy” appears in the message) will have no effect on your beliefs about the value of  $x_{39831}$  (whether “price” appears). More formally, this can be written  $p(x_{2087}|y) = p(x_{2087}|y, x_{39831})$ . (Note that this is *not* the same as saying that  $x_{2087}$  and  $x_{39831}$  are independent, which would have been written “ $p(x_{2087}) = p(x_{2087}|x_{39831})$ ”; rather, we are only assuming that  $x_{2087}$  and  $x_{39831}$  are conditionally independent *given*  $y$ .)

We now have:

$$\begin{aligned} p(x_1, \dots, x_{50000}|y) &= p(x_1|y)p(x_2|y, x_1)p(x_3|y, x_1, x_2) \cdots p(x_{50000}|y, x_1, \dots, x_{49999}) \\ &= p(x_1|y)p(x_2|y)p(x_3|y) \cdots p(x_{50000}|y) \\ &= \prod_{j=1}^d p(x_j|y) \end{aligned}$$

The first equality simply follows from the usual properties of probabilities, and the second equality used the NB assumption. We note that even though

---

<sup>2</sup>Actually, rather than looking through an English dictionary for the list of all English words, in practice it is more common to look through our training set and encode in our feature vector only the words that occur at least once there. Apart from reducing the number of words modeled and hence reducing our computational and space requirements, this also has the advantage of allowing us to model/include as a feature many words that may appear in your email (such as “cs229”) but that you won’t find in a dictionary. Sometimes (as in the homework), we also exclude the very high frequency words (which will be words like “the,” “of,” “and”; these high frequency, “content free” words are called **stop words**) since they occur in so many documents and do little to indicate whether an email is spam or non-spam.

the Naive Bayes assumption is an extremely strong assumptions, the resulting algorithm works well on many problems.

Our model is parameterized by  $\phi_{j|y=1} = p(x_j = 1|y = 1)$ ,  $\phi_{j|y=0} = p(x_j = 1|y = 0)$ , and  $\phi_y = p(y = 1)$ . As usual, given a training set  $\{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$ , we can write down the joint likelihood of the data:

$$\mathcal{L}(\phi_y, \phi_{j|y=0}, \phi_{j|y=1}) = \prod_{i=1}^n p(x^{(i)}, y^{(i)}).$$

Maximizing this with respect to  $\phi_y, \phi_{j|y=0}$  and  $\phi_{j|y=1}$  gives the maximum likelihood estimates:

$$\begin{aligned}\phi_{j|y=1} &= \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}} \\ \phi_{j|y=0} &= \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}} \\ \phi_y &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\}}{n}\end{aligned}$$

In the equations above, the “ $\wedge$ ” symbol means “and.” The parameters have a very natural interpretation. For instance,  $\phi_{j|y=1}$  is just the fraction of the spam ( $y = 1$ ) emails in which word  $j$  does appear.

Having fit all these parameters, to make a prediction on a new example with features  $x$ , we then simply calculate

$$\begin{aligned}p(y = 1|x) &= \frac{p(x|y = 1)p(y = 1)}{p(x)} \\ &= \frac{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1)}{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1) + \left(\prod_{j=1}^d p(x_j|y = 0)\right) p(y = 0)},\end{aligned}$$

and pick whichever class has the higher posterior probability.

Lastly, we note that while we have developed the Naive Bayes algorithm mainly for the case of problems where the features  $x_j$  are binary-valued, the generalization to where  $x_j$  can take values in  $\{1, 2, \dots, k_j\}$  is straightforward. Here, we would simply model  $p(x_j|y)$  as multinomial rather than as Bernoulli. Indeed, even if some original input attribute (say, the living area of a house, as in our earlier example) were continuous valued, it is quite common to **discretize** it—that is, turn it into a small set of discrete values—and apply Naive Bayes. For instance, if we use some feature  $x_j$  to represent living area, we might discretize the continuous values as follows:

Living area (sq. feet)	< 400	400-800	800-1200	1200-1600	>1600
$x_i$	1	2	3	4	5

Thus, for a house with living area 890 square feet, we would set the value of the corresponding feature  $x_j$  to 3. We can then apply the Naive Bayes algorithm, and model  $p(x_j|y)$  with a multinomial distribution, as described previously. When the original, continuous-valued attributes are not well-modeled by a multivariate normal distribution, discretizing the features and using Naive Bayes (instead of GDA) will often result in a better classifier.

### 4.2.1 Laplace smoothing

The Naive Bayes algorithm as we have described it will work fairly well for many problems, but there is a simple change that makes it work much better, especially for text classification. Let's briefly discuss a problem with the algorithm in its current form, and then talk about how we can fix it.

Consider spam/email classification, and let's suppose that, we are in the year of 20xx, after completing CS229 and having done excellent work on the project, you decide around May 20xx to submit work you did to the NeurIPS conference for publication.<sup>3</sup> Because you end up discussing the conference in your emails, you also start getting messages with the word "neurips" in it. But this is your first NeurIPS paper, and until this time, you had not previously seen any emails containing the word "neurips"; in particular "neurips" did not ever appear in your training set of spam/non-spam emails. Assuming that "neurips" was the 35000th word in the dictionary, your Naive Bayes spam filter therefore had picked its maximum likelihood estimates of the parameters  $\phi_{35000|y}$  to be

$$\begin{aligned}\phi_{35000|y=1} &= \frac{\sum_{i=1}^n 1\{x_{35000}^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}} = 0 \\ \phi_{35000|y=0} &= \frac{\sum_{i=1}^n 1\{x_{35000}^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}} = 0\end{aligned}$$

I.e., because it has never seen "neurips" before in either spam or non-spam training examples, it thinks the probability of seeing it in either type of email is zero. Hence, when trying to decide if one of these messages containing

---

<sup>3</sup>NeurIPS is one of the top machine learning conferences. The deadline for submitting a paper is typically in May-June.

“neurips” is spam, it calculates the class posterior probabilities, and obtains

$$\begin{aligned} p(y = 1|x) &= \frac{\prod_{j=1}^d p(x_j|y = 1)p(y = 1)}{\prod_{j=1}^d p(x_j|y = 1)p(y = 1) + \prod_{j=1}^d p(x_j|y = 0)p(y = 0)} \\ &= \frac{0}{0}. \end{aligned}$$

This is because each of the terms “ $\prod_{j=1}^d p(x_j|y)$ ” includes a term  $p(x_{35000}|y) = 0$  that is multiplied into it. Hence, our algorithm obtains 0/0, and doesn’t know how to make a prediction.

Stating the problem more broadly, it is statistically a bad idea to estimate the probability of some event to be zero just because you haven’t seen it before in your finite training set. Take the problem of estimating the mean of a multinomial random variable  $z$  taking values in  $\{1, \dots, k\}$ . We can parameterize our multinomial with  $\phi_j = p(z = j)$ . Given a set of  $n$  independent observations  $\{z^{(1)}, \dots, z^{(n)}\}$ , the maximum likelihood estimates are given by

$$\phi_j = \frac{\sum_{i=1}^n 1\{z^{(i)} = j\}}{n}.$$

As we saw previously, if we were to use these maximum likelihood estimates, then some of the  $\phi_j$ ’s might end up as zero, which was a problem. To avoid this, we can use **Laplace smoothing**, which replaces the above estimate with

$$\phi_j = \frac{1 + \sum_{i=1}^n 1\{z^{(i)} = j\}}{k + n}.$$

Here, we’ve added 1 to the numerator, and  $k$  to the denominator. Note that  $\sum_{j=1}^k \phi_j = 1$  still holds (check this yourself!), which is a desirable property since the  $\phi_j$ ’s are estimates for probabilities that we know must sum to 1. Also,  $\phi_j \neq 0$  for all values of  $j$ , solving our problem of probabilities being estimated as zero. Under certain (arguably quite strong) conditions, it can be shown that the Laplace smoothing actually gives the optimal estimator of the  $\phi_j$ ’s.

Returning to our Naive Bayes classifier, with Laplace smoothing, we therefore obtain the following estimates of the parameters:

$$\begin{aligned} \phi_{j|y=1} &= \frac{1 + \sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{2 + \sum_{i=1}^n 1\{y^{(i)} = 1\}} \\ \phi_{j|y=0} &= \frac{1 + \sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{2 + \sum_{i=1}^n 1\{y^{(i)} = 0\}} \end{aligned}$$

(In practice, it usually doesn't matter much whether we apply Laplace smoothing to  $\phi_y$  or not, since we will typically have a fair fraction each of spam and non-spam messages, so  $\phi_y$  will be a reasonable estimate of  $p(y = 1)$  and will be quite far from 0 anyway.)

## 4.2.2 Event models for text classification

To close off our discussion of generative learning algorithms, let's talk about one more model that is specifically for text classification. While Naive Bayes as we've presented it will work well for many classification problems, for text classification, there is a related model that does even better.

In the specific context of text classification, Naive Bayes as presented uses the what's called the **Bernoulli event model** (or sometimes **multi-variate Bernoulli event model**). In this model, we assumed that the way an email is generated is that first it is randomly determined (according to the class priors  $p(y)$ ) whether a spammer or non-spammer will send you your next message. Then, the person sending the email runs through the dictionary, deciding whether to include each word  $j$  in that email independently and according to the probabilities  $p(x_j = 1|y) = \phi_{j|y}$ . Thus, the probability of a message was given by  $p(y) \prod_{j=1}^d p(x_j|y)$ .

Here's a different model, called the **Multinomial event model**. To describe this model, we will use a different notation and set of features for representing emails. We let  $x_j$  denote the identity of the  $j$ -th word in the email. Thus,  $x_j$  is now an integer taking values in  $\{1, \dots, |V|\}$ , where  $|V|$  is the size of our vocabulary (dictionary). An email of  $d$  words is now represented by a vector  $(x_1, x_2, \dots, x_d)$  of length  $d$ ; note that  $d$  can vary for different documents. For instance, if an email starts with "A NeurIPS . . .," then  $x_1 = 1$  ("a" is the first word in the dictionary), and  $x_2 = 35000$  (if "neurips" is the 35000th word in the dictionary).

In the multinomial event model, we assume that the way an email is generated is via a random process in which spam/non-spam is first determined (according to  $p(y)$ ) as before. Then, the sender of the email writes the email by first generating  $x_1$  from some multinomial distribution over words ( $p(x_1|y)$ ). Next, the second word  $x_2$  is chosen independently of  $x_1$  but from the same multinomial distribution, and similarly for  $x_3, x_4$ , and so on, until all  $d$  words of the email have been generated. Thus, the overall probability of a message is given by  $p(y) \prod_{j=1}^d p(x_j|y)$ . Note that this formula looks like the one we had earlier for the probability of a message under the Bernoulli event model, but that the terms in the formula now mean very different things. In particular  $x_j|y$  is now a multinomial, rather than a Bernoulli distribution.

The parameters for our new model are  $\phi_y = p(y)$  as before,  $\phi_{k|y=1} = p(x_j = k|y = 1)$  (for any  $j$ ) and  $\phi_{k|y=0} = p(x_j = k|y = 0)$ . Note that we have assumed that  $p(x_j|y)$  is the same for all values of  $j$  (i.e., that the distribution according to which a word is generated does not depend on its position  $j$  within the email).

If we are given a training set  $\{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$  where  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_{d_i}^{(i)})$  (here,  $d_i$  is the number of words in the  $i$ -training example), the likelihood of the data is given by

$$\begin{aligned} \mathcal{L}(\phi_y, \phi_{k|y=0}, \phi_{k|y=1}) &= \prod_{i=1}^n p(x^{(i)}, y^{(i)}) \\ &= \prod_{i=1}^n \left( \prod_{j=1}^{d_i} p(x_j^{(i)}|y; \phi_{k|y=0}, \phi_{k|y=1}) \right) p(y^{(i)}; \phi_y). \end{aligned}$$

Maximizing this yields the maximum likelihood estimates of the parameters:

$$\begin{aligned} \phi_{k|y=1} &= \frac{\sum_{i=1}^n \sum_{j=1}^{d_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 1\}}{\sum_{i=1}^n 1\{y^{(i)} = 1\} d_i} \\ \phi_{k|y=0} &= \frac{\sum_{i=1}^n \sum_{j=1}^{d_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 0\}}{\sum_{i=1}^n 1\{y^{(i)} = 0\} d_i} \\ \phi_y &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\}}{n}. \end{aligned}$$

If we were to apply Laplace smoothing (which is needed in practice for good performance) when estimating  $\phi_{k|y=0}$  and  $\phi_{k|y=1}$ , we add 1 to the numerators and  $|V|$  to the denominators, and obtain:

$$\begin{aligned} \phi_{k|y=1} &= \frac{1 + \sum_{i=1}^n \sum_{j=1}^{d_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 1\}}{|V| + \sum_{i=1}^n 1\{y^{(i)} = 1\} d_i} \\ \phi_{k|y=0} &= \frac{1 + \sum_{i=1}^n \sum_{j=1}^{d_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 0\}}{|V| + \sum_{i=1}^n 1\{y^{(i)} = 0\} d_i}. \end{aligned}$$

While not necessarily the very best classification algorithm, the Naive Bayes classifier often works surprisingly well. It is often also a very good “first thing to try,” given its simplicity and ease of implementation.

# Chapter 5

## Kernel Methods

### 5.1 Feature maps

Recall that in our discussion about linear regression, we considered the problem of predicting the price of a house (denoted by  $y$ ) from the living area of the house (denoted by  $x$ ), and we fit a linear function of  $x$  to the training data. What if the price  $y$  can be more accurately represented as a *non-linear* function of  $x$ ? In this case, we need a more expressive family of models than linear models.

We start by considering fitting cubic functions  $y = \theta_3 x^3 + \theta_2 x^2 + \theta_1 x + \theta_0$ . It turns out that we can view the cubic function as a linear function over the a different set of feature variables (defined below). Concretely, let the function  $\phi : \mathbb{R} \rightarrow \mathbb{R}^4$  be defined as

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} \in \mathbb{R}^4. \quad (5.1)$$

Let  $\theta \in \mathbb{R}^4$  be the vector containing  $\theta_0, \theta_1, \theta_2, \theta_3$  as entries. Then we can rewrite the cubic function in  $x$  as:

$$\theta_3 x^3 + \theta_2 x^2 + \theta_1 x + \theta_0 = \theta^T \phi(x)$$

Thus, a cubic function of the variable  $x$  can be viewed as a linear function over the variables  $\phi(x)$ . To distinguish between these two sets of variables, in the context of kernel methods, we will call the “original” input value the input **attributes** of a problem (in this case,  $x$ , the living area). When the



original input is mapped to some new set of quantities  $\phi(x)$ , we will call those new quantities the **features** variables. (Unfortunately, different authors use different terms to describe these two things in different contexts.) We will call  $\phi$  a **feature map**, which maps the attributes to the features.

## 5.2 LMS (least mean squares) with features

We will derive the gradient descent algorithm for fitting the model  $\theta^T \phi(x)$ . First recall that for ordinary least square problem where we were to fit  $\theta^T x$ , the batch gradient descent update is (see the first lecture note for its derivation):

$$\begin{aligned}\theta &:= \theta + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x^{(i)} \\ &:= \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)}) x^{(i)}.\end{aligned}\tag{5.2}$$

Let  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$  be a feature map that maps attribute  $x$  (in  $\mathbb{R}^d$ ) to the features  $\phi(x)$  in  $\mathbb{R}^p$ . (In the motivating example in the previous subsection, we have  $d = 1$  and  $p = 4$ .) Now our goal is to fit the function  $\theta^T \phi(x)$ , with  $\theta$  being a vector in  $\mathbb{R}^p$  instead of  $\mathbb{R}^d$ . We can replace all the occurrences of  $x^{(i)}$  in the algorithm above by  $\phi(x^{(i)})$  to obtain the new update:

$$\theta := \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)})\tag{5.3}$$

Similarly, the corresponding stochastic gradient descent update rule is

$$\theta := \theta + \alpha (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)})\tag{5.4}$$

## 5.3 LMS with the kernel trick

The gradient descent update, or stochastic gradient update above becomes computationally expensive when the features  $\phi(x)$  is high-dimensional. For example, consider the direct extension of the feature map in equation (5.1) to high-dimensional input  $x$ : suppose  $x \in \mathbb{R}^d$ , and let  $\phi(x)$  be the vector that

contains all the monomials of  $x$  with degree  $\leq 3$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_1^2 \\ x_1x_2 \\ x_1x_3 \\ \vdots \\ x_2x_1 \\ \vdots \\ x_1^3 \\ x_1^2x_2 \\ \vdots \end{bmatrix}. \quad (5.5)$$

The dimension of the features  $\phi(x)$  is on the order of  $d^3$ .<sup>1</sup> This is a prohibitively long vector for computational purpose — when  $d = 1000$ , each update requires at least computing and storing a  $1000^3 = 10^9$  dimensional vector, which is  $10^6$  times slower than the update rule for ordinary least squares updates (5.2).

It may appear at first that such  $d^3$  runtime per update and memory usage are inevitable, because the vector  $\theta$  itself is of dimension  $p \approx d^3$ , and we may need to update every entry of  $\theta$  and store it. However, we will introduce the kernel trick with which we will not need to store  $\theta$  explicitly, and the runtime can be significantly improved.

For simplicity, we assume the initialize the value  $\theta = 0$ , and we focus on the iterative update (5.3). The main observation is that at any time,  $\theta$  can be represented as a linear combination of the vectors  $\phi(x^{(1)}), \dots, \phi(x^{(n)})$ . Indeed, we can show this inductively as follows. At initialization,  $\theta = 0 = \sum_{i=1}^n 0 \cdot \phi(x^{(i)})$ . Assume at some point,  $\theta$  can be represented as

$$\theta = \sum_{i=1}^n \beta_i \phi(x^{(i)}) \quad (5.6)$$

---

<sup>1</sup>Here, for simplicity, we include all the monomials with repetitions (so that, e.g.,  $x_1x_2x_3$  and  $x_2x_3x_1$  both appear in  $\phi(x)$ ). Therefore, there are totally  $1 + d + d^2 + d^3$  entries in  $\phi(x)$ .

for some  $\beta_1, \dots, \beta_n \in \mathbb{R}$ . Then we claim that in the next round,  $\theta$  is still a linear combination of  $\phi(x^{(1)}), \dots, \phi(x^{(n)})$  because

$$\begin{aligned}
 \theta &:= \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)}) \\
 &= \sum_{i=1}^n \beta_i \phi(x^{(i)}) + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)}) \\
 &= \sum_{i=1}^n \underbrace{(\beta_i + \alpha (y^{(i)} - \theta^T \phi(x^{(i)})))}_{\text{new } \beta_i} \phi(x^{(i)})
 \end{aligned} \tag{5.7}$$

You may realize that our general strategy is to implicitly represent the  $p$ -dimensional vector  $\theta$  by a set of coefficients  $\beta_1, \dots, \beta_n$ . Towards doing this, we derive the update rule of the coefficients  $\beta_1, \dots, \beta_n$ . Using the equation above, we see that the new  $\beta_i$  depends on the old one via

$$\beta_i := \beta_i + \alpha (y^{(i)} - \theta^T \phi(x^{(i)})) \tag{5.8}$$

Here we still have the old  $\theta$  on the RHS of the equation. Replacing  $\theta$  by  $\theta = \sum_{j=1}^n \beta_j \phi(x^{(j)})$  gives

$$\forall i \in \{1, \dots, n\}, \beta_i := \beta_i + \alpha \left( y^{(i)} - \sum_{j=1}^n \beta_j \phi(x^{(j)})^T \phi(x^{(i)}) \right)$$

We often rewrite  $\phi(x^{(j)})^T \phi(x^{(i)})$  as  $\langle \phi(x^{(j)}), \phi(x^{(i)}) \rangle$  to emphasize that it's the inner product of the two feature vectors. Viewing  $\beta_i$ 's as the new representation of  $\theta$ , we have successfully translated the batch gradient descent algorithm into an algorithm that updates the value of  $\beta$  iteratively. It may appear that at every iteration, we still need to compute the values of  $\langle \phi(x^{(j)}), \phi(x^{(i)}) \rangle$  for all pairs of  $i, j$ , each of which may take roughly  $O(p)$  operation. However, two important properties come to rescue:

1. We can pre-compute the pairwise inner products  $\langle \phi(x^{(j)}), \phi(x^{(i)}) \rangle$  for all pairs of  $i, j$  before the loop starts.
2. For the feature map  $\phi$  defined in (5.5) (or many other interesting feature maps), computing  $\langle \phi(x^{(j)}), \phi(x^{(i)}) \rangle$  can be efficient and does not

necessarily require computing  $\phi(x^{(i)})$  explicitly. This is because:

$$\begin{aligned}
\langle \phi(x), \phi(z) \rangle &= 1 + \sum_{i=1}^d x_i z_i + \sum_{i,j \in \{1, \dots, d\}} x_i x_j z_i z_j + \sum_{i,j,k \in \{1, \dots, d\}} x_i x_j x_k z_i z_j z_k \\
&= 1 + \sum_{i=1}^d x_i z_i + \left( \sum_{i=1}^d x_i z_i \right)^2 + \left( \sum_{i=1}^d x_i z_i \right)^3 \\
&= 1 + \langle x, z \rangle + \langle x, z \rangle^2 + \langle x, z \rangle^3
\end{aligned} \tag{5.9}$$

Therefore, to compute  $\langle \phi(x), \phi(z) \rangle$ , we can first compute  $\langle x, z \rangle$  with  $O(d)$  time and then take another constant number of operations to compute  $1 + \langle x, z \rangle + \langle x, z \rangle^2 + \langle x, z \rangle^3$ .

As you will see, the inner products between the features  $\langle \phi(x), \phi(z) \rangle$  are essential here. We define the **Kernel** corresponding to the feature map  $\phi$  as a function that maps  $\mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  satisfying: <sup>2</sup>

$$K(x, z) \triangleq \langle \phi(x), \phi(z) \rangle \tag{5.10}$$

To wrap up the discussion, we write the down the final algorithm as follows:

- 
1. Compute all the values  $K(x^{(i)}, x^{(j)}) \triangleq \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$  using equation (5.9) for all  $i, j \in \{1, \dots, n\}$ . Set  $\beta := 0$ .
  2. **Loop:**

$$\forall i \in \{1, \dots, n\}, \beta_i := \beta_i + \alpha \left( y^{(i)} - \sum_{j=1}^n \beta_j K(x^{(i)}, x^{(j)}) \right) \tag{5.11}$$

Or in vector notation, letting  $K$  be the  $n \times n$  matrix with  $K_{ij} = K(x^{(i)}, x^{(j)})$ , we have

$$\beta := \beta + \alpha(\vec{y} - K\beta)$$

---

With the algorithm above, we can update the representation  $\beta$  of the vector  $\theta$  efficiently with  $O(n)$  time per update. Finally, we need to show that

---

<sup>2</sup>Recall that  $\mathcal{X}$  is the space of the input  $x$ . In our running example,  $\mathcal{X} = \mathbb{R}^d$

the knowledge of the representation  $\beta$  suffices to compute the prediction  $\theta^T \phi(x)$ . Indeed, we have

$$\theta^T \phi(x) = \sum_{i=1}^n \beta_i \phi(x^{(i)})^T \phi(x) = \sum_{i=1}^n \beta_i K(x^{(i)}, x) \quad (5.12)$$

You may realize that fundamentally all we need to know about the feature map  $\phi(\cdot)$  is encapsulated in the corresponding kernel function  $K(\cdot, \cdot)$ . We will expand on this in the next section.

## 5.4 Properties of kernels

In the last subsection, we started with an explicitly defined feature map  $\phi$ , which induces the kernel function  $K(x, z) \triangleq \langle \phi(x), \phi(z) \rangle$ . Then we saw that the kernel function is so intrinsic so that as long as the kernel function is defined, the whole training algorithm can be written entirely in the language of the kernel without referring to the feature map  $\phi$ , so can the prediction of a test example  $x$  (equation (5.12).)

Therefore, it would be tempted to define other kernel function  $K(\cdot, \cdot)$  and run the algorithm (5.11). Note that the algorithm (5.11) does not need to explicitly access the feature map  $\phi$ , and therefore we only need to ensure the existence of the feature map  $\phi$ , but do not necessarily need to be able to explicitly write  $\phi$  down.

What kinds of functions  $K(\cdot, \cdot)$  can correspond to some feature map  $\phi$ ? In other words, can we tell if there is some feature mapping  $\phi$  so that  $K(x, z) = \phi(x)^T \phi(z)$  for all  $x, z$ ?

If we can answer this question by giving a precise characterization of valid kernel functions, then we can completely change the interface of selecting feature maps  $\phi$  to the interface of selecting kernel function  $K$ . Concretely, we can pick a function  $K$ , verify that it satisfies the characterization (so that there exists a feature map  $\phi$  that  $K$  corresponds to), and then we can run update rule (5.11). The benefit here is that we don't have to be able to compute  $\phi$  or write it down analytically, and we only need to know its existence. We will answer this question at the end of this subsection after we go through several concrete examples of kernels.

Suppose  $x, z \in \mathbb{R}^d$ , and let's first consider the function  $K(\cdot, \cdot)$  defined as:

$$K(x, z) = (x^T z)^2.$$

We can also write this as

$$\begin{aligned}
 K(x, z) &= \left( \sum_{i=1}^d x_i z_i \right) \left( \sum_{j=1}^d x_j z_j \right) \\
 &= \sum_{i=1}^d \sum_{j=1}^d x_i x_j z_i z_j \\
 &= \sum_{i,j=1}^d (x_i x_j) (z_i z_j)
 \end{aligned}$$

Thus, we see that  $K(x, z) = \langle \phi(x), \phi(z) \rangle$  is the kernel function that corresponds to the feature mapping  $\phi$  given (shown here for the case of  $d = 3$ ) by

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}.$$

Revisiting the computational efficiency perspective of kernel, note that whereas calculating the high-dimensional  $\phi(x)$  requires  $O(d^2)$  time, finding  $K(x, z)$  takes only  $O(d)$  time—linear in the dimension of the input attributes.

For another related example, also consider  $K(\cdot, \cdot)$  defined by

$$\begin{aligned}
 K(x, z) &= (x^T z + c)^2 \\
 &= \sum_{i,j=1}^d (x_i x_j) (z_i z_j) + \sum_{i=1}^d (\sqrt{2c} x_i) (\sqrt{2c} z_i) + c^2.
 \end{aligned}$$

(Check this yourself.) This function  $K$  is a kernel function that corresponds

to the feature mapping (again shown for  $d = 3$ )

$$\phi(x) = \begin{bmatrix} x_1x_1 \\ x_1x_2 \\ x_1x_3 \\ x_2x_1 \\ x_2x_2 \\ x_2x_3 \\ x_3x_1 \\ x_3x_2 \\ x_3x_3 \\ \sqrt{2c}x_1 \\ \sqrt{2c}x_2 \\ \sqrt{2c}x_3 \\ c \end{bmatrix},$$

and the parameter  $c$  controls the relative weighting between the  $x_i$  (first order) and the  $x_ix_j$  (second order) terms.

More broadly, the kernel  $K(x, z) = (x^T z + c)^k$  corresponds to a feature mapping to an  $\binom{d+k}{k}$  feature space, corresponding of all monomials of the form  $x_{i_1}x_{i_2}\dots x_{i_k}$  that are up to order  $k$ . However, despite working in this  $O(d^k)$ -dimensional space, computing  $K(x, z)$  still takes only  $O(d)$  time, and hence we never need to explicitly represent feature vectors in this very high dimensional feature space.

**Kernels as similarity metrics.** Now, let's talk about a slightly different view of kernels. Intuitively, (and there are things wrong with this intuition, but nevermind), if  $\phi(x)$  and  $\phi(z)$  are close together, then we might expect  $K(x, z) = \phi(x)^T \phi(z)$  to be large. Conversely, if  $\phi(x)$  and  $\phi(z)$  are far apart—say nearly orthogonal to each other—then  $K(x, z) = \phi(x)^T \phi(z)$  will be small. So, we can think of  $K(x, z)$  as some measurement of how similar are  $\phi(x)$  and  $\phi(z)$ , or of how similar are  $x$  and  $z$ .

Given this intuition, suppose that for some learning problem that you're working on, you've come up with some function  $K(x, z)$  that you think might be a reasonable measure of how similar  $x$  and  $z$  are. For instance, perhaps you chose

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right).$$

This is a reasonable measure of  $x$  and  $z$ 's similarity, and is close to 1 when  $x$  and  $z$  are close, and near 0 when  $x$  and  $z$  are far apart. Does there exist

a feature map  $\phi$  such that the kernel  $K$  defined above satisfies  $K(x, z) = \phi(x)^T \phi(z)$ ? In this particular example, the answer is yes. This kernel is called the **Gaussian kernel**, and corresponds to an infinite dimensional feature mapping  $\phi$ . We will give a precise characterization about what properties a function  $K$  needs to satisfy so that it can be a valid kernel function that corresponds to some feature map  $\phi$ .

**Necessary conditions for valid kernels.** Suppose for now that  $K$  is indeed a valid kernel corresponding to some feature mapping  $\phi$ , and we will first see what properties it satisfies. Now, consider some finite set of  $n$  points (not necessarily the training set)  $\{x^{(1)}, \dots, x^{(n)}\}$ , and let a square,  $n$ -by- $n$  matrix  $K$  be defined so that its  $(i, j)$ -entry is given by  $K_{ij} = K(x^{(i)}, x^{(j)})$ . This matrix is called the **kernel matrix**. Note that we've overloaded the notation and used  $K$  to denote both the kernel function  $K(x, z)$  and the kernel matrix  $K$ , due to their obvious close relationship.

Now, if  $K$  is a valid kernel, then  $K_{ij} = K(x^{(i)}, x^{(j)}) = \phi(x^{(i)})^T \phi(x^{(j)}) = \phi(x^{(j)})^T \phi(x^{(i)}) = K(x^{(j)}, x^{(i)}) = K_{ji}$ , and hence  $K$  must be symmetric. Moreover, letting  $\phi_k(x)$  denote the  $k$ -th coordinate of the vector  $\phi(x)$ , we find that for any vector  $z$ , we have

$$\begin{aligned}
 z^T K z &= \sum_i \sum_j z_i K_{ij} z_j \\
 &= \sum_i \sum_j z_i \phi(x^{(i)})^T \phi(x^{(j)}) z_j \\
 &= \sum_i \sum_j z_i \sum_k \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j \\
 &= \sum_k \sum_i \sum_j z_i \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j \\
 &= \sum_k \left( \sum_i z_i \phi_k(x^{(i)}) \right)^2 \\
 &\geq 0.
 \end{aligned}$$

The second-to-last step uses the fact that  $\sum_{i,j} a_i a_j = (\sum_i a_i)^2$  for  $a_i = z_i \phi_k(x^{(i)})$ . Since  $z$  was arbitrary, this shows that  $K$  is positive semi-definite ( $K \geq 0$ ).

Hence, we've shown that if  $K$  is a valid kernel (i.e., if it corresponds to some feature mapping  $\phi$ ), then the corresponding kernel matrix  $K \in \mathbb{R}^{n \times n}$  is symmetric positive semidefinite.



**Sufficient conditions for valid kernels.** More generally, the condition above turns out to be not only a necessary, but also a sufficient, condition for  $K$  to be a valid kernel (also called a Mercer kernel). The following result is due to Mercer.<sup>3</sup>

**Theorem (Mercer).** Let  $K : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$  be given. Then for  $K$  to be a valid (Mercer) kernel, it is necessary and sufficient that for any  $\{x^{(1)}, \dots, x^{(n)}\}$ ,  $(n < \infty)$ , the corresponding kernel matrix is symmetric positive semi-definite.

Given a function  $K$ , apart from trying to find a feature mapping  $\phi$  that corresponds to it, this theorem therefore gives another way of testing if it is a valid kernel. You'll also have a chance to play with these ideas more in problem set 2.

In class, we also briefly talked about a couple of other examples of kernels. For instance, consider the digit recognition problem, in which given an image (16x16 pixels) of a handwritten digit (0-9), we have to figure out which digit it was. Using either a simple polynomial kernel  $K(x, z) = (x^T z)^k$  or the Gaussian kernel, SVMs were able to obtain extremely good performance on this problem. This was particularly surprising since the input attributes  $x$  were just 256-dimensional vectors of the image pixel intensity values, and the system had no prior knowledge about vision, or even about which pixels are adjacent to which other ones. Another example that we briefly talked about in lecture was that if the objects  $x$  that we are trying to classify are strings (say,  $x$  is a list of amino acids, which strung together form a protein), then it seems hard to construct a reasonable, “small” set of features for most learning algorithms, especially if different strings have different lengths. However, consider letting  $\phi(x)$  be a feature vector that counts the number of occurrences of each length- $k$  substring in  $x$ . If we're considering strings of English letters, then there are  $26^k$  such strings. Hence,  $\phi(x)$  is a  $26^k$  dimensional vector; even for moderate values of  $k$ , this is probably too big for us to efficiently work with. (e.g.,  $26^4 \approx 460000$ .) However, using (dynamic programming-ish) string matching algorithms, it is possible to efficiently compute  $K(x, z) = \phi(x)^T \phi(z)$ , so that we can now implicitly work in this  $26^k$ -dimensional feature space, but without ever explicitly computing feature vectors in this space.

---

<sup>3</sup>Many texts present Mercer's theorem in a slightly more complicated form involving  $L^2$  functions, but when the input attributes take values in  $\mathbb{R}^d$ , the version given here is equivalent.

**Application of kernel methods:** We've seen the application of kernels to linear regression. In the next part, we will introduce the support vector machines to which kernels can be directly applied. dwell too much longer on it here. In fact, the idea of kernels has significantly broader applicability than linear regression and SVMs. Specifically, if you have any learning algorithm that you can write in terms of only inner products  $\langle x, z \rangle$  between input attribute vectors, then by replacing this with  $K(x, z)$  where  $K$  is a kernel, you can “magically” allow your algorithm to work efficiently in the high dimensional feature space corresponding to  $K$ . For instance, this kernel trick can be applied with the perceptron to derive a kernel perceptron algorithm. Many of the algorithms that we'll see later in this class will also be amenable to this method, which has come to be known as the “kernel trick.”

# Chapter 6

## Support Vector Machines

This set of notes presents the Support Vector Machine (SVM) learning algorithm. SVMs are among the best (and many believe are indeed the best) “off-the-shelf” supervised learning algorithms. To tell the SVM story, we’ll need to first talk about margins and the idea of separating data with a large “gap.” Next, we’ll talk about the optimal margin classifier, which will lead us into a digression on Lagrange duality. We’ll also see kernels, which give a way to apply SVMs efficiently in very high dimensional (such as infinite-dimensional) feature spaces, and finally, we’ll close off the story with the SMO algorithm, which gives an efficient implementation of SVMs.

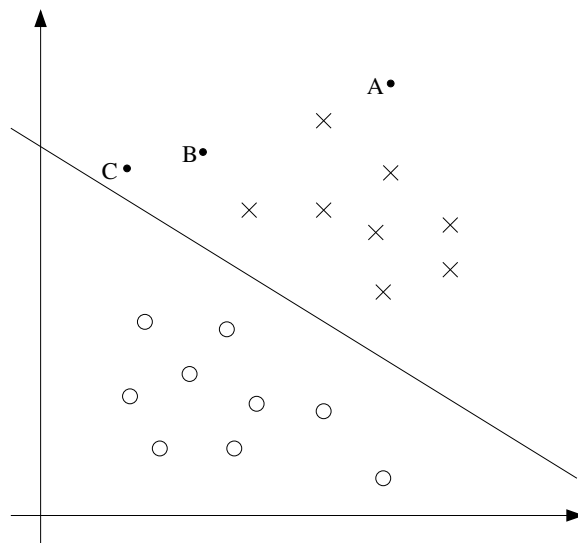
### 6.1 Margins: Intuition

We’ll start our story on SVMs by talking about margins. This section will give the intuitions about margins and about the “confidence” of our predictions; these ideas will be made formal in Section 6.3.

Consider logistic regression, where the probability  $p(y = 1|x; \theta)$  is modeled by  $h_\theta(x) = g(\theta^T x)$ . We then predict “1” on an input  $x$  if and only if  $h_\theta(x) \geq 0.5$ , or equivalently, if and only if  $\theta^T x \geq 0$ . Consider a positive training example ( $y = 1$ ). The larger  $\theta^T x$  is, the larger also is  $h_\theta(x) = p(y = 1|x; \theta)$ , and thus also the higher our degree of “confidence” that the label is 1. Thus, informally we can think of our prediction as being very confident that  $y = 1$  if  $\theta^T x \gg 0$ . Similarly, we think of logistic regression as confidently predicting  $y = 0$ , if  $\theta^T x \ll 0$ . Given a training set, again informally it seems that we’d have found a good fit to the training data if we can find  $\theta$  so that  $\theta^T x^{(i)} \gg 0$  whenever  $y^{(i)} = 1$ , and  $\theta^T x^{(i)} \ll 0$  whenever  $y^{(i)} = 0$ , since this would reflect a very confident (and correct) set of classifications for all the

training examples. This seems to be a nice goal to aim for, and we'll soon formalize this idea using the notion of functional margins.

For a different type of intuition, consider the following figure, in which x's represent positive training examples, o's denote negative training examples, a decision boundary (this is the line given by the equation  $\theta^T x = 0$ , and is also called the **separating hyperplane**) is also shown, and three points have also been labeled A, B and C.



Notice that the point A is very far from the decision boundary. If we are asked to make a prediction for the value of  $y$  at A, it seems we should be quite confident that  $y = 1$  there. Conversely, the point C is very close to the decision boundary, and while it's on the side of the decision boundary on which we would predict  $y = 1$ , it seems likely that just a small change to the decision boundary could easily have caused our prediction to be  $y = 0$ . Hence, we're much more confident about our prediction at A than at C. The point B lies in-between these two cases, and more broadly, we see that if a point is far from the separating hyperplane, then we may be significantly more confident in our predictions. Again, informally we think it would be nice if, given a training set, we manage to find a decision boundary that allows us to make all correct and confident (meaning far from the decision boundary) predictions on the training examples. We'll formalize this later using the notion of geometric margins.

## 6.2 Notation (option reading)

To make our discussion of SVMs easier, we'll first need to introduce a new notation for talking about classification. We will be considering a linear classifier for a binary classification problem with labels  $y$  and features  $x$ . From now, we'll use  $y \in \{-1, 1\}$  (instead of  $\{0, 1\}$ ) to denote the class labels. Also, rather than parameterizing our linear classifier with the vector  $\theta$ , we will use parameters  $w, b$ , and write our classifier as

$$h_{w,b}(x) = g(w^T x + b).$$

Here,  $g(z) = 1$  if  $z \geq 0$ , and  $g(z) = -1$  otherwise. This “ $w, b$ ” notation allows us to explicitly treat the intercept term  $b$  separately from the other parameters. (We also drop the convention we had previously of letting  $x_0 = 1$  be an extra coordinate in the input feature vector.) Thus,  $b$  takes the role of what was previously  $\theta_0$ , and  $w$  takes the role of  $[\theta_1 \dots \theta_d]^T$ .

Note also that, from our definition of  $g$  above, our classifier will directly predict either 1 or  $-1$  (cf. the perceptron algorithm), without first going through the intermediate step of estimating  $p(y = 1)$  (which is what logistic regression does).

## 6.3 Functional and geometric margins (option reading)

Let's formalize the notions of the functional and geometric margins. Given a training example  $(x^{(i)}, y^{(i)})$ , we define the **functional margin** of  $(w, b)$  with respect to the training example as

$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b).$$

Note that if  $y^{(i)} = 1$ , then for the functional margin to be large (i.e., for our prediction to be confident and correct), we need  $w^T x^{(i)} + b$  to be a large positive number. Conversely, if  $y^{(i)} = -1$ , then for the functional margin to be large, we need  $w^T x^{(i)} + b$  to be a large negative number. Moreover, if  $y^{(i)}(w^T x^{(i)} + b) > 0$ , then our prediction on this example is correct. (Check this yourself.) Hence, a large functional margin represents a confident and a correct prediction.

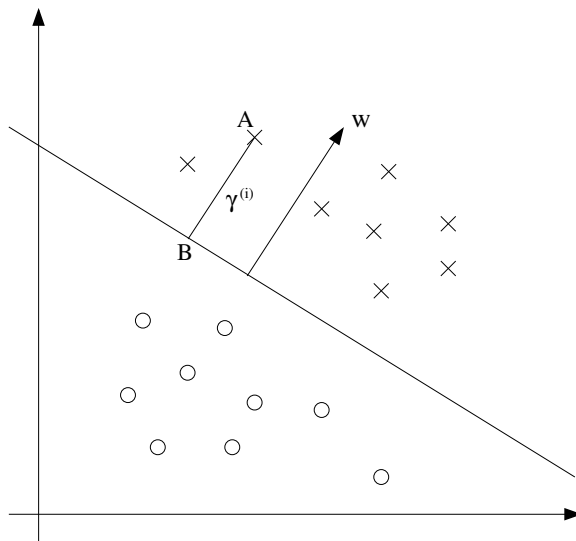
For a linear classifier with the choice of  $g$  given above (taking values in  $\{-1, 1\}$ ), there's one property of the functional margin that makes it not a very good measure of confidence, however. Given our choice of  $g$ , we note that

if we replace  $w$  with  $2w$  and  $b$  with  $2b$ , then since  $g(w^T x + b) = g(2w^T x + 2b)$ , this would not change  $h_{w,b}(x)$  at all. I.e.,  $g$ , and hence also  $h_{w,b}(x)$ , depends only on the sign, but not on the magnitude, of  $w^T x + b$ . However, replacing  $(w, b)$  with  $(2w, 2b)$  also results in multiplying our functional margin by a factor of 2. Thus, it seems that by exploiting our freedom to scale  $w$  and  $b$ , we can make the functional margin arbitrarily large without really changing anything meaningful. Intuitively, it might therefore make sense to impose some sort of normalization condition such as that  $\|w\|_2 = 1$ ; i.e., we might replace  $(w, b)$  with  $(w/\|w\|_2, b/\|w\|_2)$ , and instead consider the functional margin of  $(w/\|w\|_2, b/\|w\|_2)$ . We'll come back to this later.

Given a training set  $S = \{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$ , we also define the function margin of  $(w, b)$  with respect to  $S$  as the smallest of the functional margins of the individual training examples. Denoted by  $\hat{\gamma}$ , this can therefore be written:

$$\hat{\gamma} = \min_{i=1, \dots, n} \hat{\gamma}^{(i)}.$$

Next, let's talk about **geometric margins**. Consider the picture below:



The decision boundary corresponding to  $(w, b)$  is shown, along with the vector  $w$ . Note that  $w$  is orthogonal (at  $90^\circ$ ) to the separating hyperplane. (You should convince yourself that this must be the case.) Consider the point at A, which represents the input  $x^{(i)}$  of some training example with label  $y^{(i)} = 1$ . Its distance to the decision boundary,  $\gamma^{(i)}$ , is given by the line segment AB.

How can we find the value of  $\gamma^{(i)}$ ? Well,  $w/\|w\|$  is a unit-length vector pointing in the same direction as  $w$ . Since A represents  $x^{(i)}$ , we therefore

find that the point  $B$  is given by  $x^{(i)} - \gamma^{(i)} \cdot w / \|w\|$ . But this point lies on the decision boundary, and all points  $x$  on the decision boundary satisfy the equation  $w^T x + b = 0$ . Hence,

$$w^T \left( x^{(i)} - \gamma^{(i)} \frac{w}{\|w\|} \right) + b = 0.$$

Solving for  $\gamma^{(i)}$  yields

$$\gamma^{(i)} = \frac{w^T x^{(i)} + b}{\|w\|} = \left( \frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|}.$$

This was worked out for the case of a positive training example at  $A$  in the figure, where being on the “positive” side of the decision boundary is good. More generally, we define the geometric margin of  $(w, b)$  with respect to a training example  $(x^{(i)}, y^{(i)})$  to be

$$\gamma^{(i)} = y^{(i)} \left( \left( \frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right).$$

Note that if  $\|w\| = 1$ , then the functional margin equals the geometric margin—this thus gives us a way of relating these two different notions of margin. Also, the geometric margin is invariant to rescaling of the parameters; i.e., if we replace  $w$  with  $2w$  and  $b$  with  $2b$ , then the geometric margin does not change. This will in fact come in handy later. Specifically, because of this invariance to the scaling of the parameters, when trying to fit  $w$  and  $b$  to training data, we can impose an arbitrary scaling constraint on  $w$  without changing anything important; for instance, we can demand that  $\|w\| = 1$ , or  $|w_1| = 5$ , or  $|w_1 + b| + |w_2| = 2$ , and any of these can be satisfied simply by rescaling  $w$  and  $b$ .

Finally, given a training set  $S = \{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$ , we also define the geometric margin of  $(w, b)$  with respect to  $S$  to be the smallest of the geometric margins on the individual training examples:

$$\gamma = \min_{i=1, \dots, n} \gamma^{(i)}.$$

## 6.4 The optimal margin classifier (option reading)

Given a training set, it seems from our previous discussion that a natural desideratum is to try to find a decision boundary that maximizes the (geometric) margin, since this would reflect a very confident set of predictions

on the training set and a good “fit” to the training data. Specifically, this will result in a classifier that separates the positive and the negative training examples with a “gap” (geometric margin).

For now, we will assume that we are given a training set that is linearly separable; i.e., that it is possible to separate the positive and negative examples using some separating hyperplane. How will we find the one that achieves the maximum geometric margin? We can pose the following optimization problem:

$$\begin{aligned} \max_{\gamma, w, b} \quad & \gamma \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq \gamma, \quad i = 1, \dots, n \\ & \|w\| = 1. \end{aligned}$$

I.e., we want to maximize  $\gamma$ , subject to each training example having functional margin at least  $\gamma$ . The  $\|w\| = 1$  constraint moreover ensures that the functional margin equals to the geometric margin, so we are also guaranteed that all the geometric margins are at least  $\gamma$ . Thus, solving this problem will result in  $(w, b)$  with the largest possible geometric margin with respect to the training set.

If we could solve the optimization problem above, we’d be done. But the “ $\|w\| = 1$ ” constraint is a nasty (non-convex) one, and this problem certainly isn’t in any format that we can plug into standard optimization software to solve. So, let’s try transforming the problem into a nicer one. Consider:

$$\begin{aligned} \max_{\hat{\gamma}, w, b} \quad & \frac{\hat{\gamma}}{\|w\|} \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}, \quad i = 1, \dots, n \end{aligned}$$

Here, we’re going to maximize  $\hat{\gamma}/\|w\|$ , subject to the functional margins all being at least  $\hat{\gamma}$ . Since the geometric and functional margins are related by  $\gamma = \hat{\gamma}/\|w\|$ , this will give us the answer we want. Moreover, we’ve gotten rid of the constraint  $\|w\| = 1$  that we didn’t like. The downside is that we now have a nasty (again, non-convex) objective  $\frac{\hat{\gamma}}{\|w\|}$  function; and, we still don’t have any off-the-shelf software that can solve this form of an optimization problem.

Let’s keep going. Recall our earlier discussion that we can add an arbitrary scaling constraint on  $w$  and  $b$  without changing anything. This is the key idea we’ll use now. We will introduce the scaling constraint that the functional margin of  $w, b$  with respect to the training set must be 1:

$$\hat{\gamma} = 1.$$



Since multiplying  $w$  and  $b$  by some constant results in the functional margin being multiplied by that same constant, this is indeed a scaling constraint, and can be satisfied by rescaling  $w, b$ . Plugging this into our problem above, and noting that maximizing  $\hat{\gamma}/\|w\| = 1/\|w\|$  is the same thing as minimizing  $\|w\|^2$ , we now have the following optimization problem:

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2}\|w\|^2 \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, n \end{aligned}$$

We've now transformed the problem into a form that can be efficiently solved. The above is an optimization problem with a convex quadratic objective and only linear constraints. Its solution gives us the **optimal margin classifier**. This optimization problem can be solved using commercial quadratic programming (QP) code.<sup>1</sup>

While we could call the problem solved here, what we will instead do is make a digression to talk about Lagrange duality. This will lead us to our optimization problem's dual form, which will play a key role in allowing us to use kernels to get optimal margin classifiers to work efficiently in very high dimensional spaces. The dual form will also allow us to derive an efficient algorithm for solving the above optimization problem that will typically do much better than generic QP software.

## 6.5 Lagrange duality (optional reading)

Let's temporarily put aside SVMs and maximum margin classifiers, and talk about solving constrained optimization problems.

Consider a problem of the following form:

$$\begin{aligned} \min_w \quad & f(w) \\ \text{s.t.} \quad & h_i(w) = 0, \quad i = 1, \dots, l. \end{aligned}$$

Some of you may recall how the method of Lagrange multipliers can be used to solve it. (Don't worry if you haven't seen it before.) In this method, we define the **Lagrangian** to be

$$\mathcal{L}(w, \beta) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$$

---

<sup>1</sup>You may be familiar with linear programming, which solves optimization problems that have linear objectives and linear constraints. QP software is also widely available, which allows convex quadratic objectives and linear constraints.

Here, the  $\beta_i$ 's are called the **Lagrange multipliers**. We would then find and set  $\mathcal{L}$ 's partial derivatives to zero:

$$\frac{\partial \mathcal{L}}{\partial w_i} = 0; \quad \frac{\partial \mathcal{L}}{\partial \beta_i} = 0,$$

and solve for  $w$  and  $\beta$ .

In this section, we will generalize this to constrained optimization problems in which we may have inequality as well as equality constraints. Due to time constraints, we won't really be able to do the theory of Lagrange duality justice in this class,<sup>2</sup> but we will give the main ideas and results, which we will then apply to our optimal margin classifier's optimization problem.

Consider the following, which we'll call the **primal** optimization problem:

$$\begin{aligned} \min_w \quad & f(w) \\ \text{s.t.} \quad & g_i(w) \leq 0, \quad i = 1, \dots, k \\ & h_i(w) = 0, \quad i = 1, \dots, l. \end{aligned}$$

To solve it, we start by defining the **generalized Lagrangian**

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Here, the  $\alpha_i$ 's and  $\beta_i$ 's are the Lagrange multipliers. Consider the quantity

$$\theta_{\mathcal{P}}(w) = \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta).$$

Here, the " $\mathcal{P}$ " subscript stands for "primal." Let some  $w$  be given. If  $w$  violates any of the primal constraints (i.e., if either  $g_i(w) > 0$  or  $h_i(w) \neq 0$  for some  $i$ ), then you should be able to verify that

$$\theta_{\mathcal{P}}(w) = \max_{\alpha, \beta: \alpha_i \geq 0} f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w) \quad (6.1)$$

$$= \infty. \quad (6.2)$$

Conversely, if the constraints are indeed satisfied for a particular value of  $w$ , then  $\theta_{\mathcal{P}}(w) = f(w)$ . Hence,

$$\theta_{\mathcal{P}}(w) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{otherwise.} \end{cases}$$

---

<sup>2</sup>Readers interested in learning more about this topic are encouraged to read, e.g., R. T. Rockafeller (1970), *Convex Analysis*, Princeton University Press.

Thus,  $\theta_{\mathcal{P}}$  takes the same value as the objective in our problem for all values of  $w$  that satisfies the primal constraints, and is positive infinity if the constraints are violated. Hence, if we consider the minimization problem

$$\min_w \theta_{\mathcal{P}}(w) = \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta),$$

we see that it is the same problem (i.e., and has the same solutions as) our original, primal problem. For later use, we also define the optimal value of the objective to be  $p^* = \min_w \theta_{\mathcal{P}}(w)$ ; we call this the **value** of the primal problem.

Now, let's look at a slightly different problem. We define

$$\theta_{\mathcal{D}}(\alpha, \beta) = \min_w \mathcal{L}(w, \alpha, \beta).$$

Here, the “ $\mathcal{D}$ ” subscript stands for “dual.” Note also that whereas in the definition of  $\theta_{\mathcal{P}}$  we were optimizing (maximizing) with respect to  $\alpha, \beta$ , here we are minimizing with respect to  $w$ .

We can now pose the **dual** optimization problem:

$$\max_{\alpha, \beta: \alpha_i \geq 0} \theta_{\mathcal{D}}(\alpha, \beta) = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta).$$

This is exactly the same as our primal problem shown above, except that the order of the “max” and the “min” are now exchanged. We also define the optimal value of the dual problem's objective to be  $d^* = \max_{\alpha, \beta: \alpha_i \geq 0} \theta_{\mathcal{D}}(\alpha, \beta)$ .

How are the primal and the dual problems related? It can easily be shown that

$$d^* = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta) \leq \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta) = p^*.$$

(You should convince yourself of this; this follows from the “max min” of a function always being less than or equal to the “min max.”) However, under certain conditions, we will have

$$d^* = p^*,$$

so that we can solve the dual problem in lieu of the primal problem. Let's see what these conditions are.

Suppose  $f$  and the  $g_i$ 's are convex,<sup>3</sup> and the  $h_i$ 's are affine.<sup>4</sup> Suppose further that the constraints  $g_i$  are (strictly) feasible; this means that there exists some  $w$  so that  $g_i(w) < 0$  for all  $i$ .

---

<sup>3</sup>When  $f$  has a Hessian, then it is convex if and only if the Hessian is positive semi-definite. For instance,  $f(w) = w^T w$  is convex; similarly, all linear (and affine) functions are also convex. (A function  $f$  can also be convex without being differentiable, but we won't need those more general definitions of convexity here.)

<sup>4</sup>I.e., there exists  $a_i, b_i$ , so that  $h_i(w) = a_i^T w + b_i$ . “Affine” means the same thing as linear, except that we also allow the extra intercept term  $b_i$ .

Under our above assumptions, there must exist  $w^*, \alpha^*, \beta^*$  so that  $w^*$  is the solution to the primal problem,  $\alpha^*, \beta^*$  are the solution to the dual problem, and moreover  $p^* = d^* = \mathcal{L}(w^*, \alpha^*, \beta^*)$ . Moreover,  $w^*, \alpha^*$  and  $\beta^*$  satisfy the **Karush-Kuhn-Tucker (KKT) conditions**, which are as follows:

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, d \quad (6.3)$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, l \quad (6.4)$$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k \quad (6.5)$$

$$g_i(w^*) \leq 0, \quad i = 1, \dots, k \quad (6.6)$$

$$\alpha^* \geq 0, \quad i = 1, \dots, k \quad (6.7)$$

Moreover, if some  $w^*, \alpha^*, \beta^*$  satisfy the KKT conditions, then it is also a solution to the primal and dual problems.

We draw attention to Equation (6.5), which is called the KKT **dual complementarity** condition. Specifically, it implies that if  $\alpha_i^* > 0$ , then  $g_i(w^*) = 0$ . (I.e., the “ $g_i(w) \leq 0$ ” constraint is **active**, meaning it holds with equality rather than with inequality.) Later on, this will be key for showing that the SVM has only a small number of “support vectors”; the KKT dual complementarity condition will also give us our convergence test when we talk about the SMO algorithm.

## 6.6 Optimal margin classifiers: the dual form (option reading)

*Note: The equivalence of optimization problem (6.8) and the optimization problem (6.12), and the relationship between the primary and dual variables in equation (6.10) are the most important take home messages of this section.*

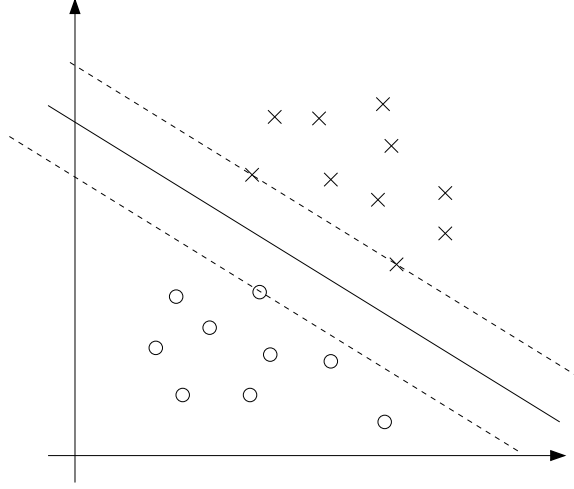
Previously, we posed the following (primal) optimization problem for finding the optimal margin classifier:

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, n \end{aligned} \quad (6.8)$$

We can write the constraints as

$$g_i(w) = -y^{(i)}(w^T x^{(i)} + b) + 1 \leq 0.$$

We have one such constraint for each training example. Note that from the KKT dual complementarity condition, we will have  $\alpha_i > 0$  only for the training examples that have functional margin exactly equal to one (i.e., the ones corresponding to constraints that hold with equality,  $g_i(w) = 0$ ). Consider the figure below, in which a maximum margin separating hyperplane is shown by the solid line.



The points with the smallest margins are exactly the ones closest to the decision boundary; here, these are the three points (one negative and two positive examples) that lie on the dashed lines parallel to the decision boundary. Thus, only three of the  $\alpha_i$ 's—namely, the ones corresponding to these three training examples—will be non-zero at the optimal solution to our optimization problem. These three points are called the **support vectors** in this problem. The fact that the number of support vectors can be much smaller than the size the training set will be useful later.

Let's move on. Looking ahead, as we develop the dual form of the problem, one key idea to watch out for is that we'll try to write our algorithm in terms of only the inner product  $\langle x^{(i)}, x^{(j)} \rangle$  (think of this as  $(x^{(i)})^T x^{(j)}$ ) between points in the input feature space. The fact that we can express our algorithm in terms of these inner products will be key when we apply the kernel trick.

When we construct the Lagrangian for our optimization problem we have:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1]. \quad (6.9)$$

Note that there're only " $\alpha_i$ " but no " $\beta_i$ " Lagrange multipliers, since the problem has only inequality constraints.

Let's find the dual form of the problem. To do so, we need to first minimize  $\mathcal{L}(w, b, \alpha)$  with respect to  $w$  and  $b$  (for fixed  $\alpha$ ), to get  $\theta_{\mathcal{D}}$ , which we'll do by setting the derivatives of  $\mathcal{L}$  with respect to  $w$  and  $b$  to zero. We have:

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} = 0$$

This implies that

$$w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}. \quad (6.10)$$

As for the derivative with respect to  $b$ , we obtain

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i y^{(i)} = 0. \quad (6.11)$$

If we take the definition of  $w$  in Equation (6.10) and plug that back into the Lagrangian (Equation 6.9), and simplify, we get

$$\mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)} - b \sum_{i=1}^n \alpha_i y^{(i)}.$$

But from Equation (6.11), the last term must be zero, so we obtain

$$\mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)}.$$

Recall that we got to the equation above by minimizing  $\mathcal{L}$  with respect to  $w$  and  $b$ . Putting this together with the constraints  $\alpha_i \geq 0$  (that we always had) and the constraint (6.11), we obtain the following dual optimization problem:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle. \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0, \end{aligned} \quad (6.12)$$

You should also be able to verify that the conditions required for  $p^* = d^*$  and the KKT conditions (Equations 6.3–6.7) to hold are indeed satisfied in

our optimization problem. Hence, we can solve the dual in lieu of solving the primal problem. Specifically, in the dual problem above, we have a maximization problem in which the parameters are the  $\alpha_i$ 's. We'll talk later about the specific algorithm that we're going to use to solve the dual problem, but if we are indeed able to solve it (i.e., find the  $\alpha$ 's that maximize  $W(\alpha)$  subject to the constraints), then we can use Equation (6.10) to go back and find the optimal  $w$ 's as a function of the  $\alpha$ 's. Having found  $w^*$ , by considering the primal problem, it is also straightforward to find the optimal value for the intercept term  $b$  as

$$b^* = -\frac{\max_{i:y^{(i)}=-1} w^{*T} x^{(i)} + \min_{i:y^{(i)}=1} w^{*T} x^{(i)}}{2}. \quad (6.13)$$

(Check for yourself that this is correct.)

Before moving on, let's also take a more careful look at Equation (6.10), which gives the optimal value of  $w$  in terms of (the optimal value of)  $\alpha$ . Suppose we've fit our model's parameters to a training set, and now wish to make a prediction at a new point input  $x$ . We would then calculate  $w^T x + b$ , and predict  $y = 1$  if and only if this quantity is bigger than zero. But using (6.10), this quantity can also be written:

$$w^T x + b = \left( \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \right)^T x + b \quad (6.14)$$

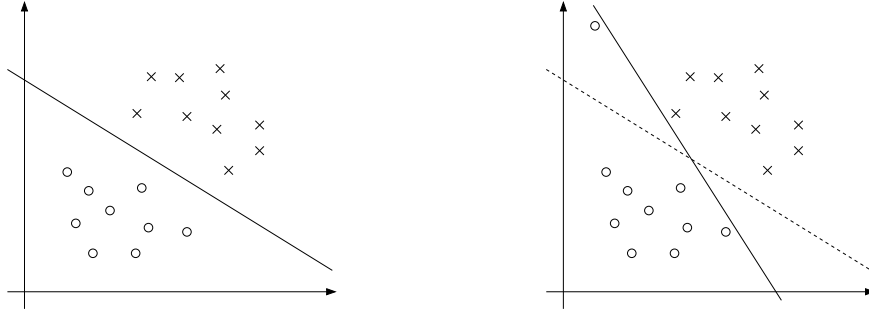
$$= \sum_{i=1}^n \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b. \quad (6.15)$$

Hence, if we've found the  $\alpha_i$ 's, in order to make a prediction, we have to calculate a quantity that depends only on the inner product between  $x$  and the points in the training set. Moreover, we saw earlier that the  $\alpha_i$ 's will all be zero except for the support vectors. Thus, many of the terms in the sum above will be zero, and we really need to find only the inner products between  $x$  and the support vectors (of which there is often only a small number) in order to calculate (6.15) and make our prediction.

By examining the dual form of the optimization problem, we gained significant insight into the structure of the problem, and were also able to write the entire algorithm in terms of only inner products between input feature vectors. In the next section, we will exploit this property to apply the kernels to our classification problem. The resulting algorithm, **support vector machines**, will be able to efficiently learn in very high dimensional spaces.

## 6.7 Regularization and the non-separable case (optional reading)

The derivation of the SVM as presented so far assumed that the data is linearly separable. While mapping data to a high dimensional feature space via  $\phi$  does generally increase the likelihood that the data is separable, we can't guarantee that it always will be so. Also, in some cases it is not clear that finding a separating hyperplane is exactly what we'd want to do, since that might be susceptible to outliers. For instance, the left figure below shows an optimal margin classifier, and when a single outlier is added in the upper-left region (right figure), it causes the decision boundary to make a dramatic swing, and the resulting classifier has a much smaller margin.



To make the algorithm work for non-linearly separable datasets as well as be less sensitive to outliers, we reformulate our optimization (using  $\ell_1$  **regularization**) as follows:

$$\begin{aligned} \min_{\gamma, w, b} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\ & \xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

Thus, examples are now permitted to have (functional) margin less than 1, and if an example has functional margin  $1 - \xi_i$  (with  $\xi > 0$ ), we would pay a cost of the objective function being increased by  $C\xi_i$ . The parameter  $C$  controls the relative weighting between the twin goals of making the  $\|w\|^2$  small (which we saw earlier makes the margin large) and of ensuring that most examples have functional margin at least 1.



As before, we can form the Lagrangian:

$$\mathcal{L}(w, b, \xi, \alpha, r) = \frac{1}{2}w^T w + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [y^{(i)}(x^T w + b) - 1 + \xi_i] - \sum_{i=1}^n r_i \xi_i.$$

Here, the  $\alpha_i$ 's and  $r_i$ 's are our Lagrange multipliers (constrained to be  $\geq 0$ ). We won't go through the derivation of the dual again in detail, but after setting the derivatives with respect to  $w$  and  $b$  to zero as before, substituting them back in, and simplifying, we obtain the following dual form of the problem:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0, \end{aligned}$$

As before, we also have that  $w$  can be expressed in terms of the  $\alpha_i$ 's as given in Equation (6.10), so that after solving the dual problem, we can continue to use Equation (6.15) to make our predictions. Note that, somewhat surprisingly, in adding  $\ell_1$  regularization, the only change to the dual problem is that what was originally a constraint that  $0 \leq \alpha_i$  has now become  $0 \leq \alpha_i \leq C$ . The calculation for  $b^*$  also has to be modified (Equation 6.13 is no longer valid); see the comments in the next section/Platt's paper.

Also, the KKT dual-complementarity conditions (which in the next section will be useful for testing for the convergence of the SMO algorithm) are:

$$\alpha_i = 0 \Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1 \quad (6.16)$$

$$\alpha_i = C \Rightarrow y^{(i)}(w^T x^{(i)} + b) \leq 1 \quad (6.17)$$

$$0 < \alpha_i < C \Rightarrow y^{(i)}(w^T x^{(i)} + b) = 1. \quad (6.18)$$

Now, all that remains is to give an algorithm for actually solving the dual problem, which we will do in the next section.

## 6.8 The SMO algorithm (optional reading)

The SMO (sequential minimal optimization) algorithm, due to John Platt, gives an efficient way of solving the dual problem arising from the derivation

of the SVM. Partly to motivate the SMO algorithm, and partly because it's interesting in its own right, let's first take another digression to talk about the coordinate ascent algorithm.

### 6.8.1 Coordinate ascent

Consider trying to solve the unconstrained optimization problem

$$\max_{\alpha} W(\alpha_1, \alpha_2, \dots, \alpha_n).$$

Here, we think of  $W$  as just some function of the parameters  $\alpha_i$ 's, and for now ignore any relationship between this problem and SVMs. We've already seen two optimization algorithms, gradient ascent and Newton's method. The new algorithm we're going to consider here is called **coordinate ascent**:

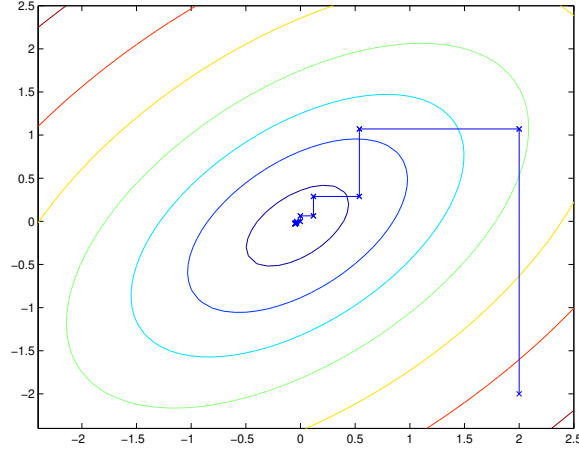
```

Loop until convergence: {
    For  $i = 1, \dots, n$ , {
         $\alpha_i := \arg \max_{\hat{\alpha}_i} W(\alpha_1, \dots, \alpha_{i-1}, \hat{\alpha}_i, \alpha_{i+1}, \dots, \alpha_n)$ .
    }
}

```

Thus, in the innermost loop of this algorithm, we will hold all the variables except for some  $\alpha_i$  fixed, and reoptimize  $W$  with respect to just the parameter  $\alpha_i$ . In the version of this method presented here, the inner-loop reoptimizes the variables in order  $\alpha_1, \alpha_2, \dots, \alpha_n, \alpha_1, \alpha_2, \dots$ . (A more sophisticated version might choose other orderings; for instance, we may choose the next variable to update according to which one we expect to allow us to make the largest increase in  $W(\alpha)$ .)

When the function  $W$  happens to be of such a form that the “arg max” in the inner loop can be performed efficiently, then coordinate ascent can be a fairly efficient algorithm. Here's a picture of coordinate ascent in action:



The ellipses in the figure are the contours of a quadratic function that we want to optimize. Coordinate ascent was initialized at  $(2, -2)$ , and also plotted in the figure is the path that it took on its way to the global maximum. Notice that on each step, coordinate ascent takes a step that's parallel to one of the axes, since only one variable is being optimized at a time.

### 6.8.2 SMO

We close off the discussion of SVMs by sketching the derivation of the SMO algorithm.

Here's the (dual) optimization problem that we want to solve:

$$\max_{\alpha} \quad W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle. \quad (6.19)$$

$$\text{s.t.} \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n \quad (6.20)$$

$$\sum_{i=1}^n \alpha_i y^{(i)} = 0. \quad (6.21)$$

Let's say we have set of  $\alpha_i$ 's that satisfy the constraints (6.20-6.21). Now, suppose we want to hold  $\alpha_2, \dots, \alpha_n$  fixed, and take a coordinate ascent step and reoptimize the objective with respect to  $\alpha_1$ . Can we make any progress? The answer is no, because the constraint (6.21) ensures that

$$\alpha_1 y^{(1)} = - \sum_{i=2}^n \alpha_i y^{(i)}.$$

Or, by multiplying both sides by  $y^{(1)}$ , we equivalently have

$$\alpha_1 = -y^{(1)} \sum_{i=2}^n \alpha_i y^{(i)}.$$

(This step used the fact that  $y^{(1)} \in \{-1, 1\}$ , and hence  $(y^{(1)})^2 = 1$ .) Hence,  $\alpha_1$  is exactly determined by the other  $\alpha_i$ 's, and if we were to hold  $\alpha_2, \dots, \alpha_n$  fixed, then we can't make any change to  $\alpha_1$  without violating the constraint (6.21) in the optimization problem.

Thus, if we want to update some subset of the  $\alpha_i$ 's, we must update at least two of them simultaneously in order to keep satisfying the constraints. This motivates the SMO algorithm, which simply does the following:

- Repeat till convergence {
1. Select some pair  $\alpha_i$  and  $\alpha_j$  to update next (using a heuristic that tries to pick the two that will allow us to make the biggest progress towards the global maximum).
  2. Reoptimize  $W(\alpha)$  with respect to  $\alpha_i$  and  $\alpha_j$ , while holding all the other  $\alpha_k$ 's ( $k \neq i, j$ ) fixed.
- }

To test for convergence of this algorithm, we can check whether the KKT conditions (Equations 6.16-6.18) are satisfied to within some *tol*. Here, *tol* is the convergence tolerance parameter, and is typically set to around 0.01 to 0.001. (See the paper and pseudocode for details.)

The key reason that SMO is an efficient algorithm is that the update to  $\alpha_i$ ,  $\alpha_j$  can be computed very efficiently. Let's now briefly sketch the main ideas for deriving the efficient update.

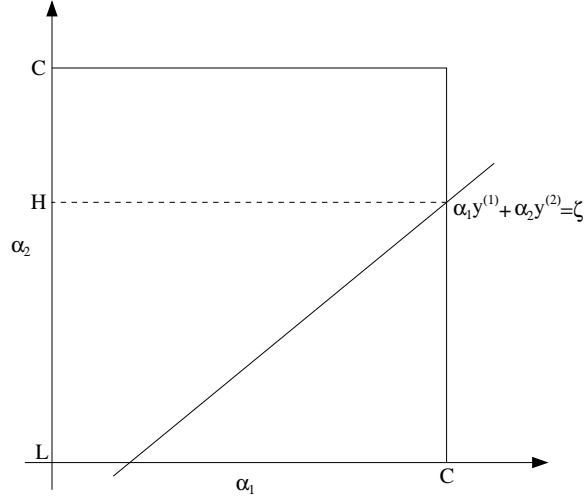
Let's say we currently have some setting of the  $\alpha_i$ 's that satisfy the constraints (6.20-6.21), and suppose we've decided to hold  $\alpha_3, \dots, \alpha_n$  fixed, and want to reoptimize  $W(\alpha_1, \alpha_2, \dots, \alpha_n)$  with respect to  $\alpha_1$  and  $\alpha_2$  (subject to the constraints). From (6.21), we require that

$$\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = - \sum_{i=3}^n \alpha_i y^{(i)}.$$

Since the right hand side is fixed (as we've fixed  $\alpha_3, \dots, \alpha_n$ ), we can just let it be denoted by some constant  $\zeta$ :

$$\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = \zeta. \tag{6.22}$$

We can thus picture the constraints on  $\alpha_1$  and  $\alpha_2$  as follows:



From the constraints (6.20), we know that  $\alpha_1$  and  $\alpha_2$  must lie within the box  $[0, C] \times [0, C]$  shown. Also plotted is the line  $\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = \zeta$ , on which we know  $\alpha_1$  and  $\alpha_2$  must lie. Note also that, from these constraints, we know  $L \leq \alpha_2 \leq H$ ; otherwise,  $(\alpha_1, \alpha_2)$  can't simultaneously satisfy both the box and the straight line constraint. In this example,  $L = 0$ . But depending on what the line  $\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = \zeta$  looks like, this won't always necessarily be the case; but more generally, there will be some lower-bound  $L$  and some upper-bound  $H$  on the permissible values for  $\alpha_2$  that will ensure that  $\alpha_1, \alpha_2$  lie within the box  $[0, C] \times [0, C]$ .

Using Equation (6.22), we can also write  $\alpha_1$  as a function of  $\alpha_2$ :

$$\alpha_1 = (\zeta - \alpha_2 y^{(2)}) y^{(1)}.$$

(Check this derivation yourself; we again used the fact that  $y^{(1)} \in \{-1, 1\}$  so that  $(y^{(1)})^2 = 1$ .) Hence, the objective  $W(\alpha)$  can be written

$$W(\alpha_1, \alpha_2, \dots, \alpha_n) = W((\zeta - \alpha_2 y^{(2)}) y^{(1)}, \alpha_2, \dots, \alpha_n).$$

Treating  $\alpha_3, \dots, \alpha_n$  as constants, you should be able to verify that this is just some quadratic function in  $\alpha_2$ . I.e., this can also be expressed in the form  $a\alpha_2^2 + b\alpha_2 + c$  for some appropriate  $a$ ,  $b$ , and  $c$ . If we ignore the “box” constraints (6.20) (or, equivalently, that  $L \leq \alpha_2 \leq H$ ), then we can easily maximize this quadratic function by setting its derivative to zero and solving. We'll let  $\alpha_2^{new, unclipped}$  denote the resulting value of  $\alpha_2$ . You should also be able to convince yourself that if we had instead wanted to maximize  $W$  with respect to  $\alpha_2$  but subject to the box constraint, then we can find the resulting value optimal simply by taking  $\alpha_2^{new, unclipped}$  and “clipping” it to lie in the

$[L, H]$  interval, to get

$$\alpha_2^{new} = \begin{cases} H & \text{if } \alpha_2^{new,unclipped} > H \\ \alpha_2^{new,unclipped} & \text{if } L \leq \alpha_2^{new,unclipped} \leq H \\ L & \text{if } \alpha_2^{new,unclipped} < L \end{cases}$$

Finally, having found the  $\alpha_2^{new}$ , we can use Equation (6.22) to go back and find the optimal value of  $\alpha_1^{new}$ .

There're a couple more details that are quite easy but that we'll leave you to read about yourself in Platt's paper: One is the choice of the heuristics used to select the next  $\alpha_i$ ,  $\alpha_j$  to update; the other is how to update  $b$  as the SMO algorithm is run.

# Part II

## Deep Learning

# Chapter 7

## Deep Learning

We now begin our study of deep learning. In this set of notes, we give an overview of neural networks, discuss vectorization and discuss training neural networks with backpropagation.

### 7.1 Supervised Learning with Non-linear Models

In the supervised learning setting (predicting  $y$  from the input  $x$ ), suppose our model/hypothesis is  $h_\theta(x)$ . In the past lectures, we have considered the cases when  $h_\theta(x) = \theta^\top x$  (in linear regression or logistic regression) or  $h_\theta(x) = \theta^\top \phi(x)$  (where  $\phi(x)$  is the feature map). A commonality of these two models is that they are linear in the parameters  $\theta$ . Next we will consider learning general family of models that are **non-linear in both** the parameters  $\theta$  and the inputs  $x$ . The most common non-linear models are neural networks, which we will define starting from the next section. For this section, it suffices to think  $h_\theta(x)$  as an abstract non-linear model.<sup>1</sup>

Suppose  $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$  are the training examples. For simplicity, we start with the case where  $y^{(i)} \in \mathbb{R}$  and  $h_\theta(x) \in \mathbb{R}$ .

**Cost/loss function.** We define the least square cost function for the  $i$ -th example  $(x^{(i)}, y^{(i)})$  as

$$J^{(i)}(\theta) = \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2 \quad (7.1)$$

---

<sup>1</sup>If a concrete example is helpful, perhaps think about the model  $h_\theta(x) = \theta_1^2 x_1^2 + \theta_2^2 x_2^2 + \dots + \theta_d^2 x_d^2$  in this subsection, even though it's not a neural network.



and define the mean-square cost function for the dataset as

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n J^{(i)}(\theta) \quad (7.2)$$

which is same as in linear regression except that we introduce a constant  $1/n$  in front of the cost function to be consistent with the convention. Note that multiplying the cost function with a scalar will not change the local minima or global minima of the cost function. Also note that the underlying parameterization for  $h_{\theta}(x)$  is different from the case of linear regression, even though the form of the cost function is the same mean-squared loss. Throughout the notes, we use the words “loss” and “cost” interchangeably.

**Optimizers (SGD).** Commonly, people use gradient descent (GD), stochastic gradient (SGD), or their variants to optimize the loss function  $J(\theta)$ . GD’s update rule can be written as<sup>2</sup>

$$\theta := \theta - \alpha \nabla_{\theta} J(\theta) \quad (7.3)$$

where  $\alpha > 0$  is often referred to as the learning rate or step size. Next, we introduce a version of the SGD (Algorithm 1), which is lightly different from that in the first lecture notes.

---

**Algorithm 1** Stochastic Gradient Descent

---

- 1: Hyperparameter: learning rate  $\alpha$ , number of total iteration  $n_{\text{iter}}$ .
- 2: Initialize  $\theta$  randomly.
- 3: **for**  $i = 1$  to  $n_{\text{iter}}$  **do**
- 4:     Sample  $j$  uniformly from  $\{1, \dots, n\}$ , and update  $\theta$  by

$$\theta := \theta - \alpha \nabla_{\theta} J^{(j)}(\theta) \quad (7.4)$$


---

Oftentimes computing the gradient of  $B$  examples simultaneously for the parameter  $\theta$  can be faster than computing  $B$  gradients separately due to hardware parallelization. Therefore, a mini-batch version of SGD is most

---

<sup>2</sup>Recall that, as defined in the previous lecture notes, we use the notation “ $a := b$ ” to denote an operation (in a computer program) in which we *set* the value of a variable  $a$  to be equal to the value of  $b$ . In other words, this operation overwrites  $a$  with the value of  $b$ . In contrast, we will write “ $a = b$ ” when we are asserting a statement of fact, that the value of  $a$  is equal to the value of  $b$ .

commonly used in deep learning, as shown in Algorithm 2. There are also other variants of the SGD or mini-batch SGD with slightly different sampling schemes.

---

**Algorithm 2** Mini-batch Stochastic Gradient Descent

---

- 1: Hyperparameters: learning rate  $\alpha$ , batch size  $B$ , # iterations  $n_{\text{iter}}$ .
- 2: Initialize  $\theta$  randomly
- 3: **for**  $i = 1$  to  $n_{\text{iter}}$  **do**
- 4:     Sample  $B$  examples  $j_1, \dots, j_B$  (without replacement) uniformly from  $\{1, \dots, n\}$ , and update  $\theta$  by

$$\theta := \theta - \frac{\alpha}{B} \sum_{k=1}^B \nabla_{\theta} J^{(j_k)}(\theta) \quad (7.5)$$


---

With these generic algorithms, a typical deep learning model is learned with the following steps. 1. Define a neural network parametrization  $h_{\theta}(x)$ , which we will introduce in Section 7.2, and 2. write the backpropagation algorithm to compute the gradient of the loss function  $J^{(j)}(\theta)$  efficiently, which will be covered in Section 7.3, and 3. run SGD or mini-batch SGD (or other gradient-based optimizers) with the loss function  $J(\theta)$ .

## 7.2 Neural Networks

Neural networks refer to broad type of non-linear models/parametrizations  $h_{\theta}(x)$  that involve combinations of matrix multiplications and other entry-wise non-linear operations. We will start small and slowly build up a neural network, step by step.

**A Neural Network with a Single Neuron.** Recall the housing price prediction problem from before: given the size of the house, we want to predict the price. We will use it as a running example in this subsection.

Previously, we fit a straight line to the graph of size vs. housing price. Now, instead of fitting a straight line, we wish to prevent negative housing prices by setting the absolute minimum price as zero. This produces a “kink” in the graph as shown in Figure 7.1. How do we represent such a function with a single kink as  $h_{\theta}(x)$  with unknown parameter? (After doing so, we can invoke the machinery in Section 7.1.)

We define a parameterized function  $h_\theta(x)$  with input  $x$ , parameterized by  $\theta$ , which outputs the price of the house  $y$ . Formally,  $h_\theta : x \rightarrow y$ . Perhaps one of the simplest parametrization would be

$$h_\theta(x) = \max(wx + b, 0), \text{ where } \theta = (w, b) \in \mathbb{R}^2 \quad (7.6)$$

Here  $h_\theta(x)$  returns a single value:  $(wx+b)$  or zero, whichever is greater. In the context of neural networks, the function  $\max\{t, 0\}$  is called a ReLU (pronounced “ray-lu”), or rectified linear unit, and often denoted by  $\text{ReLU}(t) \triangleq \max\{t, 0\}$ .

Generally, a one-dimensional non-linear function that maps  $\mathbb{R}$  to  $\mathbb{R}$  such as ReLU is often referred to as an **activation function**. The model  $h_\theta(x)$  is said to have a single neuron partly because it has a single non-linear activation function. (We will discuss more about why a non-linear activation is called neuron.)

When the input  $x \in \mathbb{R}^d$  has multiple dimensions, a neural network with a single neuron can be written as

$$h_\theta(x) = \text{ReLU}(w^\top x + b), \text{ where } w \in \mathbb{R}^d, b \in \mathbb{R}, \text{ and } \theta = (w, b) \quad (7.7)$$

The term  $b$  is often referred to as the “bias”, and the vector  $w$  is referred to as the weight vector. Such a neural network has 1 layer. (We will define what multiple layers mean in the sequel.)

**Stacking Neurons.** A more complex neural network may take the single neuron described above and “stack” them together such that one neuron passes its output as input into the next neuron, resulting in a more complex function.

Let us now deepen the housing prediction example. In addition to the size of the house, suppose that you know the number of bedrooms, the zip code and the wealth of the neighborhood. Building neural networks is analogous to Lego bricks: you take individual bricks and stack them together to build complex structures. The same applies to neural networks: we take individual neurons and stack them together to create complex neural networks.

Given these features (size, number of bedrooms, zip code, and wealth), we might then decide that the price of the house depends on the maximum family size it can accommodate. Suppose the family size is a function of the size of the house and number of bedrooms (see Figure 7.2). The zip code may provide additional information such as how walkable the neighborhood is (i.e., can you walk to the grocery store or do you need to drive everywhere). Combining the zip code with the wealth of the neighborhood may predict

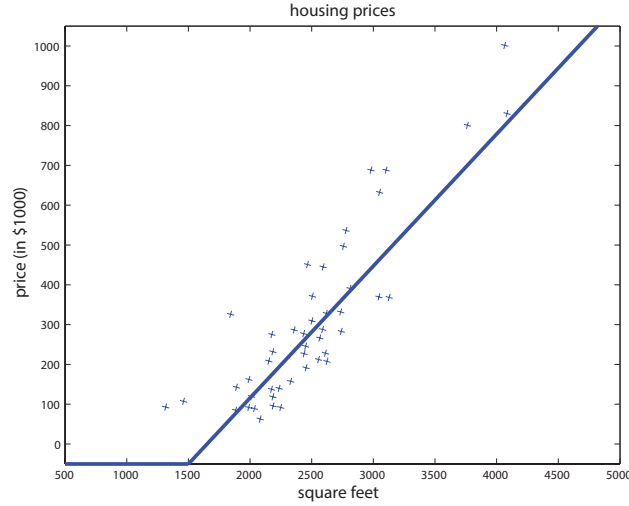


Figure 7.1: Housing prices with a “kink” in the graph.

the quality of the local elementary school. Given these three derived features (family size, walkable, school quality), we may conclude that the price of the home ultimately depends on these three features.

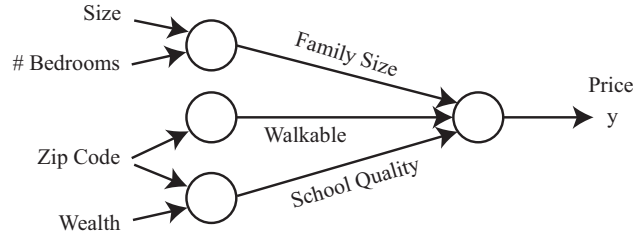


Figure 7.2: Diagram of a small neural network for predicting housing prices.

Formally, the input to a neural network is a set of input features  $x_1, x_2, x_3, x_4$ . We denote the intermediate variables for “family size”, “walkable”, and “school quality” by  $a_1, a_2, a_3$  (these  $a_i$ ’s are often referred to as “hidden units” or “hidden neurons”). We represent each of the  $a_i$ ’s as a neural network with a single neuron with a subset of  $x_1, \dots, x_4$  as inputs. Then as in Figure 7.1, we will have the parameterization:

$$\begin{aligned} a_1 &= \text{ReLU}(\theta_1 x_1 + \theta_2 x_2 + \theta_3) \\ a_2 &= \text{ReLU}(\theta_4 x_3 + \theta_5) \\ a_3 &= \text{ReLU}(\theta_6 x_3 + \theta_7 x_4 + \theta_8) \end{aligned}$$

where  $(\theta_1, \dots, \theta_8)$  are parameters. Now we represent the final output  $h_\theta(x)$  as another linear function with  $a_1, a_2, a_3$  as inputs, and we get<sup>3</sup>

$$h_\theta(x) = \theta_9 a_1 + \theta_{10} a_2 + \theta_{11} a_3 + \theta_{12} \quad (7.8)$$

where  $\theta$  contains all the parameters  $(\theta_1, \dots, \theta_{12})$ .

Now we represent the output as a quite complex function of  $x$  with parameters  $\theta$ . Then you can use this parametrization  $h_\theta$  with the machinery of Section 7.1 to learn the parameters  $\theta$ .

**Inspiration from Biological Neural Networks.** As the name suggests, artificial neural networks were inspired by biological neural networks. The hidden units  $a_1, \dots, a_m$  correspond to the neurons in a biological neural network, and the parameters  $\theta_i$ 's correspond to the synapses. However, it's unclear how similar the modern deep artificial neural networks are to the biological ones. For example, perhaps not many neuroscientists think biological neural networks could have 1000 layers, while some modern artificial neural networks do (we will elaborate more on the notion of layers.) Moreover, it's an open question whether human brains update their neural networks in a way similar to the way that computer scientists learn artificial neural networks (using backpropagation, which we will introduce in the next section.).

**Two-layer Fully-Connected Neural Networks.** We constructed the neural network in equation (7.8) using a significant amount of prior knowledge/belief about how the “family size”, “walkable”, and “school quality” are determined by the inputs. We implicitly assumed that we know the family size is an important quantity to look at and that it can be determined by only the “size” and “# bedrooms”. Such a prior knowledge might not be available for other applications. It would be more flexible and general to have a generic parameterization. A simple way would be to write the intermediate variable  $a_1$  as a function of all  $x_1, \dots, x_4$ :

$$a_1 = \text{ReLU}(w_1^\top x + b_1), \text{ where } w_1 \in \mathbb{R}^4 \text{ and } b_1 \in \mathbb{R} \quad (7.9)$$

$$a_2 = \text{ReLU}(w_2^\top x + b_2), \text{ where } w_2 \in \mathbb{R}^4 \text{ and } b_2 \in \mathbb{R}$$

$$a_3 = \text{ReLU}(w_3^\top x + b_3), \text{ where } w_3 \in \mathbb{R}^4 \text{ and } b_3 \in \mathbb{R}$$

We still define  $h_\theta(x)$  using equation (7.8) with  $a_1, a_2, a_3$  being defined as above. Thus we have a so-called **fully-connected neural network** as

---

<sup>3</sup>Typically, for multi-layer neural network, at the end, near the output, we don't apply ReLU, especially when the output is not necessarily a positive number.

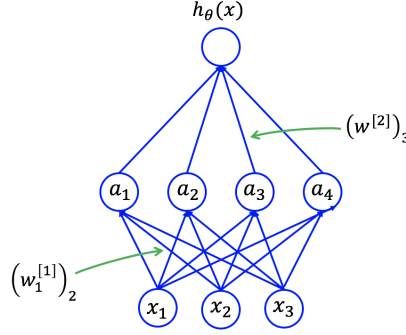


Figure 7.3: Diagram of a two-layer fully connected neural network. Each edge from node  $x_i$  to node  $a_j$  indicates that  $a_j$  depends on  $x_i$ . The edge from  $x_i$  to  $a_j$  is associated with the weight  $(w_j^{[1]})_i$  which denotes the  $i$ -th coordinate of the vector  $w_j^{[1]}$ . The activation  $a_j$  can be computed by taking the ReLU of the weighted sum of  $x_i$ 's with the weights being the weights associated with the incoming edges, that is,  $a_j = \text{ReLU}(\sum_{i=1}^d (w_j^{[1]})_i x_i)$ .

visualized in the dependency graph in Figure 7.3 because all the intermediate variables  $a_i$ 's depend on all the inputs  $x_i$ 's.

For full generality, a two-layer fully-connected neural network with  $m$  hidden units and  $d$  dimensional input  $x \in \mathbb{R}^d$  is defined as

$$\forall j \in [1, \dots, m], \quad z_j = w_j^{[1]\top} x + b_j^{[1]} \text{ where } w_j^{[1]} \in \mathbb{R}^d, b_j^{[1]} \in \mathbb{R} \quad (7.10)$$

$$a_j = \text{ReLU}(z_j),$$

$$a = [a_1, \dots, a_m]^\top \in \mathbb{R}^m$$

$$h_\theta(x) = w^{[2]\top} a + b^{[2]} \text{ where } w^{[2]} \in \mathbb{R}^m, b^{[2]} \in \mathbb{R}, \quad (7.11)$$

Note that by default the vectors in  $\mathbb{R}^d$  are viewed as column vectors, and in particular  $a$  is a column vector with components  $a_1, a_2, \dots, a_m$ . The indices  $^{[1]}$  and  $^{[2]}$  are used to distinguish two sets of parameters: the  $w_j^{[1]}$ 's (each of which is a vector in  $\mathbb{R}^d$ ) and  $w^{[2]}$  (which is a vector in  $\mathbb{R}^m$ ). We will have more of these later.

**Vectorization.** Before we introduce neural networks with more layers and more complex structures, we will simplify the expressions for neural networks with more matrix and vector notations. Another important motivation of

vectorization is the speed perspective in the implementation. In order to implement a neural network efficiently, one must be careful when using for loops. The most natural way to implement equation (7.10) in code is perhaps to use a for loop. In practice, the dimensionalities of the inputs and hidden units are high. As a result, code will run very slowly if you use for loops. Leveraging the parallelism in GPUs is/was crucial for the progress of deep learning.

This gave rise to *vectorization*. Instead of using for loops, vectorization takes advantage of matrix algebra and highly optimized numerical linear algebra packages (e.g., BLAS) to make neural network computations run quickly. Before the deep learning era, a for loop may have been sufficient on smaller datasets, but modern deep networks and state-of-the-art datasets will be infeasible to run with for loops.

We vectorize the two-layer fully-connected neural network as below. We define a weight matrix  $W^{[1]}$  in  $\mathbb{R}^{m \times d}$  as the concatenation of all the vectors  $w_j^{[1]}$ 's in the following way:

$$W^{[1]} = \begin{bmatrix} - & w_1^{[1]\top} & - \\ - & w_2^{[1]\top} & - \\ & \vdots & \\ - & w_m^{[1]\top} & - \end{bmatrix} \in \mathbb{R}^{m \times d} \quad (7.12)$$

Now by the definition of matrix vector multiplication, we can write  $z = [z_1, \dots, z_m]^\top \in \mathbb{R}^m$  as

$$\underbrace{\begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix}}_{z \in \mathbb{R}^{m \times 1}} = \underbrace{\begin{bmatrix} - & w_1^{[1]\top} & - \\ - & w_2^{[1]\top} & - \\ & \vdots & \\ - & w_m^{[1]\top} & - \end{bmatrix}}_{W^{[1]} \in \mathbb{R}^{m \times d}} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}}_{x \in \mathbb{R}^{d \times 1}} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ \vdots \\ b_m^{[1]} \end{bmatrix}}_{b^{[1]} \in \mathbb{R}^{m \times 1}} \quad (7.13)$$

Or succinctly,

$$z = W^{[1]}x + b^{[1]} \quad (7.14)$$

We remark again that a vector in  $\mathbb{R}^d$  in this notes, following the conventions previously established, is automatically viewed as a column vector, and can also be viewed as a  $d \times 1$  dimensional matrix. (Note that this is different from numpy where a vector is viewed as a row vector in broadcasting.)

Computing the activations  $a \in \mathbb{R}^m$  from  $z \in \mathbb{R}^m$  involves an element-wise non-linear application of the ReLU function, which can be computed in parallel efficiently. Overloading ReLU for element-wise application of ReLU (meaning, for a vector  $t \in \mathbb{R}^d$ ,  $\text{ReLU}(t)$  is a vector such that  $\text{ReLU}(t)_i = \text{ReLU}(t_i)$ ), we have

$$a = \text{ReLU}(z) \quad (7.15)$$

Define  $W^{[2]} = [w^{[2]\top}] \in \mathbb{R}^{1 \times m}$  similarly. Then, the model in equation (7.11) can be summarized as

$$\begin{aligned} a &= \text{ReLU}(W^{[1]}x + b^{[1]}) \\ h_\theta(x) &= W^{[2]}a + b^{[2]} \end{aligned} \quad (7.16)$$

Here  $\theta$  consists of  $W^{[1]}, W^{[2]}$  (often referred to as the weight matrices) and  $b^{[1]}, b^{[2]}$  (referred to as the biases). The collection of  $W^{[1]}, b^{[1]}$  is referred to as the first layer, and  $W^{[2]}, b^{[2]}$  the second layer. The activation  $a$  is referred to as the hidden layer. A two-layer neural network is also called one-hidden-layer neural network.

**Multi-layer fully-connected neural networks.** With this succinct notations, we can stack more layers to get a deeper fully-connected neural network. Let  $r$  be the number of layers (weight matrices). Let  $W^{[1]}, \dots, W^{[r]}, b^{[1]}, \dots, b^{[r]}$  be the weight matrices and biases of all the layers. Then a multi-layer neural network can be written as

$$\begin{aligned} a^{[1]} &= \text{ReLU}(W^{[1]}x + b^{[1]}) \\ a^{[2]} &= \text{ReLU}(W^{[2]}a^{[1]} + b^{[2]}) \\ &\dots \\ a^{[r-1]} &= \text{ReLU}(W^{[r-1]}a^{[r-2]} + b^{[r-1]}) \\ h_\theta(x) &= W^{[r]}a^{[r-1]} + b^{[r]} \end{aligned} \quad (7.17)$$

We note that the weight matrices and biases need to have compatible dimensions for the equations above to make sense. If  $a^{[k]}$  has dimension  $m_k$ , then the weight matrix  $W^{[k]}$  should be of dimension  $m_k \times m_{k-1}$ , and the bias  $b^{[k]} \in \mathbb{R}^{m_k}$ . Moreover,  $W^{[1]} \in \mathbb{R}^{m_1 \times d}$  and  $W^{[r]} \in \mathbb{R}^{1 \times m_{r-1}}$ .

The total number of neurons in the network is  $m_1 + \dots + m_r$ , and the total number of parameters in this network is  $(d+1)m_1 + (m_1+1)m_2 + \dots + (m_{r-1}+1)m_r$ .



Sometimes for notational consistency we also write  $a^{[0]} = x$ , and  $a^{[r]} = h_\theta(x)$ . Then we have simple recursion that

$$a^{[k]} = \text{ReLU}(W^{[k]}a^{[k-1]} + b^{[k]}), \forall k = 1, \dots, r-1 \quad (7.18)$$

Note that this would have been true for  $k = r$  if there were an additional ReLU in equation (7.17), but often people like to make the last layer linear (aka without a ReLU) so that negative outputs are possible and it's easier to interpret the last layer as a linear model. (More on the interpretability at the “connection to kernel method” paragraph of this section.)

**Other activation functions.** The activation function ReLU can be replaced by many other non-linear function  $\sigma(\cdot)$  that maps  $\mathbb{R}$  to  $\mathbb{R}$  such as

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (\text{sigmoid}) \quad (7.19)$$

$$\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad (\text{tanh}) \quad (7.20)$$

**Why do we not use the identity function for  $\sigma(z)$ ?** That is, why not use  $\sigma(z) = z$ ? Assume for sake of argument that  $b^{[1]}$  and  $b^{[2]}$  are zeros. Suppose  $\sigma(z) = z$ , then for two-layer neural network, we have that

$$h_\theta(x) = W^{[2]}a^{[1]} \quad (7.21)$$

$$= W^{[2]}\sigma(z^{[1]}) \quad \text{by definition} \quad (7.22)$$

$$= W^{[2]}z^{[1]} \quad \text{since } \sigma(z) = z \quad (7.23)$$

$$= W^{[2]}W^{[1]}x \quad \text{from Equation (7.13)} \quad (7.24)$$

$$= \tilde{W}x \quad \text{where } \tilde{W} = W^{[2]}W^{[1]} \quad (7.25)$$

Notice how  $W^{[2]}W^{[1]}$  collapsed into  $\tilde{W}$ .

This is because applying a linear function to another linear function will result in a linear function over the original input (i.e., you can construct a  $\tilde{W}$  such that  $\tilde{W}x = W^{[2]}W^{[1]}x$ ). This loses much of the representational power of the neural network as often times the output we are trying to predict has a non-linear relationship with the inputs. Without non-linear activation functions, the neural network will simply perform linear regression.

**Connection to the Kernel Method.** In the previous lectures, we covered the concept of feature maps. Recall that the main motivation for feature

maps is to represent functions that are non-linear in the input  $x$  by  $\theta^\top \phi(x)$ , where  $\theta$  are the parameters and  $\phi(x)$ , the feature map, is a handcrafted function non-linear in the raw input  $x$ . The performance of the learning algorithms can significantly depends on the choice of the feature map  $\phi(x)$ . Oftentimes people use domain knowledge to design the feature map  $\phi(x)$  that suits the particular applications. The process of choosing the feature maps is often referred to as **feature engineering**.

We can view deep learning as a way to automatically learn the right feature map (sometimes also referred to as “the representation”) as follows. Suppose we denote by  $\beta$  the collection of the parameters in a fully-connected neural networks (equation (7.17)) except those in the last layer. Then we can abstract right  $a^{[r-1]}$  as a function of the input  $x$  and the parameters in  $\beta$ :  $a^{[r-1]} = \phi_\beta(x)$ . Now we can write the model as

$$h_\theta(x) = W^{[r]} \phi_\beta(x) + b^{[r]} \quad (7.26)$$

When  $\beta$  is fixed, then  $\phi_\beta(\cdot)$  can viewed as a feature map, and therefore  $h_\theta(x)$  is just a linear model over the features  $\phi_\beta(x)$ . However, we will train the neural networks, both the parameters in  $\beta$  and the parameters  $W^{[r]}, b^{[r]}$  are optimized, and therefore we are not learning a linear model in the feature space, but also learning a good feature map  $\phi_\beta(\cdot)$  itself so that it’s possible to predict accurately with a linear model on top of the feature map. Therefore, deep learning tends to depend less on the domain knowledge of the particular applications and requires often less feature engineering. The penultimate layer  $a^{[r]}$  is often (informally) referred to as the learned features or representations in the context of deep learning.

In the example of house price prediction, a fully-connected neural network does not need us to specify the intermediate quantity such “family size”, and may automatically discover some useful features in the last penultimate layer (the activation  $a^{[r-1]}$ ), and use them to linearly predict the housing price. Often the feature map / representation obtained from one datasets (that is, the function  $\phi_\beta(\cdot)$  can be also useful for other datasets, which indicates they contain essential information about the data. However, oftentimes, the neural network will discover complex features which are very useful for predicting the output but may be difficult for a human to understand or interpret. This is why some people refer to neural networks as a *black box*, as it can be difficult to understand the features it has discovered.

## 7.3 Backpropagation

In this section, we introduce backpropagation or auto-differentiation, which computes the gradient of the loss  $\nabla J^{(j)}(\theta)$  efficiently. We will start with an informal theorem that states that as long as a real-valued function  $f$  can be efficiently computed/evaluated by a differentiable network or circuit, then its gradient can be efficiently computed in a similar time. We will then show how to do this concretely for fully-connected neural networks.

Because the formality of the general theorem is not the main focus here, we will introduce the terms with informal definitions. By a differentiable circuit or a differentiable network, we mean a composition of a sequence of differentiable arithmetic operations (additions, subtraction, multiplication, divisions, etc) and elementary differentiable functions (ReLU, exp, log, sin, cos, etc.). Let the size of the circuit be the total number of such operations and elementary functions. We assume that each of the operations and functions, and their derivatives or partial derivatives can be computed in  $O(1)$  time in the computer.

**Theorem 7.3.1:** *[backpropagation or auto-differentiation, informally stated] Suppose a differentiable circuit of size  $N$  computes a real-valued function  $f : \mathbb{R}^\ell \rightarrow \mathbb{R}$ . Then, the gradient  $\nabla f$  can be computed in time  $O(N)$ , by a circuit of size  $O(N)$ .<sup>4</sup>*

We note that the loss function  $J^{(j)}(\theta)$  for  $j$ -th example can be indeed computed by a sequence of operations and functions involving additions, subtraction, multiplications, and non-linear activations. Thus the theorem suggests that we should be able to compute the  $\nabla J^{(j)}(\theta)$  in a similar time to that for computing  $J^{(j)}(\theta)$  itself. This does not only apply to the fully-connected neural network introduced in the Section 7.2, but also many other types of neural networks.

In the rest of the section, we will showcase how to compute the gradient of the loss efficiently for fully-connected neural networks using backpropagation. Even though auto-differentiation or backpropagation is implemented in all the deep learning packages such as tensorflow and pytorch, understanding it is very helpful for gaining insights into the working of deep learning.

---

<sup>4</sup>We note if the output of the function  $f$  does not depend on some of the input coordinates, then we set by default the gradient w.r.t that coordinate to zero. Setting to zero does not count towards the total runtime here in our accounting scheme. This is why when  $N \leq \ell$ , we can compute the gradient in  $O(N)$  time, which might be potentially even less than  $\ell$ .

### 7.3.1 Preliminary: chain rule

We first recall the chain rule in calculus. Suppose the variable  $J$  depends on the variables  $\theta_1, \dots, \theta_p$  via the intermediate variable  $g_1, \dots, g_k$ :

$$g_j = g_j(\theta_1, \dots, \theta_p), \forall j \in \{1, \dots, k\} \quad (7.27)$$

$$J = J(g_1, \dots, g_k) \quad (7.28)$$

Here we overload the meaning of  $g_j$ 's: they denote both the intermediate variables but also the functions used to compute the intermediate variables. Then, by the chain rule, we have that  $\forall i$ ,

$$\frac{\partial J}{\partial \theta_i} = \sum_{j=1}^k \frac{\partial J}{\partial g_j} \frac{\partial g_j}{\partial \theta_i} \quad (7.29)$$

For the ease of invoking the chain rule in the following subsections in various ways, we will call  $J$  the output variable,  $g_1, \dots, g_k$  intermediate variables, and  $\theta_1, \dots, \theta_p$  the input variable in the chain rule.

### 7.3.2 One-neuron neural networks

**Simplifying notations:** In the rest of the section, we will consider a generic input  $x$  and compute the gradient of  $h_\theta(x)$  w.r.t  $\theta$ . For simplicity, we use  $o$  as a shorthand for  $h_\theta(x)$  ( $o$  stands for *output*). For simplicity, with slight abuse of notation, we use  $J = \frac{1}{2}(y - o)^2$  to denote the loss function. (Note that this overrides the definition of  $J$  as the total loss in Section 7.1.) Our goal is to compute the derivative of  $J$  w.r.t the parameter  $\theta$ .

We first consider the neural network with one neuron defined in equation (7.7). Recall that we compute the loss function via the following sequential steps:

$$z = w^\top x + b \quad (7.30)$$

$$o = \text{ReLU}(z) \quad (7.31)$$

$$J = \frac{1}{2}(y - o)^2 \quad (7.32)$$

By the chain rule with  $J$  as the output variable,  $o$  as the intermediate variable, and  $w_i$  the input variable, we have that

$$\frac{\partial J}{\partial w_i} = \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial w_i} \quad (7.33)$$

Invoking the chain rule with  $o$  as the output variable,  $z$  as the intermediate variable, and  $w_i$  the input variable, we have that

$$\frac{\partial o}{\partial w_i} = \frac{\partial o}{\partial z} \cdot \frac{\partial z}{\partial w_i}$$

Combining the equation above with equation (7.33), we have

$$\begin{aligned} \frac{\partial J}{\partial w_i} &= \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial z} \cdot \frac{\partial z}{\partial w_i} = (o - y) \cdot 1\{z \geq 0\} \cdot x_i \\ &\quad (\text{because } \frac{\partial J}{\partial o} = (o - y) \text{ and } \frac{\partial o}{\partial z} = 1\{z \geq 0\} \text{ and } \frac{\partial z}{\partial w_i} = x_i) \end{aligned}$$

Here, the key is that we reduce the computation of  $\frac{\partial J}{\partial w_i}$  to the computation of three simpler more “local” objects  $\frac{\partial J}{\partial o}$ ,  $\frac{\partial o}{\partial z}$ , and  $\frac{\partial z}{\partial w_i}$ , which are much simpler to compute because  $J$  directly depends on  $o$  via equation (7.32),  $o$  directly depends on  $a$  via equation (7.31), and  $z$  directly depends on  $w_i$  via equation (7.30). Note that in a vectorized form, we can also write

$$\nabla_w J = (o - y) \cdot 1\{z \geq 0\} \cdot x$$

Similarly, we compute the gradient w.r.t  $b$  by

$$\begin{aligned} \frac{\partial J}{\partial b} &= \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial z} \cdot \frac{\partial z}{\partial b} = (o - y) \cdot 1\{z \geq 0\} \\ &\quad (\text{because } \frac{\partial J}{\partial o} = (o - y) \text{ and } \frac{\partial o}{\partial z} = 1\{z \geq 0\} \text{ and } \frac{\partial z}{\partial b} = 1) \end{aligned}$$

### 7.3.3 Two-layer neural networks: a low-level unpacked computation

**Note:** this subsection derives the derivatives with low-level notations to help you build up intuition on backpropagation. If you are looking for a clean formula, or you are familiar with matrix derivatives, then feel free to jump to the next subsection directly.

Now we consider the two-layer neural network defined in equation (7.11). We compute the loss  $J$  by following sequence of operations

$$\begin{aligned} \forall j \in [1, \dots, m], \quad & z_j = w_j^{[1]\top} x + b_j^{[1]} \text{ where } w_j^{[1]} \in \mathbb{R}^d, b_j^{[1]} \\ & a_j = \text{ReLU}(z_j), \\ & a = [a_1, \dots, a_m]^\top \in \mathbb{R}^m \\ & o = w^{[2]\top} a + b^{[2]} \text{ where } w^{[2]} \in \mathbb{R}^m, b^{[2]} \in \mathbb{R} \\ & J = \frac{1}{2}(y - o)^2 \end{aligned} \tag{7.34}$$

We will use  $(w^{[2]})_\ell$  to denote the  $\ell$ -th coordinate of  $w^{[2]}$ , and  $(w_j^{[1]})_\ell$  to denote the  $\ell$ -coordinate of  $w_j^{[1]}$ . (We will avoid using these cumbersome notations once we figure out how to write everything in matrix and vector forms.)

By invoking chain rule with  $J$  as the output variable,  $o$  as intermediate variable, and  $(w^{[2]})_\ell$  as the input variable, we have

$$\begin{aligned}\frac{\partial J}{\partial (w^{[2]})_\ell} &= \frac{\partial J}{\partial o} \frac{\partial o}{\partial (w^{[2]})_\ell} \\ &= (o - y) \frac{\partial o}{\partial (w^{[2]})_\ell} \\ &= (o - y) a_\ell\end{aligned}$$

It's more challenging to compute  $\frac{\partial J}{\partial (w_j^{[1]})_\ell}$ . Towards computing it, we first invoke the chain rule with  $J$  as the output variable,  $z_j$  as the intermediate variable, and  $(w_j^{[1]})_\ell$  as the input variable.

$$\begin{aligned}\frac{\partial J}{\partial (w_j^{[1]})_\ell} &= \frac{\partial J}{\partial z_j} \cdot \frac{\partial z_j}{\partial (w_j^{[1]})_\ell} \\ &= \frac{\partial J}{\partial z_j} \cdot x_\ell \quad \text{(because } \frac{\partial z_j}{\partial (w_j^{[1]})_\ell} = x_\ell \text{.)}\end{aligned}$$

Thus, it suffices to compute the  $\frac{\partial J}{\partial z_j}$ . We invoke the chain rule with  $J$  as the output variable,  $a_j$  as the intermediate variable, and  $z_j$  as the input variable,

$$\begin{aligned}\frac{\partial J}{\partial z_j} &= \frac{\partial J}{\partial a_j} \frac{\partial a_j}{\partial z_j} \\ &= \frac{\partial J}{\partial a_j} 1\{z_j \geq 0\}\end{aligned}$$

Now it suffices to compute  $\frac{\partial J}{\partial a_j}$ , and we invoke the chain rule with  $J$  as the output variable,  $o$  as the intermediate variable, and  $a_j$  as the input variable,

$$\begin{aligned}\frac{\partial J}{\partial a_j} &= \frac{\partial J}{\partial o} \frac{\partial o}{\partial a_j} \\ &= (o - y) \cdot (w^{[2]})_j\end{aligned}$$

Now combining the equations above, we obtain

$$\frac{\partial J}{\partial (w_j^{[1]})_\ell} = (o - y) \cdot (w^{[2]})_j 1\{z_j \geq 0\} x_\ell$$

Next we gauge the runtime of computing these partial derivatives. Let  $p$  denotes the total number of parameters in the network. We note that  $p \geq md$  where  $m$  is the number of hidden units and  $d$  is the input dimension. For every  $j$  and  $\ell$ , to compute  $\frac{\partial J}{\partial (w_j^{[1]})_\ell}$ , apparently we need to compute at least the output  $o$ , which takes at least  $p \geq md$  operations. Therefore at the first glance computing a single gradient takes at least  $md$  time, and the total time to compute the derivatives w.r.t to all the parameters is at least  $(md)^2$ , which is inefficient.

However, the key of the backpropagation is that for different choices of  $\ell$ , the formulas above for computing  $\frac{\partial J}{\partial (w_j^{[1]})_\ell}$  share many terms, such as,  $(o - y)$ ,  $(w^{[2]})_j$  and  $1\{z_j \geq 0\}$ . This suggests that we can re-organize the computation to leverage the shared computation.

It turns out the crucial shared quantities in these formulas are  $\frac{\partial J}{\partial o}$ ,  $\frac{\partial J}{\partial z_1}, \dots, \frac{\partial J}{\partial z_m}$ . We now write the following formulas to compute the gradients efficiently in Algorithm 3.

---

**Algorithm 3** Backpropagation for two-layer neural networks

---

- 1: Compute the values of  $z_1, \dots, z_m$ ,  $a_1, \dots, a_m$  and  $o$  as in the definition of neural network (equation (7.34)).
- 2: Compute  $\frac{\partial J}{\partial o} = (o - y)$ .
- 3: Compute  $\frac{\partial J}{\partial z_j}$  for  $j = 1, \dots, m$  by

$$\frac{\partial J}{\partial z_j} = \frac{\partial J}{\partial o} \frac{\partial o}{\partial a_j} \frac{\partial a_j}{\partial z_j} = \frac{\partial J}{\partial o} \cdot (w^{[2]})_j \cdot 1\{z_j \geq 0\} \quad (7.35)$$

- 4: Compute  $\frac{\partial J}{\partial (w_j^{[1]})_\ell}$ ,  $\frac{\partial J}{\partial b_j^{[1]}}$ ,  $\frac{\partial J}{\partial (w^{[2]})_j}$ , and  $\frac{\partial J}{\partial b^{[2]}}$  by

$$\begin{aligned} \frac{\partial J}{\partial (w_j^{[1]})_\ell} &= \frac{\partial J}{\partial z_j} \cdot \frac{\partial z_j}{\partial (w_j^{[1]})_\ell} = \frac{\partial J}{\partial z_j} \cdot x_\ell \\ \frac{\partial J}{\partial b_j^{[1]}} &= \frac{\partial J}{\partial z_j} \cdot \frac{\partial z_j}{\partial b_j^{[1]}} = \frac{\partial J}{\partial z_j} \\ \frac{\partial J}{\partial (w^{[2]})_j} &= \frac{\partial J}{\partial o} \frac{\partial o}{\partial (w^{[2]})_j} = \frac{\partial J}{\partial o} \cdot a_j \\ \frac{\partial J}{\partial b^{[2]}} &= \frac{\partial J}{\partial o} \frac{\partial o}{\partial b^{[2]}} = \frac{\partial J}{\partial o} \end{aligned}$$


---

### 7.3.4 Two-layer neural network with vector notation

As we have done before in the definition of neural networks, the equations for backpropagation becomes much cleaner with proper matrix notation. Here we state the algorithm first and also provide a cleaner proof via matrix calculus.

Let

$$\begin{aligned}\delta^{[2]} &\triangleq \frac{\partial J}{\partial o} \in \mathbb{R} \\ \delta^{[1]} &\triangleq \frac{\partial J}{\partial z} \in \mathbb{R}^m\end{aligned}\tag{7.36}$$

Here we note that when  $A$  is a real-valued variable,<sup>5</sup> and  $B$  is a vector or matrix variable, then  $\frac{\partial A}{\partial B}$  denotes the collection of the partial derivatives with the same shape as  $B$ .<sup>6</sup> In other words, if  $B$  is a matrix of dimension  $m \times d$ , then  $\frac{\partial A}{\partial B}$  is a matrix in  $\mathbb{R}^{m \times d}$  with  $\frac{\partial A}{\partial B_{ij}}$  as the  $ij$ th-entry. Let  $v \odot w$  denote the entry-wise product of two vectors  $v$  and  $w$  of the same dimension. Now we are ready to describe backpropagation in Algorithm 4.

---

**Algorithm 4** Back-propagation for two-layer neural networks in vectorized notations.

---

- 1: Compute the values of  $z \in \mathbb{R}^m$ ,  $a \in \mathbb{R}^m$ , and  $o$
- 2: Compute  $\delta^{[2]} = (o - y) \in \mathbb{R}$
- 3: Compute  $\delta^{[1]} = (o - y) \cdot W^{[2]\top} \odot 1\{z \geq 0\} \in \mathbb{R}^{m \times 1}$
- 4: Compute

$$\begin{aligned}\frac{\partial J}{\partial W^{[2]}} &= \delta^{[2]} a^\top \in \mathbb{R}^{1 \times m} \\ \frac{\partial J}{\partial b^{[2]}} &= \delta^{[2]} \in \mathbb{R} \\ \frac{\partial J}{\partial W^{[1]}} &= \delta^{[1]} x^\top \in \mathbb{R}^{m \times d} \\ \frac{\partial J}{\partial b^{[1]}} &= \delta^{[1]} \in \mathbb{R}^m\end{aligned}$$


---

<sup>5</sup>We will avoid using the notation  $\frac{\partial A}{\partial B}$  for  $A$  that is not a real-valued variable.

<sup>6</sup>If you are familiar with the notion of total derivatives, we note that the dimensionality here is different from that for total derivatives.



**Derivation using the chain rule for matrix multiplication.** To have a succinct derivation of the backpropagation algorithm in Algorithm 4 without working with the complex indices, we state the extensions of the chain rule in vectorized notations. It requires more knowledge of matrix calculus to state the most general result, and therefore we will introduce a few special cases that are most relevant for deep learning. Suppose  $J$  is a real-valued output variable,  $z \in \mathbb{R}^m$  is the intermediate variable and  $W \in \mathbb{R}^{m \times d}, u \in \mathbb{R}^d$  are the input variables. Suppose they satisfy:

$$\begin{aligned} z &= Wu + b, \text{ where } W \in \mathbb{R}^{m \times d} \\ J &= J(z) \end{aligned} \tag{7.37}$$

Then we can compute  $\frac{\partial J}{\partial u}$  and  $\frac{\partial J}{\partial W}$  by:

$$\frac{\partial J}{\partial u} = W^\top \frac{\partial J}{\partial z} \tag{7.38}$$

$$\frac{\partial J}{\partial W} = \frac{\partial J}{\partial z} \cdot u^\top \tag{7.39}$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial z} \tag{7.40}$$

We can verify the dimensionality is indeed compatible because  $\frac{\partial J}{\partial z} \in \mathbb{R}^m$ ,  $W^\top \in \mathbb{R}^{d \times m}$ ,  $\frac{\partial J}{\partial u} \in \mathbb{R}^d$ ,  $\frac{\partial J}{\partial W} \in \mathbb{R}^{m \times d}$ ,  $u^\top \in \mathbb{R}^{1 \times d}$ .

Here the chain rule in equation (7.38) only works for the special cases where  $z = Wu$ . Another useful case is the following:

$$\begin{aligned} a &= \sigma(z), \text{ where } \sigma \text{ is an element-wise activation, } z, a \in \mathbb{R}^d \\ J &= J(a) \end{aligned}$$

Then, we have that

$$\frac{\partial J}{\partial z} = \frac{\partial J}{\partial a} \odot \sigma'(z) \tag{7.41}$$

where  $\sigma'(\cdot)$  is the element-wise derivative of the activation function  $\sigma$ , and  $\odot$  is element-wise product of two vectors of the same dimensionality.

Using equation (7.38), (7.39), and (7.41), we can verify the correctness of Algorithm 4. Indeed, using the notations in the two-layer neural network

$$\begin{aligned} \frac{\partial J}{\partial z} &= \frac{\partial J}{\partial a} \odot \text{ReLU}'(z) && \left( \begin{array}{l} \text{by invoking equation (7.41) with setting} \\ J \leftarrow J, a \leftarrow a, z \leftarrow a, \sigma \leftarrow \text{ReLU.} \end{array} \right) \\ &= (o - y)W^{[2]\top} \odot \text{ReLU}'(z) && \left( \begin{array}{l} \text{by invoking equation (7.38) with setting} \\ J \leftarrow J, z \leftarrow o, W \leftarrow W^{[2]}, u \leftarrow a, b \leftarrow b^{[2]} \end{array} \right) \end{aligned}$$

Therefore,  $\delta^{[1]} = \frac{\partial J}{\partial z}$ , and we verify the correctness of Line 3 in Algorithm 4. Similarly, let's verify the third equation in Line 4,

$$\begin{aligned} \frac{\partial J}{\partial W^{[1]}} &= \frac{\partial J}{\partial z} \cdot x^\top && \left( \begin{array}{l} \text{by invoking equation (7.39) with setting} \\ J \leftarrow J, z \leftarrow z, W \leftarrow W^{[1]}, u \leftarrow x, b \leftarrow b^{[1]} \end{array} \right) \\ &= \delta^{[1]} x^\top && \text{(because we have proved } \delta^{[1]} = \frac{\partial J}{\partial z} \text{)} \end{aligned}$$

### 7.3.5 Multi-layer neural networks

In this section, we will derive the backpropagation algorithms for the model defined in (7.17). Recall that we have

$$\begin{aligned} a^{[1]} &= \text{ReLU}(W^{[1]}x + b^{[1]}) \\ a^{[2]} &= \text{ReLU}(W^{[2]}a^{[1]} + b^{[2]}) \\ &\dots \\ a^{[r-1]} &= \text{ReLU}(W^{[r-1]}a^{[r-2]} + b^{[r-1]}) \\ a^{[r]} &= z^{[r]} = W^{[r]}a^{[r-1]} + b^{[r]} \\ J &= \frac{1}{2}(a^{[r]} - y)^2 \end{aligned}$$

Here we define both  $a^{[r]}$  and  $z^{[r]}$  as  $h_\theta(x)$  for notational simplicity.

Define

$$\delta^{[k]} = \frac{\partial J}{\partial z^{[k]}} \tag{7.42}$$

The backpropagation algorithm computes  $\delta^{[k]}$ 's from  $k = r$  to 1, and computes  $\frac{\partial J}{\partial W^{[k]}}$  from  $\delta^{[k]}$  as described in Algorithm 5.

## 7.4 Vectorization Over Training Examples

As we discussed in Section 7.1, in the implementation of neural networks, we will leverage the parallelism across the multiple examples. This means that we will need to write the forward pass (the evaluation of the outputs) of the neural network and the backward pass (backpropagation) for multiple training examples in matrix notation.

**The basic idea.** The basic idea is simple. Suppose you have a training set with three examples  $x^{(1)}, x^{(2)}, x^{(3)}$ . The first-layer activations for each

---

**Algorithm 5** Back-propagation for multi-layer neural networks.

---

- 1: Compute and store the values of  $a^{[k]}$ 's and  $z^{[k]}$ 's for  $k = 1, \dots, r - 1$ , and  $J$ . ▷ This is often called the “forward pass”
- 2: Compute  $\delta^{[r]} = \frac{\partial J}{\partial z^{[r]}} = (z^{[r]} - o)$ .
- 3: **for**  $k = r - 1$  to 1 **do**
- 4:     Compute

$$\delta^{[k]} = \frac{\partial J}{\partial z^{[k]}} = \left( W^{[k+1]^\top \delta^{[k+1]} \right) \odot \text{ReLU}'(z^{[k]})$$

- 5:     Compute

$$\begin{aligned} \frac{\partial J}{\partial W^{[k+1]}} &= \delta^{[k+1]} a^{[k]^\top} \\ \frac{\partial J}{\partial b^{[k+1]}} &= \delta^{[k+1]} \end{aligned}$$


---

example are as follows:

$$\begin{aligned} z^{[1](1)} &= W^{[1]}x^{(1)} + b^{[1]} \\ z^{[1](2)} &= W^{[1]}x^{(2)} + b^{[1]} \\ z^{[1](3)} &= W^{[1]}x^{(3)} + b^{[1]} \end{aligned}$$

Note the difference between square brackets  $[\cdot]$ , which refer to the layer number, and parenthesis  $(\cdot)$ , which refer to the training example number. Intuitively, one would implement this using a for loop. It turns out, we can vectorize these operations as well. First, define:

$$X = \begin{bmatrix} \begin{array}{c} | \\ x^{(1)} \\ | \end{array} & \begin{array}{c} | \\ x^{(2)} \\ | \end{array} & \begin{array}{c} | \\ x^{(3)} \\ | \end{array} \end{bmatrix} \in \mathbb{R}^{d \times 3} \quad (7.43)$$

Note that we are stacking training examples in columns and *not* rows. We can then combine this into a single unified formulation:

$$Z^{[1]} = \begin{bmatrix} \begin{array}{c} | \\ z^{[1](1)} \\ | \end{array} & \begin{array}{c} | \\ z^{[1](2)} \\ | \end{array} & \begin{array}{c} | \\ z^{[1](3)} \\ | \end{array} \end{bmatrix} = W^{[1]}X + b^{[1]} \quad (7.44)$$

You may notice that we are attempting to add  $b^{[1]} \in \mathbb{R}^{4 \times 1}$  to  $W^{[1]}X \in \mathbb{R}^{4 \times 3}$ . Strictly following the rules of linear algebra, this is not allowed. In

practice however, this addition is performed using *broadcasting*. We create an intermediate  $\tilde{b}^{[1]} \in \mathbb{R}^{4 \times 3}$ :

$$\tilde{b}^{[1]} = \begin{bmatrix} | & | & | \\ b^{[1]} & b^{[1]} & b^{[1]} \\ | & | & | \end{bmatrix} \quad (7.45)$$

We can then perform the computation:  $Z^{[1]} = W^{[1]}X + \tilde{b}^{[1]}$ . Often times, it is not necessary to explicitly construct  $\tilde{b}^{[1]}$ . By inspecting the dimensions in (7.44), you can assume  $b^{[1]} \in \mathbb{R}^{4 \times 1}$  is correctly broadcast to  $W^{[1]}X \in \mathbb{R}^{4 \times 3}$ .

The matricization approach as above can easily generalize to multiple layers, with one subtlety though, as discussed below.

**Complications/Subtlety in the Implementation.** All the deep learning packages or implementations put the data points in the rows of a data matrix. (If the data point itself is a matrix or tensor, then the data are concentrated along the zero-th dimension.) However, most of the deep learning papers use a similar notation to these notes where the data points are treated as column vectors.<sup>7</sup> There is a simple conversion to deal with the mismatch: in the implementation, all the columns become row vectors, row vectors become column vectors, all the matrices are transposed, and the orders of the matrix multiplications are flipped. In the example above, using the row major convention, the data matrix is  $X \in \mathbb{R}^{3 \times d}$ , the first layer weight matrix has dimensionality  $d \times m$  (instead of  $m \times d$  as in the two layer neural net section), and the bias vector  $b^{[1]} \in \mathbb{R}^{1 \times m}$ . The computation for the hidden activation becomes

$$Z^{[1]} = XW^{[1]} + b^{[1]} \in \mathbb{R}^{3 \times m} \quad (7.46)$$

---

<sup>7</sup>The instructor suspects that this is mostly because in mathematics we naturally multiply a matrix to a vector on the left hand side.

# Bibliography

- Mikhail Belkin, Daniel Hsu, and Ji Xu. Two models of double descent for weak features. *SIAM Journal on Mathematics of Data Science*, 2(4):1167–1180, 2020.
- Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. 2019.
- Preetum Nakkiran. More data can hurt for linear regression: Sample-wise double descent. 2019.
- Preetum Nakkiran, Prayaag Venkat, Sham Kakade, and Tengyu Ma. Optimal regularization can mitigate double descent. 2020.
- Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and Experiment*, 2021 (12):124003, 2021.
- Manfred Opper. Statistical mechanics of learning: Generalization. *The handbook of brain theory and neural networks*, pages 922–925, 1995.
- Manfred Opper. Learning to generalize. *Frontiers of Life*, 3(part 2):763–775, 2001.