

Notes on Machine Learning
from Andrew Ng's Coursera Course

Simon Zahn

August 29, 2016

Contents

Preface	ii
1 Introduction	1
1.1 Supervised Learning	1
1.2 Unsupervised Learning	2
2 Linear Regression	3
2.1 Univariate Linear Regression	3
2.1.1 The Hypothesis Function	3
2.1.2 The Cost Function	4
2.1.3 Gradient Descent	5
2.2 Multivariate Linear Regression	6
2.2.1 Gradient Descent for Multiple Variables	7
2.2.2 Feature Scaling	7
2.2.3 Tips for Gradient Descent	8
2.2.4 Polynomial Regression	9
2.3 Vectorized Equations	9
2.4 The Normal Equation	11
2.4.1 Normal Equation Noninvertibility	11
3 Logistic Regression	12
3.1 Binary Classification	12
3.2 Hypothesis Representation	14
3.2.1 Interpretation of the Logistic Hypothesis Function	14
3.3 Decision Boundary	15
A Notation	16
B Linear Algebra Review	17

Preface

Hello, this is a preface.

Chapter 1

Introduction

Two definitions of Machine Learning are offered at the start of the course:

Arthur Samuel ‘The field of study that gives computers the ability to learn without being explicitly programmed.’

Tom Mitchell ‘A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks T , as measured by P , improves with experience E .’

1.1 Supervised Learning

In supervised learning, we have a data set and already know what the correct output should be, knowing that there is an existing relationship between the input and output. There are two types of supervised learning: **classification**, and **regression**.

In a regression problem, we’re trying to predict results with a continuous output; i.e. we are trying to map input variables to some continuous function. In a classification problem, we are instead trying to predict results with a discrete output. Let’s look at some examples:

Regression If we’re trying to predict the price of a house given data about the house such as its size, location, etc., this is a regression problem.

Regression Given a picture of a person, try to predict his/her age.

Classification Given a picture of a person, try to predict if he/she is of high school, college, or graduate age.

Classification As a bank, decide whether or not to give a loan to a potential borrower.

1.2 Unsupervised Learning

Unsupervised learning allows us to approach problems with little or no idea about what our results should look like. We can try and derive structure from data where we don't necessarily know the effect of the variables. We can derive this structure by clustering the data based on relationships among the variables in the data. With unsupervised learning, there is no feedback based on the prediction results.

One example of unsupervised learning is clustering. Imagine we take 1000 essays written by US college students. We can try and automatically group these essays into a smaller number that are somehow similar or related by different variables; word frequency, sentence length, page count, etc.

Here is an example of unsupervised learning that isn't clustering: the cocktail party algorithm. This is a way to find structure in messy data, such as the identification of individual voices and music from a mesh of sounds at a cocktail party.¹

¹For more information about auditory filtering, look at Wikipedia's [Cocktail Party Effect](#) article.

Chapter 2

Linear Regression

In regression problems, we take a variable (or multiple variables) as input, and try to fit the output to a continuous expected result function.

2.1 Univariate Linear Regression

In univariate linear regression, we want to predict a single output value \hat{y} from a single input value x . Since this is supervised learning, we already have an idea about what the input/output relationship should look like.

2.1.1 The Hypothesis Function

Imagine we have a problem where the input is x and the output is y . In order to do machine learning, there should exist a relationship (a pattern) between the input and output variables. Let's say this function is $y = f(x)$. In this situation, f is known as the target function. However, this function f is unknown to us, so we need to try and guess what it is. To do that, we form a *hypothesis* function $h(x)$ that approximates the unknown $f(x)$.

For single variable linear regression, our hypothesis function takes two parameters: θ_0 and θ_1 . As such, we often write it as $h_\theta(x)$, and it takes the form

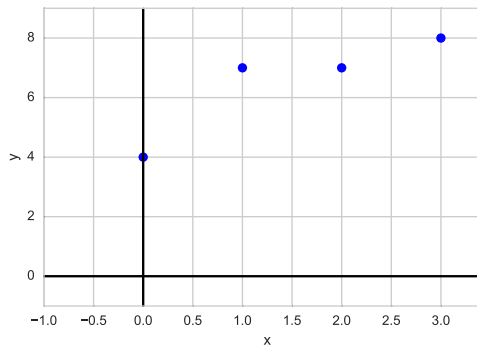
$$\hat{y} = h_\theta(x) = \theta_0 + \theta_1 x \quad (2.1)$$

Note that this is the equation of a straight line ($y = mx + b$). We're trying to find the values for θ_0 and θ_1 to get our estimated output \hat{y} . In other words, we're trying to determine the function h_θ that maps our input data (the x 's) to our output data (the y 's).

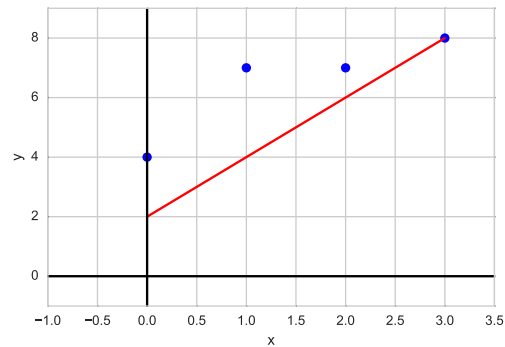
Suppose we have the following set of training data:

Input x	Output y
0	4
1	7
2	7
3	8

We can plot these points, as shown in Figure 2.1a. Let's make a random guess at our hypothesis function: $\theta_0 = 2$ and $\theta_1 = 2$, making our hypothesis function $h_\theta(x) = 2 + 2x$, as shown in Figure 2.1b.



(a) Plotting our example points.



(b) Plotting our example points.

Using this hypothesis function for $x = 1$, we have $\hat{y} = h_\theta(1) = 2 + 2 \cdot 1 = 4$. In this case, $\hat{y} = 4$, but $y = 7$, so maybe this isn't the best fit hypothesis.

2.1.2 The Cost Function

The cost function,¹ is a function used for parameter estimation, where the input to the cost function is some function of the difference between estimated and the true values for an instance of data. In this case, we can use the cost function to measure the accuracy of our hypothesis function.

The cost function looks at something similar to an average² of all the results of the hypothesis with inputs from the x 's compared to the actual output y 's. We define our cost function as follows:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 \quad (2.2)$$

This is known as the **mean squared error**. If we set \bar{x} equal to the mean of the squares all the $h_\theta(x_i) - y_i$, then the cost function is just the mean of \bar{x} . The term $\frac{1}{2m}$ is merely a convenience for the computation of gradient descent, which we'll see very shortly.

¹The cost function can also be called the loss function.

²It's actually something a bit fancier than a standard average.

2.1.3 Gradient Descent

We now have our hypothesis function defined, as well as a way of measuring how well it fits the data. Now, we have to estimate the parameters in the hypothesis function, and that's where gradient descent comes in.

Let's graph our cost function as a function of the parameter estimates. This can be somewhat confusing, as we are moving up to a higher level of abstraction. We are not graphing x and y itself, but the parameter range of our hypothesis function and the cost resulting from selecting particular sets of parameters. We put θ_0 on the x -axis, and θ_1 on the y -axis, with the cost function on the vertical z -axis.

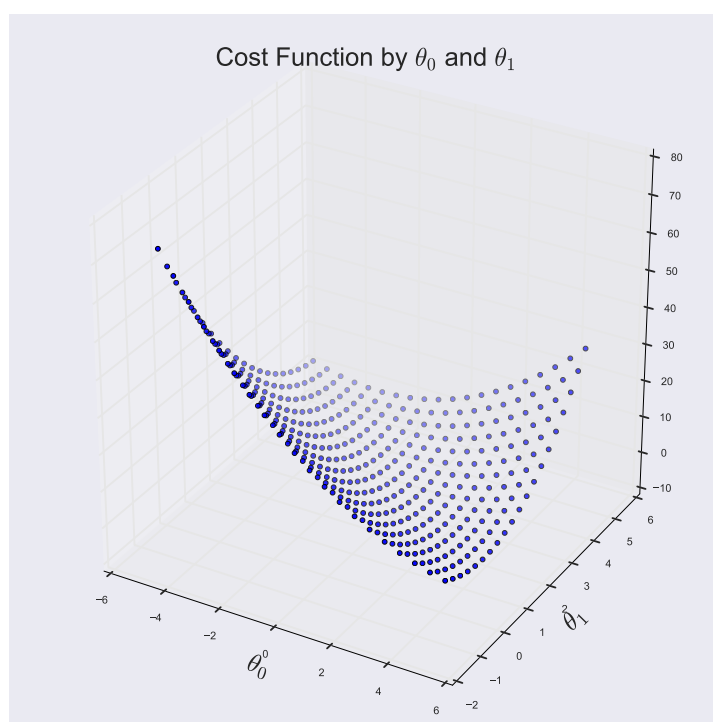


Figure 2.2: Plot of the cost function $J(\theta_0, \theta_1)$ using our hypothesis $h_\theta(x)$.

Our goal is to take the parameters θ_0 and θ_1 for when the cost function is at its minimum. We can calculate this value by taking the derivative of the cost function, which gives us direction of the steepest gradient to move towards. Take a step in that direction, and repeat. The step size is determined by the parameter α , which is called the **learning rate**. The gradient descent algorithm is:

Repeat until convergence:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (2.3)$$

where $j = 0, 1$ represents the feature index number.

Gradient Descent for Linear Regression

When specifically applied to the case of univariate linear regression, we can derive another form of the gradient descent equation. If we substitute our actual hypothesis function and cost function, we can modify the equation to

Repeat until convergence: {

$$\begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m ((h_{\theta}(x_i) - y_i) x_i) \end{aligned} \quad (2.4)$$

}

where m is the size of the training set, θ_0 is a constant that will be changing simultaneously with θ_1 , and x_1, y_i are values of the given training set. Note that we have separated out the two cases for θ_j into separate equations for θ_0 and θ_1 , and that for θ_1 we are multiplying x_i at the end due to the derivative.

2.2 Multivariate Linear Regression

Let's start by looking at some sample housing data with multiple features.

Size (feet ²)	# of Bedrooms	# of Floors	Age (years)	Price (in 1000's of \$)
x_1	x_2	x_3	x_4	y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
560	1	1	12	155

In this, we can introduce some notation:

- The variables x_1, x_2 , etc. are the features.
- The variable y is the output variable.
- The number of input features is denoted n . In this example, $n = 4$.
- m specifies the number of training examples (rows). Here, $m = 5$.
- $x^{(i)}$ is the input (features) of the i^{th} training example. So $x^{(2)}$ is the column vector $[1416, 3, 2, 40, 232]$.

- $x_j^{(i)}$ is feature j in the i^{th} training example. Here, $x_1^{(4)} = 852$.

At this point, we can define the multivariable form of the hypothesis function for linear regression:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \cdots + \theta_n x_n \quad (2.5)$$

For convenience of notation, we will define $x_0 = 1$ for all feature vectors ($x_0^{(i)} = 1$). So now, if we include x_0 , our hypothesis function takes the form:

$$h_{\theta}(x) = \sum_{i=0}^n \theta_i x_i \quad (2.6)$$

Now, we can also write the x values and θ values as vectors:

$$\vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \text{and} \quad \vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

In vector notation, this is

$$h_{\theta}(x) = \vec{\theta}^T \vec{x} \quad (2.7)$$

where we transpose $\vec{\theta}$ into a row vector so we're able to take the inner product.

Now that we have our vector $\vec{\theta} \in \mathbb{R}^{n+1}$, the cost function is

$$J(\vec{\theta}) = \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 \quad (2.8)$$

2.2.1 Gradient Descent for Multiple Variables

Using our expanded hypothesis and cost functions, the gradient descent algorithm becomes:

Repeat until convergence: {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right) \cdot x_j^{(i)} \quad \text{for } j = 0, 1, \dots, n \quad (2.9)$$

}

2.2.2 Feature Scaling

When features are in very different ranges, it can slow down gradient descent dramatically (and also mess up our machine learning algorithms!), because θ will descend quickly on small ranges and slowly on large ranges, and so will oscillate inefficiently down to

the minimum. The way to prevent this is to ensure that all the ranges are roughly the same, ideally:

$$-1 \leq x \leq 1$$

Two techniques to accomplish this are **feature scaling** and **mean normalization**. Feature scaling involved dividing the input values by the range (max value minus the min value) of the input variable, resulting in a new range of just 1.

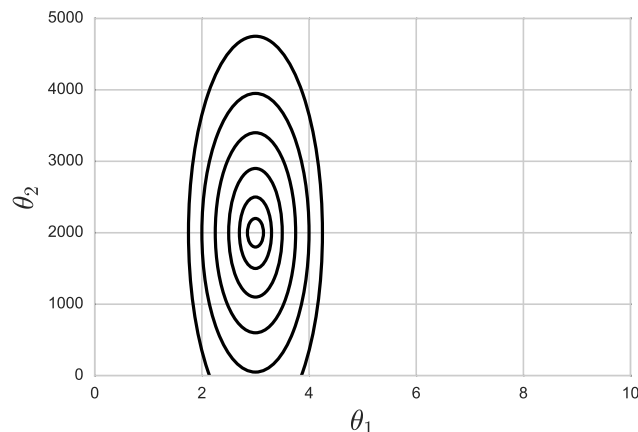


Figure 2.3: When one feature is on a much larger scale than the other, the plot of the cost function will be stretched out in the direction of the larger feature. Here, imagine that θ_1 is the number of bedrooms a house has, and θ_2 is the size in square feet.

Mean normalization involves subtracting the mean value for an input variable from the value for that input variable, resulting in a new mean of zero. To implement both of these simultaneously, use the following formula:

$$x_i := \frac{x_i - \mu_i}{s_i} \quad (2.10)$$

where μ_i is the average value of x_i , and s_i can either be the range ($x_{\max} - x_{\min}$) or the standard deviation.

2.2.3 Tips for Gradient Descent

Here are some of Professor Andrew Ng's tips on implementing gradient descent.

1. **Plot $J(\theta)$.** If you plot $J(\theta)$ as the ordinate and the number of iterations as the abscissa,³ the graph should be steadily decreasing with increasing number of iterations. If $J(\theta)$ ever increases, then α is probably too large.
2. If $J(\theta)$ decreases by less than E in one iteration, where E is some very small number, such as 10^{-3} , then you can declare convergence.

³On a Cartesian coordinate plane, the ordinate is the y -axis and the abscissa is the x -axis.

- For sufficiently small α , $J(\theta)$ should decrease with every iteration. To choose α , try a range of values for α with threefold increases, such as:

$$\cdots \rightarrow 0.001 \rightarrow 0.003 \rightarrow 0.01 \rightarrow 0.03 \rightarrow 0.1 \rightarrow 0.3 \rightarrow 1 \rightarrow 3 \rightarrow \cdots$$

- Sometimes, it's better to define new features instead of using the ones given. For example, if we have a house with features frontage⁴ and depth,⁵ you can combine these into a new feature called area, which is how much land the house sits on.

2.2.4 Polynomial Regression

The form of the hypothesis doesn't necessarily need to be linear if that doesn't fit the data well. We can change the behavior or curve of our hypothesis function by making it quadratic, cubic, square root, or some other form.

For example, if our hypothesis function is $h_\theta(x) = \theta_0 + \theta_1 x_1$, we can create additional features based on x_1 , to get the quadratic function $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$, or the cubic function $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$.

When thinking about nonlinear features, it is important to keep in mind that features scaling becomes even more essential than it was for linear regression. If x_1 has a range of 1 to 1000, the x_1^2 has a range of 1 to 1,000,000.

2.3 Vectorized Equations

Let's revisit our housing example from §2.2. Recall that we looked at the following example data, and we'll add an extra column for x_0 that always takes a value of one:

x_0	Size (feet ²) x_1	# Bedrooms x_2	# Floors x_3	Age (years) x_4	Price (in \$1000's) y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178
1	560	1	1	12	155

From this, we construct a matrix X that contains all of the features from the training data, and a vector \vec{y} of all the output data.

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \\ 1 & 560 & 1 & 1 & 12 \end{bmatrix} \quad \text{and} \quad \vec{y} = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \\ 155 \end{bmatrix}$$

⁴The width of the land in the front of the house.

⁵The width of the land on the side of the house.

Here, X is a $m \times (n + 1)$ matrix, and \vec{y} is a m -dimensional vector.

Let's go through this again, but this time in full abstraction. Say we have m examples $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$ and each $x^{(i)}$ has n features. Then, we have an $(n + 1)$ -dimensional feature vector:

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad (2.11)$$

The matrix X , which is also called the **design matrix**, is constructed by taking the transpose of each vector $x^{(i)}$. Each feature vector $x^{(i)}$ becomes a row in the design matrix. Just as previously, the output vector \vec{y} is obtained by taking all the labels and stacking them up into an m -dimensional vector, and the vector $\vec{\theta}$ is created from stacking all of the parameters for the hypothesis function.

$$X = \begin{bmatrix} \text{---} & \text{---} & (x^{(1)})^\top & \text{---} & \text{---} \\ \text{---} & \text{---} & (x^{(2)})^\top & \text{---} & \text{---} \\ \text{---} & \text{---} & (x^{(3)})^\top & \text{---} & \text{---} \\ \text{---} & \text{---} & (x^{(4)})^\top & \text{---} & \text{---} \\ \text{---} & \text{---} & \vdots & \text{---} & \text{---} \\ \text{---} & \text{---} & (x^{(m)})^\top & \text{---} & \text{---} \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ y^{(4)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \end{bmatrix} \quad (2.12)$$

Think back to the start of this section when we separated our table into the design matrix X and the output vector y . The design matrix is simply the data as stored in a table put into a matrix.

In our matrix notation for multivariate regression, the hypothesis function takes the form

$$h_\theta(X) = X\vec{\theta} \quad (2.13)$$

This will always work since X is an $m \times n$ matrix, and $\vec{\theta}$ is an $n \times 1$ vector. In a similar fashion, the cost function in matrix notation is

$$J(\vec{\theta}) = \frac{1}{2m} (X\vec{\theta} - \vec{y})^\top (X\vec{\theta} - \vec{y}) \quad (2.14)$$

The gradient descent rule can be expressed as

$$\vec{\theta} := \vec{\theta} - \alpha \nabla J(\theta) \quad (2.15)$$

There ∇ is the gradient (vector derivative) operator. If we solve this out using our vectorized hypothesis function, we get

$$\vec{\theta} := \vec{\theta} - \frac{\alpha}{m} X^\top (X\vec{\theta} - \vec{y}) \quad (2.16)$$

2.4 The Normal Equation

The normal equation is a method of solving for the optimal θ analytically, that is, without iteration. From calculus, if we want to find the minimum of a quadratic equation, we set the derivative equal to zero, and solve. We can apply the same logic to the cost function. If we take the partial derivative $\partial/\partial\theta_j J(\theta)$ and set this equal to zero for all values of j , we'll analytically solve for the minimum.

The derivation of the normal equation is fairly involved from a linear algebra perspective, so at this point just take it as a fact:

$$\vec{\theta} = (X^T X)^{-1} X^T \vec{y} \quad (2.17)$$

When deciding whether to use gradient descent or the normal equation, consider the following:

Gradient Descent	Normal Equation
Need to choose α	No need to choose α
Needs many iterations	No need to iterate
$O(kn^2)$	$O(n^3)$, need to calculate $X^T X$
Works well when n is large	Slow if n is very large

With the normal equation, computing the inverse has complexity $O(n^3)$. If we have a large number of features, this will cause the normal equation to perform slowly. In practice, when n exceeds 10,000, it would probably be a good idea to use gradient descent.

2.4.1 Normal Equation Noninvertibility

When implementing the normal equation, sometimes the matrix $X^T X$ is noninvertible. The common causes are:

- Redundant features, where two or more features are linearly dependent
- Too many features (i.e. $m \leq n$). In this case, delete some features or use regularization (which we'll get to later)

We typically avoid this problem by coding a pseudoinverse, instead of taking the actual inverse.

Chapter 3

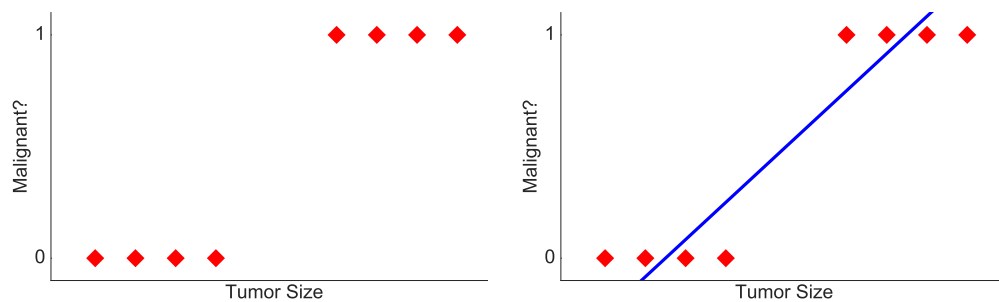
Logistic Regression

Now we turn away from regression to classification problems. Don't be confused by the name 'logistic regression,' it's actually just named for the mathematical function and it's a common approach to classification.

3.1 Binary Classification

In classification problems, instead of our output being continuous, we expect it to fall into discrete classes. We'll start with the simplest case: binary classification. Here, we have our output variable $y \in \{0, 1\}$. Typically, we take 0 as the negative class and 1 as the positive class.

Consider the following example: we have a sample of eight tumors, and we want to determine if they're malignant based on the tumor size. These are plotted in Figure 3.1a. One thing we can do is assume a linear relationship with hypothesis $h_{\theta}(\vec{x}) = \vec{\theta}^T \vec{x}$. This shown in Figure 3.1b.



(a) Plot of tumors by size.

(b) Plot of tumors by size with linear regression line $y = \frac{x-1}{6} - 0.25$.

To try and make predictions, we can threshold the output at $h_{\theta}(x) = 0.5$, and then:

- If $h_{\theta}(x) \geq 0.5$, then predict $y = 1$

- If $h_{\theta}(x) < 0.5$, then predict $y = 0$

and you can see this in Figure 3.2.

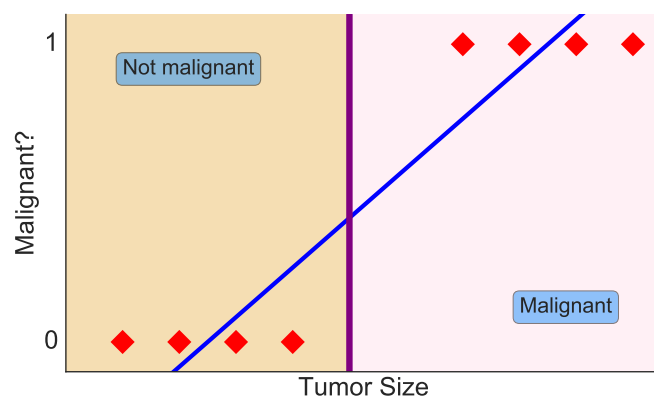


Figure 3.2: Linear regression plotted with classification regions.

In this example, it would seem like linear regression is a good classifier. However, what if we add a new data point for a large tumor. Suddenly, our results look like this

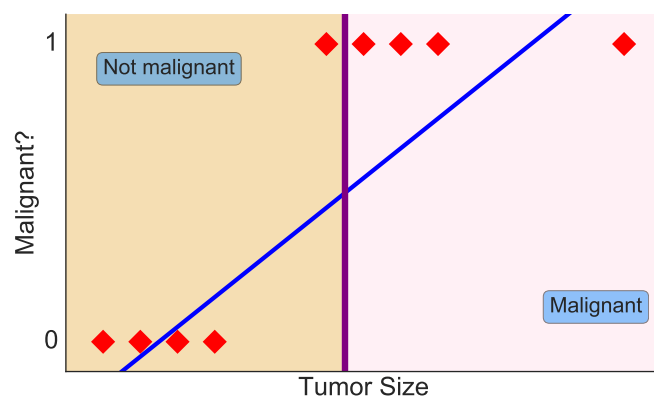


Figure 3.3: Linear regression plotted with classification regions after a new data point is added. Notice how one of the malignant tumors is now being misclassified as benign.

and now, we have a malignant tumor being misclassified as benign. Ergo, maybe linear regression isn't the best way to build a binary classifier.

3.2 Hypothesis Representation

In linear regression, our hypothesis was $h_{\theta}(\vec{x}) = \vec{\theta}^T \vec{x}$. For logistic regression, we want our hypothesis to satisfy $0 \leq h_{\theta}(x) \leq 1$. To do this, we use the sigmoid function.¹ To make this work, we modify our hypothesis to be

$$h_{\theta}(x) = g(\vec{\theta}^T \vec{x}) \quad (3.1)$$

where the function $g(z)$ is defined as

$$g(z) = \frac{1}{1 + e^{-z}} \quad (3.2)$$

Thus, to get the hypothesis function using the sigmoid function, just set $z = \vec{\theta}^T \vec{x}$.

The sigmoid function, shown in Figure 3.4, maps any real number onto the interval $(0, 1)$. This makes it immensely useful for transforming an arbitrary function for use with classification.

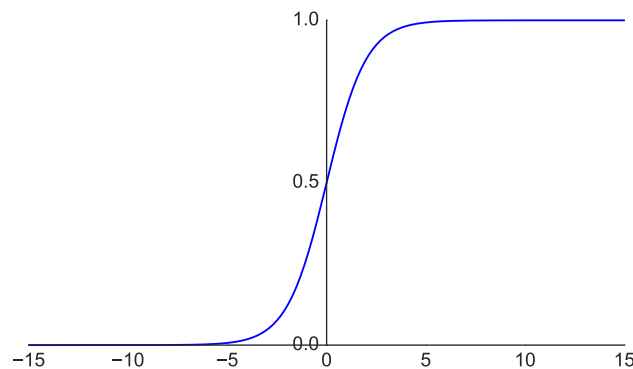


Figure 3.4: Sample plot of the sigmoid function.

3.2.1 Interpretation of the Logistic Hypothesis Function

When examining the hypothesis function output for logistic regression, we interpret $h_{\theta}(x)$ as the estimated probability that $y = 1$ on an input example x . For example, let's revisit the tumor size question from above. We have

$$\vec{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$$

If our hypothesis $h_{\theta}(x) = 0.7$, then we can tell the patient that there is a 70% chance of the tumor being malignant.

¹This is also called the logistic function, and is the namesake for logistic regression.

Slightly more formally, we interpret $h_\theta(x)$ as:²

$$h_\theta(x) = P(y = 1|x; \theta) \tag{3.3}$$

Thus, by the rules of probability:

$$P(y = 0|x; \theta) + P(y = 1|x; \theta) = 1 \tag{3.4}$$

$$P(y = 0|x; \theta) = 1 - P(y = 1|x; \theta) \tag{3.5}$$

3.3 Decision Boundary

²This is read as "the probability that $y = 1$, given x , parameterized by θ ."

Appendix A

Notation

$\hat{}$ predicted output
 \hat{y} predicted output of the variable y

\mathbb{Z} the set of integers; zahlen is the German word for numbers
 \oplus the earth
 \odot the sun
 c speed of light
 G Newton's gravitational constant
 C circumference
 S distance

Appendix B

Linear Algebra Review