

# Homework 2

## Question One

### Least mean square 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:

$\theta_j := \theta_j - \alpha \partial / \partial \theta_j J(\theta)$ . (This update is simultaneously performed for all values of  $j = 0, \dots, n$ .)

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:

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

For a single training example, this gives the update rule:

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

Update rules:

Consider the objective function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  that takes any multi-dimensional vector  $x = [x_1, x_2, \dots, x_d]^T$  as its input. The gradient of  $f(x)$  with respect to  $x$  is defined by the vector of partial derivatives:

$$\nabla_x f(x) = [\partial f(x) / \partial x_1, \partial f(x) / \partial x_2, \dots, \partial f(x) / \partial x_d]^T.$$

To keep our notation compact we may use the notation  $\nabla f(x)$  and  $\nabla_x f(x)$  interchangeably when there is no ambiguity about which parameters we are optimizing over. Each element  $\partial f(x)/\partial x_i$  of the gradient indicates the rate of change for  $f$  at the point  $x$  with respect to the input  $x_i$  only. To measure the rate of change of  $f$  in any direction that is represented by a unit vector  $u$ , in multivariate calculus, we define the directional derivative of  $f$  at  $x$  in the direction of  $u$  as:

$$D_u f(x) = \lim_{h \rightarrow 0} \frac{f(x+hu) - f(x)}{h}, D_u f(x) = \lim_{h \rightarrow 0} \frac{f(x+hu) - f(x)}{h},$$

which can be rewritten according to the chain rule as

$$D_u f(x) = \nabla f(x) \cdot u. D_u f(x) = \nabla f(x) \cdot u.$$

Since  $D_u f(x)$  gives the rates of change of  $f$  at the point  $x$  in all possible directions, to minimize  $f$ , we are interested in finding the direction where  $f$  can be reduced fastest. Thus, we can minimize the directional derivative  $D_u f(x)$  with respect to  $u$ . Since

$$D_u f(x) = \|\nabla f(x)\| \cdot \|u\| \cdot \cos(\theta) = \|\nabla f(x)\| \cdot \cos(\theta) D_u f(x) = \|\nabla f(x)\| \cdot \|u\| \cdot \cos(\theta) = \|\nabla f(x)\| \cdot \cos(\theta), \text{ where:}$$

$\theta$  is the angle between  $\nabla f(x)$  and  $u$ , the minimum value of  $\cos(\theta)$  is  $-1$  when  $\theta = \pi$ . Therefore,  $D_u f(x)$  is minimized when  $u$  is at the opposite direction of the gradient  $\nabla f(x)$ . Now we can iteratively reduce the value of  $f$  with the following gradient descent update:

$$x := x - \eta \nabla f(x), x := x - \eta \nabla f(x),$$

where the positive scalar  $\eta$  is called the learning rate or step size.

## Question Two

When faced with a regression problem, why might linear regression, and specifically why might the least-squares cost function  $J$ , be a reasonable choice?

Let us consider a set of probabilistic assumptions, under which least-squares regression is derived as a very natural algorithm. Again let us assume that the target variables and the inputs are related via the equation  $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$ , 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  $q(i)$  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 " $q(i) \sim N(0, \sigma^2)$ ." I.e., the density of  $q(i)$  is given by:

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

This implies that:

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

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 as:

$y(i) | x(i); \theta \sim N(\theta^T x^{(i)}, \sigma^2)$ . 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  $p(\sim y | X; \theta)$ . This quantity is typically viewed a function of  $\sim y$  (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:  $L(\theta) = L(\theta; X, \sim y) = p(\sim y | X; \theta)$ . Note that by the independence assumption on the  $q(i)$ 's (and hence also the  $y(i)$ 's given the  $x(i)$ 's), this can also be written:

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

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  $\ell(\theta)$ :

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

Hence, maximizing  $\ell(\theta)$  gives the same answer as minimizing  $\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$ , 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.

### Question Three

First, the derivative of  $\ln(x)$ :  $d/dx \ln(x) = 1/x$

Second, the (very elegant) derivative of the sigmoid:  $d\sigma(z)/dz = \sigma(z)(1-\sigma(z))$

Finally, the chain rule of derivatives. Suppose we are computing the derivative of a composite function  $f(x) = u(v(x))$ . The derivative of  $f(x)$  is the derivative of  $u(x)$  with respect to  $v(x)$  times the derivative of  $v(x)$  with respect to  $x$ :

$$d f / d x = d u / d v \cdot d v / d x$$

First, we want to know the derivative of the loss function with respect to a single weight  $w_j$  (we'll need to compute it for each weight, and for the bias):

$$\begin{aligned} \partial LL(w,b) / \partial w_j &= \partial \partial w_j [-y \log \sigma(w \cdot x + b) + (1-y) \log(1-\sigma(w \cdot x + b))] \\ &= - [ \partial / \partial w_j y \log \sigma(w \cdot x + b) + \partial / \partial w_j (1-y) \log[1-\sigma(w \cdot x + b)] \end{aligned}$$

Next, using the chain rule, and relying on the derivative of  $\log$ :

$$\partial LL(w,b) / \partial w_j = - y / \sigma(w \cdot x + b) \partial / \partial w_j \sigma(w \cdot x + b) - (1-y) / (1-\sigma(w \cdot x + b)) \partial / \partial w_j (1-\sigma(w \cdot x + b))$$

Rearranging terms:

$$\partial LL(w,b) / \partial w_j = - [ y / \sigma(w \cdot x + b) - (1-y) / (1-\sigma(w \cdot x + b)) ] \partial / \partial w_j \sigma(w \cdot x + b)$$

And now plugging in the derivative of the sigmoid, and using the chain rule one more time, we end up with Eq:

$$\begin{aligned} \partial LL(w,b) / \partial w_j &= - [ y - \sigma(w \cdot x + b) ] / \sigma(w \cdot x + b) [1 - \sigma(w \cdot x + b)] \sigma(w \cdot x + b) [1 - \sigma(w \cdot x + b)] \partial (w \cdot x + b) / \partial w_j \\ &= - [ y - \sigma(w \cdot x + b) ] \sigma(w \cdot x + b) [1 - \sigma(w \cdot x + b)] \sigma(w \cdot x + b) [1 - \sigma(w \cdot x + b)] x_j \\ &= - [ y - \sigma(w \cdot x + b) ] x_j \\ &= [ \sigma(w \cdot x + b) - y ] x_j \end{aligned}$$

#### Question Four

X <sub>1</sub>	X <sub>2</sub>	y
0	0	2
0	1	3
1	0	3
1	1	4

#### Gradient Descent Iteration #1

Let's start with the initial values of 0 for both coefficients.

$$\theta_0=0$$

$$\theta_1=0$$

$$\theta_2=0$$

$$y = \theta_0 + \theta_1 * x_1 + \theta_2 * x_2$$

$$y = 0 + 0 * x_1 + 0 * x_2$$

We can calculate the error for a prediction as follows:

$$\text{error} = p(i) - y(i)$$

Where  $p(i)$  is the prediction for the  $i$ 'th instance in our dataset and  $y(i)$  is the  $i$ 'th output variable for the instance in the dataset.

We can now calculate the predicted value for  $y$  using our starting point coefficients for the first training instance:

$$X_1=0, x_2=0, y=2$$

$$p(i) = 0 + 0 * 0 + 0 * 0$$

$$p(i) = 0$$

Using the predicted output, we can calculate our error:

$$\text{error} = 0 - 2$$

$$\text{error} = -2$$

We can now use this error in our equation for gradient descent to update the weights. We will start with updating the intercept first, because it is easier.

We can say that  $\theta_0$  is accountable for all of the error. This is to say that updating the weight will use just the error as the gradient. We can calculate the update for the  $\theta_0$  coefficient as follows:

$$\theta_0(t+1) = \theta_0(t) - \text{alpha} * \text{error}$$

Where  $\theta_0(t+1)$  is the updated version of the coefficient we will use on the next training instance,  $\theta_0(t)$  is the current value for  $\theta_0$  alpha is our learning rate and error is the error we calculate for the training instance. Let's use a small learning rate of 0.5 and plug the values into the equation to work out what the new and slightly optimized value of  $\theta_0$  will be:

$$\theta_0(t+1) = 0 - 0.5 * -2$$

$$\theta_0(t+1) = 1$$

Now, let's look at updating the value for  $\theta_1$ . We use the same equation with one small change. The error is filtered by the input that caused it. We can update  $\theta_1$  using the equation:

$$\theta_1(t+1) = \theta_1(t) - \text{alpha} * \text{error} * x_1$$

Where  $\theta_1(t+1)$  is the update coefficient,  $\theta_1(t)$  is the current version of the coefficient, alpha is the same learning rate described above, error is the same error calculated above and x is the input value.

We can plug in our numbers into the equation and calculate the updated value for  $\theta_1$ :

$$\theta_1(t+1) = 0 - 0.5 * -2 * 0$$

$$\theta_1(t+1) = 0$$

Now, let's look at updating the value for  $\theta_2$ . We use the same equation with one small change. The error is filtered by the input that caused it. We can update  $\theta_2$  using the equation:

$$\theta_2(t+1) = \theta_2(t) - \alpha * \text{error} * x_2$$

Where  $\theta_2(t+1)$  is the update coefficient,  $\theta_2(t)$  is the current version of the coefficient,  $\alpha$  is the same learning rate described above, error is the same error calculated above and  $x_2$  the input value.

We can plug in our numbers into the equation and calculate the updated value for  $\theta_2$ :

$$\theta_2(t+1) = 0 - 0.5 * -2 * 0$$

$$\theta_2(t+1) = 0$$

We have just finished the first iteration of gradient descent and we have updated our weights to be  $\theta_0=1$  and  $\theta_1=0$   $\theta_2=0$

## Gradient Descent Iteration #2

$X_1=0$ ,  $x_2=1$ ,  $y=3$  and  $\theta_0=1$ ,  $\theta_1=0$  and  $\theta_2=0$

$$y = \theta_0 + \theta_1 * x_1 + \theta_2 * x_2$$

$$p(i) = 1 + 0 * 0 + 0 * 0$$

$$p(i) = 1$$

$$\text{error} = 1 - 3$$

$$\text{error} = -2$$

$$\theta_0(t+2) = \theta_{t+1}(t) - \alpha * \text{error}$$

$$\theta_0(t+2) = 1 - 0.5 * -2$$

$$\theta_0(t+2) = 2$$

$$\theta_1(t+2) = \theta_{t+1}(t) - \alpha * \text{error} * x_1$$

$$\theta_1(t+2) = 0 - 0.5 * -2 * 0$$

$$\theta_1(t+1) = 0$$

$$\theta_2(t+2) = \theta_{t+1}(t) - \alpha * \text{error} * x_2$$

$$\theta_2(t+2) = 0 - 0.5 * -2 * 1$$

$$\theta_2(t+1) = 1$$

### **Gradient Descent Iteration #3**

$$X_1=1, x_2=0, y=3 \text{ and } \theta_0=2, \theta_1=0 \text{ and } \theta_2=1$$

$$y = \theta_0 + \theta_1 * x_1 + \theta_2 * x_2$$

$$p(i) = 2 + 0 * 1 + 1 * 0$$

$$p(i) = 2$$

$$\text{error} = 2 - 3$$

$$\text{error} = -1$$

$$\theta_0(t+3) = \theta_{t+2}(t) - \alpha * \text{error}$$

$$\theta_0(t+3) = 2 - 0.5 * -1$$

$$\theta_0(t+3) = 2 + 0.5$$



$$\theta_0(t+3) = 2.5$$

$$\theta_1(t+3) = \theta_{t+2}(t) - \alpha * \text{error} * x_1$$

$$\theta_1(t+3) = 0 - 0.5 * -1 * 1$$

$$\theta_1(t+3) = 0.5$$

$$\theta_2(t+3) = \theta_{t+2}(t) - \alpha * \text{error} * x_2$$

$$\theta_2(t+3) = 1 - 0.5 * -1 * 0$$

$$\theta_2(t+3) = 1$$

#### **Gradient Descent Iteration #4**

$$X_1=1 \ x_2=1, y=4 \text{ and } \theta_0=2.5, \theta_1=0.5 \text{ and } \theta_2=1$$

$$y = \theta_0 + \theta_1 * x_1 + \theta_2 * x_2$$

$$p(i) = 2.5 + 0.5 * 1 + 1 * 1$$

$$p(i) = 4$$

$$\text{error} = 4 - 4$$

$$\text{error} = 0$$

$$\theta_0(t+4) = \theta_{t+3}(t) - \alpha * \text{error}$$

$$\theta_0(t+4) = 2.5 - 0.5 * 0$$

$$\theta_0(t+4) = 2.5$$

$$\theta_1(t+4) = \theta_{t+3}(t) - \alpha * \text{error} * x_1$$

$$\theta_1(t+4) = 0.5 - 0.5 * 0 * 1$$

$$\theta_1(t+4) = 0.5$$

$$\theta_2(t+4) = \theta_{t+3}(t) - \text{alpha} * \text{error} * x_2$$

$$\theta_2(t+4) = 1 - 0.5 * 0 * 1$$

$$\theta_2(t+4) = 1$$