

Let $W^{(1)} = \begin{bmatrix} -2 & 4 & -1 \\ 6 & 0 & -3 \end{bmatrix}$ & $\vec{b}^{(1)} = \begin{bmatrix} 0.1 \\ -2.5 \end{bmatrix}$ (13)

in $\vec{a}^{(1)} = \sigma(W^{(1)} \cdot \vec{a}^{(0)} + \vec{b}^{(1)})$ with

& $\vec{a}^{(0)} = \begin{bmatrix} 0.3 \\ 0.4 \\ 0.1 \end{bmatrix}$

Find $\vec{a}^{(1)}$.

Soln:- $W^{(1)} \cdot \vec{a}^{(0)} = \begin{bmatrix} -2 & 4 & -1 \\ 6 & 0 & -3 \end{bmatrix} \begin{bmatrix} 0.3 \\ 0.4 \\ 0.1 \end{bmatrix} = \begin{bmatrix} -0.6 + 1.6 - 0.1 \\ 1.8 - 0.3 \end{bmatrix} = \begin{bmatrix} 0.9 \\ 1.5 \end{bmatrix}$

$\therefore W^{(1)} \cdot \vec{a}^{(0)} + \vec{b}^{(1)} = \begin{bmatrix} 1.0 \\ -1.0 \end{bmatrix}$

$\Rightarrow \vec{a}^{(1)} = \begin{bmatrix} \sigma(1.0) \\ \sigma(-1.0) \end{bmatrix} = \begin{bmatrix} 0.76 \\ -0.76 \end{bmatrix}$.

More simple neural networks

AIM : Iteratively update the values of all weights & biases so that the network learns to perform tasks (like classify data) based on a set of training examples.

Training a network means finding the right weights & biases using the training example. (which are pairs of matching inputs & outputs)

Backpropagation is a classical training method. Here, we look at the first at the output layer & then calculate gradients backward towards the input layer with chain rule. (of the cost func)

In a many-layered neural network, we have

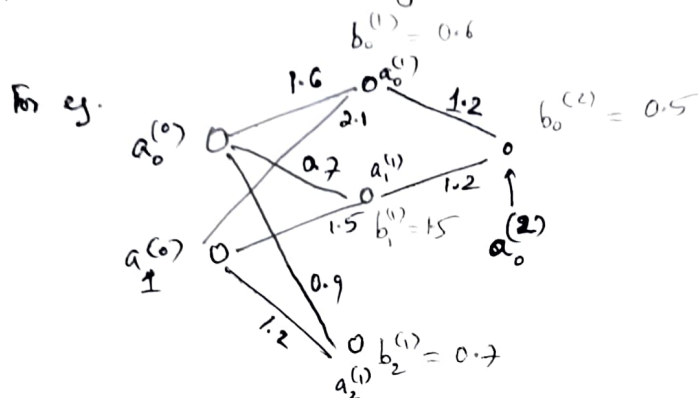
$$a^{(L)} = \sigma(w^{(L)} \cdot a^{(L-1)} + b^{(L)})$$

where L denotes the final layer.

assume the case where there is one

Let us ~~fix~~ a training example, say input is X & output is Y . X will be a vector of same size as $a^{(0)}$ & Y will be of same size as $a^{(L)}$.

~~a^(L) = predicted~~ For neural network to start prediction, we will randomly initialise ~~the~~ all the weights & biases.



Now, for these values of weights & biases, I can find ~~the~~ ^{that the error} cost incurred is

$$C = \frac{1}{2} (a_0^{(2)} - Y)^2$$

weights & biases

For the same ~~fixed~~ training example, if I think of ~~weights & biases~~ as variables, then cost will ^{be} a function which depends on the variable $a_0^{(2)}$ which now depends on ~~these~~ ^{on these} weights & biases.

$$C(a_0^{(2)}) = \frac{1}{2} (a_0^{(2)} - Y)^2$$

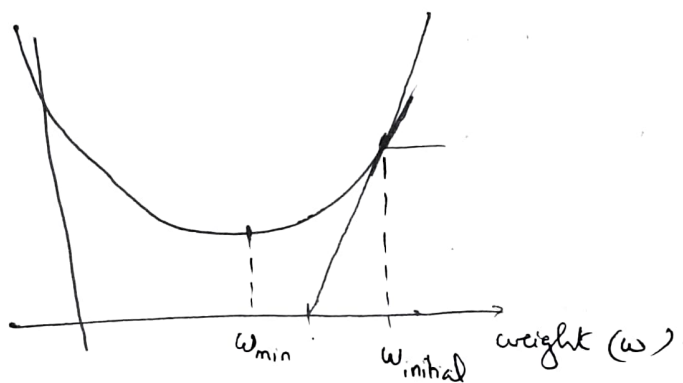
To MINIMIZE COST

~~To minimise cost~~ Suppose for the initial weights & biases, the cost value was too high, we update the weights & biases so as to minimise the cost. ~~in~~ ^(in an iterative manner)

Note that ~~this~~ ^{this} updation will change $a_0^{(2)}$ & in ~~turn~~ ^{turn} minimise C .

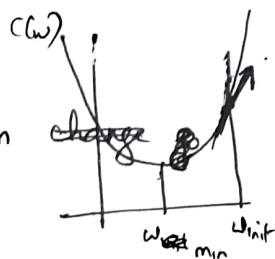
Consider a simple example $a^{(0)} \xrightarrow{w} a^{(1)}$ where $\sigma(x) = x$ & $b^{(1)} = 0$, ie, $a^{(1)} = w a^{(0)} \Rightarrow C = \frac{1}{2} (a^{(1)} - Y)^2 = \frac{1}{2} (w a^{(0)} - Y)^2$ (Parabola)

Get
(c)



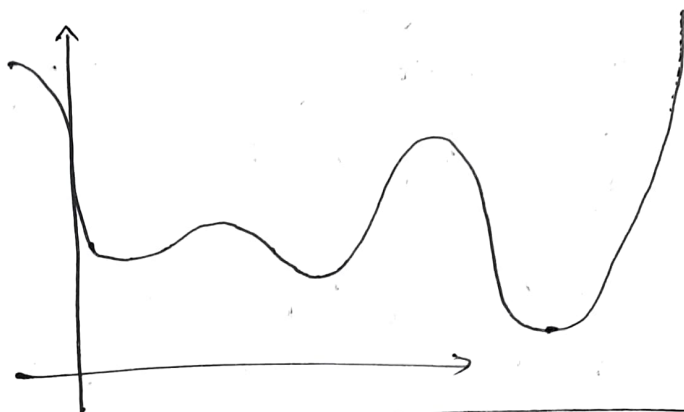
If weight is too large or too small, ^(cost) error is high & hence to 15
 at w_{min} , it has minimum cost.

If at your randomly chosen w_{init} , $\frac{\partial C}{\partial w}$ is $\neq 0$, then ~~change~~
 reduce C by going in opposite dir. of gradient.
 However, ~~the~~



A Brief Look at Challenges :-

→ Multiple local
minimas



→ Gradient exploding & vanishing problems
~~exploding~~

Chain Rule
Backpropagation in Case 1

$$a^{(0)} \rightarrow w \rightarrow a^{(1)}$$

In class, I had shown a challenge which appears in another iterative approach called Newton's Method. Here, we will use only Gradient Descent Method.

Recall $a^{(1)} = \sigma(w a^{(0)} + b^{(1)}) = \sigma(w a^{(0)} + b)$ ^{Use $w^{(1)} = w$ for simplicity.}

$$\text{Here, } C = (a^{(1)} - y)^2$$

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial w}$$

~~So~~ Implementation of gradient can be made less expensive by using intermediate variables, say

$$z^{(1)} = w^{(1)} a^{(0)} + b^{(1)} \quad \text{or} \quad z = w a^{(0)} + b$$

$$\text{i.e., } \frac{\partial C}{\partial w} = \frac{\partial C}{\partial a^{(1)}} \times \frac{\partial a^{(1)}}{\partial z} \times \frac{\partial z}{\partial w}$$

Then, $a^{(1)} = \sigma(z^{(1)})$ & hence,

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z} \frac{\partial z}{\partial w}$$

Eg 1) Consider the ~~neural network~~

$$a^{(0)} \xrightarrow{w^{(1)}} a^{(1)}$$

Let $\sigma(x) = x$ (linear activation ~~activation~~ func.) & $b^{(1)} = 0$

Forward propagation

The process of obtaining the ^{vector} ~~value~~ $a^{(1)}$ corresponding to a given set of ^(random/updated) biases & Weight matrices by considering a specific training example is called forward propagation. (FP) & thus, finding the error value

The process of propagating the error (obtained at final layer) ~~& passing~~ to previous layers so as to update the weights & bias is called backpropagation. (BP) (in a ^{desired} required manner)

The process of FP & BP are repeated until the ^{value of} error fn is 'satisfactory'.

Backpropagation basics

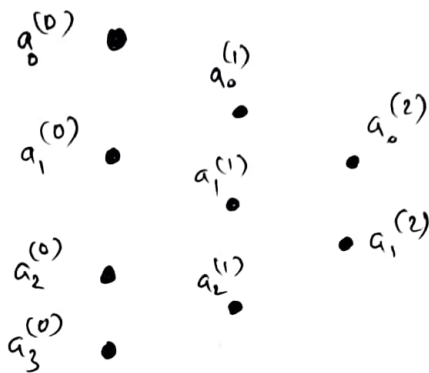
Recall that $a^{(n)} = \sigma(z^{(n)})$ is the feedforward equation at the n^{th} layer where $z^{(n)} = W^{(n)} a^{(n-1)} + b^{(n)}$

$$\text{Cost } C = \frac{1}{2} \|a^{(n)} - y\|^2$$

(For a single training example)

In order to update W & b , we need to determine

Consider an example



$0 \leq n \leq N$
 $\frac{\partial C}{\partial W^{(n)}}$ & $\frac{\partial C}{\partial b^{(n)}}$
 (other relevant gradients intermittent)

(I will also write $\frac{\partial C}{\partial W}$ & $\frac{\partial C}{\partial b}$ from now on)

Here, $a^{(2)} = W^{(2)} a^{(1)} + b^{(2)}$

$$C = \frac{1}{2} \|a^{(2)} - y\|^2 \text{ where}$$

$$a^{(2)} = \sigma(z^{(2)}) \quad \& \quad z^{(2)} = W^{(2)} a^{(1)} + b^{(2)}$$

$$a^{(1)} = \sigma(z^{(1)}) \quad \& \quad z^{(1)} = W^{(1)} a^{(0)} + b^{(1)}$$

$$\frac{\partial C}{\partial W^{(2)}} = \frac{\partial C}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial W^{(2)}} \quad \text{--- (1) \&}$$

$$\frac{\partial C}{\partial W^{(1)}} = \frac{\partial C}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial W^{(1)}} \quad \text{--- (2)}$$

Here, if $\frac{\partial C}{\partial W^{(2)}}$ has been computed, then ~~most of the~~ ^{few of those} computations in (1)

can be reused to compute $\frac{\partial C}{\partial W^{(1)}}$. Further,

$$(2) \Leftrightarrow \frac{\partial C}{\partial W^{(1)}} = \frac{\partial C}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial W^{(1)}}$$

(2) can be generalised as (for N layered network)

$$\frac{\partial C}{\partial w^{(i)}} = \frac{\partial C}{\partial a^{(N)}} \frac{\partial a^{(N)}}{\partial a^{(N-1)}} \frac{\partial a^{(N-1)}}{\partial a^{(N-2)}} \dots \frac{\partial a^{(i+1)}}{\partial a^{(i)}} \frac{\partial a^{(i)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial w^{(i)}}$$

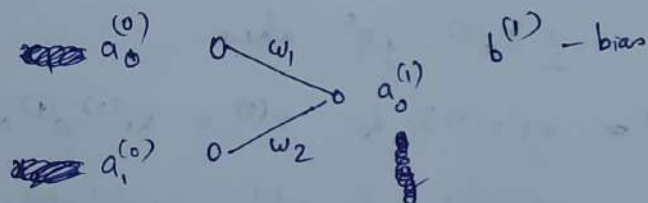
Nth layer to ith layer

$$\text{Similarly, } \frac{\partial C}{\partial w^{(i-1)}} = \frac{\partial C}{\partial a^{(N)}} \frac{\partial a^{(N)}}{\partial a^{(N-1)}} \frac{\partial a^{(N-1)}}{\partial a^{(N-2)}} \dots \frac{\partial a^{(i+1)}}{\partial a^{(i)}} \frac{\partial a^{(i)}}{\partial a^{(i-1)}} \frac{\partial a^{(i-1)}}{\partial z^{(i-1)}} \frac{\partial z^{(i-1)}}{\partial w^{(i-1)}}$$

Reusable computations from the ith layer

Now, let us do a problem involving neural networks with forward & back propagation, for single training example.

Pr 1) Consider the following neural network



Using the activation function, $\sigma(x) = \frac{1}{1+e^{-x}}$ (Sigmoid) & learning rate $\alpha = 0.1$, perform 1 step of FP & BP for the training example

$X = (x_1, x_2)$	Y
$(0.8, 0.6)$	1.0

with initial weights & biases as given here:-
 $w_1 = 0.3, w_2 = 0.5$
 $b = 0.1$

Soln:- $\sigma'(x) = -(1+e^{-x})^{-2} \times (-e^{-x}) = \frac{e^{-x}}{(1+e^{-x})^2}$

$$a^{(0)} = \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix}; W^{(1)} = \begin{bmatrix} 0.3 \\ 0.5 \end{bmatrix}; b^{(1)} = 0.1; a^{(1)} = \sigma(z^{(1)})$$

$$a^{(1)} = \sigma(z^{(1)}) \text{ where } z^{(1)} = W^{(1)} a^{(0)} + b^{(1)} = w_1 a_0^{(0)} + w_2 a_1^{(0)} + b^{(1)}$$

~~Let us use~~ forward propagation

$$z^{(1)} = \begin{bmatrix} 0.3 & 0.5 \end{bmatrix} \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix} + 0.1 = 0.64$$

$$a^{(1)} = \sigma(z^{(1)}) = 0.6547$$

$$C = \frac{1}{2} (a^{(1)} - y)^2 = \frac{1}{2} [0.6547 - 1]^2 = \frac{0.1192}{2} = 0.0596$$

Backpropagation

$$w_{i, \text{new}} = w_{i, \text{old}} - \alpha \frac{\partial C}{\partial w_i} (w_{i, \text{old}})$$

$$\begin{aligned} \frac{\partial C}{\partial w_1} &= \frac{\partial C}{\partial a^{(1)}} \times \frac{\partial a^{(1)}}{\partial z^{(1)}} \left(\frac{\partial z^{(1)}}{\partial w_1} \right) \\ &= \frac{1}{2} \times 2 \times (a^{(1)} - y) \times \sigma'(z^{(1)}) \times \frac{\partial z^{(1)}}{\partial w_1} \end{aligned}$$

$$\begin{aligned} z^{(1)} &= w_1 a_0^{(0)} + w_2 a_1^{(0)} + b^{(1)} \\ \Rightarrow \frac{\partial z^{(1)}}{\partial w_1} &= a_0^{(0)} \quad \& \quad \frac{\partial z^{(1)}}{\partial w_2} = a_1^{(0)} \\ &= 0.8 \qquad \qquad \qquad = 0.6 \end{aligned}$$

1114, $\frac{\partial C}{\partial i}$

$$W^{(1), \text{new}} = W^{(1), \text{old}} - 0.1 \times \begin{bmatrix} \frac{\partial C}{\partial w_1} (W^{(1), \text{old}}) \\ \frac{\partial C}{\partial w_2} (W^{(1), \text{old}}) \end{bmatrix}^T$$

Only diff between these two is

Let us simplify $W^{(1), \text{new}}$ & $W^{(1), \text{old}}$ notations to simply W^{new} & W^{old} . this term

$$\begin{aligned} \left. \frac{\partial C}{\partial w_1} \right|_{W^{\text{old}}} &= (0.6547 - 1) \times \frac{e^{-0.64}}{(1 + e^{-0.64})^2} \times 0.8 \\ &= -0.3453 \times 0.226 \times 0.8 = -0.0624 \end{aligned}$$

$$\frac{\partial C}{\partial w_2} = \frac{\partial C}{\partial a^{(1)}} \times \frac{\partial a^{(1)}}{\partial z^{(1)}} \times \left(\frac{\partial z^{(1)}}{\partial w_2} \right)$$

$$\Rightarrow \left. \frac{\partial C}{\partial w_2} \right|_{W^{\text{old}}} = -0.3453 \times 0.226 \times 0.6 = -0.047$$

$$\begin{aligned} \therefore W^{\text{new}} &= \cancel{W^{\text{old}}} \begin{bmatrix} 0.3 \\ 0.5 \end{bmatrix}^T - 0.1 \begin{bmatrix} -0.0624 \\ -0.047 \end{bmatrix}^T \\ &= \begin{bmatrix} 0.30624 \\ 0.5047 \end{bmatrix}^T \end{aligned}$$

We can form the foll. table

$W^{(1)}$	$a^{(0)}$	$\frac{\partial C}{\partial W^{(1)}}$	$\frac{\partial C}{\partial z^{(1)}}$	$W^{(1), new}$	$\frac{\partial C}{\partial b^{(1)}}$	$b^{(1), new}$
$\begin{bmatrix} 0.3 \\ 0.5 \end{bmatrix}^T$	$\begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix}$	0.0596	$\begin{bmatrix} -0.0624 \\ -0.047 \end{bmatrix}^T$	$\begin{bmatrix} 0.30624 \\ 0.5047 \end{bmatrix}^T$	-0.078	0.1078

$$\text{Here, } \frac{\partial C}{\partial b^{(1)}} = \frac{\partial C}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial b^{(1)}}$$

$$\Rightarrow \frac{\partial C}{\partial b^{(1)}} \bigg|_{W^{old}} = -0.3453 \times 0.226 \times 1 = -0.078$$

$$\text{hence, } b^{new} = 0.1 - 0.1 \times (-0.078) = 0.1078$$

Thus, we form the foll table (2nd step exercise)

$b^{(1)}$	$W^{(1)}$	$a^{(1)}$	C	$\nabla_{W^{(1)}} C$	$\frac{\partial C}{\partial b^{(1)}}$	$W^{(1), new}$	$b^{(1), new}$
0.1	$\begin{bmatrix} 0.3 \\ 0.5 \end{bmatrix}^T$	0.655	0.0596	$\begin{bmatrix} -0.0624 \\ -0.047 \end{bmatrix}^T$	-0.078	$\begin{bmatrix} 0.30624 \\ 0.5047 \end{bmatrix}^T$	0.1078
0.1078	$\begin{bmatrix} 0.30624 \\ 0.5047 \end{bmatrix}^T$	0.6583	0.0584	$\begin{bmatrix} -0.0615 \\ -0.0461 \end{bmatrix}^T$	-0.0768	$\begin{bmatrix} 0.3124 \\ 0.5093 \end{bmatrix}^T$	0.1155

Linear Regression problem: Background & brief overview

Let input $X = (x_1, x_2, \dots, x_n)$ & $y \in \mathbb{R}$. Let there be m training examples

Training example	X					Y
	x_1	x_2	x_3	...	x_n	
1	x_1^1	x_2^1	x_3^1	...	x_n^1	y^1
2	x_1^2	x_2^2	x_3^2	...	x_n^2	
...						
m	x_1^m	x_2^m	x_3^m	...	x_n^m	y^m

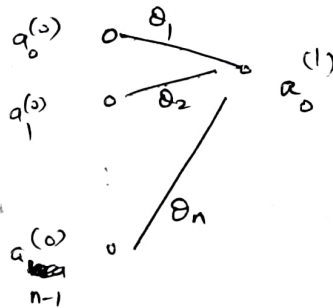
~~Let us say~~ Let us say $h_0(X) = \sum_{i=0}^n \theta_i x_i$

$\theta_0 \rightarrow$ bias

$\theta_1, \dots, \theta_n$ } weights

Let $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$

$\theta_1, \dots, \theta_n$ will be the weights of a neural network with n inputs & 1 output & θ_0 will be its bias.



Cost fn $J(\theta) = \frac{1}{2} \sum_{i=0}^m [h_0(x^{(i)}) - y^{(i)}]^2$

To minimise $J(\theta)$, we ~~choose~~ start with an initial guess of θ & use the search algorithm 'Gradient Descent' that repeatedly changes θ to make $J(\theta)$ smaller & smaller according to

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

This update is performed on all $0 \leq j \leq n$ simultaneously.

α - learning rate

Since θ is added with a vector along the -ve of the gradient of J , the new θ obtained will be reducing $J(\theta)$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}^{\text{new}} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}^{\text{old}} - \alpha \begin{bmatrix} \frac{\partial J}{\partial \theta_0}(\theta^{\text{old}}) \\ \frac{\partial J}{\partial \theta_1}(\theta^{\text{old}}) \\ \vdots \\ \frac{\partial J}{\partial \theta_n}(\theta^{\text{old}}) \end{bmatrix}$$

Computing $\frac{\partial J}{\partial \theta}$

$$J = \frac{1}{2} \sum_{i=1}^m (\theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)} - y)^2$$

$$= \frac{1}{2} \sum_{i=1}^m \left(\theta_j x_j^{(i)} + \left(\sum_{\substack{k=0 \\ k \neq j}}^n \theta_k x_k^{(i)} \right) - y \right)^2$$

$$\therefore \frac{\partial J}{\partial \theta_j}(\theta) = \frac{1}{2} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} \left[\theta_j x_j^{(i)} + \left(\sum_{\substack{k=0 \\ k \neq j}}^n \theta_k x_k^{(i)} \right) - y \right]^2$$

$$= \frac{1}{2} \sum_{i=1}^m 2 \left(\theta_j x_j^{(i)} + \left(\sum_{\substack{k=0 \\ k \neq j}}^n \theta_k x_k^{(i)} \right) - y \right) x_j^{(i)}$$

$$= \sum_{i=1}^m x_j^{(i)} (h_{\theta}(x^{(i)}) - y)$$

For ~~a single~~ training (eg. ie, $m=1$), say $\boxed{(x_1^{(1)}, \dots, x_n^{(1)}), y^{(1)}}$
a single

$$\frac{\partial J}{\partial \theta_j}(\theta) = x_j^{(1)} (h_{\theta}(x^{(1)}) - y^{(1)})$$

Algorithm in batch gradient descent

Repeat until convergence,

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} - \alpha \begin{bmatrix} \sum_{i=1}^m x_{\theta_0}^{(i)} (h_{\theta}(x^{(i)}) - y^{(i)}) \\ \vdots \\ \sum_{i=1}^m x_n^{(i)} (h_{\theta}(x^{(i)}) - y^{(i)}) \end{bmatrix}$$

Stochastic gradient descent

Here, we don't subtract the whole $\frac{\partial J}{\partial \theta_j}$ from θ_j , but a part, namely that corr. to each training example one at a time

Algorithm

for $i=1$ to m

$$\left\{ \begin{array}{l} \theta_j = \theta_j - \alpha x_j^{(i)} (h_{\theta}(x^{(i)}) - y^{(i)}) \end{array} \right.$$

ie,

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \begin{bmatrix} \theta_0 - \alpha x_0^{(1)} (h_{\theta}(x^{(1)}) - y^{(1)}) \\ \theta_1 - \alpha x_1^{(1)} (h_{\theta}(x^{(1)}) - y^{(1)}) \\ \vdots \\ \theta_n - \alpha x_n^{(1)} (h_{\theta}(x^{(1)}) - y^{(1)}) \end{bmatrix}$$

First training example's component from $\frac{\partial J}{\partial \theta_0}$

At $i=1$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \begin{bmatrix} \theta_0 - \alpha x_0^{(1)} (h_{\theta}(x^{(1)}) - y^{(1)}) \\ \theta_1 - \alpha x_1^{(1)} (h_{\theta}(x^{(1)}) - y^{(1)}) \\ \vdots \\ \theta_n - \alpha x_n^{(1)} (h_{\theta}(x^{(1)}) - y^{(1)}) \end{bmatrix} = \theta^{\text{new}}$$

↓

At $i=2$

$$\begin{bmatrix} \theta_0^{\text{new}} - \alpha x_0^{(2)} (h_{\theta^{\text{new}}}(x^{(2)}) - y^{(2)}) \\ \theta_1^{\text{new}} - \alpha x_1^{(2)} (h_{\theta^{\text{new}}}(x^{(2)}) - y^{(2)}) \\ \vdots \\ \theta_n^{\text{new}} - \alpha x_n^{(2)} (h_{\theta^{\text{new}}}(x^{(2)}) - y^{(2)}) \end{bmatrix}$$

↓

and so on

When $m=1$, stochastic G.D = batch G.D

Eg 1) Consider the neural network

$$a^{(0)} \xrightarrow{\omega} a^{(1)}$$

Let the activation function be linear, i.e., $\sigma(x) = x$ & $b^{(1)} = 0$.

Initialise $\omega = 0.8$ & perform 3 steps of FP & BP for the training example $x = 1.5$ | $y = 0.5$

Soln :- $a^{(1)} = \sigma(a^{(0)}\omega^{(1)} + b^{(1)}) = 0.8 a^{(0)}$

1st step
of FP

$$C = (a^{(1)} - y)^2$$

Putting $a^{(0)} = 1.5$, we get $a^{(1)} = 0.8 \times 1.5 = 1.2$
Then, $C = (1.2 - 0.5)^2 = 0.49$.

1st step of
BP

~~$$\frac{\partial C}{\partial \omega} = \frac{\partial C}{\partial a^{(1)}} \times \frac{\partial a^{(1)}}{\partial \omega} = 2(a^{(1)} - y) \times \frac{\partial a^{(1)}}{\partial \omega}$$~~

~~Since we want to update ω , write $a^{(1)} = a^{(0)}\omega = 1.5\omega$~~



MANIPAL INSTITUTE OF TECHNOLOGY

MANIPAL

(A Constituent unit of MAHE, Manipal)

~~1st step of BP~~

We now have the entries of the 1st row of the following table

$a^{(0)}$	w_{old}	b	$a^{(1)}$	Cost	$\left. \frac{\partial C}{\partial w} \right _{w=w_{old}}$	$w_{new} = w_{old} - \eta \frac{\partial C}{\partial w}(w_{old})$
1.5	0.8	0	1.2	0.49	$2(1.5 - 0.5)$	

To determine $\frac{\partial C}{\partial w}(w_{old})$, note that $\frac{\partial C}{\partial a^{(1)}} = 2(a^{(1)} - y)$ & $\frac{\partial a^{(1)}}{\partial w} = a^{(0)}$

$$\therefore \frac{\partial C}{\partial w}(0.8) = 2(1.2 - 0.5) \times 1.5 = 2.1$$

$$\eta = 0.1 \Rightarrow w_{new} = w_{old} - \eta \frac{\partial C}{\partial w}(0.8)$$

$$= 0.8 - 0.1 \times 2.1 = 0.59$$

This completes the entries of 1st row. Thus, we have performed 1 round of FP & BP, & adjusted weights. Now, let us see the change in the cost due to change in w . This leads to 2nd step of FP

2nd step of FP

$$a^{(1)} = w a^{(0)} \quad \& \quad C = (a^{(1)} - y)^2 \quad \text{where } w = 0.59 \quad (\text{instead of } 0.8)$$

$$\Rightarrow a^{(1)} = 0.59 \times 1.5 = 0.885 \quad \& \quad C = 0.148225$$

$$\text{Now, } \frac{\partial C}{\partial w}(w_{old}) = \frac{\partial C}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial w} = 2(a^{(1)} - 0.5) \times 1.5 = 1.155$$

$$\Rightarrow w_{new} = 0.59 - 0.1 \times (1.155) = 0.4745$$

$$\& \quad C = (0.885 - 0.5)^2 = 0.148225$$

This is the 2nd row of above table (shown in next page)

Also, $w_{\text{new}} = w_{\text{old}} - 0.1 \times \frac{\partial C}{\partial w}(w_{\text{old}})$.

(18)

$$= 0.59 - 0.1 \times 1.155$$

$$= 0.4745$$

The following table can be constructed.

	$a^{(0)}$	w_{old}	b	$a^{(1)}$	C	$\frac{\partial C}{\partial w}(w_{\text{old}})$ $= 3(1.5w_{\text{old}} - 0.5)$	$w_{\text{new}} = w_{\text{old}} - 0.1 \frac{\partial C}{\partial w}(w_{\text{old}})$
Step 1	1.5	0.8	0	1.2	0.49	2.1	0.59
Step 2	1.5	0.59	0	0.885	0.148	1.155	0.4745
Step 3	1.5	0.4745	0	0.71175	0.045	0.63525	0.410975
							\vdots
							0.333

Repeated iteration eventually converges to 0.333 & this will be the ideal weight for which the cost will be minimum, ~~Note that in this simple case, $C = (w_{\text{old}} - y)$~~ For this training example (1.5, 0.5). Now, we note that

$$C = (1.5w - 0.5)^2 \text{ \& }$$

$$\frac{\partial C}{\partial w} = 2(1.5w - 0.5) = 0 \text{ \& } \frac{d^2 C}{dw^2} = 3 > 0$$

already tells you to choose $w = \frac{1}{3}$, as C is minimum (global) in this case. However, in general, for vector valued inputs & outputs, the computations involved with equating partial derivatives to zero & evaluation Hessian are too expensive when compared to (stochastic) gradient descent algorithm.