

1.

For counter = 1 upto Niter		
	Choose an integer $k$ uniformly at random from $\{1, 2, 3, \dots, N\}$	
	$\mathbf{x}^{\{k\}}$ is current training data point	
	$\mathbf{a}^{[1]} = \mathbf{x}^{\{k\}}$	
	For $l = 2$ upto $L$	
		$\mathbf{z}^{[l]} = \mathbf{W}^{[l]}\mathbf{a}^{[l-1]} + \mathbf{b}^{[l]}$
		$\mathbf{a}^{[l]} = \sigma\left(\mathbf{z}^{[l]}\right)$
		$\mathbf{D}^{[l]} = \text{diag}\left(\sigma'\left(\mathbf{z}^{[l]}\right)\right)$
	end	
	$\delta^{[L]} = \mathbf{D}^{[L]}\left(\mathbf{a}^{[L]} - \mathbf{y}\left(\mathbf{x}^{\{k\}}\right)\right)$	
	For $l = L - 1$ downto 2	
		$\delta^{[l]} = \mathbf{D}^{[l]}\left(\mathbf{W}^{[l+1]}\right)^T \delta^{[l+1]}$
	end	
	For $l = L$ downto 2	
		$\mathbf{W}^{[l]} \rightarrow \mathbf{W}^{[l]} - \eta \delta^{[l]} \mathbf{a}^{[l-1]T}$
		$\mathbf{b}^{[l]} \rightarrow \mathbf{b}^{[l]} - \eta \delta^{[l]}$
	end	
	end	

The picture<sup>1</sup> above is the pseudo code of SGD.

Now, we just follow the steps to calculate  $\theta^1$ .

Take  $a^{[1]} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ .

For  $l = 2$  to 2, so we only need to deal with the case  $l = 2$ .

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} = \begin{pmatrix} 5 & 6 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} + 4 = 21.$$

$$a^{[2]} = \sigma(z^{[2]}) = \sigma(21), \text{ where } \sigma(x) = \frac{1}{1 + e^{-x}}, \text{ the sigmoid function.}$$

$$D = \text{diag}(\sigma'(z^{[2]})) = \sigma(21)(1 - \sigma(21))$$

$$\delta^{[2]} = \sigma(21)(1 - \sigma(21))(\sigma(21) - 3).$$

The next for loop start from  $l = 1 < 2$ , so we can ignore it.

For  $l = 2$  to 2, again, we only need to deal with the case  $l = 2$ .

Substitute  $W^{[2]} - \eta \delta^{[2]}(a^{[1]})^T = \begin{pmatrix} 5 & 6 \end{pmatrix} - \eta \cdot \sigma(21)(1 - \sigma(21))(\sigma(21) - 3) \begin{pmatrix} 1 & 2 \end{pmatrix}$  for  $W^{[2]}$ .

On the other hand, substitute  $b^{[2]} - \eta \delta^{[2]} = 4 - \eta \cdot \sigma(21)(1 - \sigma(21))(\sigma(21) - 3)$  for  $b^{[2]}$ .

Hence, we have  $\theta^1 = \begin{pmatrix} 4 - \eta \cdot k & 5 - \eta \cdot k & 6 - 2\eta \cdot k \end{pmatrix}$ , where  $k = \sigma(21)(1 - \sigma(21))(\sigma(21) - 3)$ .

<sup>1</sup>The picture is captured from

2. (a) Recall that  $\sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$ .

$$\begin{aligned}
 \sigma'(x) &= -\frac{1}{(1+e^{-x})^2} \cdot (-e^{-x}) \\
 &= \frac{e^{-x}}{1+2e^{-x}+e^{-2x}} \\
 &= \frac{e^x}{e^{2x}+2e^x+1} \\
 &= \frac{1}{e^x+1} \cdot \frac{e^x}{e^x+1} \\
 &= \sigma(x)(1-\sigma(x)) \\
 &= -(\sigma(x))^2 + \sigma(x) \\
 &= -\sigma^2 + \sigma
 \end{aligned}$$

$$\begin{aligned}
 \sigma''(x) &= -2\sigma\sigma' + \sigma' \\
 &= (-2\sigma+1)(-\sigma^2+\sigma) \\
 &= \sigma(1-\sigma)(1-2\sigma)
 \end{aligned}$$

$$\begin{aligned}
 \sigma'''(x) &= (6\sigma^2-6\sigma+1) \cdot \sigma' \\
 &= (6\sigma^2-6\sigma+1)\sigma(1-\sigma)
 \end{aligned}$$

(b) Recall that  $\tanh x = \frac{e^{2x}-1}{e^{2x}+1}$ .

Observe that the denominator is similar.

$$\left\{ \begin{array}{lcl} \sigma(2x) & = & \frac{e^{2x}}{e^{2x}+1} \\ 1-\sigma(2x) & = & \frac{1}{e^{2x}+1} \end{array} \right.$$

Hence, we have

$$\tanh x = 2\sigma(2x) - 1$$

3. I guess that the next questions may be

- (a) Is the SGD convergent?
- (b) If so, the point it converges to must be the global minima?
- (c) If it is not the global minima, how do we fix the "stop"?

Actually, I have already heard something about these questions and I just leave a brief reference answer here.

For (a), yes, but I have never proved it.

For (b), no, it might be saddle point, or in general case, most of them are just saddle points <sup>2</sup>.

However, the graph of the Coss function may be rugged, so the problem we met more frequently is we can't even touch the saddle point, unless we do the for loop with plenty of times.

Thus, also answer (c) simultaneously, we can use the concept of momentum <sup>3</sup> in Physics, like substitute the gradient we minus now with the "general" one which considering the former or the sum of all gradients we have now or change  $\eta$  in each  $l$ , for example, divide it by the norm of the gradient <sup>4</sup>.

---

<sup>2</sup>See [What Happened When Gradient is almost zero?](#)

<sup>3</sup>See [【機器學習2021】類神經網路訓練不起來怎麼辦\(二\)：批次\(batch\) 與動量\(momentum\) 22:37](#)

<sup>4</sup>See [【機器學習2021】類神經網路訓練不起來怎麼辦\(三\)：自動調整學習速率\(Learning Rate\) 00:00](#)