

⇒ Week-5

⇒ Cost function for neural nets

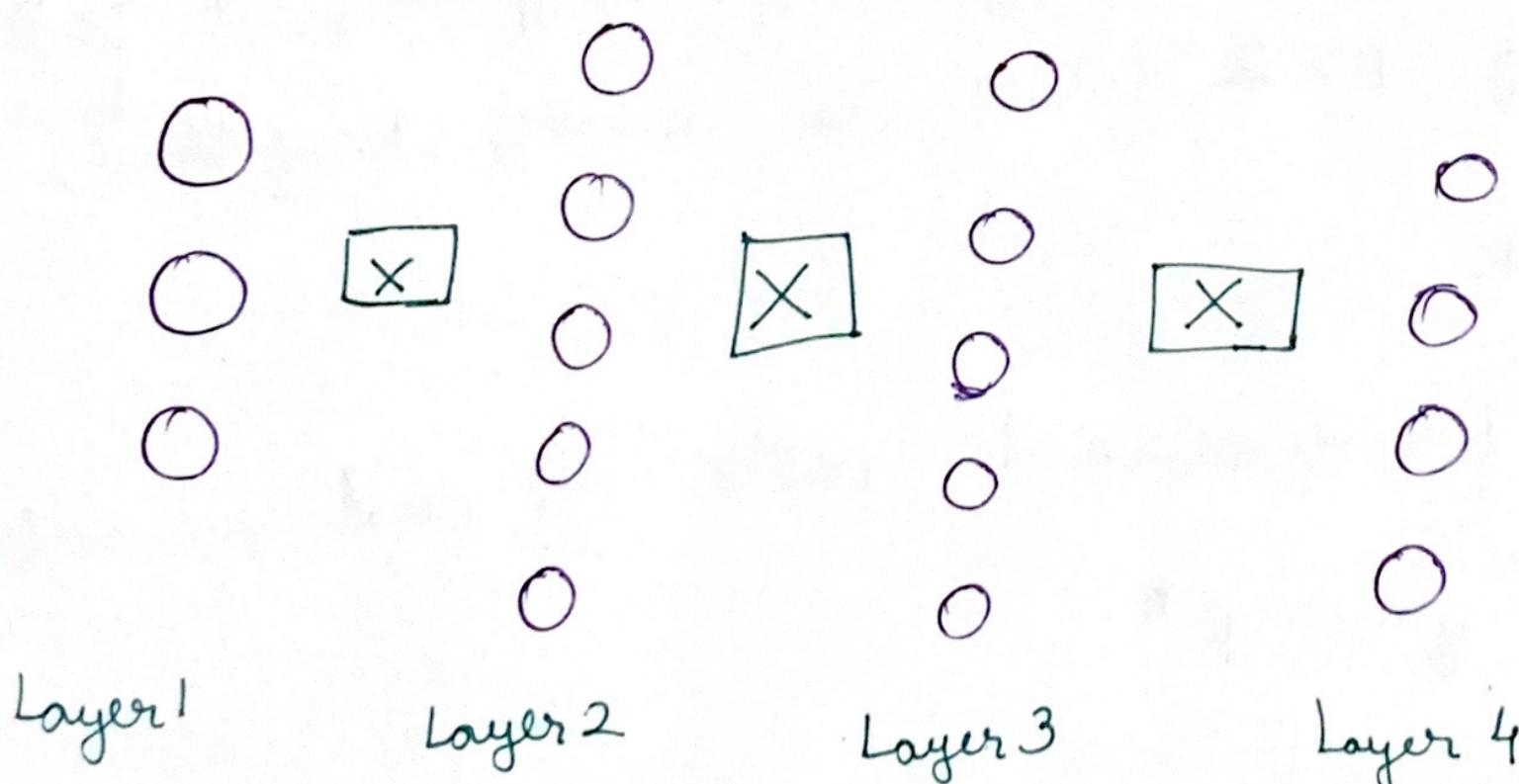
* Lets define a few variables that we can use:

→ L = total number of layers in the network

→ s_l = number of units (not counting bias unit) in layer \cancel{L}

→ K = no. of output units/classes

eg?



→ \boxed{X} : denotes that every ^{neuron} ~~neuron~~ before of ^{the} layer which is before the 'mark' is connected to each neuron of the layer present after the 'mark'

→ Here, $L = 4$, $s_1 = 3$, $s_2 = 5$, $s_3 = 5$, $s_4 = 4$
 $K = 4$

⇒ Types of classification problems with neural nets

→ 1st type : Binary classification

* 1 output (0 or 1)

* $R = 1$

* $D_L = 1$

→ 2nd type : Multiclass classification

* K distinct classifications . (Typically $K \geq 3$)

* → If $K=2$ (Just use binary classification)

* $D_L = K$

* $y \rightarrow k$ dimensional vector of real numbers

or $y \in R^K$

eg: $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$

Pedestrian

Car

Bike

Truck

\Rightarrow Cost function

- * Recall that the cost funcⁿ for regularized logistic regression was :

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(h_\theta(x^i)) + (1-y^i) \log(1-h_\theta(x^i))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

- * For neural nets, our cost function is a generalization of the eqⁿ above. Instead of one output we generate k outputs.

$h_\theta(x) \in R^{*K}$; ~~(not R^K)~~
 $(h_\theta(x))_i = i^{\text{th}} \text{ output}$

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^i \log(h_\theta(x^i))_k + (1-y_k^i) \log(1-(h_\theta(x^i))_k) \right] + \frac{\lambda}{2m} \sum_{L=1}^{L-1} \sum_{i=1}^{S_L} \sum_{j=1}^{S_{L+1}} (\theta_{j,i}^L)^2$$



... Dayaaaaaaa..

→ There are 2 halves to the neural net logistic regression

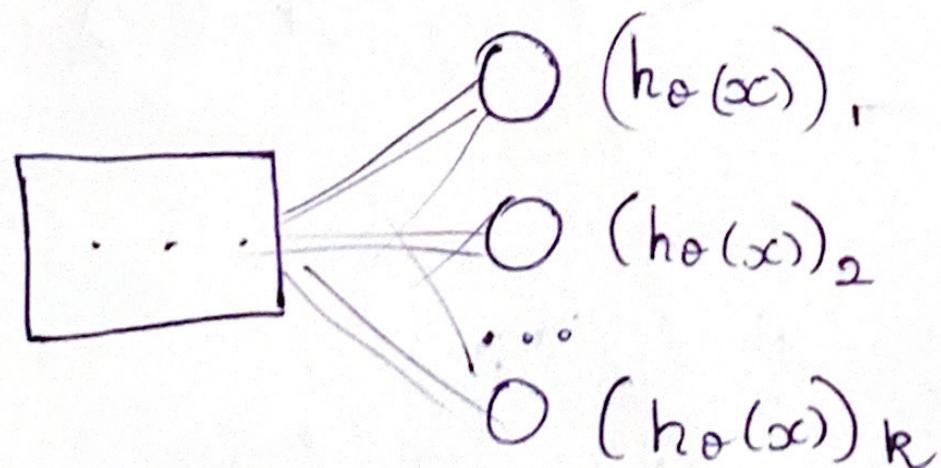
cost function



1st

$$\rightarrow \text{1st half: } -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^R y_k^i \log(h_0(x^i))_k + (1-y_k^i) \log(1-h_0(x^i))_k \right]$$

* This is just saying for each training example (i.e. 1 to m - the first summation), calculate the logistic regression for every output neuron in the output layer (i.e. 1 to R - the second summation). It's averaging the summation over $m \rightarrow$ total training examples



TL;DR

TL;DR: The triple double sum simply adds up the logistic regression costs calculated for each cell in the output layer.

$$\rightarrow \text{2nd half: } \frac{\lambda}{2m} \sum_{L=1}^{L-1} \sum_{i=1}^{s_L} \sum_{j=1}^{s_{L+1}} (\theta_{ji})^2$$

The triple sum simply adds up the squares of all
 the individual θ s in the entire network.

* Note the i in the triple sum does NOT refer
 to training sample i .

* As similar to logistic regression, the bias terms
 are not included in the triple sum.

⇒ Backpropagation algorithm : Algorithm used to minimize the cost function of our neural network.

(just like what we were doing with gradient descent in logistic & linear regression)

* Our goal is to compute :

$$\min_{\theta} J(\theta)$$

* To compute this [↑] we can use algorithms such as gradient descent & advance ~~of descent~~ optimization algorithms, which require 2 things of the foll. things :

1) $J(\theta)$ → i.e. the cost function itself! Just calculate it manually & code it!

2) $\frac{\partial J(\theta)}{\partial \theta_i}$

$\delta \theta_i^L$ → This is NOT AT ALL easy!

- * θ is indeed in 3 dimensions because we have separate parameter values for each mode in each layer going to each mode in the following layer.
- * Remember $\frac{\partial}{\partial \theta_{ij}^L} J(\theta)$ is a REAL number (not a vector or a matrix) ~~across all~~.
- \Rightarrow Mechanics behind back propagation
- * For each node, we can calculate (δ_j^l) - this is the error of node j in layer L .
- * $a_j^L \rightarrow$ activation of node j in layer L .
 - Totally theoretical value, so we can expect some error compared to the "real" value
 - The delta term captures this error!

* But the problem is, "What is this "real" or ideal value that'd give us optimum results) & how do we calculate it?"

The NN is a totally artificial construct. So, the only "real" value we have is our actual classification (our y value) — so that's where we start!

* Looking at the sample on pg-1 (Week 5)
is looking at the fourth (or output) layer,

$$\boxed{\delta_j^4 = a_j^4 - y_j}$$

= [Activation of]
[the unit] - [Actual value (as observed)
in training sample]

$$= h_0(x)_j - y_j$$

* $\delta^4 = a^4 - y$ → Vectorized implementation

Vector of errors in 4th layer Vector of activation values for the 4th layer

* The idea behind backpropagation is that once you know what the error in your ^{output} nodes is, you backpropagate this error to further calculate error in nodes present in the previous layer(s), by using the below formula (derived later):

* $\delta^l = ((o^l)^T \delta^{l+1}) * g'(z^l)$

Element wise multiply

* → If $g(z) \equiv \text{Sigmoid func}^n = 1/(1+e^{-z})$

then $g'(z) = g(z) * (1-g(z))$ [Derived later]

∴ $\delta^l = ((o^l)^T \delta^{l+1}) * a^l * (1-a^l)$

* Using this we can calculate $\frac{\partial J(\theta)}{\partial \theta_{i,j}^l}$ in the foll. way: (ignoring the regularization term)
 will derive it later

$$\boxed{\frac{\partial J(\theta)}{\partial \theta_{i,j}^l} = a_j^l s_i^{l+1}}$$

→ For m training samples we can just calculate average (mean) of all m samples:

$$\therefore \boxed{\frac{\partial J(\theta)}{\partial \theta_{i,j}^l} = \frac{1}{m} \sum_{t=1}^m a_j^{t+l} s_i^{t+l+1}}$$

for each weight and
 This'll be calculated & will be fed into advanced optimization algs. to perform gradient descent. e.g.: In Matlab, we use fminunc

* Intuitively the gradient $\frac{\partial J(\theta)}{\partial \theta_{i,j}^l}$

how much the weight $\theta_{i,j}^l$ needs to change (increase or decrease) in order to minimize $J(\theta)$.

* Now, once we know (upon upon see :-), we can finally put together these mechanics of backprop, we can finally put together into a formal backprop algorithm:

→

Given training set $\{(x^1, y^1), \dots, (x^m, y^m)\}$

Set $\Delta_{ij}^L = 0 \quad \forall i, j$

For training sample $t = 1 \text{ to } m$:

Set $a^t = x^t$

Perform forward propagation to compute a^t for $l = 2, 3, \dots, L$

Using y^t , compute $\delta^L = a^L - y^t$

Compute $\delta^{L-1}, \delta^{L-2}, \dots, \delta^2$

$\Delta_{ij}^L := \Delta_{ij}^L + a_j^{t+1} \delta_i^{L+1} \quad \forall i, j, l$

$\Delta_{ij}^L := \frac{1}{m} \Delta_{ij}^L + \lambda \alpha_{ij}^L \quad \text{if } j \neq 0$

$\Delta_{ii}^L := \frac{1}{m} \Delta_{ii}^L \quad \text{if } j = 0$

Now, let's discuss what each term does.

* $\Delta_{ij}^L = 0 \quad \forall i, j, l$: Used to compute $\frac{\partial J(\theta)}{\partial \theta_{ij}}$
Represents θ_{ij}

Just initialize a matrix full of 0s.

* For each training sample $t = 1 \dots m$:

1) Set $a^1 := x^t$

$$\begin{matrix} a^1 & a^2 & & a^{L-1} & \\ 0 & 0 & & 0 & a^L \\ 0 & \text{X} & 0 & \dots & 0 & \text{X} \\ 0 & 0 & & & 0 & 0 \end{matrix}$$

2) Perform forward propagation to compute

a^l for $l = 2, 3, \dots, L$

→ For one training sample (x, y) :

Forward propagation:

$$\begin{aligned} a^1 &= x \\ z^2 &= \theta^1 a^1 \\ a^2 &= g(z^2) \\ &\vdots \\ z^L &= \theta^{L-1} a^{L-1} \end{aligned}$$

3) Using y^t , compute $\delta^L = a^L - y^t$



→ $L = \text{Total number of layers}$

→ $a^L = \text{Vector of outputs of the activation units for the last layer}$

→ Error values for last layer are simply b/w our theoretical value & actual one.

4) Compute $\delta^{L-1}, \delta^{L-2}, \dots, \delta^2$



→ Use
$$\delta^L = ((\theta^L)^T \delta^{L+1}) * a^L * (1-a^L)$$

5) $\Delta_{ij}^L := \Delta_{ij}^L + a_{ij}^L \delta_i^{L+1}$



→
$$\Delta^L := \Delta^L + \delta^{L+1} (a^L)^T$$
 → Vectorized implementation

* Up

→ The 'D' matrix



* The 'D' : Calculates mean of all the values of $\delta \Delta_{iz}^L$ for all training samples m.

* $D_{iz}^L := \frac{1}{m} (\Delta_{iz}^L + \lambda \Delta_{iz}^L)$ if $z \neq 0$

$\Delta_{iz}^L := \frac{1}{m} \Delta_{iz}^L$ if $z = 0$ } No regularization for bias-terms!



⇒ Exploration of derivatives used in backpropagation

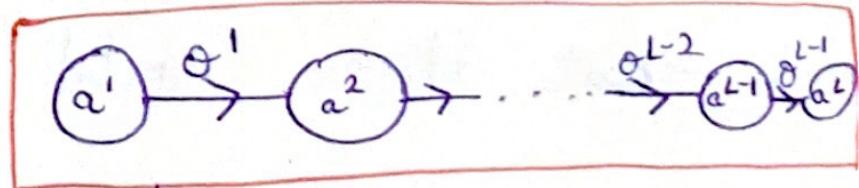
* We know that for a logistic regression classifier (which is what all of the output neurons in a neural network are), we use the cost funcⁿ:

$$J(\theta) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

∴ apply this over the 'K' output neurons,

and for all 'm' samples. In the foll. sample

say our NN looks like this



* The eqⁿ to compute the partial derivatives of the theta terms in the output ~~neuron~~ :

$$\frac{\partial J(\theta)}{\partial \theta^{(L-1)}} = \frac{\partial J(\theta)}{\partial a^L} \frac{\partial a^L}{\partial z^L} \frac{\partial z^L}{\partial \theta^{L-1}}$$

[Chain rule]

Notice we are computing this for a single neuron & we've dropped i completely for simplification.

* The eqⁿ to compute partial derivatives of the theta terms in the [lost] hidden layer neurons (i.e. neurons are of layer L-1) :

$$\frac{\partial J(\theta)}{\partial \theta^{L-2}} = \frac{\partial J(\theta)}{\partial a^L} \frac{\partial a^L}{\partial z^L} \frac{\partial z^L}{\partial a^{L-1}} \frac{\partial a^{L-1}}{\partial z^{L-1}} \frac{\partial z^{L-1}}{\partial \theta^{L-2}}$$

* Clearly, the 2 eqⁿs share some pieces in common, so a delta term (δ^L) can be used for the common pieces b/w the output layer & the hidden layer immediately before it (with the possibility that there could be many hidden layer if we wanted) :

$$\delta^L = \frac{\partial J(\theta)}{\partial a^L} \frac{\partial a^L}{\partial z^L}$$

* And we can go ahead & use another delta term (δ^{L-1}) for the pieces that would be shared by the first hidden layer & a hidden layer before that, if we had one.

$$\delta^{L-1} = \frac{\partial J(o)}{\partial a^L} \frac{\partial a^L}{\partial z^L} \frac{\partial z^L}{\partial a^{L-1}} \frac{\partial a^{L-1}}{\partial z^{L-1}}$$

$$\therefore \boxed{\delta^{L-1} = \delta^L \frac{\partial z^L}{\partial a^{L-1}} \frac{\partial a^{L-1}}{\partial z^{L-1}}}$$

* With these delta terms, our equations become:

$$\boxed{\frac{\partial J(o)}{\partial o^{L-1}} = \delta^L \frac{\partial z^L}{\partial o^{L-1}}}$$

$$\boxed{\frac{\partial J(o)}{\partial o^{L-2}} = \delta^{L-1} \frac{\partial z^{L-1}}{\partial o^{L-2}}}$$

⇒ Now, let's calculate these derivatives

→ Deriving the sigmoid gradient function

$$\sigma(x) \equiv \text{Sigmoid function} = \frac{1}{1+e^{-x}}$$

$$\sigma'(x) = \frac{d}{dx} (1+e^{-x})^{-1}$$

$$= -1 (1+e^{-x})^{-2} \frac{d}{dx} (1+e^{-x})$$

$$= \frac{-1}{(1+e^{-x})^2} e^{-x} (-1)$$

$$\therefore \sigma'(x) = \frac{-e^{-x}}{(1+e^{-x})^2} = \left(\frac{1}{1+e^{-x}}\right) \left(-\frac{e^{-x}}{1+e^{-x}}\right)$$

$$= \left(\frac{1}{1+e^{-x}}\right) \left(\frac{\cancel{1+e^{-x}}}{\cancel{1+e^{-x}}} - \frac{1}{1+e^{-x}}\right)$$

$$\therefore \boxed{\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))}$$

→ Let's start with the output layer :

$$\frac{\partial J(\theta)}{\partial \theta^{L-1}} = \delta^L \frac{\partial z^L}{\partial \theta^{L-1}}$$

$$\delta^L = \frac{\partial J(\theta)}{\partial a^L} \frac{\partial a^L}{\partial z^L}$$

* Given $J(\theta) = -y \log(a^L) - (1-y) \log(1-a^L)$

$$\delta \hat{a}^L = k_0(\theta)$$



$$\therefore \frac{\partial J(\theta)}{\partial a^L} = -\frac{y}{a^L} + \frac{(1-y)}{1-a^L}$$

* Given $a = g(z)$ where $g = \frac{1}{1+e^{-z}}$

$$\therefore \frac{\partial a^L}{\partial z^L} = a^L (1-a^L)$$

{ Using the derivation on prev. page }

BB

* Substituting values of $\frac{\partial a^L}{\partial z^L}$ & $\frac{\partial J(o)}{\partial a^L}$ in
the eqⁿ of δ^L :

$$\delta^L = \frac{\partial J(o)}{\partial a^L} \quad \frac{\partial a^L}{\partial z^L}$$

$$= \left(\frac{1-y}{1-a^L} - \frac{y}{a^L} \right) \left(a^L \quad (1-a^L) \right)$$

$$= \frac{a^L(1-y) - y(1-a^L)}{a^L(1-a^L)} \cdot [a^L(1-a^L)]$$

$$= a^L - \cancel{ya^L} - y + \cancel{a^Ly}$$

∴ $\boxed{\delta^L = a^L - y}$

* Given $z^L = \phi^{L-1} a^{L-1}$

∴
$$\frac{\partial z^L}{\partial \phi^{L-1}} = a^{L-1}$$

→ Putting it all together for the final layer

$$\frac{\partial J(\phi)}{\partial \phi^{L-1}} = \delta^L \frac{\partial z^L}{\partial \phi^{L-1}}$$

∴
$$\frac{\partial J(\phi)}{\partial \phi^{L-1}} = (\alpha^L - y) (a^{L-1})$$

→ Let's continue on for the hidden layer :

$$\frac{\partial J(\theta)}{\partial \theta^{L-2}} = \delta^{L-1} \frac{\partial z^{L-1}}{\partial \theta^{L-2}}$$

* $\frac{\partial z^L}{\partial a^{L-1}} = \theta^{L-1} \quad \left\{ \text{As } z^L = \theta^{L-1} \cdot a^{L-1} \right\}$

* $\frac{\partial a^{L-1}}{\partial z^{L-1}} = a^{L-1} (1 - a^{L-1}) \quad \left\{ \text{As } \sigma'(z) = \sigma(z)(1 - \sigma(z)) \right\}$

$$\therefore \delta^{L-1} = \delta^L \frac{\partial z^{L-1}}{\partial a^{L-1}} \frac{\partial a^{L-1}}{\partial z^{L-1}}$$

$$\therefore \boxed{\delta^{L-1} = \delta^L \theta^{L-1} a^{L-1} (1 - a^{L-1})}$$

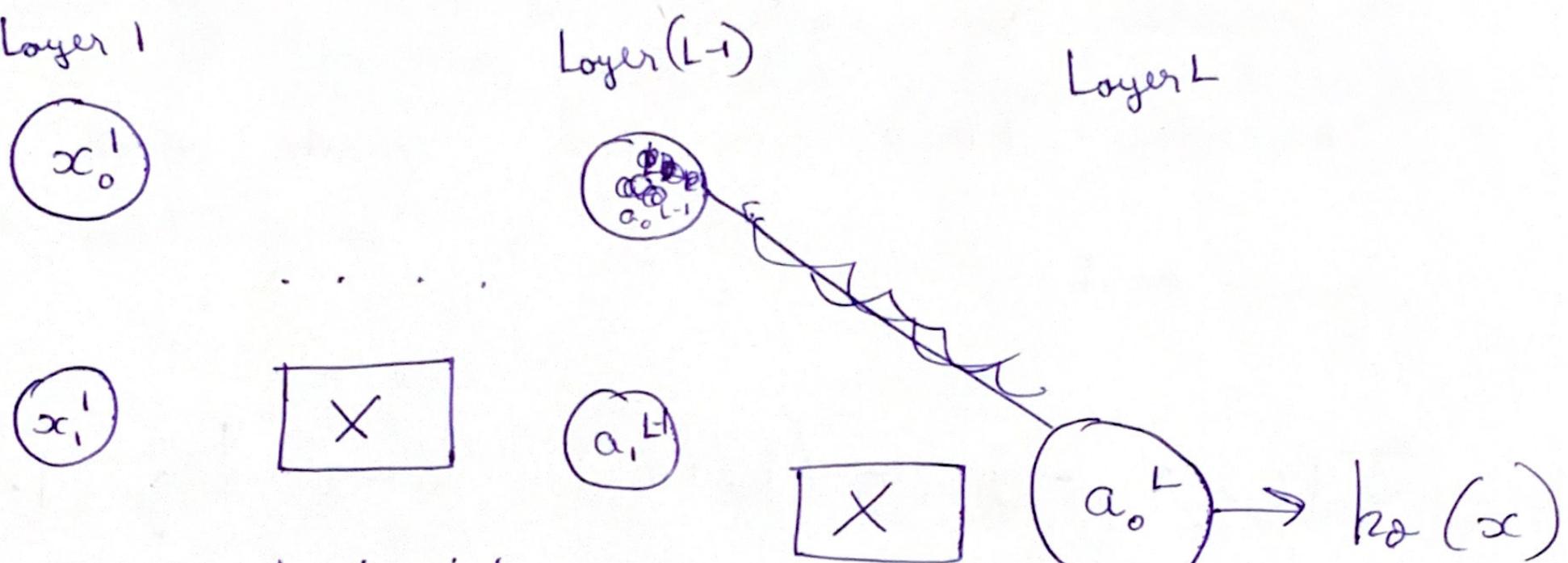
$$\therefore \frac{\partial J(\theta)}{\partial \theta^{L-2}} = (\delta^{L-1}) \frac{\partial z^{L-1}}{\partial \theta^{L-2}} \quad \left\{ \text{Use } z^{L-1} = \theta^{L-2} a^{L-2} \right\}$$

$$= \left(\delta^L \frac{\partial z^L}{\partial a^{L-1}} \frac{\partial a^{L-1}}{\partial z^{L-1}} \right) a^{L-2} \quad \left\{ \begin{array}{l} \text{Use } \\ z^L = \theta^{L-1} a^{L-1} \\ \text{and } \\ \theta^{L-1} = \theta^{L-2} a^{L-2} \end{array} \right.$$

$$\therefore \boxed{\frac{\partial J(\theta)}{\partial \theta^{L-2}} = \left[(a^{L-2} \delta^L) (a^{L-1}) (1 - a^{L-1}) \right] (a^{L-2})}$$

has fully connected layers

- * Supp. our neural net \hat{y} looks something like this :



For any arbitrary node j in layer L

$$a_j^L = g(z_j^L) = g(\theta_{j,0}^{L-1} a_0^{L-1} + \theta_{j,1}^{L-1} a_1^{L-1} + \dots + \theta_{j,s_{L-1}}^{L-1} a_{s_{L-1}}^{L-1})$$

where s_{L-1} = no. of nodes in layer (L-1)

a_k^{L-1} = k^{th} mode/neuron of layer (L-1)

* Since we have only one output node in our network, we can drop the using similar arguments as earlier (in the single neuron connected NN) we can write our cost function as :

$$J(\theta) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

* Calculating $\frac{\delta J(\theta)}{\delta \theta_{ijg_L^{L-1}}}$:

$$\frac{\delta J(\theta)}{\delta \theta_{ijg_L^{L-1}}} = \frac{\delta J(\theta)}{\delta a_g^L} \cdot \frac{\delta a_g^L}{\delta z_{g_L^L}} \cdot \frac{\delta z_{g_L^L}}{\delta \theta_{ijg_L^{L-1}}}$$

* For prev. layer or [last] hidden layer :

$$\frac{\delta J(\theta)}{\delta \theta_{ijg_L^{L-2}}} = \frac{\delta J(\theta)}{\delta a}$$