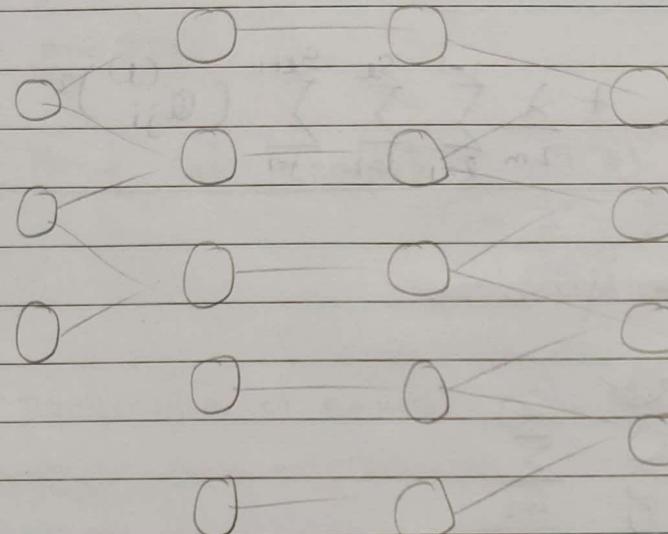


## Week 5

-- @ cost function -- --

Focus → Application of Neural networks on classification problems.

### \* Neural Network (classification)



layer 1      layer 2      layer 3      layer 4

Training data  
 $\rightarrow \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$L =$  total no. of layers in network.  
 $L=4$

$S_L =$  no. of units in layer L  
 (not counting bias unit)

$S_1 = 3, S_2 = 5, S_3 = 5, S_4 = S_L = 4$

• Binary classification:

$$y = 0 / y = 1$$

| 1 output unit

$$h_\theta(x) \in \mathbb{R}$$

$$S_L = 1$$

$$k = 1$$

• Multiclass classification ( $k$  classes)  
 $(k \geq 3)$

$$\Rightarrow y \in \mathbb{R}^k \text{ e.g. } \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

,  $k$  output units.

$$h_\theta(x) \in \mathbb{R}^k$$

$$S_L = k$$

## \* Cost function

- logistic reg<sup>n</sup>.

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

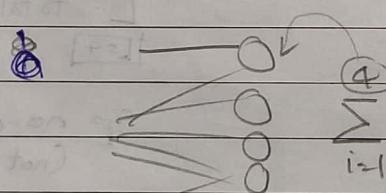
- neural network.

$h_{\theta}(x) \in \mathbb{R}^k$   $\xrightarrow{\text{k dimensional vector}}$  selects  $i^{\text{th}}$  element of vector  
 $(h_{\theta}(x))_i = i^{\text{th}}$  output

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

$(K=1)$  if we have

binary classification problem.



- Basically the logistic regression algorithm's cost function but summing that cost function over each of my four output units in turn.
- In regularization term, basically we are summing  $\sum_{i,j} \theta_{ij}$

Doesn't make much difference except we don't sum over terms corresponding bias values ( $i=0$ )

$$\sum_i \sum_j$$

Ques Suppose we try to minimize  $J(\theta)$ , we need to supply code to compute?

→  $J(\theta)$  and partial derivative terms  $\frac{\partial}{\partial \theta_{ij}}$   
for every  $i, j, l$ .

- @ Back propagation Algorithm
- \* Back propagation algorithm

Gradient computation

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

Need code to compute  $\rightarrow J(\theta)$

$$\rightarrow \frac{\partial}{\partial \theta_{ij}} J(\theta)$$

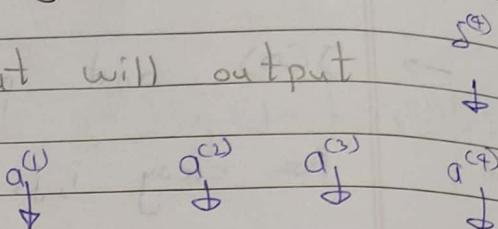
→ parameter of neural network  $\theta_{ij}^{(l)} \in \mathbb{R}$

## The flow of algorithm

- Given one training example  $(x, y)$ :

- Forward propagation:  $\rightarrow$  what will output

$$a^{(1)} = x$$



$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$

add bias term

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)})$$

$$z^{(4)} = \Theta^{(3)} a^{(3)}$$

$$a^{(4)} = h_\theta(x) = g(z^{(4)})$$

- Gradient computation: Backpropagation algorithm

$\hookrightarrow$  To compute partial derivatives

Intuition:  $\delta_j^{(l)}$  = "error" of node  $j$  in layer  $l$

$a_j^{(l)}$  = "activation"

For each output unit (layer  $L=4$ )

$$\delta_j^{(4)} = a_j^{(4)} - y_j$$

label value from training example.

$(h_\theta(x))_j$

$$\delta^{(4)} = a^{(4)} - y \rightarrow \text{(vectorize Implementation)}$$

$\rightarrow$  Dimension = # output units ( $C$ )

\* → Element wise multiplication of those two vectors.  
 $g^{(2)}$  → Deactivation of activation function  $g$ , elevated at the input values given by  $a^{(2)}$

**Sadguru**  
 Page No: 58  
 Date: / /

$$\delta^{(3)} = (\Theta^{(2)})^T \delta^{(4)} \cdot * g'(z^{(2)})$$

$$\delta^{(2)} = (\Theta^{(1)})^T \delta^{(3)} \cdot * g'(z^{(1)})$$

(2 PTO for detail)

$$a^{(3)} \cdot * (1 - a^{(3)})$$

$$a^{(2)} \cdot * (1 - a^{(2)})$$

No  $\delta^{(0)}$  term → no errors, we don't want to change input.

Back propagating  
 errors from output layer to layer 3.

$$\delta^{(2)} \leftarrow \delta^{(3)} \leftarrow \delta^{(4)}$$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

(Ignoring  $\lambda$  / if  $\lambda=0$  /  
 Ignoring regularization term  
 ↳ Fix later)

\* Back propagation algorithm.

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

Set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ ) — capital  $\Delta$ .  
 (will be used to compute  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$ )

$\Delta$ 's are going to be used as accumulators  
 that will slowly add things in order to  
 compute partial derivatives.

loop:

For  $i = 1$  to  $m$  ---  $i$ th iteration corresponds to  $(x^{(i)}, y^{(i)})^T$   
 Set  $a^{(0)} = x^{(i)}$  training set.

Perform forward propagation to compute  $a^{(l)}$  for  
 $l = 2, 3, \dots, L$

Using  $y^{(i)}$ , compute  $\delta^{(L)} = a^{(L)} - y^{(i)}$

compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$  -- back prop

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l+1)} \delta_i^{(l+1)}$$

$$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \quad \text{if } j \neq 0 \quad \text{Two cases}$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j = 0$$

corresponds to bias term

$$\frac{\partial J(\theta)}{\partial \Theta_{ij}^{(l)}} = D_{ij}^{(l)} \quad \text{This will be used in gradient descent}$$

que

 $(x^{(1)}, y^{(1)})$ ,  $(x^{(2)}, y^{(2)})$ 

FP = forward propagation / BP = Backward propagation

Sequence of operation.

→ FP using  $x^{(1)}$  followed by BP using  $y^{(1)}$ .

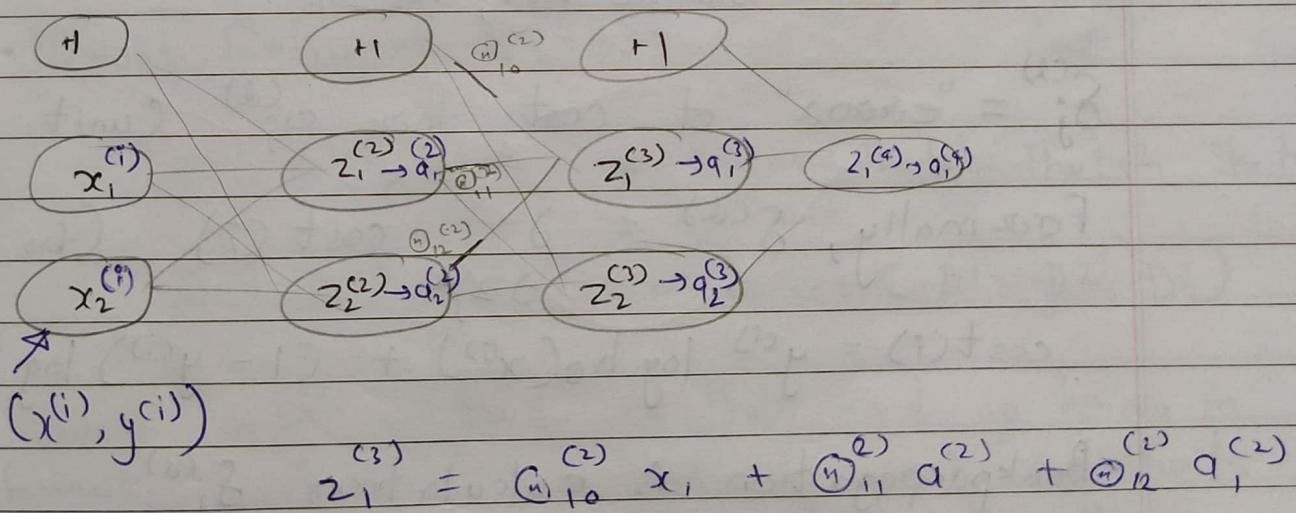
Then

FP using  $x^{(2)}$  →  $y^{(2)}$ 

★ Forward Propagation In this vid :- How to use back propog. to compute derivatives of your loss fun.

→ C ★ Back Propagation intuition - - -

• Forward propagation.



$$z_1^{(1)} = w_{11}^{(1)} x_1^{(i)} + w_{12}^{(1)} x_2^{(i)} + b_1^{(1)}$$

At a time we use only one example, FP it. BP it.

Sadguru  
Page No: 61  
Date: / /

- What is back propagation doing?

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{m_l} \sum_{j=1}^{s_{l+1}} (\theta_{ji}^{(l)})^2$$

Focussing on a single example  $x^{(i)}, y^{(i)}$ , the case of 1 output unit, and ignoring regularization ( $\lambda=0$ )

$$\text{cost}(i) = y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log h_\theta(x^{(i)})$$

(Think of  $\text{cost}(i) \approx (h_\theta(x^{(i)}) - y^{(i)})^2$ )

↳ Measures how well is neural network doing on example?

Just to clarify intuition.

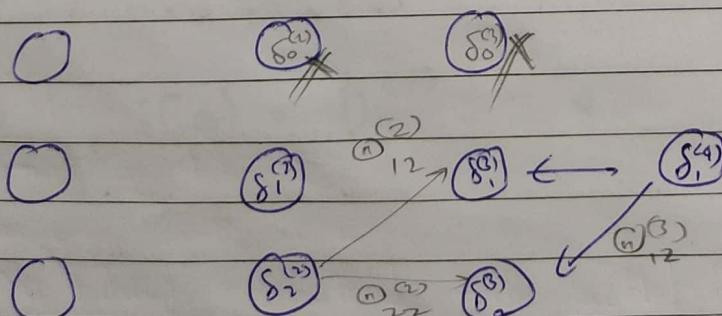
Works in similar way, not same.

$\delta_j^{(l)}$  = "error" of cost for  $\theta_{ji}^{(l)}$  (unit  $j$  in layer  $l$ )

Formally,  $\delta^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$  (for  $j \geq 0$ ), where

$$\text{cost}(i) = y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log h_\theta(x^{(i)})$$

• Backpropagation is calculation  $\delta_j^{(l)}$   
↳ quote error of activation value.



$$\delta_1^{(4)} = y^{(i)} - a_1^{(4)}$$

$$\delta_2^{(2)} = \Theta_{12}^{(2)} \delta_1^{(3)} + \Theta_{22}^{(2)} \delta_2^{(3)}$$

$$\delta_2^{(3)} = \Theta_{12}^{(3)} \delta_1^{(4)}$$

Implementation Note Unrolling --

Implementation  
detail of

unrolling parameters

from matrices

into vectors,

which we

need in

order to

use advanced

optimization  
routine.

Advanced optimization ~~for~~  $\rightarrow$  logistic regn.

output = cost function & derivatives

function [jVal, gradient] = costFunction(theta)

$\rightarrow \mathbb{R}^{n+1}$

Input

$\rightarrow \mathbb{R}^m$  vector

Advanced optimization algorithm optTheta = fminunc(@costFunction, initialTheta, options)

• Neural Network (L=4):

$\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$  - matrices (Theta1, Theta2, Theta3)

$D^{(1)}, D^{(2)}, D^{(3)}$  - matrices (D1, D2, D3)

"unroll" into vectors so they end up being in  
format suitable for passing in to as the theta  
off for getting out for gradient there  
 $\rightarrow$  so, we could pass into fminunc (vector)

e.g. input layer, hidden layer, output layer.

$$S_1 = 10, S_2 = 10, S_3 = 1$$

$$\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$$

$$\mathbb{R}^{(S_{l+1}) \times (S_l)}$$

$$D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$$

$\alpha$

$\alpha$

$\alpha$

$\alpha$

$\alpha$

$\alpha$

$\alpha$

$\alpha$

$$\text{thetaVec} = [\Theta^{(1)}(:, :); \Theta^{(2)}(:, :); \Theta^{(3)}(:, :)];$$

Vector

Matrix's

$$DVec = [D^{(1)}(:, :); D^{(2)}(:, :); D^{(3)}(:, :)];$$

$$\Theta^{(1)} = \text{reshape}(\text{thetaVec}(1:110), 10, 11);$$

$$\Theta^{(2)} = \text{reshape}(\text{thetaVec}(111:220), 10, 11);$$

$$\Theta^{(3)} = \text{reshape}(\text{thetaVec}(221:231), 1, 11);$$

We Just converted the  $\Theta$ [Matrix] to  $\Theta$ [Vector]

Ques :  $D_1$  is  $10 \times 6$  matrix &  $D_2$  is a  $1 \times 11$  matrix.  
You set  $DVec = [D1(:, 1); D2(:)]$ ;  
How to get  $D_2$  back from  $DVec$ .  
 $\rightarrow \text{reshape}(DVec(61:71), 1, 11)$

Sadguru  
Page No: 64  
Date: / /

## • Learning Algorithm.

We have to perform ~~Optimization~~  
because we are going to use it.

• Have initial parameter  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$

• Unroll to get  $\Theta$  (long vector) initialTheta to pass to

$fminunc(@costFunction, initialTheta, options)$

# Demo: >>

>> Theta1 = ones(10, 11)

>> Theta2 = 2 \* ones(10, 11)

>> Theta3 = 3 \* ones(10, 11)

>> thetaVec = [Theta1(:); Theta2(:); Theta3(:)];

>> size(thetaVec) = 231, 1

>> thetaVec =

1
2
3
3

>> reshape(thetaVec(1:110), 10, 11)

>> reshape(thetaVec(111:220), 10, 11)

>> reshape(thetaVec(221:231), 1, 11)

Advantage of Matrix representation → Advantage of Vector are

easier to do forward prop ↴ Easier to use advanced optimization algorithm  
& back prop. ↴ because they tend to think our parameters are unrolled in big vector

## # Implementing cost function.

(2) function [JVal, gradientVec] = costFunction(thetaVec)

$\theta_1, \theta_2, \theta_3$  will be in matrix form From ThetaVec, get  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$  using reshape

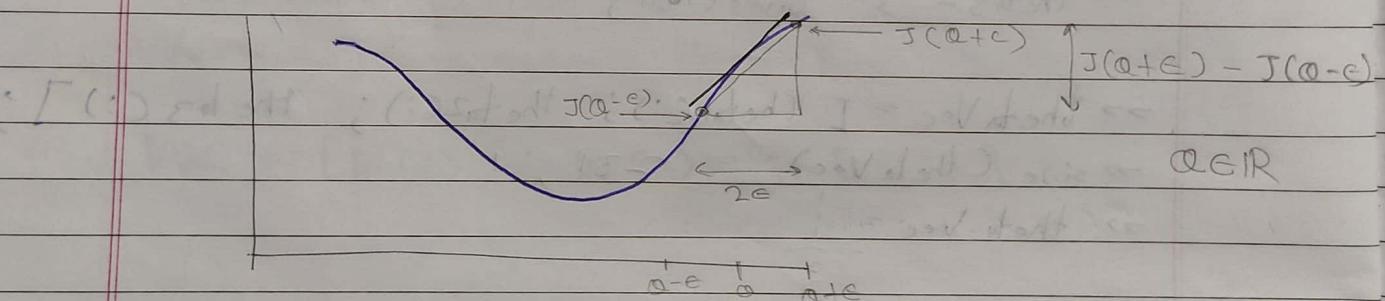
Use forward prop / back prop to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\theta)$ .

I think its calculating DV<sub>ec</sub> Unroll  $D^{(1)}, D^{(2)}, D^{(3)}$  to get gradient Vec.

## --@ Gradient Checking

→ Gradient Checking → Checking whether it is running right or not.

## # Numerical estimation of gradients :-



Hence we want slope of /

we estimate slope of both lines to be sum.

double sided difference.

one sided difference

$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

$$\frac{J(\theta + \epsilon) - J(\theta)}{\epsilon}$$

If we take  $\epsilon$

$$\epsilon = 10^{-4} \text{ (In range of } 10^{-4})$$

very small, it is actually definition of derivative.

Ques

$$J(\theta) = \theta^3, \text{ let } \theta=1, \epsilon=0.01$$

$\frac{\partial J(\theta)}{\partial \theta} = 3$ , what you get by approximation formula  $\approx 3.0001$

Previously we considered  $\theta$  as scalar parameter

Now we will be considering Moving on more general case.

### \* Parameter Vector $\theta$ .

Maths

•  $\theta \in \mathbb{R}^n$  -- e.g.  $\theta$  is "unrolled" version of  $(\theta^{(1)}, \theta^{(2)}, \theta^{(3)})$

•  $\theta = \theta_1, \theta_2, \dots, \theta_n$

$\Rightarrow$

$$\frac{\partial J(\theta)}{\partial \theta_1} \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial J(\theta)}{\partial \theta_2} \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial J(\theta)}{\partial \theta_n} \approx \frac{J(\theta_1, \theta_2, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_n - \epsilon)}{2\epsilon}$$

Program

for i=1:n — dimension of  $\theta$ .

thetaPlus = theta;

thetaPlus(i) = thetaPlus(i) + EPSILON;

thetaMinus = theta;

thetaMinus(i) = thetaMinus(i) - EPSILON;

$$\text{gradApprox}(i) = \frac{(J(\theta_{\text{plus}}) - J(\theta_{\text{minus}}))}{2 * \text{EPSILON}}$$

end;

i We implement for loop to calculate top <sup>partial</sup> derivatives of the cost function for respect to every parameter in that network. and

P.T.O.

Gradient = slope.

Sadguru

Page No:

67

Date:

/ /

computes slowly

numerically computed derivative

computes tasks

Backprop test

Check that good Approx  $\approx D_{Vec}$

back prop

This must be equal, at least for few decimals. To make sure everything is

- ii We can then take the gradient that we got from back prop.

•  $D_{Vec}$  was derivative we get from back prop.

• Backprop was relatively efficient way to compute a derivative or a partial derivative of a cost fun with respect to all our parameters.

Now we can use  $D_{Vec}$  in advanced optimization algorithm

#### \* Implementation Note:-

1) Implement back prop to compute  $D_{Vec}$  (controlled  $D^{(0)}, D^{(1)}, D^{(2)}$ ).

2) Implement numerical gradient check to compute good Approx.

3) Make sure they give similar values. (1.82)

4) Turn off gradient checking. Using back prop code for learning.

Just move on, everything is ok, go on using back prop to train model.

Disable gradient checking code before training classifier

computes slow

Ques

Ans → We use back prop over numerical gradient computation during learning.

because

numerical gradient algo is slow.

## --@ Random initialization

Before running any algo, gradient descent or any other advanced optimization, we need to pick some initial value for parameters theta.

Even advanced optimization algo, assumes we will pass initialTheta. After initializing, we will go to gradient descent to minimize  $J(\theta)$ .

We initialize  $\theta = [0, 0]$  in regression (even in logistic). This will not work in neural network.

$\rightarrow \text{optTheta} = \text{fminunc}(@\text{costFunction}, \underline{\text{initialTheta}}, \text{options})$

## \* Zero Initialization

Random youtube video:-

Bias → initialize all to zero ( $c_0$ )

Weights →

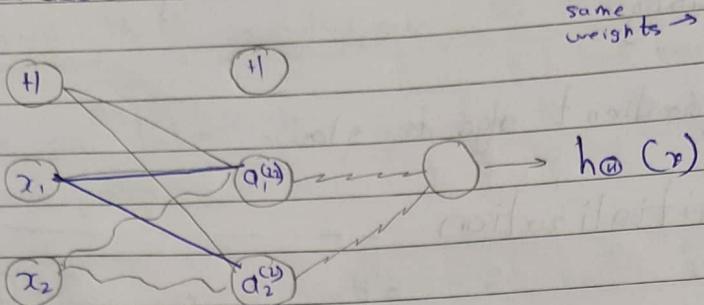
- Can't initialize to zero → then all gradients will be zero.

- Can't initialize all weights to same value.  
→ all hidden units will behave same.

Soln →  $W_{ij}^{(k)}$  from  $U[-b, b]$ , where  $b = \sqrt{6}$

$H_x$  = no of units in layer  $k$ .  
 $\sqrt{H_x + H_{x-1}}$

## ★ Zero initialization



If we initialize  $\Theta_{ij}^{(l)} = 0$ , for all  $i, j, l$ .

$$\rightarrow a_1^{(c2)} = a_2^{(c2)}$$

Because outgoing weights are same.

$$\rightarrow \delta_1^{(c2)} = \delta_2^{(c2)}$$

This means

$$\frac{\partial J(\Theta)}{\partial \Theta_{01}^{(c1)}} = \frac{\partial J(\Theta)}{\partial \Theta_{02}^{(c1)}}$$

Even after one gradient descent update

i.e.  $\Theta_{01}^{(c1)} := \alpha \frac{\partial J(\Theta)}{\partial \Theta_{01}^{(c1)}}$

$$\Theta_{02}^{(c1)} = \alpha \frac{\partial J(\Theta)}{\partial \Theta_{02}^{(c1)}}$$

It will remain same.

i.e.  $\Theta_{01}^{(c1)} = \Theta_{02}^{(c1)}$  --- (non-zero value)

(It changes from 0 to non-zero, but it will be equal)

Even after one another update, both values will be same  
similarly  $\Theta_{01}^{(c1)} = \Theta_{12}^{(c1)}$  &  $\Theta_{31}^{(c1)} = \Theta_{32}^{(c1)}$

So, even after one iteration,  $a_1^{(c2)} = a_2^{(c2)}$   
Hidden layers will have same functions as input (checked units are same)

So, this will highly redundant representation.  
because, it will consider only one feature,  
as all  $a_{ij}^{(0)} = a_{j(i)}^{(0)} = a_{(i)j}^{(0)}$

- Therefore we go with random initialization.

### \* Random initialization: Symmetry breaking.

problem of symmetric ways  $\rightarrow$  The problem we faced in zero initialization

(Ways all being same)

Random initialization breaks symmetry

Initialize each  $a_{ij}^{(0)}$  to a random value in  $[-\varepsilon, \varepsilon]$   
(i.e.  $-\varepsilon \leq a_{ij}^{(0)} \leq \varepsilon$ )

e.g.  $\Theta_{01} = \text{rand}(10, 11) * (2 * \text{INIT-EPSILON}) - \text{INIT-EPSILON}$   
 $\hookrightarrow$  (random 10x11 matrix between 0 & 1)

$\Theta_{02} = \text{rand}(1, 11) * (2 * \text{INIT-EPSILON}) - \text{INIT-EPSILON}$

$\alpha = \text{rand}(1, 1) * (2 * \text{INIT-EPSILON}) - \text{INIT-EPSILON}$   $\rightarrow \alpha$  will only one number  
 $a_{ij}^{(0)} = \alpha$  for all  $i, j, l$

This will not work, & will fail to break symmetry.

- The idea is to randomly initialize the values to small values close to zero. b/w  $-\varepsilon$  &  $\varepsilon$

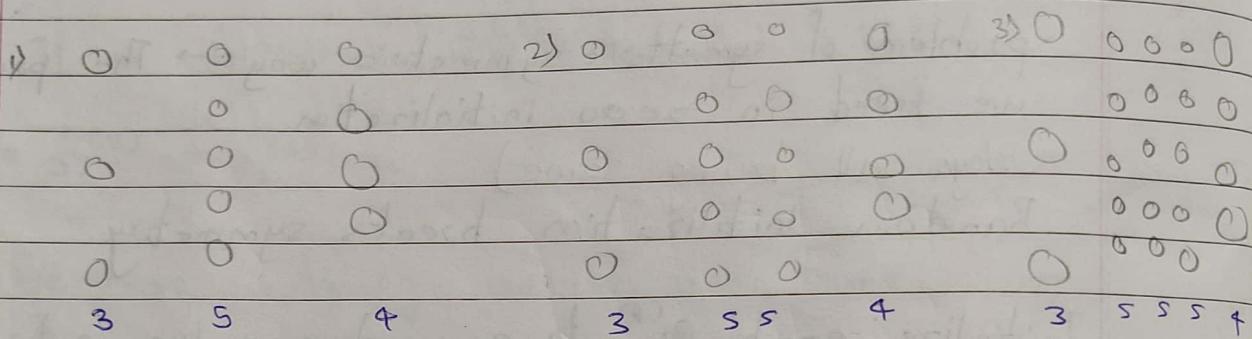
Then use gradient descent/any other advanced optimization video to minimize  $J(\Theta)$  as a function of parameters theta starting from randomly chosen initial value for parameters.

-- @ Putting it together

★ Training a neural network.

To train neural network

↳ First choose architecture.



# No. of input units: Dimension of features  $[x^1]$

# No. of output units: Number of classes no of classes

in our classification prob.

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ or } y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$y \in \{1, 2, 3, \dots, 10\}$$

~~$y = s$~~

$$y = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

# No. of hidden layers

- Reasonable default :-

- 1 hidden layer

or

- $\geq 1$  hidden layer, having same no of hidden units in every layer (usually more hidden layers, the better)  $\rightarrow$  but becomes computationally expensive.

- is most common architecture.

No. of hidden must be somewhat comparable to input layers. (3 to 4 times no of inputs)

★ Training a neural network.

- 1) Randomly initialize weights (near to zero)

- 2) Implement forward propagation to get  $h_{\Theta}(x^{(i)})$  for any  $x^{(i)}$

- 3) Implement forward propagation to code to compute cost function  $J(\Theta)$

- 4) Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \Theta_j} J(\Theta)$   
 $\frac{\partial}{\partial \Theta_j} J(\Theta)$  P.T.O. to compute

The idea is, we take  $(x^1, y^1)$  (first training ex) FP & BP  
 then we take  $(x^2, y^2)$  FP, BP  
m times.

Page No: 73  
 Date: / /

no of training examples.

for  $i = 1 : m$

{  
 Perform forward propagation and back  
 propagation using  $(x^{(i)}, y^{(i)})$

(Net activations  $a^{(l)}$ ) and delta term ( $\delta^{(l)}$ ) for  $L_{2,L}$

$$\Delta^{(d)} := \Delta^{(d)} + \delta^{(d+1)} (a^{(d)})^T$$

Compute  $\frac{\partial}{\partial w_{jk}^{(d)}} J(\theta)$

$$D_{jk}^{(d)} := \frac{1}{m} \Delta_{jk}^{(d)} + \lambda \Theta_{jk}^{(d)} - \Theta_{j0}^{(d)}$$

$$D_{jk}^{(d)} := \frac{1}{m} \Delta_{jk}^{(d)} \quad \dots \quad (If j=0)$$

$$\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ x^{(d)} & 0 & 0 \end{matrix}$$

- 5) Use gradient checking to compute  $\frac{\partial}{\partial w_{jk}^{(d)}} J(\theta)$   
 computed using back propagation vs using mathematical numerical estimate of gradient of  $J(\theta)$

[good Approx  $\approx DVec$ ]

If found right, disable gradient checking code.

6) Use gradient descent or advanced optimization method with backpropagation to try to minimize  $J(\theta)$  as a function of parameters ( $\theta$ )

BP computes gradients

$$\frac{\partial J(\theta)}{\partial \theta_{jk}^{(l)}}$$

$J(\theta)$  for neural network is non-convex

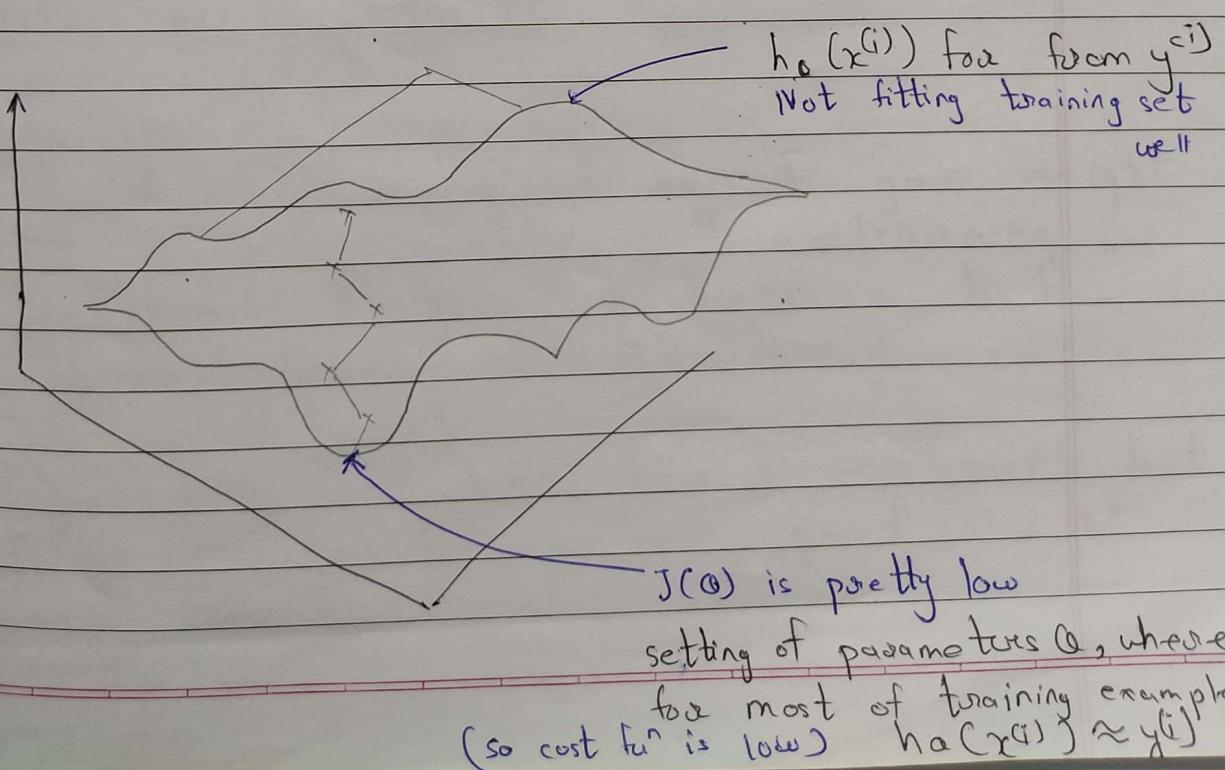
In theory algo like gradient descent / AO method can stuck in local optimum.

but

In practice gradient descent / AO finds very good local minimum.

It does not give guarantee to global minimum.

In practice this is not a huge problem.



Gradient descent - start from random point  
And it will repeatedly go downhill

Back propagation - Computes direction of Gradient

Gradient descent - Taking little steps downhill.  
until it gets global optimum  
not necessarily

Q If we use gradient descent with BP.  
to minimize  $J(\theta)$  as a fun of  $\theta$ .

How to verify the learning algo is correct

→ Plot  $J(\theta)$  as a function of number of iterations  
and make sure it is decreasing with every iteration.

(At least non-increasing)

→ Plot  $J(\theta)$  as a function of the number of iterations to make sure the parameter values are improving in classification accuracy