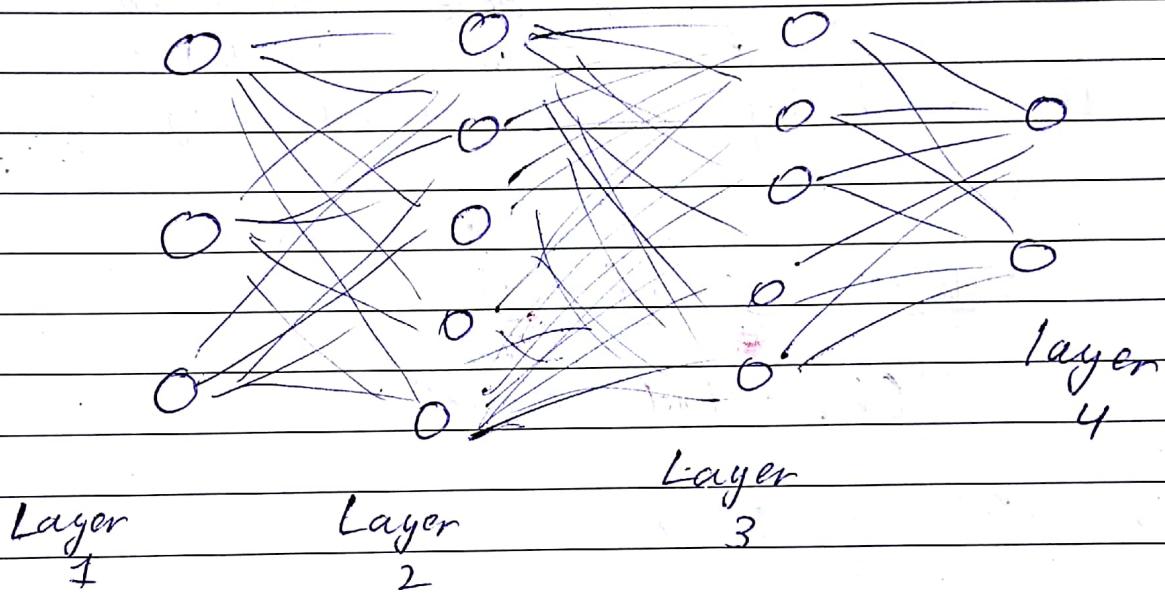


My
Wing
know
our score.

Cost function

The focus of neural networks will be predominantly on classification problems.



Expressing the entire data set:

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

L = total no. of layers in a network
(in this case $L=4$)

s_l = no. of m units (not counting bias unit) in layer l

Date: _____

In this case $s_1 = 3, s_2 = 5, s_3 = 5, s_4 = 2$

Binary classification

$y=0, y=1$ (1 output unit)

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

$$s_L = 1$$

$$\text{or } k=1$$

k denotes the number
of nodes (units) in
the output layer

Multi-class classification

(k distinct classes)

$$y \in \mathbb{R}^k \quad s_L = k$$

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

$$y = h_\theta(x) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Date: _____

Neural Network cost function:

$h_{\theta}(x) \in \mathbb{R}^k$, $(h_{\theta}(x))_i = i^{\text{th}} \text{ output}$

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k) \right] +$$

$$\frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{m_l} \sum_{j=1}^{s_{l+1}} (\theta_{ji}^{(l)})^2$$

cost of a
multi-class classification
neural network

Explanation:

$$X = \begin{pmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & \dots & x_n^{(2)} \\ \vdots & & & & \\ x_1^{(m)} & x_2^{(m)} & x_3^{(m)} & \dots & x_n^{(m)} \end{pmatrix}$$

$$\gamma = \begin{pmatrix} y_1^{(1)} & y_2^{(1)} & \dots & y_{s_L}^{(1)} \\ y_1^{(2)} & y_2^{(2)} & \dots & y_{s_L}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ y_1^{(m)} & y_2^{(m)} & \dots & y_{s_L}^{(m)} \end{pmatrix}$$

When outputting the prediction:

$$\begin{pmatrix} h(x^{(1)})_1 & h(x^{(1)})_2 & \dots & h(x^{(1)})_{s_L} \\ h(x^{(2)})_1 & h(x^{(2)})_2 & \dots & h(x^{(2)})_{s_L} \\ \vdots & \vdots & \ddots & \vdots \\ h(x^{(m)})_1 & h(x^{(m)})_2 & \dots & h(x^{(m)})_{s_L} \end{pmatrix}$$

Thus:

$$\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k$$

means to go through every row of the dataset & cycle through all the outputs (with k) & apply $\log(h_\theta(x^{(i)}))_k$ to the value $y_k^{(i)}$ that is one cost function for $y=1$

Date: _____

Since there is only one such value then this function is applied once per row.

The remaining Os in $y^{(i)}$ will then be applied the following function:

$$\sum_{i=1}^m \sum_{k=1}^K (1 - y_k^{(i)}) \log (1 - (h_\theta(x^{(i)}))_k)$$

(cost function for $y=0$)

$$\sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

simply implies
to use every
 $\theta_{ji}^{(l)}$ value possible
regularization for each layer in the neural network

Note that the values of i, j start at 1, the bias is not included within regularization

Back propagation algo

Date: _____

We wish to find $\min_{\theta} J(\theta)$

and for that we need:

$$-J(\theta)$$

$$-\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) \quad \theta_{ij}^{(l)} \in \mathbb{R}$$

Given that there is only 1 training example: (x, y) then forward propagation is executed as such:

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

bias included

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

$$a^{(2)} = g(z^{(2)})$$

$$\text{add } a_0^{(2)}$$

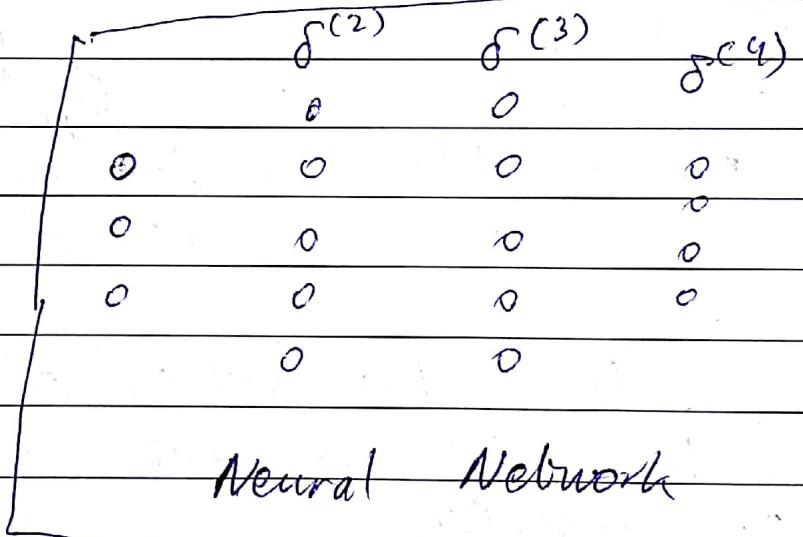
$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)})$$

$$\text{add } a_0^{(3)}$$

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

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



Backpropagation with 1 training example

Date: _____

To compute the derivatives we will use the backpropagation algo.

Intuition: compute $\delta_j^{(l)}$ = error in $a_j^{(l)}$ or error in node j in layer l

Consider this neural network example,

For each output unit (layer $l=4$):

$$\delta_j^{(4)} = a_j^{(4)} - y_j \rightarrow (h_\theta(x))_j$$

Vectorised implementation:

$$\delta^{(4)} = a^{(4)} - y$$

$$\delta^{(4)} \in \mathbb{R}^{s_L}, a^{(4)} \in \mathbb{R}^{s_L}, y \in \mathbb{R}^{s_L}$$

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

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

No $\delta^{(1)}$ as there can be no error associated with inputs
(we don't want to change them)

Date: _____

It can be proved that:

$$g'(z^{(3)}) = a^{(3)} \cdot (1 - a^{(3)})$$

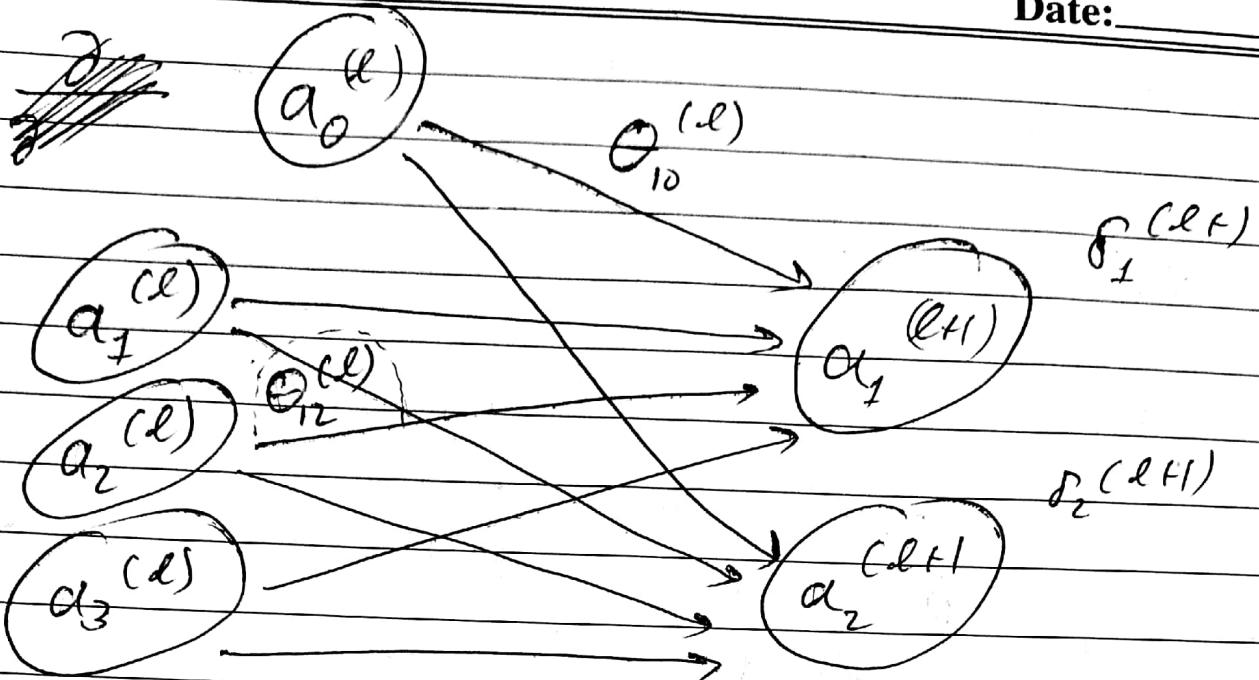
$$g'(z^{(2)}) = a^{(2)} \cdot (1 - a^{(2)})$$

This process is called back propagation because we compute the error of the output layer first ($\delta^{(n)}$ in this case) & then we keep taking steps backwards to compute $\delta^{(3)}$ then $\delta^{(2)}$ & so on

Via complicated mathematical proof it is possible to prove (with the help of $\delta^{(n)}$, $\delta^{(3)}$ & etc) that,

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

(ignoring the regularization term i.e. $\lambda=0$)



$$\frac{\partial}{\partial \theta_{12}^{(l)}} J(\theta) = a_2^{(l)} \delta_1^{(l+1)}$$

$$\frac{\partial}{\partial \theta_{10}^{(l)}} J(\theta) = a_0^{(l)} \delta_1^{(l+1)} = \delta_1^{(l+1)}$$

Note that $\theta_{0j}^{(l)}$ does not exist

thus this is not possible!

$$\frac{\partial}{\partial \theta_{0j}^{(l)}} J(\theta) = a_j^{(l)} \delta_0^{(l+1)}$$

No connection from $a_j^{(l)}$ to $\delta_0^{(l+1)}$

Backpropagation algorithm with multiple data points in training set:
 $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j)

→ This will be used to compute $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$
 capital δ
 (pronounced delta)

For $i=1$ to m :

- Set $a^{(i)} = x^{(i)}$
- Set $a^{(1)} = x^{(i)}$
- perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \dots, L$
- Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$
- Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

→ backpropagation & $\delta^{(1)}$ are non-existent as we do not wish to find the error on the inputs

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

→ on each instance of input, every $\Delta_{ij}^{(l)}$ is updated

I update per cycle for unique i, j, l & for every unique i, j, l the values are accumulated from i to m

↑ {different i 's}

$\Delta_{ij}^{(l)}$ can not exist as $\theta_{ij}^{(l)}$ is not calculated, ~~but already has the value of 1~~

Vectorised form of equation:

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

$s_{l+1} \times (s_l + 1)$

This is the case because of how θ is:

$$\theta = \begin{bmatrix} \theta_{10} & \theta_{12} & \dots \\ \theta_{20} & \theta_{21} & \dots \\ \theta_{30} & \theta_{31} & \dots \end{bmatrix}$$

$$f^{(l+1)} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \vdots & \ddots & \vdots \\ f_{s_{l+1}}^{(l+1)} & \dots & \dots \end{bmatrix}$$

$$a^{(l)} = [a_0^{(l)} \ a_1^{(l)} \ \dots \ a_{s_l}^{(l)}]$$

$$(1 \times s_{l+1})$$

$$\theta = \begin{bmatrix} \theta_{10} & \theta_{12} & \dots \\ \theta_{20} & \theta_{21} & \dots \\ \theta_{30} & \theta_{31} & \dots \end{bmatrix}$$

$$(s_{l+1} \times 1)$$

$$s_{l+1} \times (s_l + 1)$$

possible because:

$$\frac{\partial}{\partial \theta_{10}} a^{(l)} = a_0^{(l)} \delta_i^{(l+1)}$$

is possible

~~apply average &~~ add regularization

Date: _____

After the loop, execute:

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

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

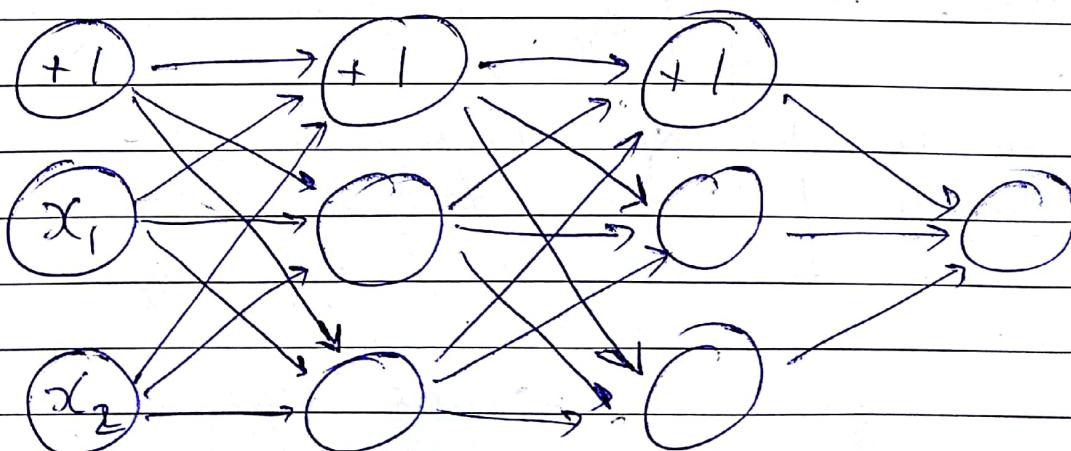
No regularization in bias term

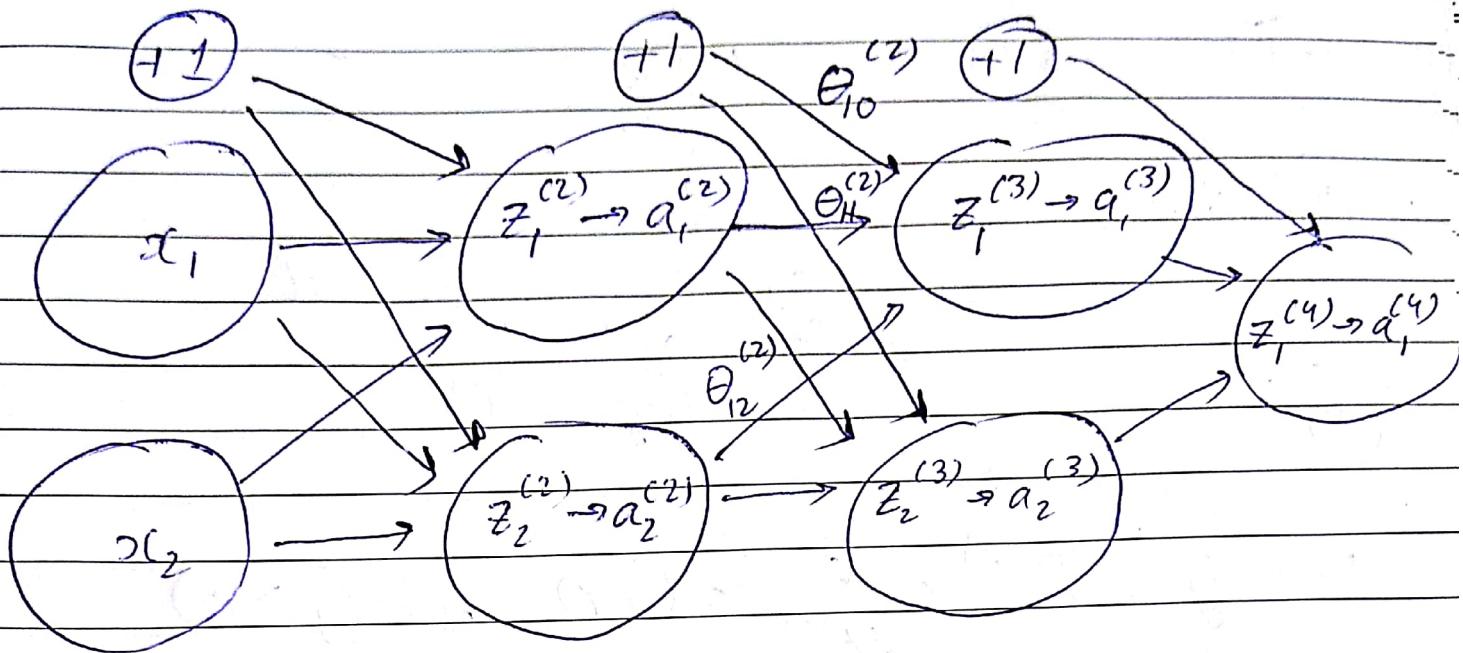
$$\left[\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)} \right]$$

With complex mathematical proof this can be proved

Back propagation Intuition

Use the following example to illustrate intuition



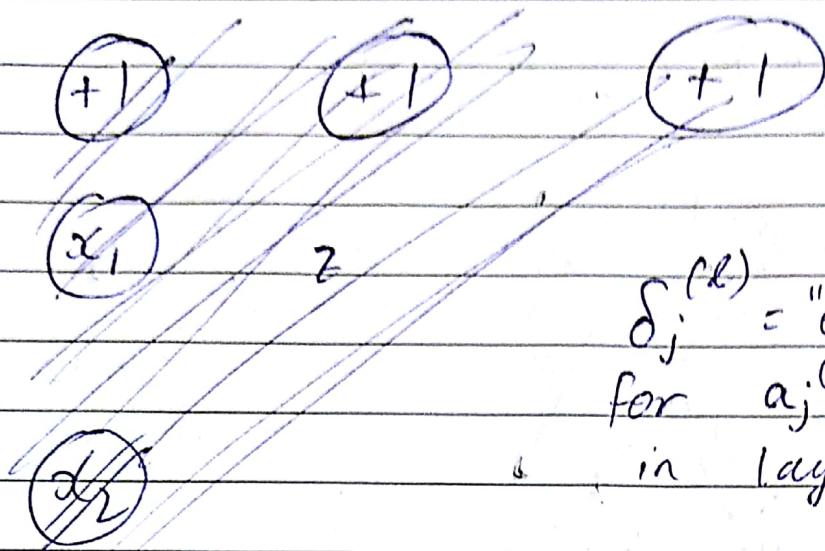


$$z_1^{(3)} = \theta_{10}^{(2)} \times 1 + \theta_{11}^{(2)} \cdot a_1^{(2)} + \theta_{12}^{(2)} \cdot a_2^{(2)}$$

Computations flow in this direction
for forward propagation

←
Computations flow in this direction
for backward propagation

~~Backward~~ Back propagation can be considered
the inverse operation of forward propagation



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

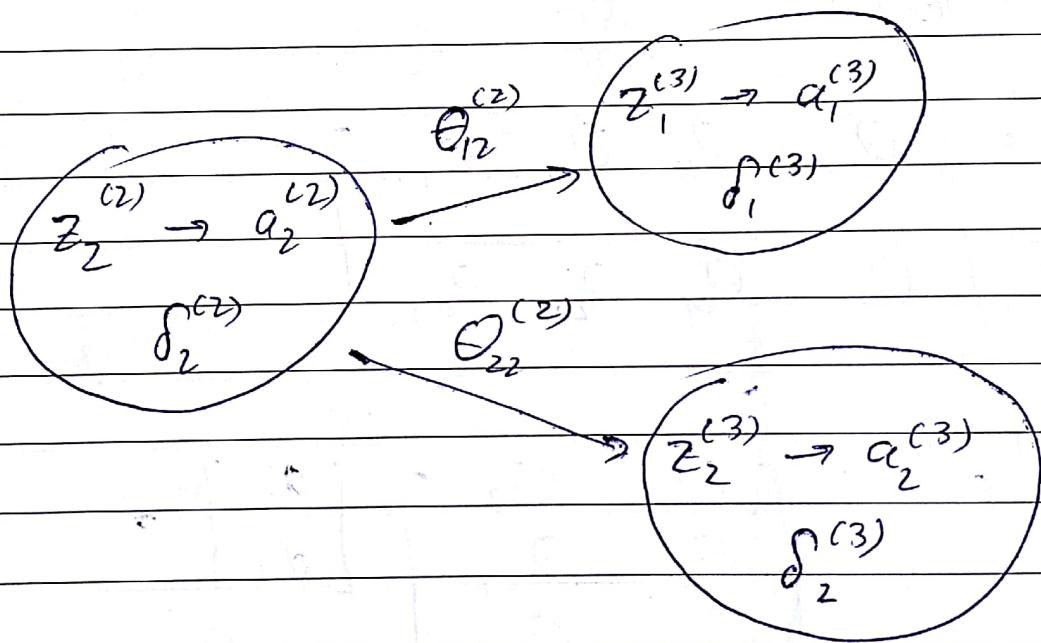
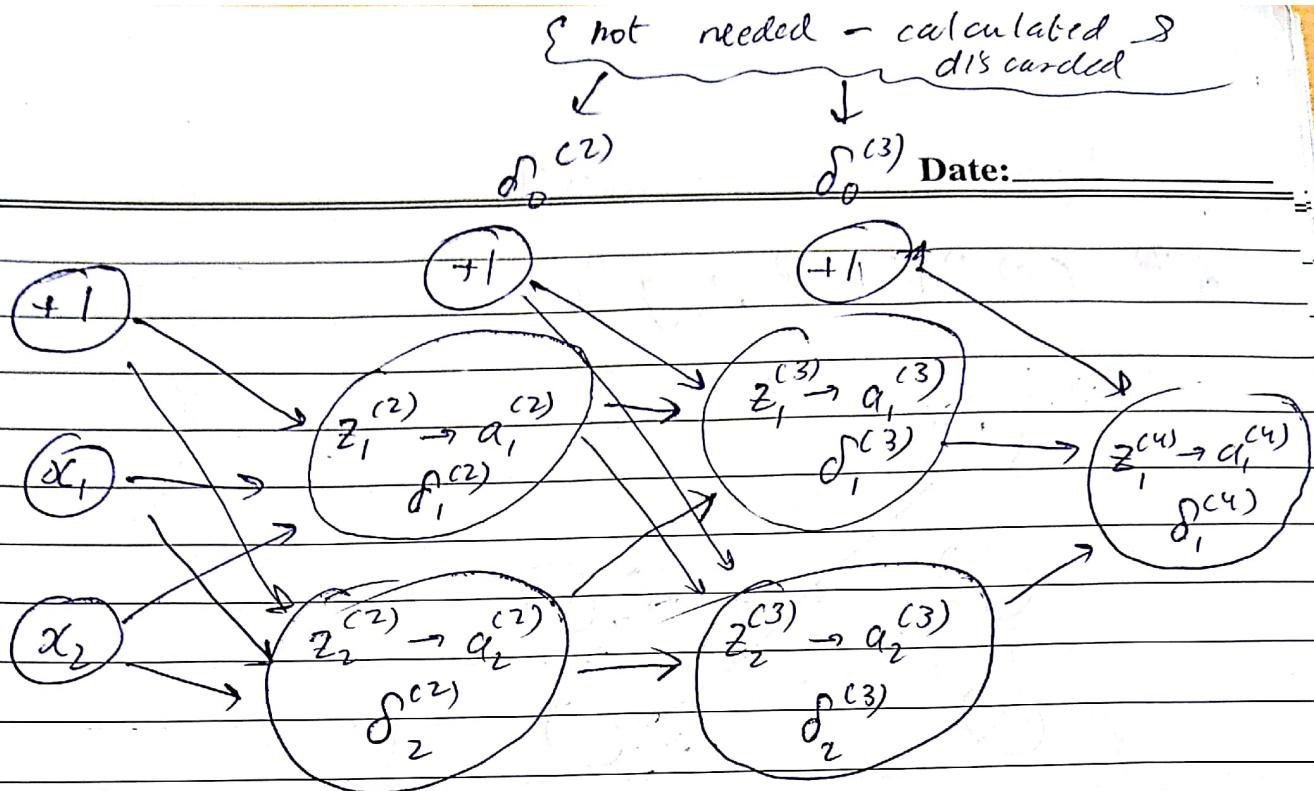
Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$ for ($j \geq 0$)

where:

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

$\frac{\partial \text{cost}(i)}{\partial z_j^{(l)}}$ → what is the minute change
in cost if a minute change
in ~~$z_j^{(l)}$~~ $z_j^{(l)}$ occurs. A minute
change in $z_j^{(l)}$ causes for the
 $h_\theta(x^{(i)})$ to change in value

$$\delta_i^{(u)} = y^{(i)} - a_i^{(u)}$$

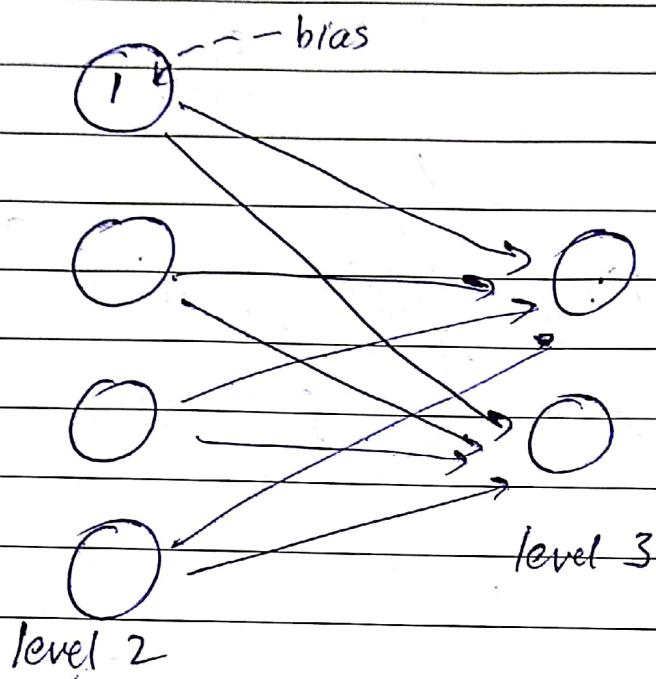


$$\delta_2^{(2)} = \theta_{12}^{(2)} \delta_1^{(3)} + \theta_{22}^{(2)} \delta_2^{(3)}$$

How $\delta^{(c)}$ are calculated :

Date: _____

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



$$\Theta^{(2)} = \begin{bmatrix} \Theta_{10} & \Theta_{11} & \Theta_{12} & \Theta_{13} \\ \Theta_{20} & \Theta_{21} & \Theta_{22} & \Theta_{23} \end{bmatrix}$$

$$(\Theta^{(2)})^T \delta^{(3)} = \begin{bmatrix} \Theta_{10} & \Theta_{20} \\ \Theta_{11} & \Theta_{21} \\ \Theta_{12} & \Theta_{22} \\ \Theta_{13} & \Theta_{23} \end{bmatrix} \begin{bmatrix} \delta_4 \\ \delta_3 \end{bmatrix} \quad (2 \times 1)$$

(4×2)

$$= \begin{bmatrix} \theta_{10} \delta_4 + \theta_{20} \delta_3 \\ \theta_{11} \delta_4 + \theta_{21} \delta_3 \\ \theta_{12} \delta_4 + \theta_{22} \delta_3 \\ \theta_{13} \delta_4 + \theta_{23} \delta_3 \end{bmatrix}$$

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

does not exist

$$= \begin{bmatrix} \theta_{10} \delta_4 + \theta_{20} \delta_3 \\ \theta_{11} \delta_4 + \theta_{21} \delta_3 \\ \theta_{12} \delta_4 + \theta_{22} \delta_3 \\ \theta_{13} \delta_4 + \theta_{23} \delta_3 \end{bmatrix} * g \begin{bmatrix} z_0 \\ z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} \delta_0^{(2)} \\ \delta_1^{(2)} \\ \delta_2^{(2)} \\ \delta_3^{(2)} \end{bmatrix}$$

Unrolling parameters

Changing matrices to vectors

function [jVal, gradient] = costFunction(theta)

optTheta = fminunc(@(costFunction,
initialTheta, options))

default behaviour is that
these are vectors

Date: _____

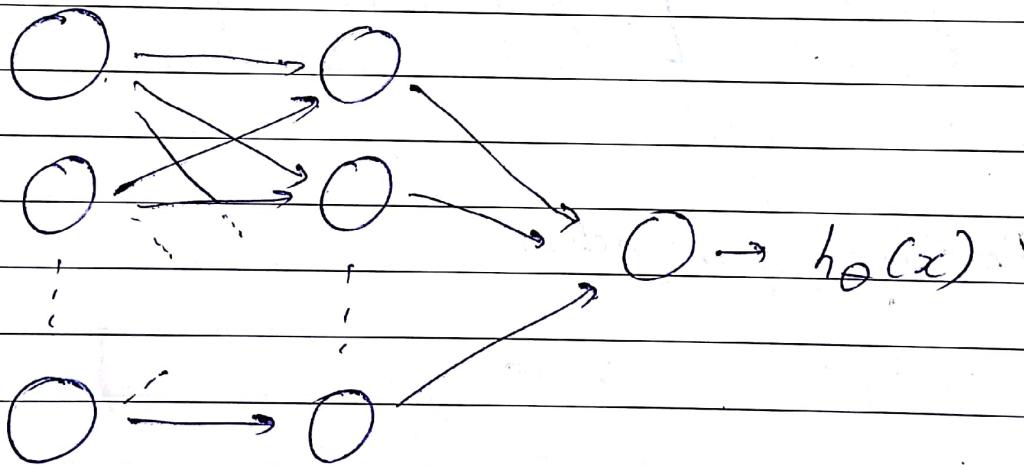
However with a 4 layer Neural Network this is the case:

- $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)} \rightarrow$ matrices
 - $D^{(1)}, D^{(2)}, D^{(3)} \rightarrow$ matrices

We must "unroll" these matrices into vectors so that they are in a format suitable for costFunction & fminunc

Example:

$$s_1 = 10, \quad s_2 = 10, \quad s_3 = 7$$



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

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

Seems like mistake is
made here: $\mathcal{O}^{(3)}$ will
not exist

Victory

In octave!

$\text{thetaVec} = [\Theta_1(:); \Theta_2(:); \Theta_3(:)];$

unroll a
matrix into
a vector

places all the unrolled
fragments together

$DVec = [D_1(:), D_2(:), D_3(:)];$

Going back from vector to matrix representations

$\Theta_1 = \text{reshape}(\text{thetaVec}(1:110), 10, 11)$

$\Theta_2 = \text{reshape}(\text{thetaVec}(111:220), 10, 11)$

$\Theta_3 = \text{reshape}(\text{thetaVec}(221:231), 1, 11)$

✓
reshape into
a 1×11 matrix

Final learning algorithm

- Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$

- Unroll to get initialTheta to pass to
 $fminunc(@costFunction, initialTheta, options)$

Date: _____

function [JVal, gradientVec] = costFunction(thetaVec)

- From thetaVec, get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ via reshape
- Use forward prop/back prop to compute $O^{(1)}, O^{(2)}, O^{(3)}$ & $J(\Theta)$
- Unroll $O^{(1)}, O^{(2)}, O^{(3)}$ to get gradientVec

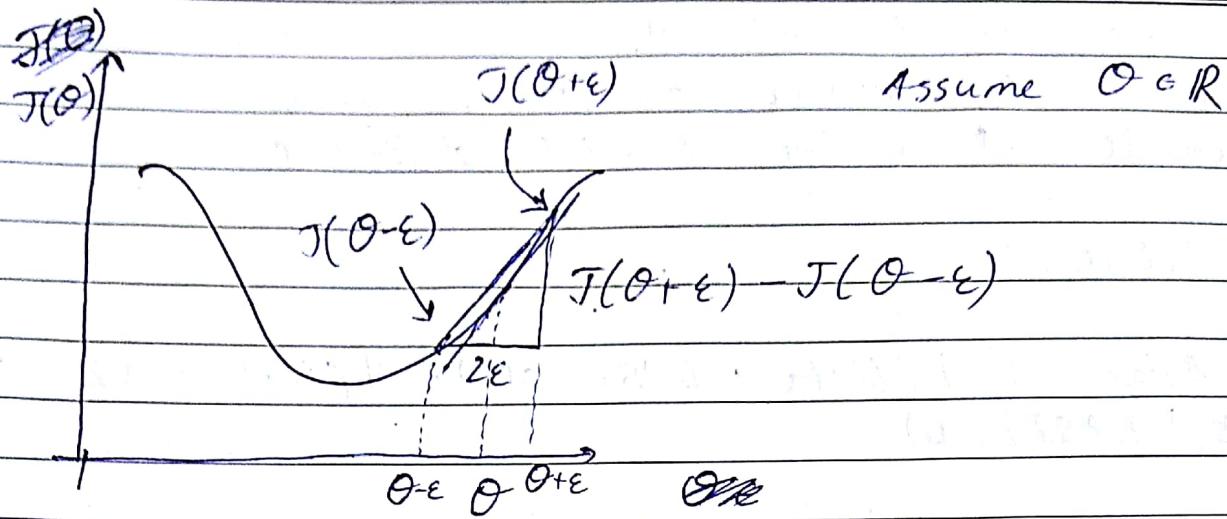
Gradient checking

Back propagation is a complex algorithm & hence subtle bugs can make their way into the ~~the~~ implementation. The cost might be decreasing on every iteration of an optimisation with the bug as well but you might end up with a ~~the~~ neural network with a higher level of error if compared to the bug free implementation.

- ↳ gradient checking eliminates this problem
- ↳ it is smart to implement gradient checking in other algorithms as well

Gradient checking is used to verify that the actual gradient of the cost function is being calculated

Date: _____



Assume $\theta \in \mathbb{R}$

We can find an approximation to $\frac{d}{d\theta} J(\theta)$ by connecting a line between $J(\theta + \epsilon)$ & $J(\theta - \epsilon)$ & finding out its gradient

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

ϵ is usually a small value $\approx 10^{-4}$

2-sided difference estimate:

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

1-sided difference estimate:

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

We use the 2-sided difference for our gradient approximation as it is more accurate than the 1-sided difference

In Octave:

$$\text{grad Aprox} = \frac{J(\theta + \text{EPSILON}) - J(\theta - \text{EPSILON})}{2 * \text{EPSILON}}$$

Previous example considered a single value θ or $\theta \in \mathbb{R}$. For $\theta \in \mathbb{R}^n$ (θ can be unrolled version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$)

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

$$\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \varepsilon_1, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \varepsilon_1, \theta_2, \theta_3, \dots, \theta_n)}{2\varepsilon}$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \varepsilon_2, \dots, \theta_n) - J(\theta_1, \theta_2 - \varepsilon_2, \dots, \theta_n)}{2\varepsilon}$$

etc

Date: _____

In Octave :

for $i = 1 : n$:

① thetaPlus = theta

② thetaPlus(i) = thetaPlus(i) + EPSILON

③ thetaMinus = theta;

④ thetaMinus(i) = thetaMinus(i) - EPSILON

⑤ gradApprox(i) = $(J(\theta_{\text{plus}}) - J(\theta_{\text{minus}})) / (2 * \text{EPSILON})$

end;

① - ② :

Copy the theta vector & add epsilon to the element of the current iteration:

ThetaPlus is refreshed with theta on each iteration so previous epsilon is removed

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i + \epsilon \\ \vdots \\ \theta_n \end{bmatrix}$$

③ - ④ : Copy the theta Vector & subtract epsilon from the element of the current iteration:

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i - \epsilon \\ \vdots \\ \theta_n \end{bmatrix}$$

Date: _____

Once all gradApprox values are computed
check that gradApprox \approx DVec

If the difference between gradApprox & DVec is in decimal places then
this increases the confidence that
backprop computed correctly.

Implementation note:

- Implement backprop to compute DVec (unrolled $O^{(1)}, O^{(2)}, O^{(3)}$)
- Implement numerical gradient check to compute gradApprox
- Make sure they give similar values
- Turn off gradient checking
- Use backprop for learning

→ This is done because computing the gradient with gradient checking is computationally expensive as compared to using backprop

Important:

Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient on every iteration of gradient descent (or in inner loop of costFunction()) then your code will be very slow

Random Initialization

Initial value of θ

For gradient descent & advanced optimisation method, need ~~so~~ initial value of θ

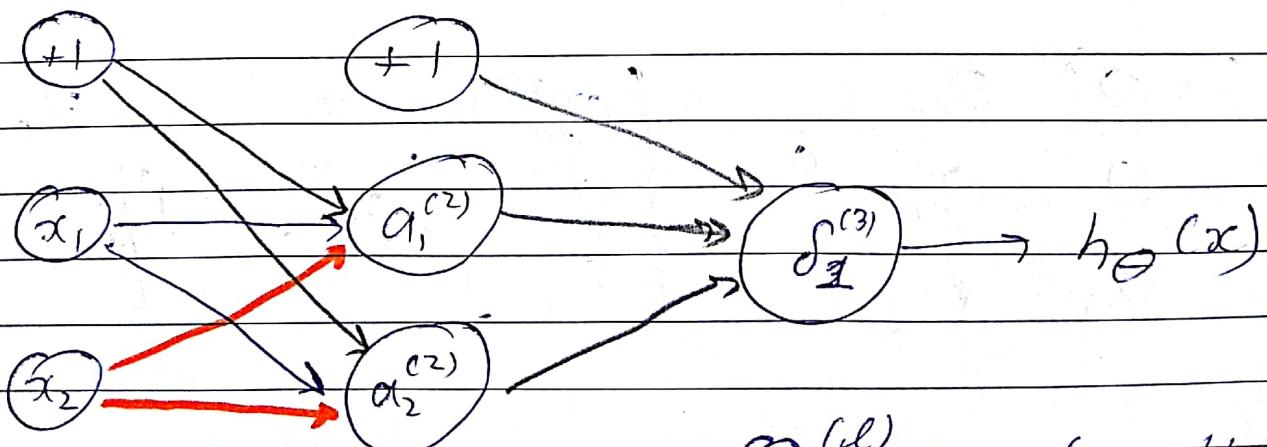
`optTheta = fminfminunc(@costFunction, initialTheta, options)`.

Consider gradient descent

Set `initialTheta = zeros(n, 1)`?

↳ worked okay in logistic regression but does not work with neural networks

Suppose all weights are initialised to 0:



$$\theta_{ij}^{(l)} = 0 \text{ for all } i, j, l$$

Date:

All the lines in pencil ~~8 lines in pen~~ are the same, all the lines in pen are the same & all in marker are the same hence: $\alpha_1^{(2)} = \alpha_2^{(2)}$

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

Thus $\delta_1^{(2)} = \delta_2^{(2)}$

$$\frac{\partial}{\partial \Theta_{ij}^{(2)}} J(\Theta) = \alpha_j^{(2)} \delta_i^{(2+1)}$$

$$\Theta^{(1)} = \begin{bmatrix} \Theta_{10} & \Theta_{11} & \Theta_{12} \\ \Theta_{20} & \Theta_{21} & \Theta_{22} \end{bmatrix}$$

$$\Theta^{(2)} = \begin{bmatrix} \Theta_{10} & \Theta_{11} & \Theta_{12} \end{bmatrix}$$

$$\delta^{(2)} = \begin{bmatrix} \Theta_{10}^{(2)} \\ \Theta_{11}^{(2)} \\ \Theta_{12}^{(2)} \end{bmatrix} \left[\begin{bmatrix} \delta_1^{(3)} \end{bmatrix} \right] * g' \begin{bmatrix} z_0^{(2)} \\ z_1^{(2)} \\ z_2^{(2)} \end{bmatrix}$$

Date: _____

$$\delta^{(2)} = \left[\begin{array}{l} \theta_{10}^{(2)} \delta_1^{(3)} \times g'(z_0^{(2)}) \\ \theta_{11}^{(2)} \delta_1^{(3)} \times g'(z_1^{(2)}) \\ \theta_{12}^{(2)} \delta_1^{(3)} \times g'(z_2^{(2)}) \end{array} \right].$$

all $z^{(2)}$ elements are equal as ~~weights~~
~~each node in layer 2 has~~
 the same exact weights, so all
 elements are equal in $\theta^{(2)}$ as well
 hence all elements in $\delta^{(2)}$ are equal:

$$\boxed{\delta_1^{(2)} = \delta_2^{(2)}}$$

Thus $\frac{\partial}{\partial \theta_{10}^{(1)}} J(\theta) = \frac{\partial}{\partial \theta_{20}^{(1)}} J(\theta)$

because $\frac{\partial}{\partial \theta_{10}^{(1)}} J(\theta) = a_0^{(1)} \delta_1^{(2)}$

$$= \frac{\partial}{\partial \theta_{20}^{(1)}} J(\theta) = a_0^{(1)} \delta_2^{(2)}$$

since $\delta_1^{(2)} = \delta_2^{(2)}$

And when θ_{10} & θ_{20} are updated
then they will be equal since initial
value, & its derivative are equal ~~for both~~ in
 θ_{10} & θ_{20}

→ Thus all the lines of the same colour
remain the same value, but will be
non-zero

After each update, parameters corresponding
to inputs going into each of 2
hidden units are identical

i.e. $a_1^{(2)} = a_2^{(2)}$ since the weights are
the same & applied to the same
input values after 1 iteration

→ As gradient descent is continued the
coloured lines will always remain
the same thus each node in ~~the~~
an ~~top~~ hidden layer computes the exact
same function

→ The layer 3 node sees the same
feature from all its inputs

Random initialization is used to get
around this problem

Random initialization:

has nothing to do with ϵ in gradient checking
Date: _____

Symmetry breaking

Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$ (i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)

E.g. \rightarrow random 10x11 matrix (between 0-81)

$\Theta_{1:1} = \text{rand}(10, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON}$

\rightarrow in $[-\epsilon, \epsilon]$

Symmetry is set not only when all weights are 0 but also when any all weights are equal; we need to break the symmetry by setting random variables

Putting it together

Training a neural network

Pick a network architecture (connectivity pattern between neurons) - # of ~~hidden~~ layers & # of nodes in these layers

No. of input units: dimensions of features
 $\mathbf{x}^{(i)}$

No. of output units: Number of classes

$y \in \{1, 2, \dots, 10\}$ where if $y = 1$

then $y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \} \text{ 10 elements total}$

Reasonable default: 1 hidden layer, or if > 1 hidden layer, have same no. of hidden units in every layer (~~more~~ usually the more the better)

→ however having more hidden units makes training more expensive

Date: _____

Having more nodes in hidden layers than input layers is usually better

Training a neural network

- 1) Randomly initialize weights
- 2) Implement forward propagation to get $h_\theta(x^{(i)})$ for any $x^{(i)}$
- 3) Implement code to compute cost function $J(\theta)$
- 4) Implement backprop to compute partial derivatives $\frac{\partial}{\partial \theta^{(l)}} J(\theta)$

for $i = 1:m$ } iterate over data rows
in $(x^{(i)}, y^{(i)})$:

Perform forward propagation &
backpropagation using example
 $(x^{(i)}, y^{(i)})$

(Get activations $a^{(l)}$ & delta terms
 $\delta^{(l)}$ for $l=2, \dots, L$)

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

Compute $\frac{\partial}{\partial \theta^{(l)}} J(\theta)$ with regularization & etc

Date: _____

5) Use gradient checking to compare
 $\frac{\partial}{\partial \theta_j} J(\theta)$ computed using backprop

vs using numerical estimate of
gradient of $J(\theta)$

→ Then disable gradient checking
code

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

$J(\theta)$ is non-convex for neural
networks & is susceptible to
local minima