

Machine Learning

Week- 3 (Classification problem)

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

m training examples.

$$\text{each } x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \xrightarrow{\text{+ } i=1 \text{ to } m} x_0^{(i)} = 1, \quad \alpha \in \mathbb{R}^{n+1}$$

$$\alpha = \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} \quad (n+1) \times 1$$

$$y \in \{0, 1\}$$

$$h_{\alpha}(x) = \frac{1}{1 + e^{-\alpha^T x}}$$

$$[\alpha_0 \ \alpha_1 \ \dots \ \alpha_n] * \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$(\alpha_0 x_0 + \alpha_1 x_1 + \dots + \alpha_n x_n)$$

→ how do I choose parameters α ?

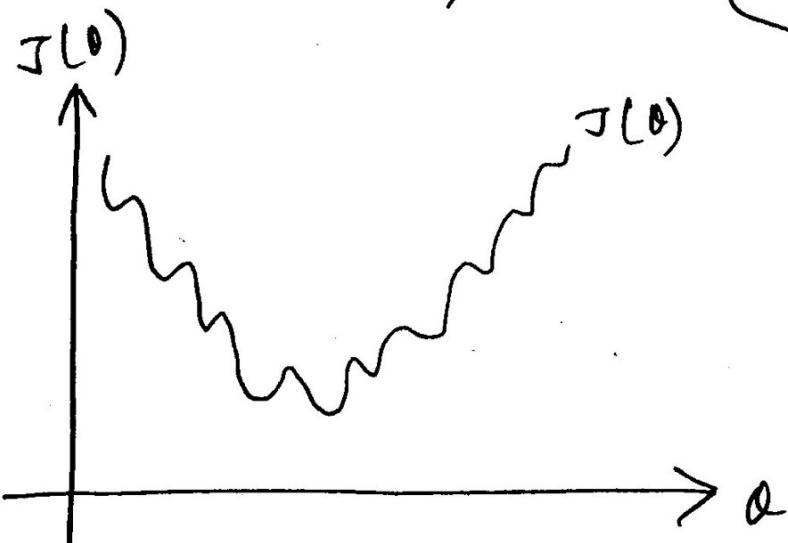
Cost function in linear regression:

$$J(\alpha) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} \left(h_{\alpha}(x^{(i)}) - y^{(i)} \right)^2$$

$$= \frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\alpha}(x^{(i)}), y)$$

If we use $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(x^{(i)}), y^{(i)})$

$$\text{Cost}(h_\theta(x), y) = \frac{1}{2} (h_\theta(x) - y)^2$$



$$\frac{1}{1 + e^{-\alpha^T x}}$$

non-convex.

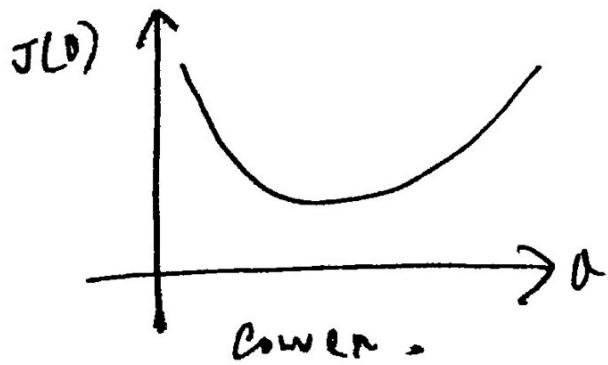
always

gradient descent will not converge to global

minimum in

a convex fn.

we want something like

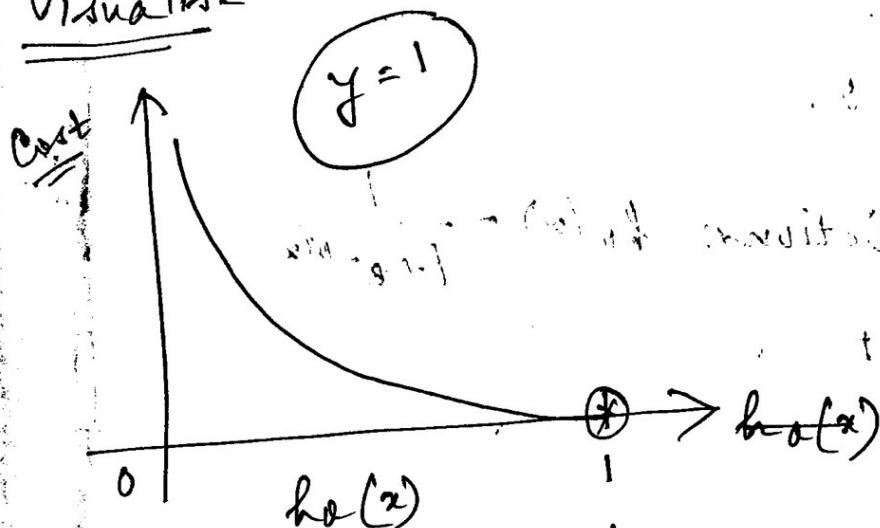


convex.

New Cost function

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0. \end{cases}$$

Visualise

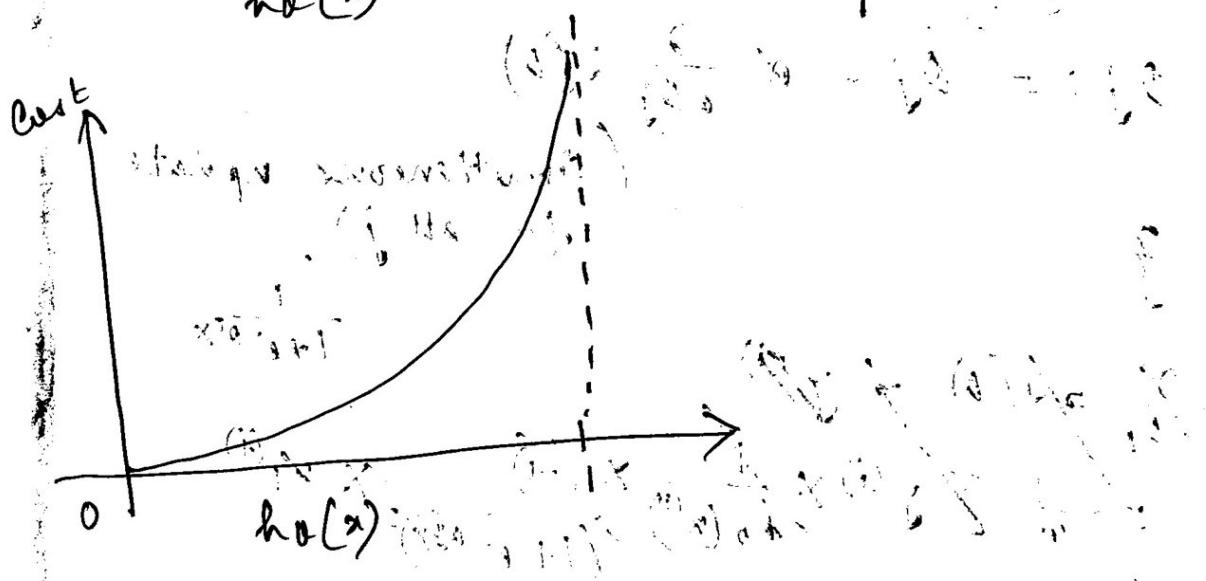


As $h_\theta(x) = 1$,

if $y = 1$,
Cost = 0.

If $h_\theta(x) = 0$,
 $\text{Cost} = -\log(1/h_\theta(x)) = 0$.

but $y = 1$,
Cost $\rightarrow \infty$.



Simplified cost function

$$\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

→ Maximum likelihood estimation

→ Convex, Convex,

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

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

\downarrow
 $\min J(\theta)$ w.r.t. θ .

To make predictions: $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$

gradient descent

Repeat {

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

(simultaneous update
for all j).
}

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \left[-\frac{1}{m} \sum_{i=1}^m y^{(i)} \times \frac{1}{h_\theta(x^{(i)})} \times \frac{(1)}{(1+e^{-\theta^T x})^2} \times x_j^{(i)} \right] \\ &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \times \frac{1}{h_\theta(x^{(i)})} \times \frac{(1)}{(1+e^{-\theta^T x})^2} \times x_j^{(i)} \end{aligned}$$

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (e^{(i)}) x_j^{(i)}$$

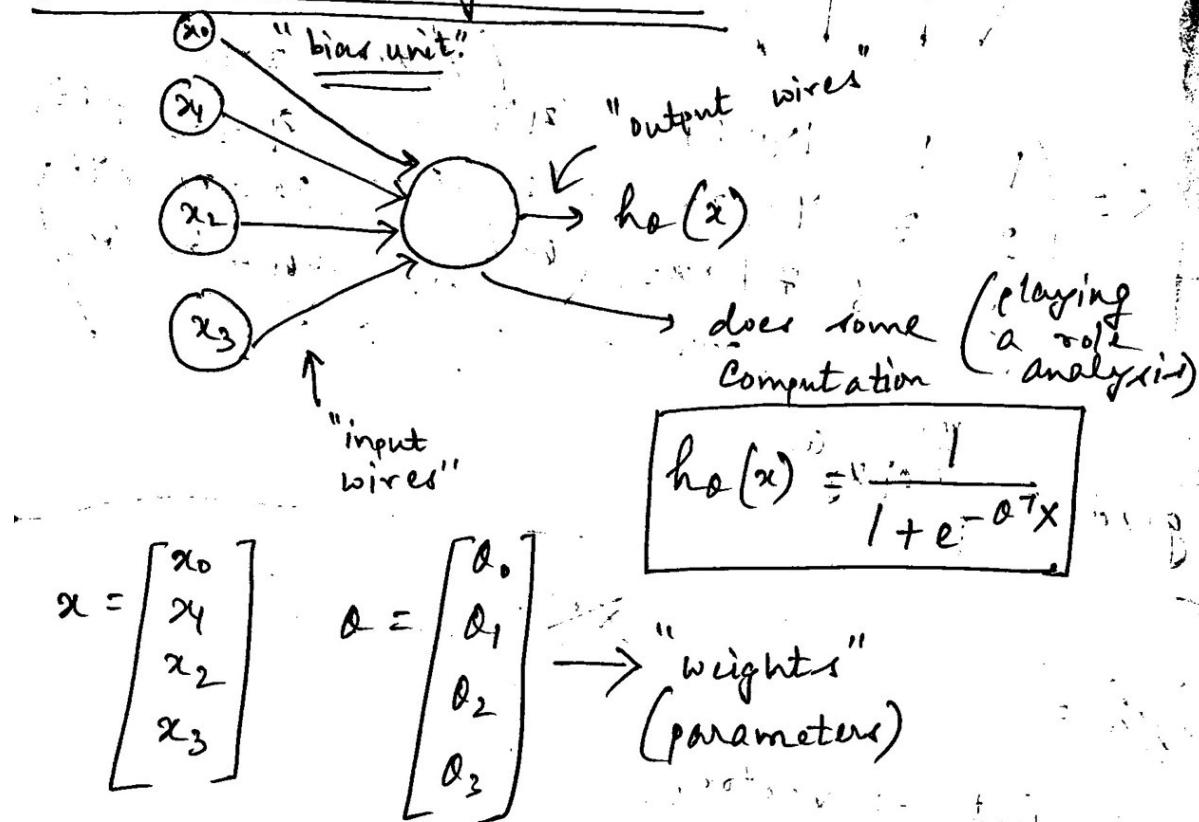
}.

Week-4 (Neural Networks)

Model representation (1)

Modelling how neurons in the brain work

Neuron model: logistic unit



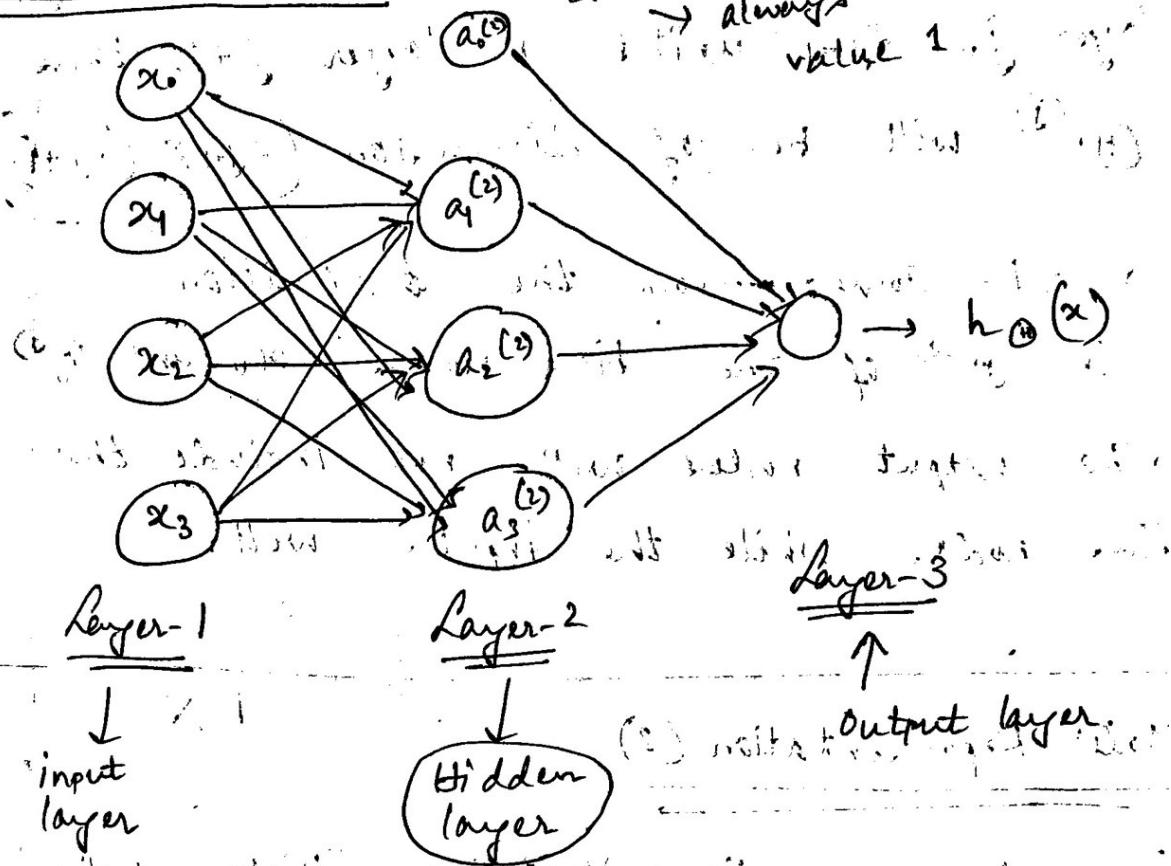
$x_0 \rightarrow$ bias unit \rightarrow always $= 1$.

Sigmoid (logistic) activation function

This above diagram represents a single neuron.

Neural network \rightarrow group of different neurons strong together.

Neural Network



$a_i^{(j)}$ = "activation" of unit i in layer j

$\Theta^{(j)}$ = matrix of weights controlling function mapping from layer j to layer $(j+1)$

$$a_1^{(2)} = g \left(\Theta_{10}^{(2)} x_0 + \Theta_{11}^{(2)} x_1 + (\Theta_{12}^{(2)} x_2 + \Theta_{13}^{(2)} x_3) \right)$$

$$a_2^{(2)} = g \left(\Theta_{20}^{(2)} x_0 + (\Theta_{21}^{(2)} x_1 + (\Theta_{22}^{(2)} x_2 + \Theta_{23}^{(2)} x_3) \right)$$

$$a_3^{(2)} = g \left(\Theta_{30}^{(2)} x_0 + (\Theta_{31}^{(2)} x_1 + (\Theta_{32}^{(2)} x_2 + \Theta_{33}^{(2)} x_3) \right)$$

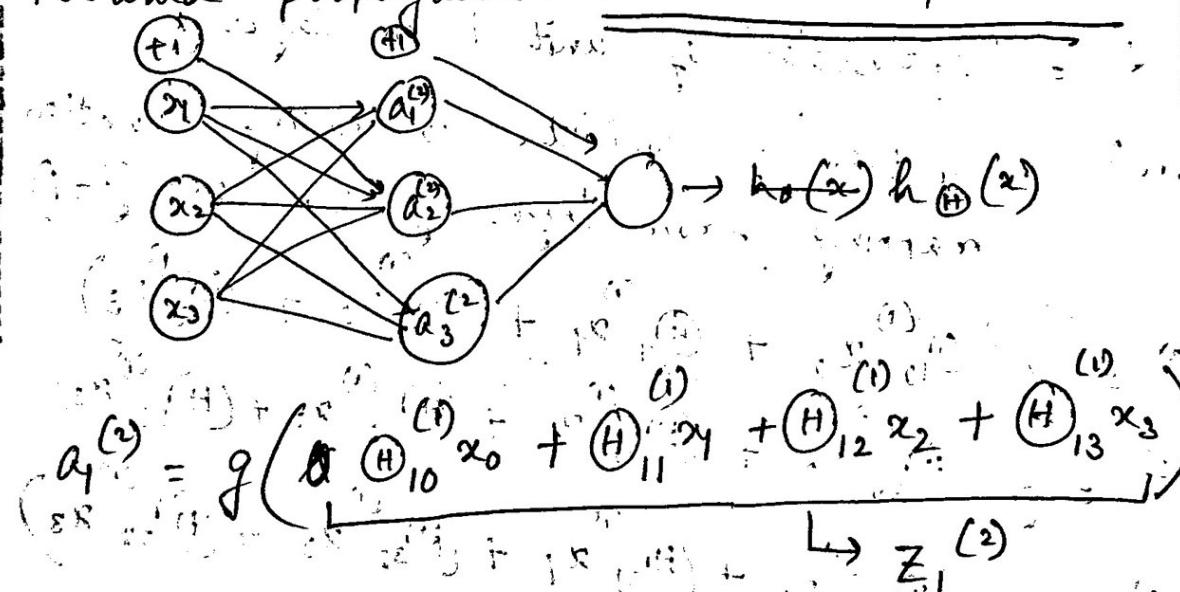
$$h_{\Theta}(x) = a_4^{(3)} = g \left(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)} \right)$$

$\Theta^{(2)} \in \mathbb{R}^{3 \times 4}$
dimension matrix.

- If network has s_j units in layer j , s_{j+1} units in layer $j+1$, then
- $\Theta^{(j)}$ will be of dimension $s_{j+1} \times s_j$
- +1 comes from the addition in $\Theta^{(1)}$ of the "bias nodes", x_0 and Θ_0 .
- * The output nodes will not include the bias nodes while the inputs will.

Model Representation (2)

Forward propagation: Vectorized implementation



$$a_1^{(2)} = g(z_1^{(2)})$$

$$a_2^{(2)} = g(z_2^{(2)})$$

$$a_3^{(2)} = g(z_3^{(2)})$$

$$\begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{bmatrix} = z^{(2)}$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

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

$$z^{(2)} = \textcircled{H}^{(1)} \otimes a^{(1)}$$

$$a^{(2)} = g(z^{(2)}) \rightarrow \text{elementwise sigmoid fn.}$$

$$\text{Add } a_0^{(2)} + 1 \rightarrow h_0 a_0^{(2)} \in \mathbb{R}^4$$

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

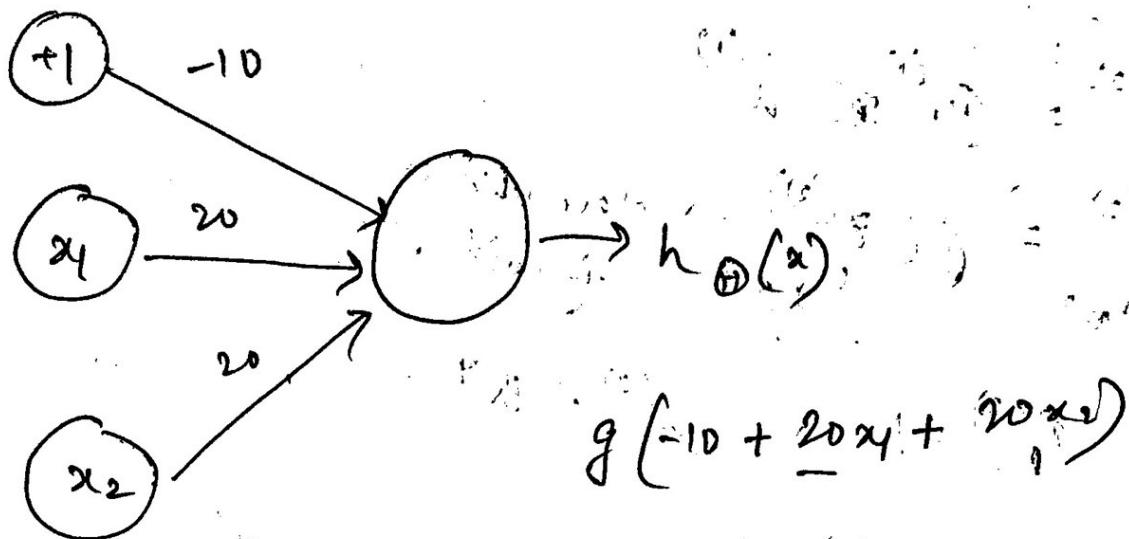
$$\underline{h_{\textcircled{H}}(x)} = a^{(3)} = g(z^{(3)})$$

This process is called forward propagation

$$a^{(3)} = g\left(\textcircled{H}_{10}^{(2)} a_0^{(2)} + \textcircled{H}_{11}^{(2)} a_1^{(2)} + \textcircled{H}_{12}^{(2)} a_2^{(2)} + \textcircled{H}_{13}^{(2)} a_3^{(2)}\right)$$

Applications

Examples and intuitions I



x_1	x_2	$h_{\theta}(x)$
0	0	≈ 0
0	1	≈ 1
1	0	≈ 1
1	1	≈ 1

A simple example of applying neural networks is by predicting x_1 AND x_2 which is the logical 'and' operator and is only true if both x_1 and x_2 are 1.

Applications & intuitions-II

(Not x_1 and Not x_2) NOR

say = 1 if and only if $x_1 = x_2 = 0$

~~$g(x_1 \wedge \neg x_1 \wedge \neg x_2)$~~

→ For negation just put a large weight in front of the number you want to negate.

(+) matrices for AND, NOR, and OR are:

$$\text{AND: } \Theta H^{(1)} = \begin{bmatrix} -30 & 20 & 20 \end{bmatrix} \rightarrow \text{AND}$$

$$H^{(1)} = \begin{bmatrix} 10 & -20 & -20 \end{bmatrix} \rightarrow \text{NOR}$$

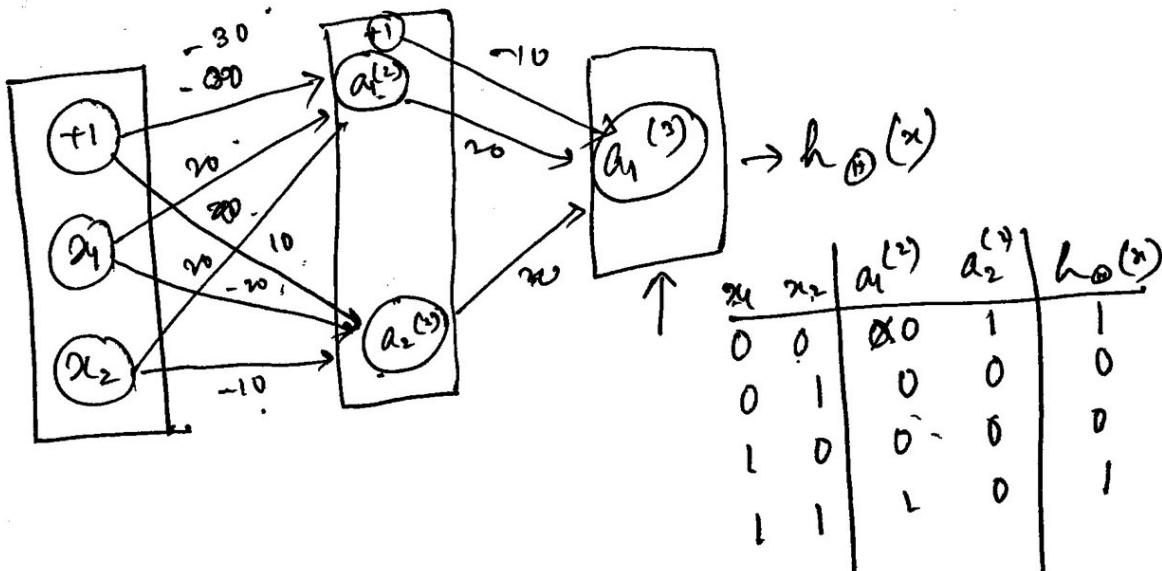
$$H^{(1)} = \begin{bmatrix} 10 & 20 & 20 \end{bmatrix} \rightarrow \text{OR}$$

We combine these to get the XNOR operation

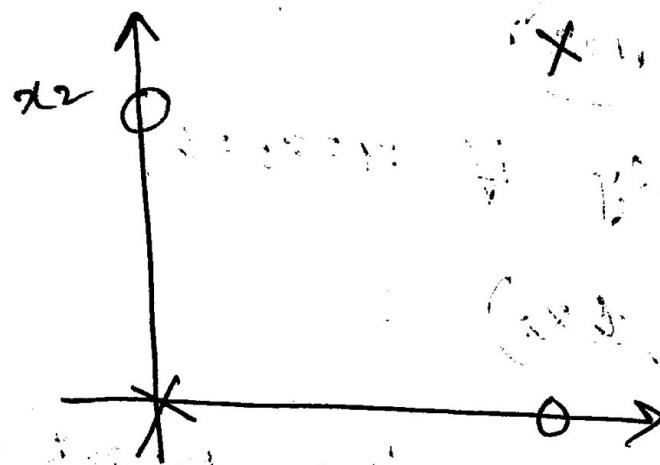
$-x_1 \times \text{NOR } x_2 \rightarrow 1 \text{ if } x_1 \& x_2 = 0 \text{ and } 2$

or

$$x_1 \& x_2 = 0$$



$x_1 \text{ XNOR } x_2$



Result if $x \rightarrow$ positive examples
 $(1) \text{ (True)}$
 $x \rightarrow$ negative examples
 $(0) \text{ (False)}$

Multiclass classification

~~Output~~ $h_{\Theta}(x)$ is a vector of size Θ
the number of classes.

Size \rightarrow Associated with $\Theta^{(1)}, \Theta^{(2)} \in \mathbb{R}^{10 \times 6}$

5 hidden layers.

Vectorization

$$\underline{n=2}$$

gradient descent

$$\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_1 := \theta_1 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ \theta_2 := \theta_2 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)} \end{array} \right.$$

$$\theta := \theta - \alpha \underset{\text{vector subtraction}}{\circledcirc} s$$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$

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

$$s = \begin{bmatrix} s_0 \\ s_1 \\ s_2 \end{bmatrix} \quad s_0 = \left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \right) \downarrow \mathbb{R}$$

$$u = 2v + 5w$$

$$u = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix}, v = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}, w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$\begin{bmatrix} 2v_1 + 5w_1 \\ 2v_2 + 5w_2 \\ 2v_3 + 5w_3 \end{bmatrix} \quad e^{(1)} = \begin{bmatrix} x_0^{(1)} \\ x_1^{(1)} \\ \vdots \\ x_n^{(1)} \end{bmatrix} \quad + \quad e^{(2)} = \begin{bmatrix} x_0^{(2)} \\ x_1^{(2)} \\ \vdots \\ x_n^{(2)} \end{bmatrix}$$

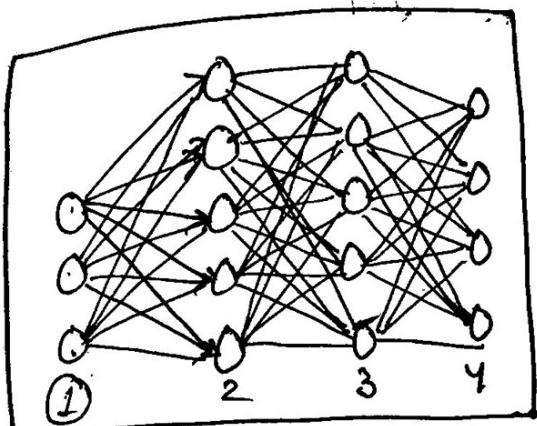
Week-5 (Neural Networks: Learning)

Cost function

Application of neural network to classification problem.

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

$\rightarrow L = \text{total no. of layers in the network}$



$\rightarrow s_L = \text{no. of units (except bias unit)}$

$$\text{in layer } l: (x_l, s_1 = 3, s_2 = 5, s_3 = 5, s_4 = 4)$$

Binary ($K=1$)

$$y = 0 \text{ or } 1$$

1 output unit

$$h_{(H)}(x) \in \mathbb{R}^1$$

$$s_L = 1 \quad K \geq 1$$

Multi-class (K classes)

$$y \in \mathbb{R}^K \quad l = (1, 0, 0, \dots, 0)$$

K output units

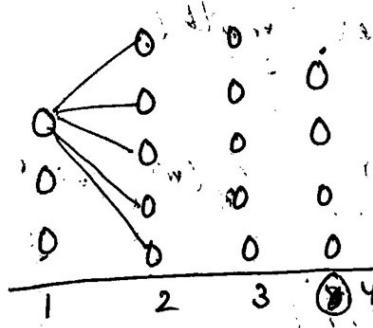
$$h_{(H)}(x) \in \mathbb{R}^K$$

$$s_L = K \quad (K \geq 3)$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(h_{\theta}(x^{(i)}) + (1 - y_i) \ln(1 - h_{\theta}(x^{(i)}))$$

Backpropagation Algorithm

$\{(x, y)\} \rightarrow$ only one training set. example



Forward propagation

$$x = a^{(1)} ; \quad (H) \in \mathbb{R}^{5 \times 4}$$

$$z^{(2)} \in \mathbb{R}^{5 \times 1}$$

~~$$z^{(2)} = (H)^{(1)} a^{(1)}$$~~

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

$$a^{(2)} \in \mathbb{R}^{6 \times 1}$$

~~$$z^{(3)} = (H)^{(2)} a^{(2)}$$~~

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

~~$$z^{(4)} = (H)^{(3)} a^{(3)}$$~~

$$a^{(4)} = g(z^{(4)})$$

Backward propagation

$\delta \rightarrow$ error

$$\delta^{(4)} = y_{\text{true}}(h_{(4)}(x)) - y$$

~~$$\text{Vectorize } \delta^{(4)} = (h_{(4)}(x) - y)$$~~

$$\delta^{(3)} = (H)^{(3)} \top \delta^{(4)} * g'(z^{(3)}) \rightarrow a^{(3)} * (1 - a^{(3)})$$

$$\delta^{(2)} = (H)^{(2)} \top \delta^{(3)} * g'(z^{(2)})$$

No $\delta^{(1)}$

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y^{(i)}_k \log h_\theta(x^{(i)})_k + (1-y^{(i)}_k) \log (1-h_\theta(x^{(i)})_k) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^L \sum_{j=1}^{s_l} \sum_{j'=1}^{s_{l+1}} (\Theta_{j,j'}^{(l)})^2$$

min $J(\Theta)$

Θ

Need to compute $\circled{1} J(\Theta)$

$$\circled{2} \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$$

\rightarrow main aim: gradient computation,
for cost minimisation.

Simple case One training example - $(x^{(0)}, y)$

Forward propagation

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

$$z^{(2)} = \Theta^{(1)} a^{(1)} \quad \Theta^{(1)} \in \mathbb{R}^{4 \times 3}$$

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

$$z^{(3)} = \Theta^{(2)} a^{(2)} \quad \Theta^{(2)} \in \mathbb{R}^{1 \times 4}$$

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

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

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

→ Gradient Computation: "backprop algo b2"

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

For each output unit, we find error:

$$\delta_1^{(4)} = \delta_1^{(4)} = a_1^{(4)} - y_1 \rightarrow (h_{\Theta}(x))_{j1}$$

Vectorize: $\underline{\delta^{(4)}} = \underline{a^{(4)} - y} \quad \underline{\delta^{(3)}} = (\Theta)^{(3)}^T \delta^{(4)} * g'(z^{(3)})$

→ use this formula: $\delta^{(3)} = (\Theta)^{(2)}^T \delta^{(2)} * g'(z^{(2)})$

$$\delta^{(2)} = (\Theta)^{(2)}^T \delta^{(1)} * g'(z^{(1)}) \rightarrow \underline{a^{(2)} * (1-a^{(2)})}$$

(No. $\delta^{(1)}$) → input has no error.

$$\frac{\partial J(\theta)}{\partial \Theta_{ij}^{(l)}} = a_j^{(l)} \delta_i^{(l+1)} \quad (\text{ignoring } \lambda)$$

this can be

proved mathematically

Remember: ① how to compute

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

→ Steps in back propagation algorithm :

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

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

→ used to compute $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$

for $t = 1$ to m accumulator

— set $a^{(1)} = x^{(t)}$

— Perform forward propagation to

— Compute $a^{(l)}$ for $l = 2, 3, \dots, L$

— Using $y^{(t)}$, compute $s^{(l)} = \underline{a^{(l)} - y^{(t)}}$

— Compute $s^{(L-1)}, s^{(L-2)}, \dots, s^{(2)}$ ~~$s^{(1)}$~~

— $\Delta_{ij}^{(l)} := \underline{\Delta_{ij}^{(l)} + a_j^{(l)} s_i^{(l+1)}}$ using
backprop

— $A^{(l)} := A^{(l)} + s^{(l+1)} (a^{(l)})^T$

→ assuming $A^{(l)}$ is
a matrix subscripted by
 (i, j)

→ accumulate the
partial derivative terms
from previous slide.

Note : $j \rightarrow$ denotes node number.

Go outside for loop.

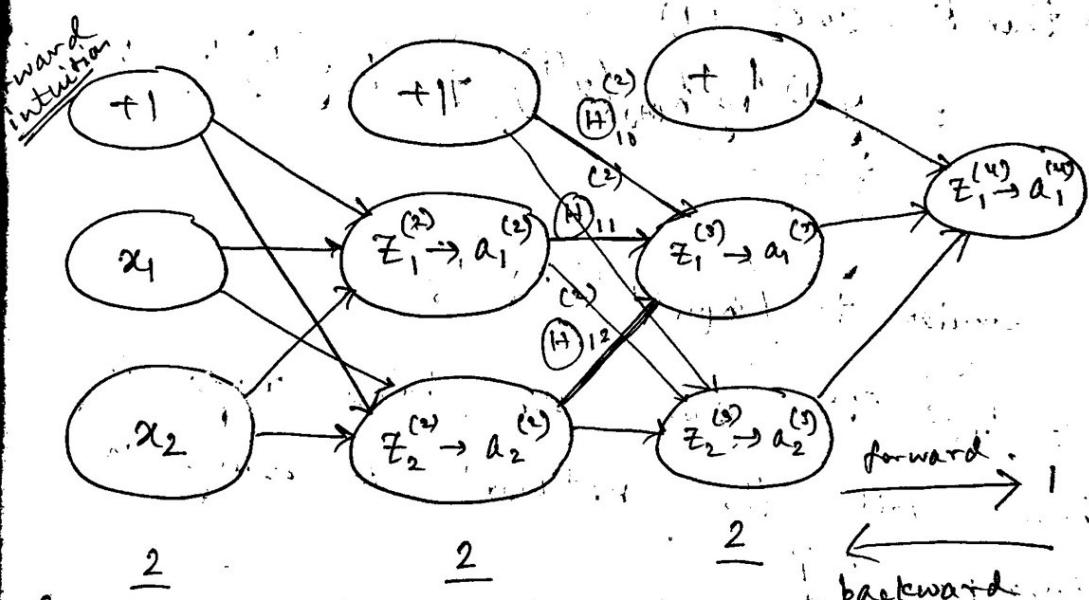
$$D_{ij}^{(l)} := \frac{1}{m} \left(\Delta_{ij}^{(l)} + \lambda (\Theta)_{ij}^{(l)} \right) \text{ if } j \neq 0$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$$

$$\frac{\partial}{\partial (\Theta)_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)},$$

can be proved mathematical

Backpropagation intuition:



(excluding bias unit)

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

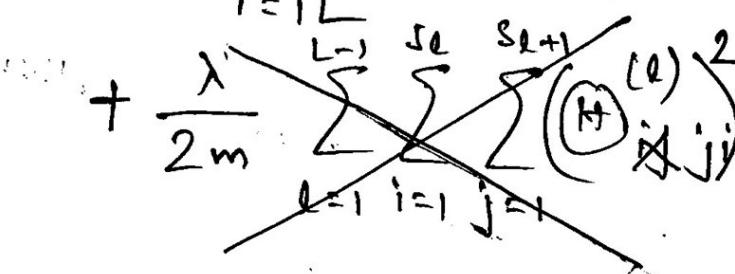
For every node calculation

Every node is a weighted sum of nodes in previous layer.

→ What is backpropagation doing?

for $K = 1$

Cost function is

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


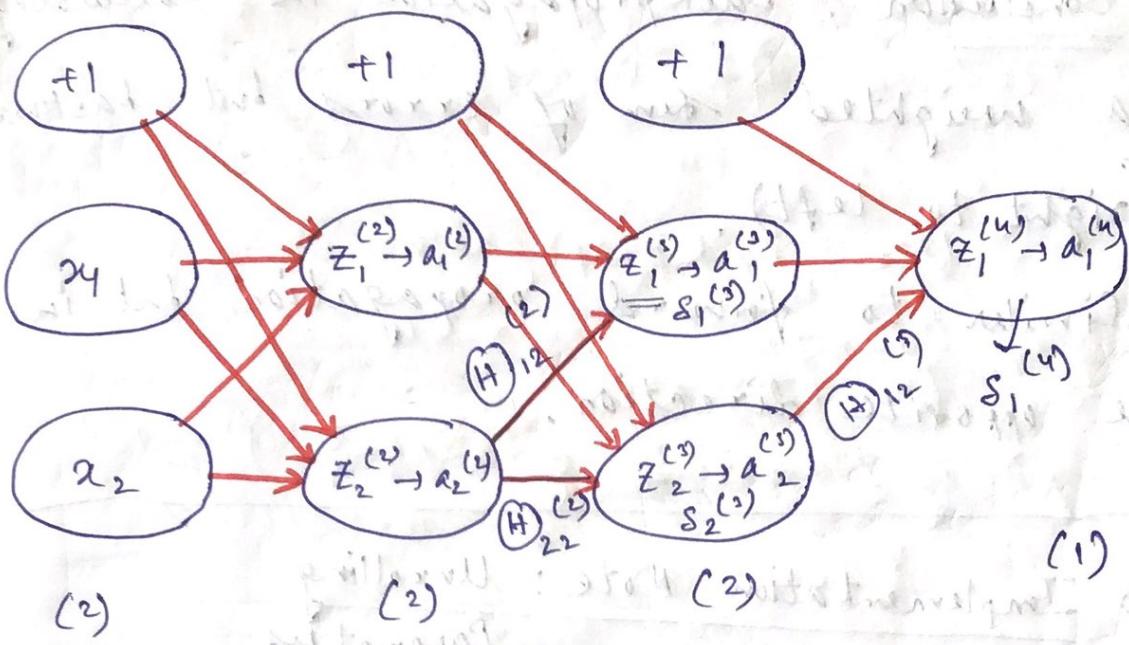
For one training example $(x^{(i)}, y^{(i)})$

→ think Cost (i)

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

↓
Similar to $(h_{\theta}(x^{(i)}) - y^{(i)})^2$

→ Cost (i) gives an idea regarding how well the network is doing on example i ?



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

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

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

$$\delta_1^{(3)} = \frac{\partial}{\partial z_1^{(3)}} \text{cost}(t)$$

\rightarrow basically $\delta_j^{(l)}$ measures the impact on of
 a change in $z_j^{(l)}$ on $\text{cost}(t)$] $t=1, \dots, m$

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

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

$$\rightarrow \underline{\underline{\delta_2^{(2)} = (\text{H})_{12}^{(2)} \times \delta_1^{(3)} + (\text{H})_{22}^{(2)} \times \delta_2^{(3)}}}$$

→ Conclusion : Backpropagation is calculating
a weighted sum of errors but backward
(right to left).

Similar to forward propagation but in
the opposite direction.

→ Implementation Note : Unrolling
Parameters.

Parameter $\Theta_{11}, \Theta_{12}, \dots$ into
a big vector so that it can be
used into advanced optimisation
algorithms.

From vectors \rightarrow reshape to get
the original matrices.

Neural net networks learning : Gradient Checking

$$\rightarrow J(\theta) = \theta^3$$

$$\begin{aligned} & \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} = \frac{(1.01)^3 - (-0.99)^3}{2 \times 0.01} \\ &= \frac{0.060002}{0.02} \approx 3.0001 \end{aligned}$$

Gradient interpretation: $\frac{d}{dx} f(x) \approx \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$
 ϵ is very small.

For $\theta \in \mathbb{R}^n$:

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

→ unrolled version of $(H^{(1)}, H^{(2)}, \theta^{(3)})$

Numerically compute of derivative

for $i = 1 : n$

$\text{thetaPlus} = \text{theta}$

$\text{thetaPlus}(i) = \text{theta}(i) + \text{Epsilon}$

$\text{thetaMinus} = \text{theta}$

$\text{thetaMinus}(i) = \text{theta}(i) - \text{Epsilon}$

$\text{gradApprox}(i) = \frac{J(\theta^+) - J(\theta^-)}{2 * \epsilon}$

End.

Check if
 $\text{gradApprox} \approx \text{above}$

from
 backprop

gradient
 checking.

Note: After gradient checking,
turn off it off. (disable
gradient checking code)

- Use backpropagation code for
Computing gradients
- ~~④~~ Numerical gradient computation
is very slow & inefficient.

→ Random Initialization # symmetry
breaking

~~x and (m, n) → initialise matrix~~
with all values $\in [0, 1]$

$$\in [-\varepsilon, \varepsilon]$$

$$\begin{aligned} &\cancel{x \in [0, 1]} \\ &\cancel{(2\varepsilon) \approx 0} \quad \cancel{r \approx 0} - r \\ &\cancel{(2\varepsilon) + r \approx 0} \quad \cancel{2\varepsilon r < 2\varepsilon} \quad \cancel{2(0.0001) \times 1} \\ &\cancel{r = 0.5} \quad \cancel{r \neq 0} \\ &r = 0.1 \end{aligned}$$

$$\begin{aligned} &r = 0.5 \quad \cancel{r = 0} \quad \cancel{0, 2\varepsilon} = \varepsilon \\ &\varepsilon = 0.001 \quad \cancel{[0, 2\varepsilon]} = \varepsilon \\ &2\varepsilon = 0.002 \quad \left(\begin{matrix} 0.001 & 0.001 \\ -0.001 & 0.001 \end{matrix} \right) \quad \cancel{[0, 2\varepsilon]} = \varepsilon \\ &2\varepsilon r = 0.002 \times 0.5 \quad \varepsilon \in [-\varepsilon, \varepsilon] \quad \frac{0.002}{-0.001} = 0.001 \\ &\therefore = 0.0010 - 0.001 \quad \frac{0.002}{0.001} = 2 \\ &\therefore = 0. \end{aligned}$$

→ Putting it together

Training a neural network

→ Pick a network architecture

(Connectivity pattern between neurons)

No. of inputs: dimension of features, $x^{(i)}$

No. of output units: Number of classes

Hidden layers ??

Reasonable default: 1 hidden layer, or

if > 1 hidden layer, have same no. of

hidden units in every layer (usually the more the better)

→ Steps in training a neural network -

① Randomly initialise weights $\theta^{(i)}$

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 backpropagation to compute partial derivatives

$$\boxed{\frac{\partial}{\partial \theta^{(i)}} J(\theta)}$$