

Week 5

Notation

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$

L - total no. of layers in network

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

Lets consider 2 classification cases:

- binary classification.

output unit $\left\{ \begin{array}{l} y = 0 \text{ or } 1 \\ S_L = 1 \text{ (only one unit in output layer)} \\ \text{or this will be written as } k=1 \end{array} \right.$ $h_\theta(x) \in \mathbb{R}$

- Multi-class classification (k classes)

output units $\left\{ \begin{array}{l} y \in \mathbb{R}^k \text{ eg. } \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \\ \text{petrol car, motorcycle, truck} \\ h_\theta(x) \in \mathbb{R}^k \\ S_L = k \text{ where } k \geq 3 \end{array} \right.$

Cost Function

Logistic regression

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

\nwarrow
since θ_0 not considered

Neural network (for k terms)

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

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

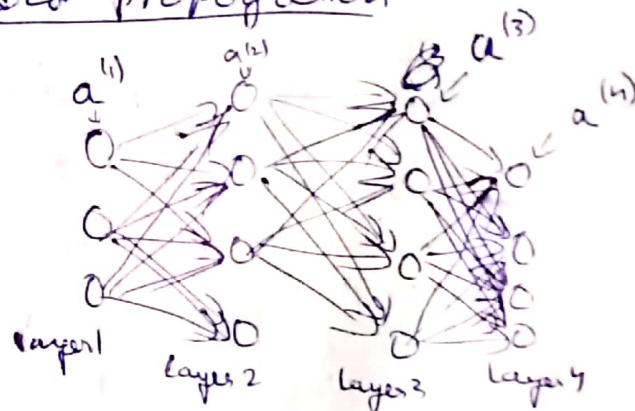
Gradient Computation

The goal is to $\min_{\theta} J(\theta)$

For gradient descent or any other advanced optimization algorithms, we need,

- $J(\theta)$ \leftarrow we can use the formula to find this
- $\frac{\partial}{\partial \theta^{(l)}} J(\theta)$ \leftarrow we will learn how to find this

Forward propagation



Given one training example (x, y) ,

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

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

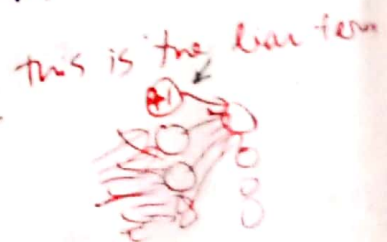
$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$

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

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

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

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



Backpropagation Algorithm

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

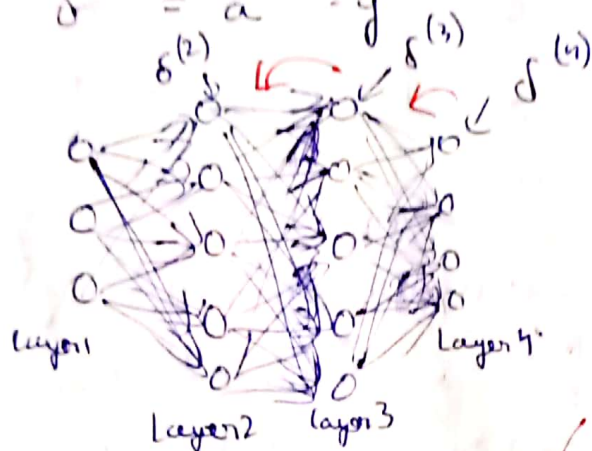
for each output unit (layer $L=4$)

$$\delta_i^{(4)} = a_i^{(4)} - y_i$$

This is nothing but prediction $(h_0(x))_i$

In vectorized form,

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



Differentiation of $g(z^{(3)})$

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} * g'(z^{(3)}) \rightarrow a^{(3)} * (1 - a^{(3)})$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} * g'(z^{(2)}) \rightarrow a^{(2)} * (1 - a^{(2)})$$

We find that (through calculus),

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_i^{(l)} \delta_j^{(l+1)} \quad \left(\begin{array}{l} \text{ignoring } \lambda; \\ \text{if } \lambda = 0 \end{array} \right)$$

Note: We don't find $\delta^{(1)}$ since it is the input layer x

The steps for backpropagation Algorithm

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

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

← This is the accumulation factor that'll be used to find $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$

For $i = 1$ to $m \leftarrow (x^{(i)}, y^{(i)})$

→ 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 $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ using backpropagation algorithm

→ $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_i^{(l)} \delta_j^{(l+1)}$

← In matrix form $\Delta^{(l)} := \Delta^{(l)} + a^{(l+1)} (a^{(l)})^T$

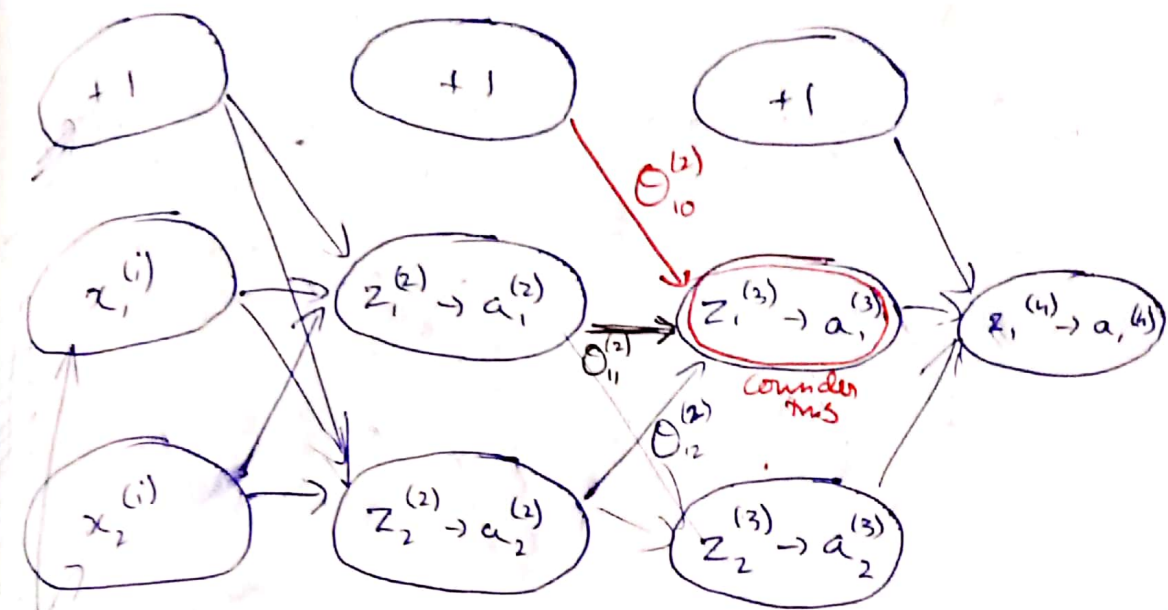
→ $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}$ if $j \neq 0$

← ~~Backpropagate~~
Gradient terms accumulator
 $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$ if $j = 0$

→ $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = D_{ij}^{(l)}$

Note: Backpropagation is neural-network terminology for minimizing our cost function

Intuition of backpropagation



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

Backpropagation is kind of opposite of this

consider cost function,

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

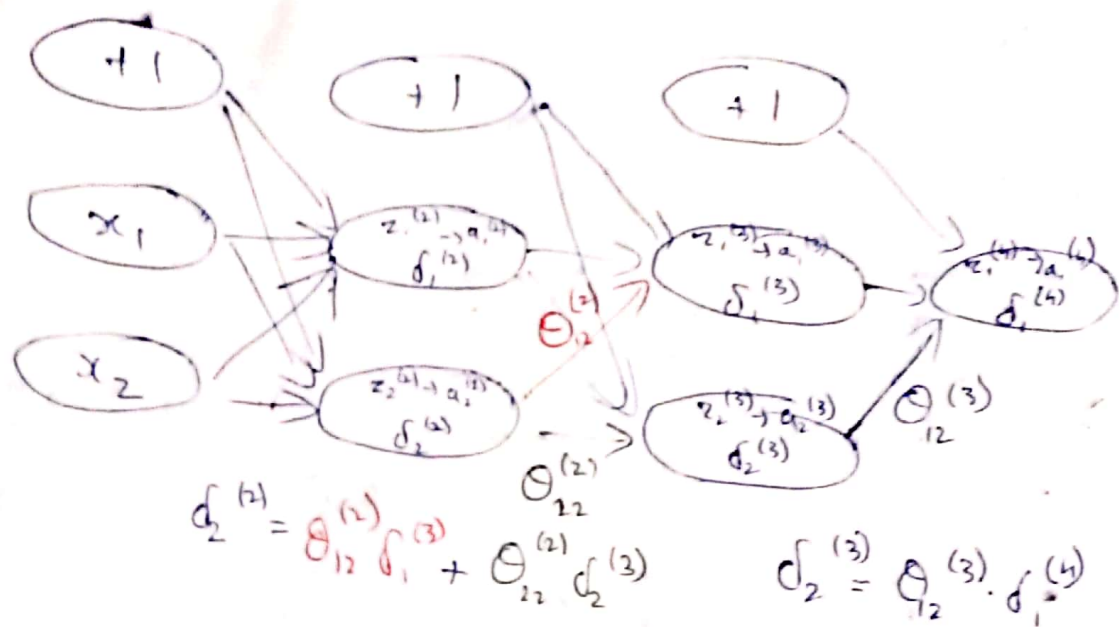
Focusing on a single example $x^{(i)}, y^{(i)}$, the case of 1 output unit, and ignoring regularization ($\lambda = 0$),

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

We can think of this as $\text{cost}(i) \approx (h_{\theta}(x^{(i)}) - y^{(i)})^2$

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



Unrolling parameters

Taking matrices and turning them into vectors $A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ $B = \begin{bmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \\ 2 & 2 & 2 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 2 \\ 2 \\ 2 \\ 2 \\ 2 \end{bmatrix}$
~~and the vector~~

Example

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

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

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

$$\text{thetaVec} = [\text{theta1}(:); \text{theta2}(:); \text{theta3}(:)]$$

$$D\text{Vec} = [D1(:); D2(:); D3(:)]$$

$$\text{theta1} = \text{reshape}(\text{thetaVec}(1:110), 10, 11)$$

$$\text{theta2} = \text{reshape}(\text{thetaVec}(111:220), 10, 11)$$

$$\text{theta3} = \text{reshape}(\text{thetaVec}(221:231), 1, 11)$$

Learning algorithm

- Have initial parameters $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$
- Unroll to get initial Theta to pass to `minimize (@costFunction, initialTheta, options)`

function [JVal, gradientVec] = costFunction(thetaVec)

From thetaVec, get $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$

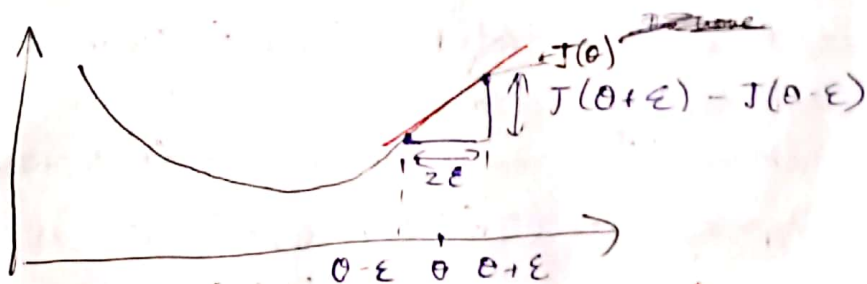
Use forward prop/back prop to compute $D^{(1)}, D^{(2)}, D^{(3)} \leftarrow J(\theta)$

Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradientVec

Gradient Checking

To make sure forward prop/back prop is working properly.

Numerical estimation of gradients



We need to find the derivative to get the slope of the line. We use the above method to approximate the derivative.

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

Use a small ~~epsilon~~ ϵ (epsilon) to get a more accurate derivative approximation (around 10^{-4})

Parameter vector θ

$\theta \in \mathbb{R}^n$ (θ is "unrolled" version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$)

$$\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$$

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

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

\vdots

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

On octave,
epsilon = $1e-4$;
for $i = 1:n$,

thetaPlus = theta;

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

thetaMinus = theta;

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

gradApprox(i) = (J(thetaPlus) - J(thetaMinus)) / (2 * EPSILON);

end;

We need to check if gradApprox \approx DVec

From
backprop

Implementation:

→ Implement backprop to compute DVec

→ Implement numerical gradient check to get gradApprox

→ Make sure they give similar values

→ Turn off gradient checking. Using backprop for learning.

Note: Be sure to turn off gradient checking code ~~in~~ after verifying if the NN is working properly, since it makes the code run slowly, when learning

Random initialization

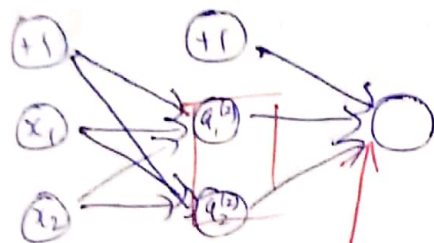
For gradient descent and advanced optimization methods, we need to set an initial value of θ

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

Consider gradient descent, can we set $\text{initialTheta} = \text{zeros}(n, 1)$? **All zeros**

No. This won't work in neural network.

Zero initialization



$$\Rightarrow \theta_{ij}^{(1)} = 0 \text{ for all } i, j, l$$

$$a_1^{(2)} = a_2^{(2)} \quad \& \quad \delta_1^{(2)} = \delta_2^{(2)}$$

Therefore $\frac{\partial}{\partial \theta_{o1}^{(1)}} J(\theta) = \frac{\partial}{\partial \theta_{o2}^{(1)}} J(\theta)$ and $\theta_{o1}^{(1)} = \theta_{o2}^{(1)}$

So even after each update, $a_1^{(2)} = a_2^{(2)}$ will remain, and the process repeats

So the neural network only sees one value and doesn't learn

So we use random initialization to perform symmetry breaking.

Random initialisation:

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

E.g.

Random 10x11 matrix with values between 0 & 1

$\text{Theta1} = \text{rand}(10, 11) * (2 * \text{INIT_EPSILON})$

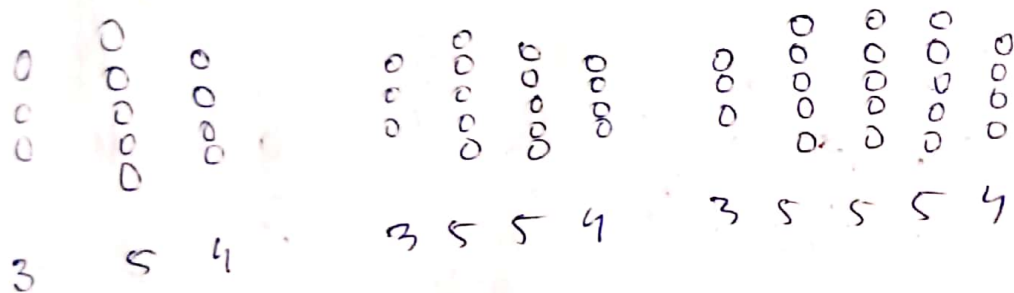
$- \text{INIT_EPSILON};$

← You'll get a value between $-\epsilon$ & ϵ

$\text{Theta2} = \text{rand}(1, 11) * (2 * \text{INIT_EPSILON})$

$- \text{INIT_EPSILON};$

Picking a Neural Network Architecture



→ No. of input units: Dimension of feature $x^{(1)}$

→ No. of output units: Number of classes

Make sure it is like $y = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \dots$

→ Reasonable default is 1 ~~hidden~~ hidden layer

→ As if no. of hidden layers > 1 , have same no. of hidden units in every layer (usually the more units, the better).

Training a neural network

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

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

for $i = 1:m$

Perform forward prop & back prop 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 $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$

5. Use gradient checking to compare

$\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$ computed using back prop vs.

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

Then disable gradient checking code

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

