

WEEK-5

Neural Networks: Representation & Learning

Week 5:

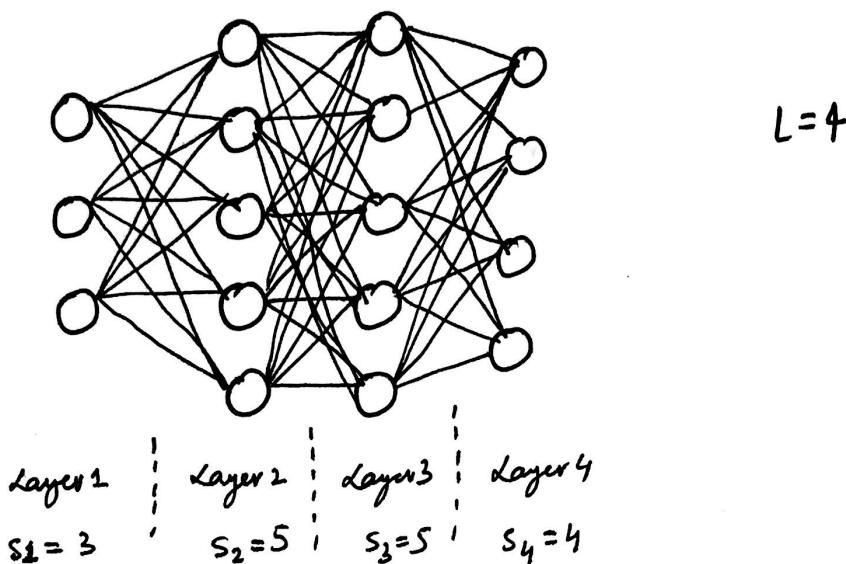
Cost function & Back Propagation

- fitting the parameters for a training set.

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

L = total no. of layers in network

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



①

Binary classification

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

$$\stackrel{\geq 0}{\rightarrow} h_{\theta}(x)$$

1 o/p unit

$$s_L = 1$$

$$K = 1$$

②

Multi-class classification

$$y \in \mathbb{R}^K$$

$$\left[\quad \right]_{K \times 1}$$

K output units.

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

$$s_L = K$$

Cost function for Neural Network:

→ Logistic regression

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

→ Neural Network:

$\underline{h_\theta(x)} \in \mathbb{R}^K$
 K -dimensional vector | $(h_\theta(x))_i^o = i^{\text{th}}$ output.

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

Back Propagation Algorithm

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

Need to compute

① $J(\theta)$

② $\frac{\partial J(\theta)}{\partial \theta_{ij}^l}$

$\theta_{ij}^l \in \mathbb{R}$

⇒ Basically partially derivative of all θ

Ex: Gradient computation: one training example (x, y) (58)

$$L=4$$

\Rightarrow forward propagation:

$$\begin{aligned} a^1 &= x \\ z^2 &= \theta^1 a^1 \\ a^2 &= g(z^2) \quad (\text{add } a_0^{(2)}) \\ z^3 &= \theta^2 a^2 \\ a^3 &= g(z^3) \quad (\text{add } a_0^{(3)}) \\ z^4 &= \theta^3(a^3) \\ a^4 &= g(z^4) \end{aligned}$$

Intuition: $\delta_j^l =$ "error" of node j in layer l

$$\begin{aligned} &\Rightarrow \left\{ \begin{array}{l} \text{for each output unit } (l=u) \\ \delta_j^4 = a_j^4 - y_j \end{array} \right. \xrightarrow{\text{hypothesis } (h_\theta(x))_j} \begin{array}{c} \dots \\ \delta_j^4 \\ \dots \end{array} \quad L=4 \\ &\Rightarrow \delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} \cdot \underbrace{g'(z^{(3)})}_{a^{(3)} \cdot (1-a^{(3)})} \\ &\Rightarrow \delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} \cdot \underbrace{g'(z^{(2)})}_{a^{(2)} \cdot (1-a^{(2)})} \end{aligned}$$

$$\boxed{\frac{\partial J(\theta)}{\partial \theta_{ij}^{(2)}} = a_{ij}^{(2)} \delta_i^{(2+1)}} \Rightarrow (\text{ignoring } x_j \text{ if } \lambda=0)$$

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

(5)

Back propagation Algorithm

① Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$

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

loop:

② for $i=1$ to m ,

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

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

end loop:

$$\left\{ \begin{array}{l} \Delta_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l-1)} + \frac{\lambda}{m} \theta_{ij}^{(l)} \quad \text{if } j \neq 0 \\ \Delta_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l-1)} \quad \text{if } j = 0 \end{array} \right\} \quad \frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = \Delta_{ij}^{(l)}$$

Backpropagation Intuition

→ what is backpropagation doing?

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

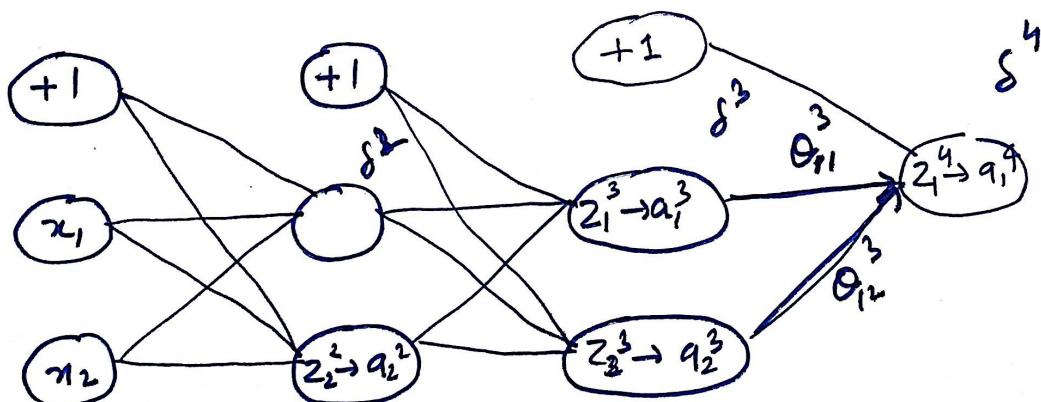
→ focussing on a single example x^i, y^i , the case of 1 output unit ($\lambda=0$) ($K=1$)

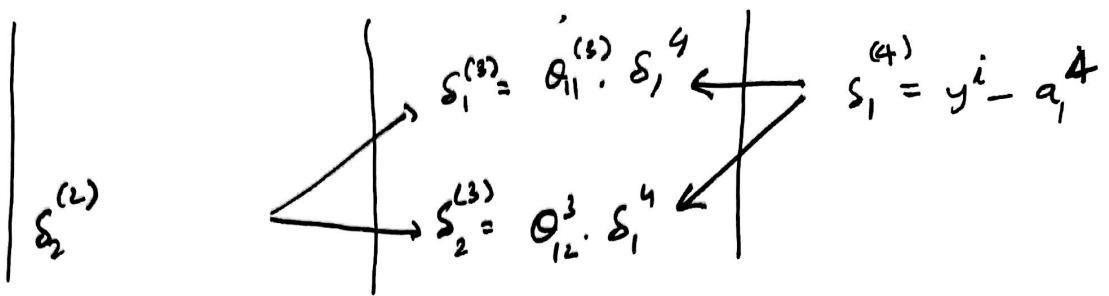
→ δ_j^l = "error" of cost for a_j^l (unit j in layer l).

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

↓
calculating the error, so it can change
the θ /activation unit value

→





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

Implementation Note: Unrolling parameters.

→ when using advanced optimization, fminunc uses initial theta as vector

function [Ejval, grad] = cost function(theta)
 ...
 opttheta = fminunc(@cost function, initial Theta, 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 vectors"

Example:

$$S_1 = 10$$

$$S_2 = 10$$

$$S_3 = 1$$

$$\begin{array}{c|c|c} \theta^1 & \theta^2 & \theta^3 \\ R^{10 \times 11} & R^{10 \times 11} & R^{1 \times 11} \\ \hline D^1 & R^{10 \times 11} \end{array}$$

~~$\theta^1 \in R^{10 \times 11}$~~

```
# thetaVec = [Theta1(:); Theta2(:); Theta3(:)];
DVec = [D1(:); D2(:)];
```

This sets and unroll all matrices into a long vector.

```
# Theta1 = reshape(thetaVec(1:110), 10, 11);
Theta2 = reshape(thetaVec(111:220), 10, 11);
Theta3 = reshape(thetaVec(221:330), 10, 11);
```

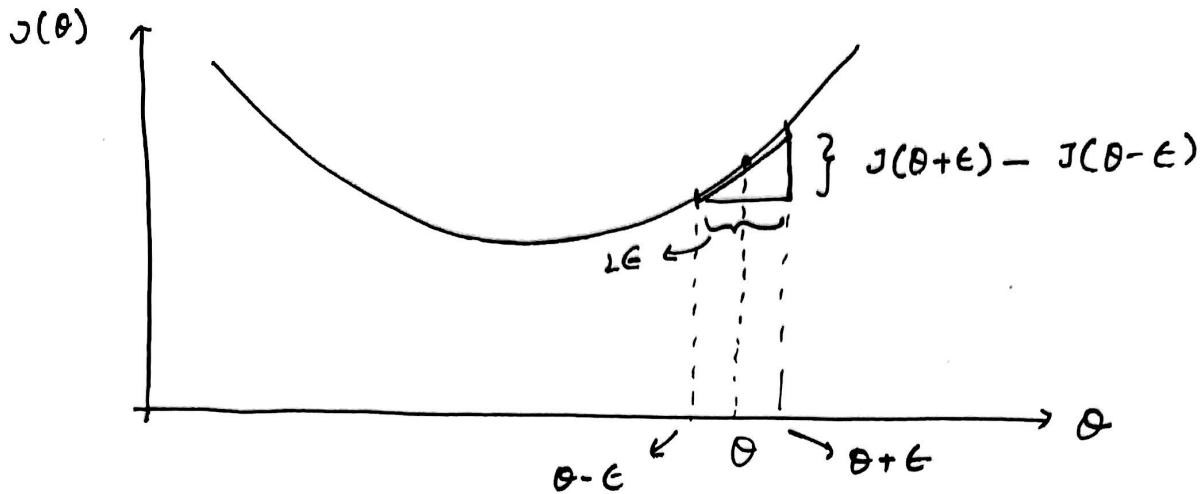
Learning Algorithm:

- ① Have initial parameters $\theta^1, \theta^2, \theta^3$
 - ② unroll to get ~~initial parameters~~ initial theta
 - ③ fminunc (@cost fn, initial Theta, 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 .
and $J(\theta)$

Gradient checking



$$\therefore \frac{\partial J(\theta)}{\partial \theta} = \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \left. \begin{array}{l} \text{gradient / slope} \\ (\text{Two-sided difference}) \end{array} \right\}$$

$\epsilon \approx 10^{-4}$ } set epsilon

Implement in Octave: grad Appox = $\frac{J(\theta + \text{Appox}) - J(\theta - \text{Appox})}{2 * \text{Appox}}$

Parameter Vector $\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n] \in \mathbb{R}^n$

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

for $i = 1:n$,

thetaminus, thetaplus = theta;

thetaplus(i) = thetaplus(i) + epsilon

thetaminus(i) = theta minus - epsilon

grad approx(i) = $\frac{J(\thetaplus(i)) - J(\thetaminus(i))}{2\epsilon}$

end;

Now, check

$$\boxed{\text{gradApprox}(i) \approx D\text{vec}}$$

from backpropagation

● Implementation Note:

- ① Implement backprop to compute Dvec (unrolled D1, D2, D3)
 - ② Implement numerical gradient check to compute gradApprox.
 - ③ Make sure they give similar values
- *④ Turn off gradient checking. Using backprop code for learning.

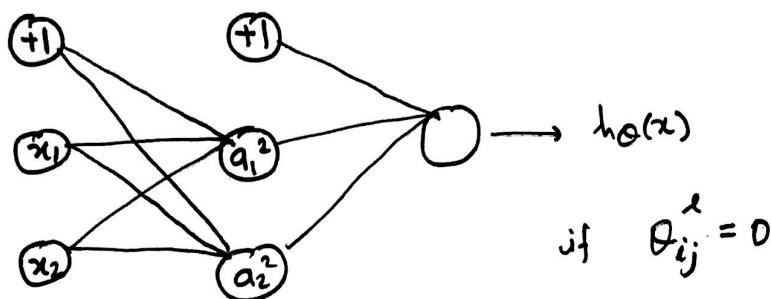
Because : If you run numerical gradient computation on every iteration of gradient descent, your code will be very slow.

Random Initialization: initial value of θ

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

- * if $\text{initialTheta} = \text{zeros}(n, 1)$
 - works well for logistic regression
 - does not work well in Neural Network

Case 1: zero initialization



$$\alpha_1^2 = \alpha_2^2, \text{ Also, } \delta_1^2 = \delta_2^2$$

$$\frac{\partial J(\theta)}{\partial \theta_1^{(i)}} = \frac{\partial J(\theta)}{\partial \theta_{02}^{(i)}}$$

$$\theta_{01}^{(i)} = \theta_{02}^{(i)}$$

if activation units are identical, then it results in redundantly complex network producing same o/p

- * After each update, parameters corresponding to inputs going into each of two hidden units are identical!

Case 2: Random initialization : [Symmetry Breaking] ***

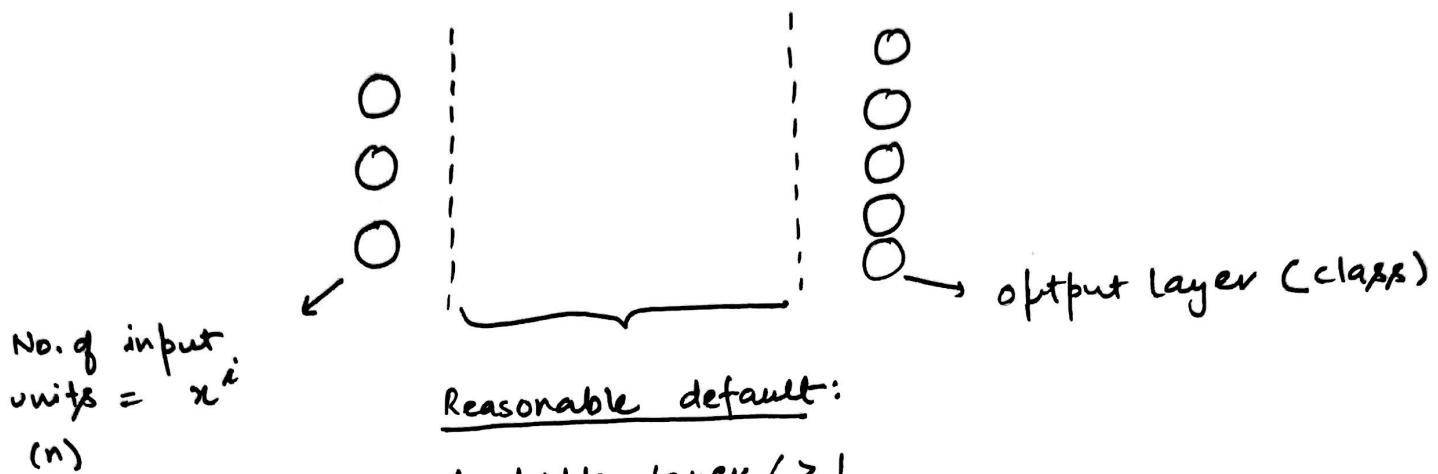
→ initialize each θ_{ij}^2 to a random value in $[-\epsilon, \epsilon]$

{es.
 $\text{Theta1} = \text{rand}(10, 11) * (2 * \text{Init_Epsilon}) - \text{init_epsilon}$ }

Summary: Putting it together

→ Training a Neural Network

- ① Pick a Network architectures (connectivity pattern between neurons)



Reasonable default:

1 hidden layer / > 1
hidden layer, have same
no. of hidden units. in
every layer [usually the more the
better] → computationally expensive]

→ Steps:

1. Randomly initialize weights.
2. Implement FP to get $h_{\theta}(x)^i$ for any x^i
3. Implement code to compute cost fn. $J(\theta)$
4. Implement Back propagation (BP) to compute
partial derivatives $\frac{\partial J(\theta)}{\partial \theta_{ij}^k}$
5. Use gradient checking to compare $\frac{\partial J(\theta)}{\partial \theta_{jk}^i}$ computed using
BP vs. using numerical estimate
6. Use Grad descent or advanced optimization method with
BP to try to minimize $J(\theta)$ as a fn. of parameters θ