

NEURAL NETWORKS & COST FUNCTION:

- m # training examples
 - L # layers in network (including I/O layers)
 - $S_l = H$ units in layer l (not counting bias unit)

\Rightarrow Binary $[1]$ or $[0]$ vs. Multi-class eg. $[1] [0]$ etc. (k -classes)

$$\hookrightarrow S_L = k \quad (=1 \text{ for binary})$$

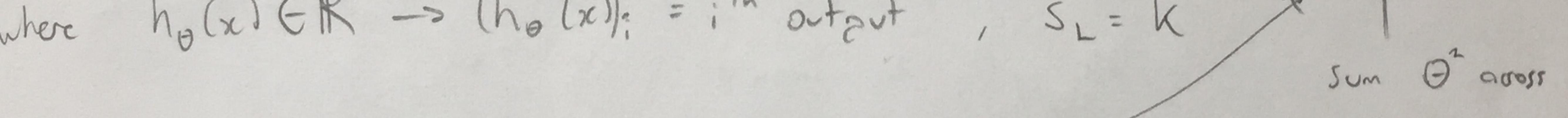
\Rightarrow Recall Cost Fmt for Logistic Regression

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

⇒ Neural Network:

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

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



sum Θ^2 across
all i, j, l vals,
excluding bias terms

Must minimize \rightarrow need gradient (partial derivatives) $\frac{\partial}{\partial \theta_{ij}^{(e)}} J(\theta)$

ex] Just one training example (x, y) : (4 layers, $L=4$)

1] Forward Propagation: $a^{(1)} = ?$

$$z^{(2)} = \Theta_a^{(1)} c_1$$

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

$\gamma^{(3)}$ $\alpha^{(2)} \quad (2)$

$$z''' = \Theta'''_a$$

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

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

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

$\alpha \rightarrow$ "activations"

2] Back Propagation:

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

⇒ for each output unit ($L=4$) $\rightarrow \delta_j^{(4)} = a_j^{(4)} - y_j$

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

↳ $\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)})$ $a^{(3)} \cdot (1-a^{(3)})$

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

No $\delta^{(1)}$ (no error in features of training set!)

⇒ Possible to prove $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)}$ (ignoring λ for now)

Back Prop

→ Algorithm: Set $\Delta_{ij}^{(l)} = 0$ for all l, i, j (use to compute $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$)

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

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

Use forward prop. to compute $a^{(l)}$ for $l=2, 3, \dots, L$

Use $y^{(i)}$ to compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

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

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)} \Rightarrow \Delta^{(l)} = \underbrace{\Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T}_{\text{VECTORIZED}}$

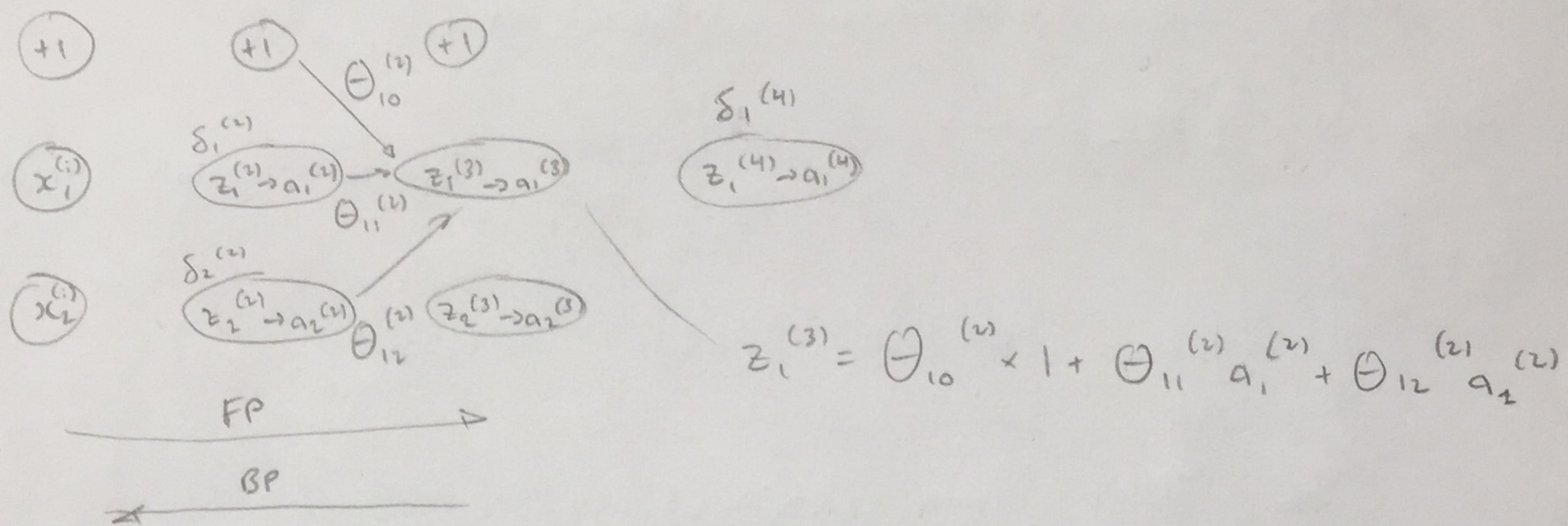
$D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \text{ if } j > 0$

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

↳ D terms = partial derivatives!

⇒ Back Propagation not easy to understand!

⇒ Forward Prop example:



⇒ Single example $(x^{(i)}, y^{(i)})$, $\lambda=0$

$$\hookrightarrow \text{cost}(i) = y^{(i)} \log [h_\theta(x^{(i)})] + (1-y^{(i)}) \log [1-h_\theta(x^{(i)})] \quad (\approx [h_\theta(x^{(i)}) - y^{(i)}]^2)$$

↳ "how well is network doing on example i?"

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

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

$$\Rightarrow \delta_1^{(4)} = y^{(i)} - a_1^{(4)} \rightarrow \text{get } \delta_1^{(3)}, \delta_1^{(2)}, \delta_2^{(3)}, \delta_2^{(2)}$$

$$\hookrightarrow \delta_2^{(2)} = \Theta_{12}^{(2)} \delta_1^{(3)} + \Theta_{22}^{(2)} \delta_{12}^{(3)} \leftarrow \text{WEIGHTED SUM!}$$

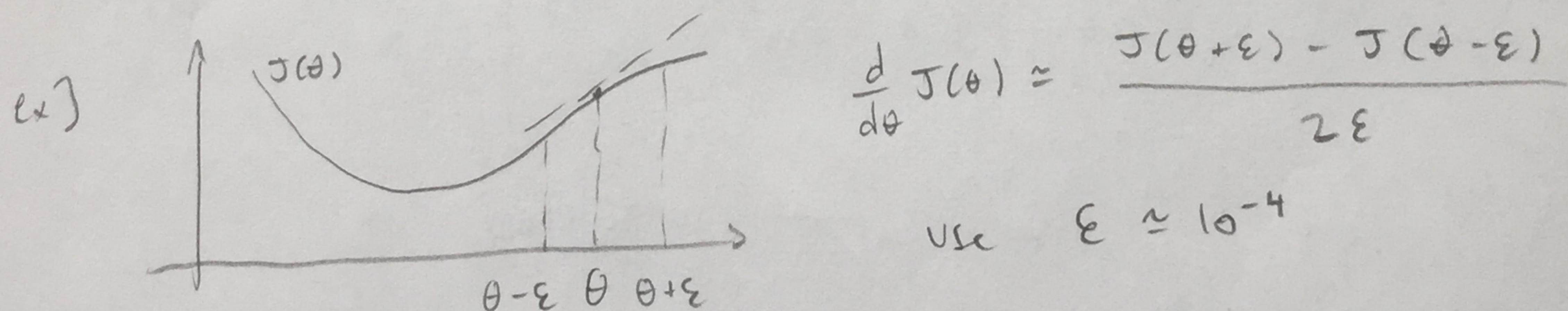
$$\delta_2^{(3)} = \Theta_{12}^{(3)} \delta_1^{(4)}$$

Back Prop Implementation:

function [Jval, gradient] = costFunction(theta) OCTAVE
 optTheta = fminunc(@costFunction, initialTheta, options)
 ⇒ for N.N. of $L=4 \rightarrow \begin{bmatrix} \Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)} \\ D^{(1)}, D^{(2)}, D^{(3)} \end{bmatrix}$ Matrices → "unroll" into vectors
 ↳ thetaVec = [Theta1(:); Theta2(:); Theta3(:)];
 DVec = same ↑
 then: Theta1 = reshape(thetaVec(1:110), 10, 11);
 Theta2 = " " " (111:220) " "
 Theta3 = " " " (221:231), 1, 11);
 ⇒ unroll Θ 's to get initialTheta for fminunc!
 ⇒ in cost function → ① reshape thetaVec to get $\Theta^{(l)}$'s
 ② use FP & BP to get $D^{(l)}$'s & $J(\theta)$
 ③ unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradientVec

Gradient Checking:

↳ Not necessarily good if $J(\theta)$ is decreasing → make sure everything O.K.



implementation: gradApprox = $(J(\text{theta} + \text{epsilon}) - J(\text{theta} - \text{epsilon})) / (2 * \text{epsilon});$

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

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

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \theta_{n-1}, \theta_n + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_{n-1}, \theta_n - \epsilon)}{2\epsilon}$$

implementation: for i=1:n

thetaPlus = theta;

thetaPlus(i) = thetaPlus(i) + ε;

thetaMinus = theta;

thetaMinus(i) = thetaMinus(i) - ε;

grad Approx(i) = $(J(\text{thetaPlus}) - J(\text{thetaMinus})) / (2\epsilon)$;

end;

⇒ check if grad Approx ≈ DVec! (up to small values)

⇒ gradient checking is just to make sure BP implementation is correct

↳ turn off gradient check after confirmation! (otherwise slow)

Random Initialization:

How to choose initialTheta? → used all zeros for multiclass regression...

⇒ Can't use all 0's for N.N.! → set $\alpha_1^{(2)} = \alpha_2^{(2)}$, $\delta_1^{(2)} = \delta_2^{(2)}$, $\Theta_{01}^{(1)} = \Theta_{02}^{(1)}$

Solution: Initialize each $\Theta_{ij}^{(l)}$ to random value in $[-\epsilon, \epsilon]$

↳ Theta1 = $\underbrace{\text{rand}(10, 11) * (2 * \epsilon)}_{\text{b/n } 0 \text{ & } 1} - \epsilon$;

⇒ each $\Theta_{ij}^{(l)}$ must be random!

Theta2 = $\text{rand}(1, 11) * (2 * \epsilon) - \epsilon$;

(break symmetry)

Putting it All Together:

① Pick N.N. Architecture (# layers, # units in each layer)

⇒ # input units = # of features ($x^{(i)}$)

⇒ # output units = # of classes (k)

⇒ reasonable default → 1 hidden layer

↳ if > 1 hidden layer, same # hidden units in every layer

⇒ more hidden units usually better, but ↑ comp. \$

② Randomly initialize weights

③ FP to set $h_\theta(x^{(i)})$ for any $x^{(i)}$

④ Compute cost function $J(\theta)$

⑤ BP to compute partial derivatives $\frac{\partial}{\partial \theta_j} J(\theta)$

for $i = 1:m$ (each $(x^{(i)}, y^{(i)})$)

FP & BP → get $a^{(l)}$ & $\delta^{(l)}$ for $l = 2 \dots L$

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

...

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

end

⑥ Use gradient checking to compare P.D. terms from BP & numerical solns.

↳ then, disable gradient checking

⑦ Use gradient descent, fminunc, etc. w/ BP to get minimum of $J(\theta)$
w/ optimal θ params.

* Note: $J(\theta)$ not convex → could get stuck @ local min

↳ still pretty good.