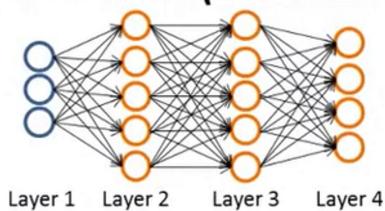


Week 5

• Cost Function and Backpropagation:

a. Cost Function:

Neural Network (Classification)



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

$L =$ total no. of layers in network $L=4$

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

$$s_1=3 \quad s_2=5 \quad s_3=5 \quad s_4=s_L=4.$$

Binary classification

$$y=0 \text{ or } 1. \quad h_{\Theta}(x) \in \mathbb{R} \quad \text{1 output layer. } s_L=1$$

Multi-class classification (K classes).

$$y \in \mathbb{R}^K. \quad \text{e.g. } \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}.$$

K output units.

$$h_{\Theta}(x) \in \mathbb{R}^K \quad s_L=k \quad (k \geq 3).$$

Neural network:

$$h_{\Theta}(x) \in \mathbb{R}^K \quad (h_{\Theta}(x))_i = i^{\text{th}} \text{ output}$$

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

sum of K outputs.

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

don't sum $i=0$.

layer $l \rightarrow$ layer $l+1$.

b. Backpropagation Algorithm:

We need to compute: $J(\Theta)$, $\frac{\partial}{\partial \Theta_{ij}} J(\Theta)$, $(\Theta_{ij})_{ij} \in \mathbb{R}$.

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

$a_j^{(l)}$ = "activation" of node j in layer l .

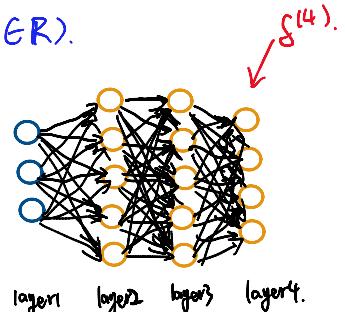
$$\delta_j^{(4)} = a_j^{(4)} - y_j. \quad (h_{\Theta}(x))_j \quad \text{verorize} \Rightarrow \delta^{(4)} = a^{(4)} - y$$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} * g'(z^{(3)}). \quad \text{element-wise.}$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} * g'(z^{(2)}). \quad a^{(2)} * (1-a^{(2)}).$$

(No $\delta^{(1)}$).

$$\frac{\partial}{\partial \Theta_{ij}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}. \quad \text{ignoring } \lambda \text{ or } \lambda=0.$$



Backpropagation algorithm

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

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j). (use to compute $\frac{\partial}{\partial \Theta^{(l)}} J(\Theta)$)

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

Set $\underline{a}^{(1)} = \underline{x}^{(i)}$

→ Perform forward propagation to compute $\underline{a}^{(l)}$ for $l = 2, 3, \dots, L$

→ Using $\underline{y}^{(i)}$, compute $\delta^{(L)} = \underline{a}^{(L)} - \underline{y}^{(i)}$

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

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

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + \underline{\delta}^{(l+1)} (\underline{a}^{(l)})^T$.

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

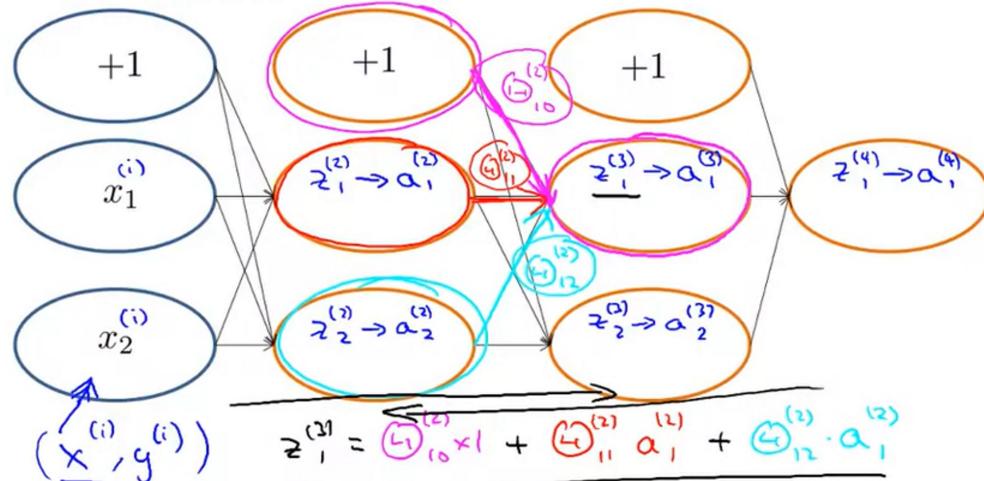
→ $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$ if $j = 0$

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

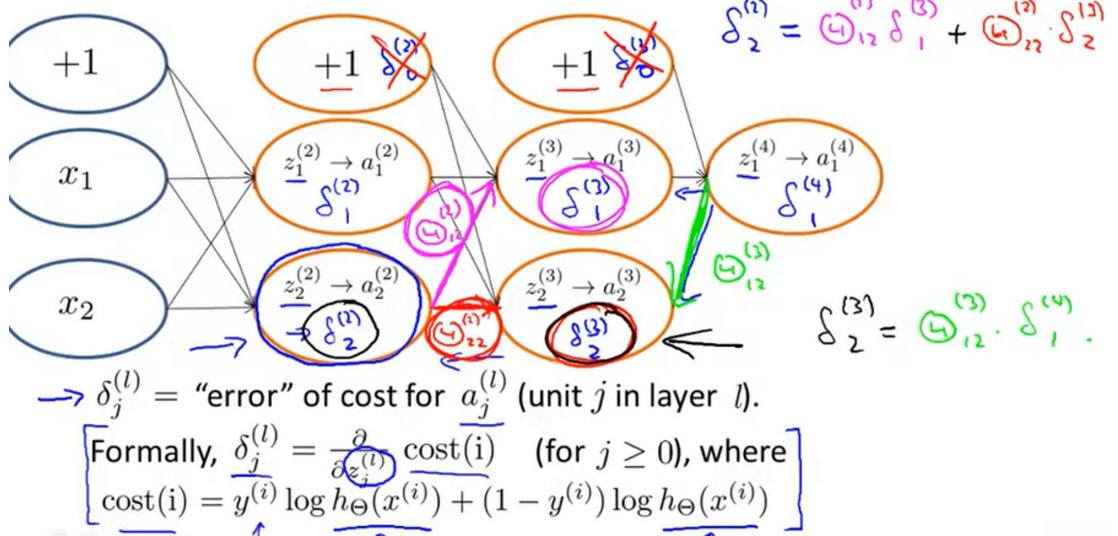
c. Backpropagation Intuition:

Look at the Forward propagation first.

Forward Propagation



Forward Propagation



• Backpropagation in Practice:

a. Implementation Note: Unrolling Parameters:

Advanced optimization

```

function [jVal, gradient] = costFunction(theta)
    ...
    ↪  $\Theta^{n+1}$  (is vector). ↪  $\Theta^{n+1}$  (vectors)
optTheta = fminunc(@costFunction, initialTheta, 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” into vectors

So we are going to unroll it into vectors.

Example

$s_1 = 10, s_2 = 10, s_3 = 1$

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

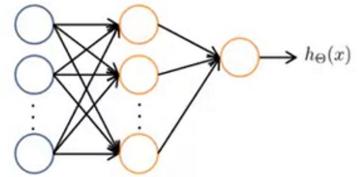
→ $D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$

→ “unroll” into vectors. $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$

→ thetaVec = [Theta1(:); Theta2(:); Theta3(:)];

→ DVec = [D1(:); D2(:); D3(:)];

Theta1 = reshape(thetaVec(1:110), 10, 11);
Theta2 = reshape(thetaVec(111:220), 10, 11);
Theta3 = reshape(thetaVec(221:231), 1, 11);



Learning Algorithm

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

```

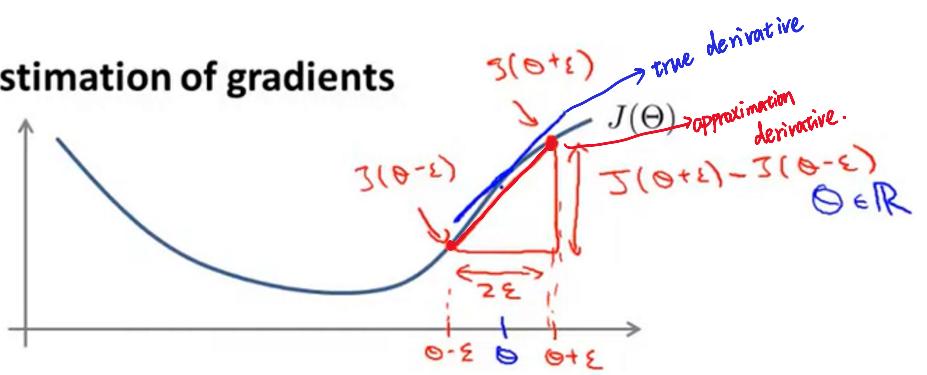
function [jval, gradientVec] = costFunction(thetaVec)
    → From thetaVec, get  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ . reshape
    → Use forward prop/back prop to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\Theta)$ .
    Unroll  $D^{(1)}, D^{(2)}, D^{(3)}$  to get gradientVec.

```

b. Gradient Checking:

Through gradient descent, it looks like J of theta is decreasing, but you might wind up with a neural network that has a higher level of error.

Numerical estimation of gradients



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

$\epsilon = 10^{-4}$ { too small : difficult to compute.
too big : bad result. }

vector θ : $\theta \in \mathbb{R}^n$

$$\begin{aligned}\frac{\partial}{\partial \theta_1} J(\theta) &\approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2 \epsilon} \\ \frac{\partial}{\partial \theta_2} J(\theta) &\approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2 \epsilon} \\ &\vdots \\ \frac{\partial}{\partial \theta_n} J(\theta) &\approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_{n-1}, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_{n-1}, \theta_n - \epsilon)}{2 \epsilon}.\end{aligned}$$

Check approximation $\approx \text{DVec}$ → from backpropagation.

Implementation Note:

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

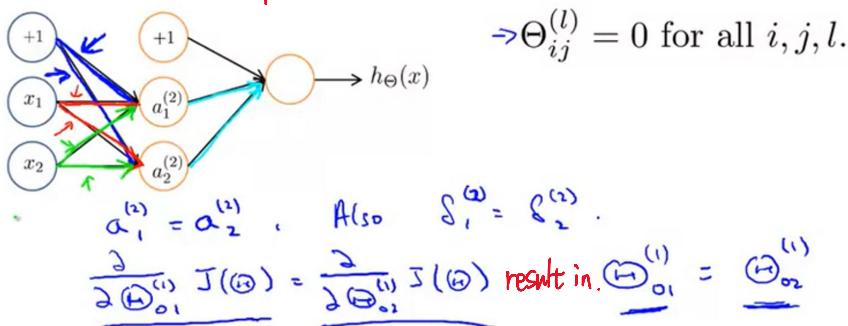
Important:

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

very slow. ($\overset{\downarrow}{g^{(1)}}, \overset{\downarrow}{g^{(2)}}, \overset{\downarrow}{g^{(3)}}$) efficient.

c. Random Initialization:

Zero initialization if we initial with D



After each update, parameters corresponding to

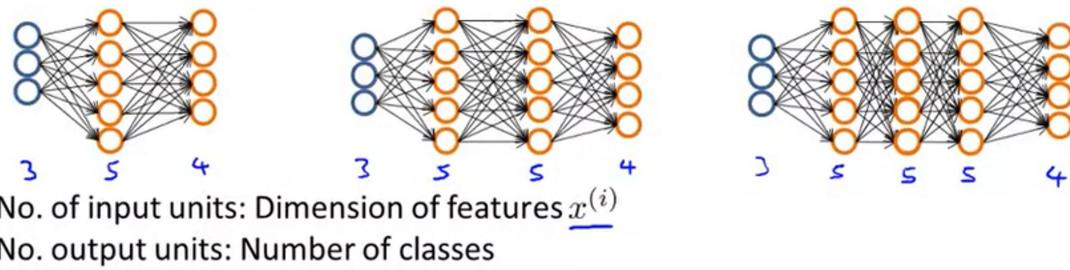
inputs going into each of two hidden units are identical.

Random initialization: Symmetry breaking

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

d. Putting It Together:

① Pick a suitable neural network:



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

$$y \in \{1, 2, 3, \dots, 103\} \quad \text{transform } y \text{ into one-hot vector.}$$

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow$$

② randomly initialize weights

③ implement forwardpropogation to get $h_{\Theta}(x^{(i)})$ for any $x^{(i)}$.

④ implement code to compute cost function $J(\Theta)$.

⑤ implement backpropogation to compute partial derivatives $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$

⑥ use gradient checking to compare $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ computed using backpropogation vs. using estimate of gradient of $J(\Theta)$.

⑦ use gradient descent or advanced optimization method with backpropogation to try to minimize $J(\Theta)$ as a function of parameters Θ .