

Backpropagation 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 $\frac{\partial}{\partial \Theta^{(l)}} J(\Theta)$)

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

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 $\underline{\delta}^{(L)} = \underline{a}^{(L)} - \underline{y}^{(i)}$

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

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

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

Backpropagation algorithm

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

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

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

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 $\underline{\delta}^{(L)} = \underline{a}^{(L)} - \underline{y}^{(i)}$

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

→ $\underline{\Delta}_{ij}^{(l)} := \underline{\Delta}_{ij}^{(l)} + \underline{a}_j^{(l)} \underline{\delta}_i^{(l+1)}$

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

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

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

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

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 $\underline{\delta}^{(L)} = \underline{a}^{(L)} - \underline{y}^{(i)}$

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

→ $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + \underline{a}_j^{(l)} \underline{\delta}_i^{(l+1)}$

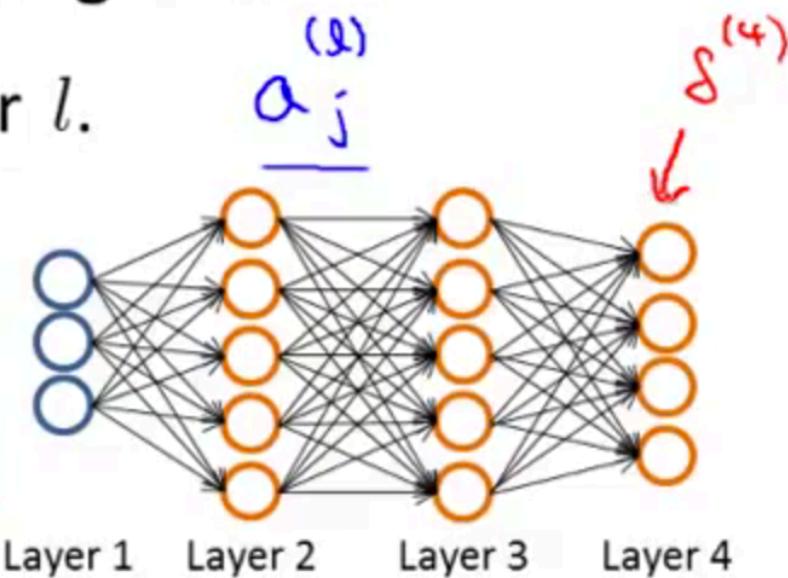
.

Gradient computation: Backpropagation algorithm

Intuition: $\underline{\delta_j^{(l)}}$ = “error” of node j in layer l .

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

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



$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} * g'(z^{(3)})$$

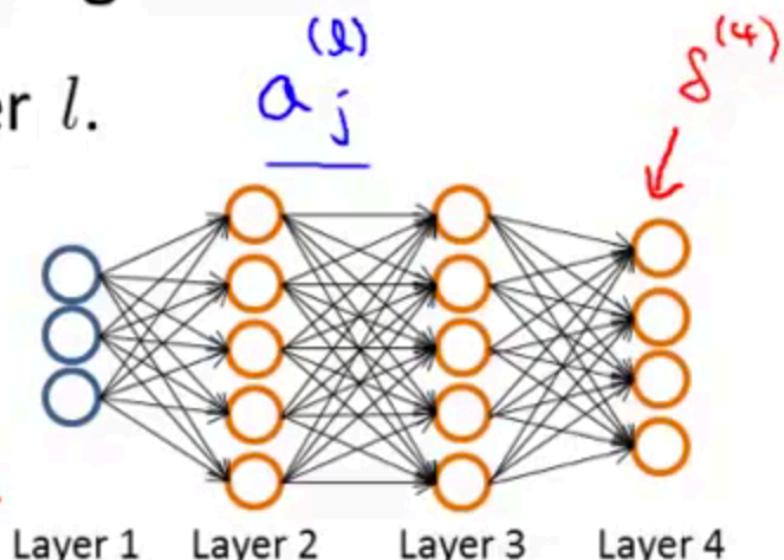
$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} * g'(z^{(2)})$$

Gradient computation: Backpropagation algorithm

Intuition: $\underline{\delta_j^{(l)}}$ = “error” of node j in layer l .

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

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



$$\rightarrow \underline{\delta^{(3)}} = (\underline{\Theta^{(3)}})^T \underline{\delta^{(4)}} * \underline{g'(z^{(3)})}$$

$$\underline{a^{(3)}} * \underline{(1 - a^{(3)})}$$

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

Logistic regression:

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

Θ_0

Neural network:

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

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

Gradient computation

Given one training example ($\underline{x}, \underline{y}$):

Forward propagation:

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

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

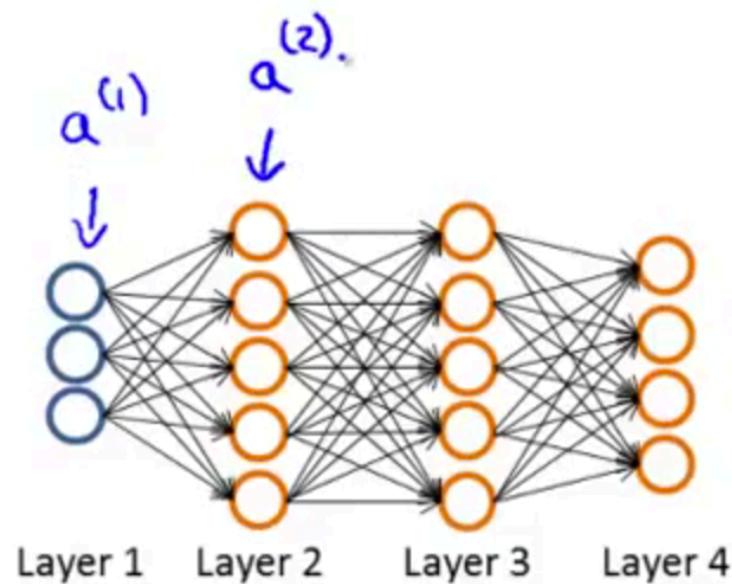
$$\rightarrow \underline{a}^{(2)} = g(\underline{z}^{(2)}) \quad (\text{add } a_0^{(2)})$$

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

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

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

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



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`.

```
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;
```

Random initialization: Symmetry breaking

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

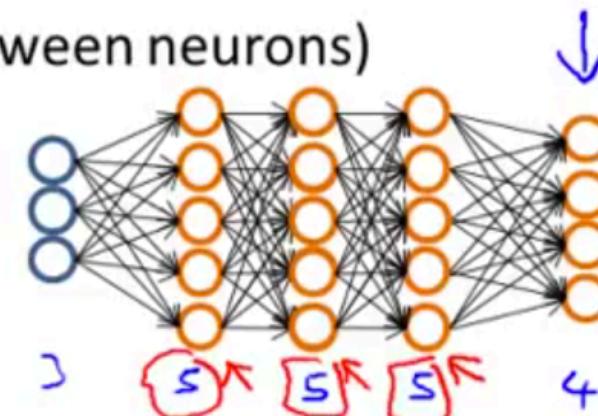
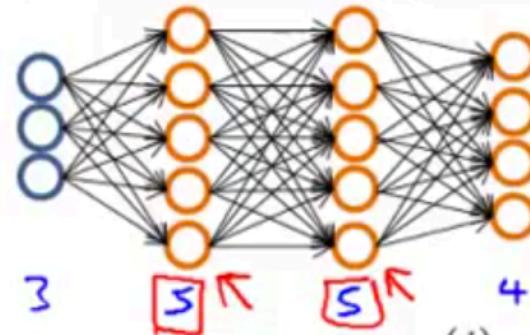
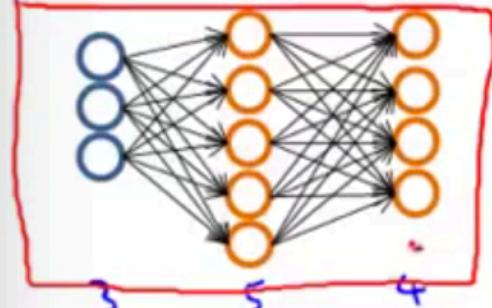
E.g.

```
Theta1 = rand(10,11)*(2*INIT_EPSILON)
        - INIT_EPSILON;
```

```
Theta2 = rand(1,11)*(2*INIT_EPSILON)
        - INIT_EPSILON;
```

Training a neural network

Pick a network architecture (connectivity pattern between neurons)



→ No. of input units: Dimension of features $x^{(i)}$

→ No. output units: Number of classes

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

$$y \in \{1, 2, 3, \dots, 10\}$$

~~$y = 5$~~

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

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$