Lecture #40

# Artificial Intelligence for Data Science

**Week-11:**
MACHINE LEARNING (Part VIII)
**Neural Network Learning**

Course Instructor:

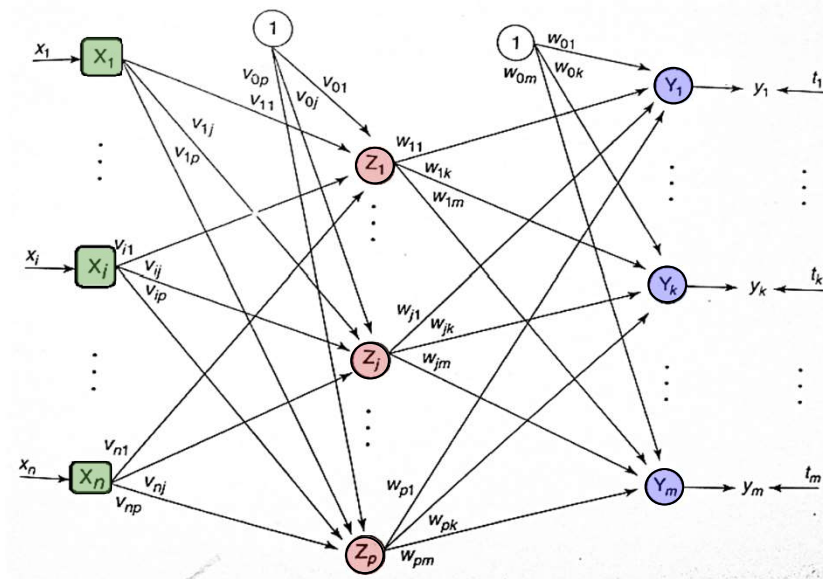**Dr. Monidipa Das**

**Assistant Professor**

**Department of Computational and Data Sciences**

**Indian Institute of Science Education and Research Kolkata, India 741246**
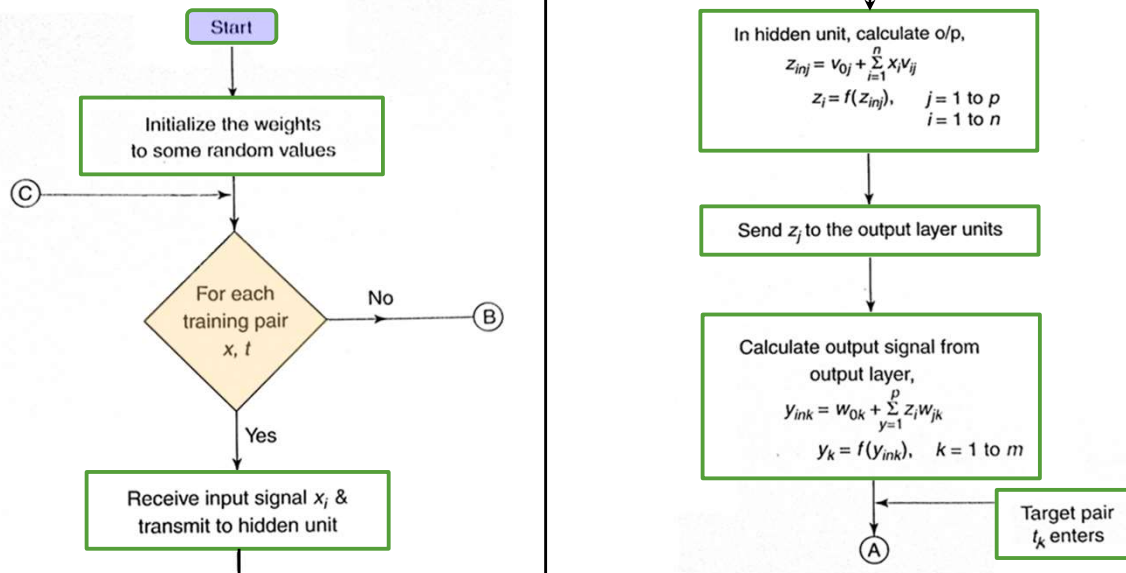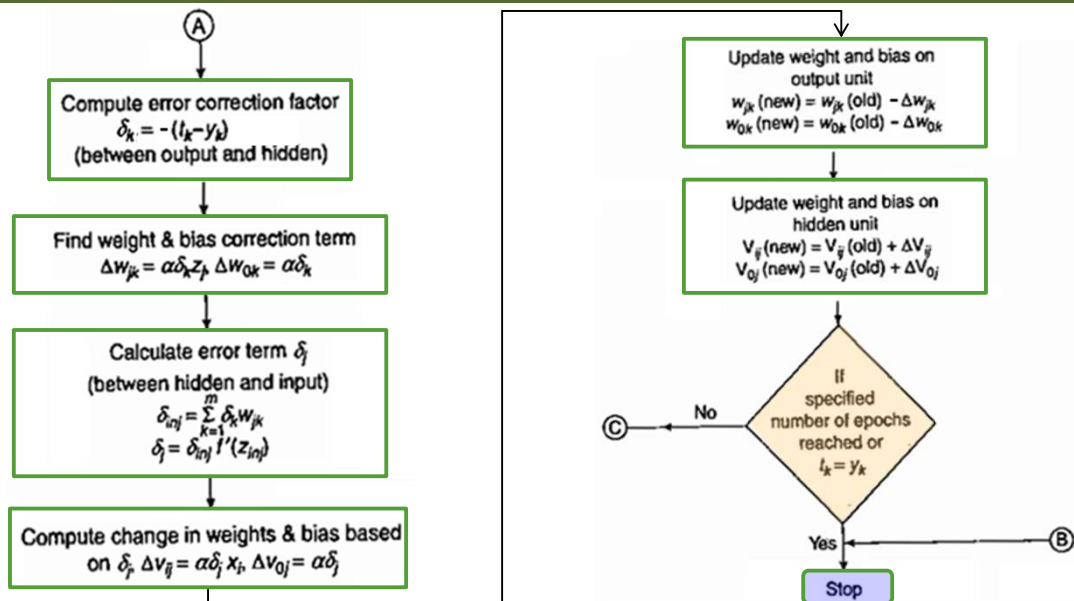
---

# Back Propagation Network (BPN)

- Architecture

# BPN Training

# BPN Training

# Neural Network

$a_i^{(j)}$ = "activation" of unit $i$ in layer $j$

$\Theta^{(j)}$ = weight matrix controlling function mapping from layer $j$ to layer $j+1$

$$a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$
$$a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$
$$a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$
$$h_\Theta(x) = a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

If network has $s_j$ units in layer $j$ *and* $s_{j+1}$ units in layer $j+1$, then $\Theta^{(j)}$ has dimension $s_{j+1} \times (s_j + 1)$ .
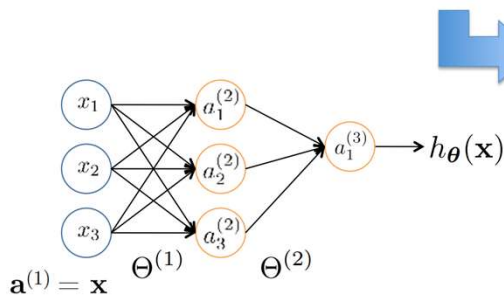
$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4} \qquad \Theta^{(2)} \in \mathbb{R}^{1 \times 4}$$

# Vectorization

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3\right) = g\left(z_1^{(2)}\right)$$
$$a_2^{(2)} = g\left(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3\right) = g\left(z_2^{(2)}\right)$$
$$a_3^{(2)} = g\left(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3\right) = g\left(z_3^{(2)}\right)$$
$$h_\Theta(\mathbf{x}) = g\left(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)}\right) = g\left(z_1^{(3)}\right)$$



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

**Feed-Forward Steps:**
$$\mathbf{z}^{(2)} = \Theta^{(1)} \mathbf{x}$$
$$\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$$
Add $a_0^{(2)} = 1$
$$\mathbf{z}^{(3)} = \Theta^{(2)} \mathbf{a}^{(2)}$$
$$h_\Theta(\mathbf{x}) = \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$$

# Forward Propagation



Feed-Forward Steps:

$$\mathbf{z}^{(2)} = \Theta^{(1)}\mathbf{x} \color{red}{= \Theta^{(1)}\boldsymbol{a}^{(1)}}$$
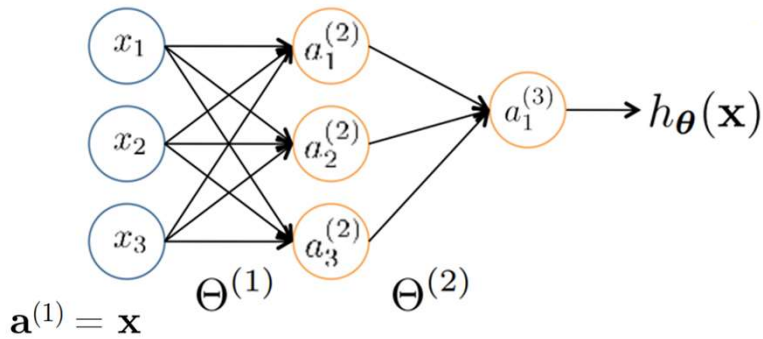
$$\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$$

$$\text{Add } a_0^{(2)} = 1$$

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

$$h_\Theta(\mathbf{x}) = \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$$

# Forward Propagation

- Given one labeled training instance $(\mathbf{x},\ y)$:

### Forward Propagation

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \Theta^{(1)}\mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$     [add $a_0^{(2)}$]
- $\mathbf{z}^{(3)} = \Theta^{(2)}\mathbf{a}^{(2)}$
- $\mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$     [add $a_0^{(3)}$]
- $\mathbf{z}^{(4)} = \Theta^{(3)}\mathbf{a}^{(3)}$
- $\mathbf{a}^{(4)} = h_\Theta(\mathbf{x}) = g(\mathbf{z}^{(4)})$

# Cost Function

**Neural Network:**

$$h_\Theta \in \mathbb{R}^K \qquad (h_\Theta(\mathbf{x}_i))_k = k^{th} \text{output of the } i\text{-th sample}$$

$$J(\Theta) = -\frac{1}{n}\left[\sum_{i=1}^{n}\sum_{k=1}^{K} y_{ik}\log\left(h_\Theta(\mathbf{x}_i)\right)_k + (1 - y_{ik})\log\left(1 - (h_\Theta(\mathbf{x}_i))_k\right)\right]$$

$k^{th}$ class:    true, predicted
not $k^{th}$ class: true, predicted

# Optimizing the Neural Network

$$J(\Theta) = -\frac{1}{n}\left[\sum_{i=1}^{n}\sum_{k=1}^{K} y_{ik}\log(h_\Theta(\mathbf{x}_i))_k + (1 - y_{ik})\log\left(1 - (h_\Theta(\mathbf{x}_i))_k\right)\right]$$

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

Solve via: $\min_\Theta J(\Theta)$

> We can use Gradient Descent (GD)

Need code to compute:
- $J(\Theta)$
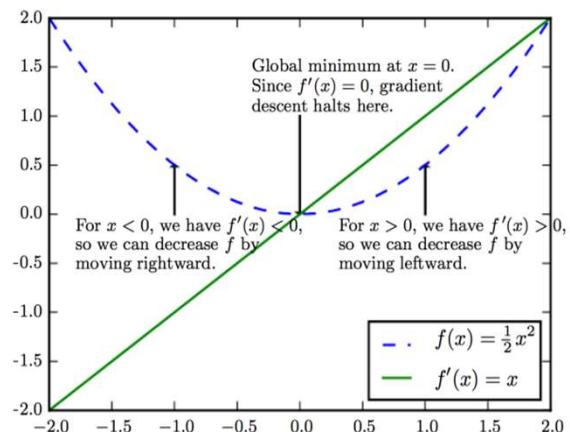- $\frac{\partial}{\partial\Theta_{ij}^{(l)}} J(\Theta)$

## How Gradient Descent (GD) works using derivatives 11

- Criterion $f(x)$ minimized by moving from current solution in direction of the negative of gradient

- For $x>0$, $f(x)$ increases with $x$ and $f'(x)>0$
- For $x<0$, $f(x)$ decreases with $x$ and $f'(x)<0$
- Use $f'(x)$ to follow function downhill
- Reduce $f(x)$ by going in direction opposite sign of derivative $f'(x)$



Global minimum at $x = 0$. Since $f'(x) = 0$, gradient descent halts here.

For $x < 0$, we have $f'(x) < 0$, so we can decrease $f$ by moving rightward.

For $x > 0$, we have $f'(x) > 0$, so we can decrease $f$ by moving leftward.

$f(x) = \frac{1}{2}x^2$

$f'(x) = x$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# Method of Gradient Descent 12

- Decrease $f$ by moving in the direction of the negative gradient
  - This is known as the method of *gradient descent*

- The method proposes a new point    $x' = x - \alpha \nabla_x f(x)$

  - where $\alpha$ is the *learning rate*, a positive scalar. Set to a small constant.

- This converges when every element of the gradient is zero
  - In practice, very close to zero

- We may be able to avoid iterative algorithm and jump to the critical point by solving the equation $\nabla_x f(x) = 0$ for $x$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# Gradient Descent Algorithm

**Gradient Descent: Algorithmic representation**

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(simultaneously update $j = 0$ and $j = 1$)

}

**Parameters**

**Cost function, which is to be minimized**

**Correct updating of parameters**

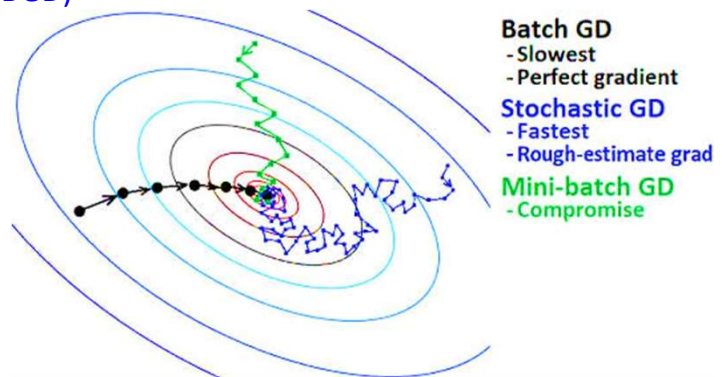$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$
$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$
$$\theta_0 := \text{temp0}$$
$$\theta_1 := \text{temp1}$$

**Incorrect updating of parameters**

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$
$$\theta_0 := \text{temp0}$$
$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$
$$\theta_1 := \text{temp1}$$

Dr. Monidipa Das, Department of CDS, IISER Kolkata
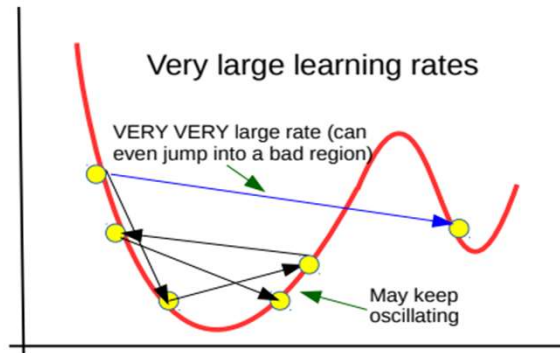
# Types of Gradient Descent

- There are **three major types** of Gradient Descent Algorithms:
  - Batch Gradient Descent (BGD)
  - Stochastic Gradient Descent (SGD)
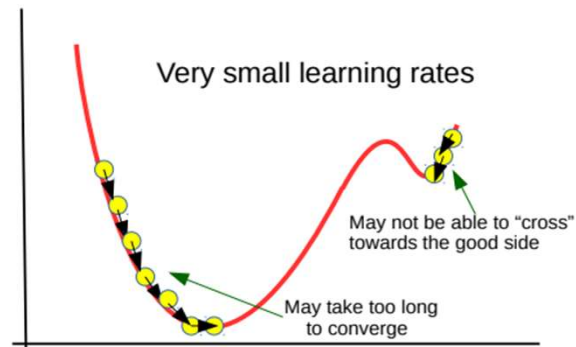  - Mini-Batch Gradient Descent (MBGD)



**Batch GD**
- Slowest
- Perfect gradient

**Stochastic GD**
- Fastest
- Rough-estimate grad

**Mini-batch GD**
- Compromise

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# Learning rate is very important



Very large learning rates

VERY VERY large rate (can even jump into a bad region)

May keep oscillating

Too large, the next point will perpetually bounce haphazardly across the bottom of the well

Very small learning rates

May not be able to "cross" towards the good side
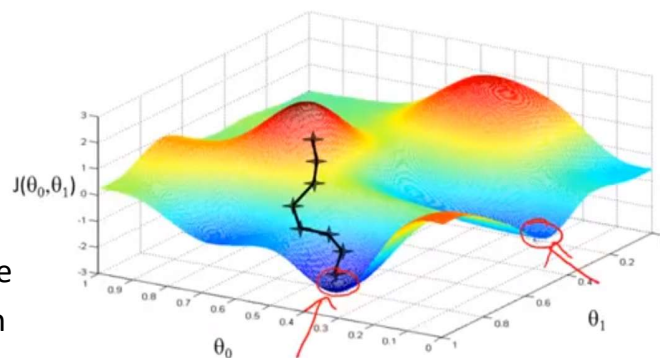
May take too long to converge

Too small learning rate will take too long

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# GD: Advantages and Disadvantages

- Advantages
  - Ease of implementation
  - Convergence guarantee
  - Scalability

- Disadvantages
  - Sensitivity to learning rate
  - Sensitivity to initialization
  - Convergence speed



$J(\theta_0,\theta_1)$

Courtesy: Andrew Ng

Dr. Monidipa Das, Department of CDS, IISER Kolkata

# Optimizing the Neural Network

$$J(\Theta) = -\frac{1}{n} \left[ \sum_{i=1}^{n} \sum_{k=1}^{K} y_{ik} \log(h_\Theta(\mathbf{x}_i))_k + (1 - y_{ik}) \log\left(1 - (h_\Theta(\mathbf{x}_i))_k\right) \right]$$

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

Solve via: $\min_{\Theta} J(\Theta)$

> $J(\Theta)$ is not convex, so GD on a neural net yields a local optimum
> - But, tends to work well in practice

Need code to compute:
- $J(\Theta)$
- $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$
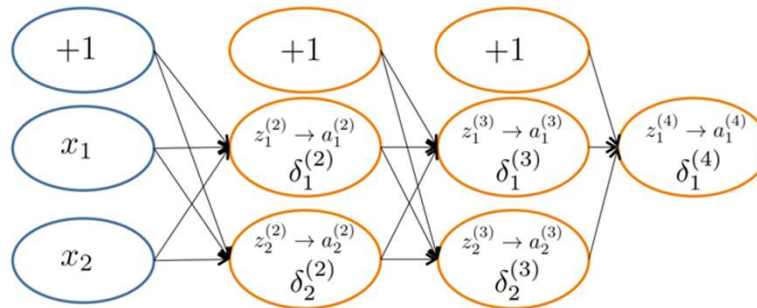
# Backpropagation

- Similar to the perceptron learning algorithm, we cycle through our examples
  - If the output of the network is correct, no changes are made
  - If there is an error, weights are adjusted to reduce the error

- The trick is to assess the blame for the error and divide it among the contributing weights
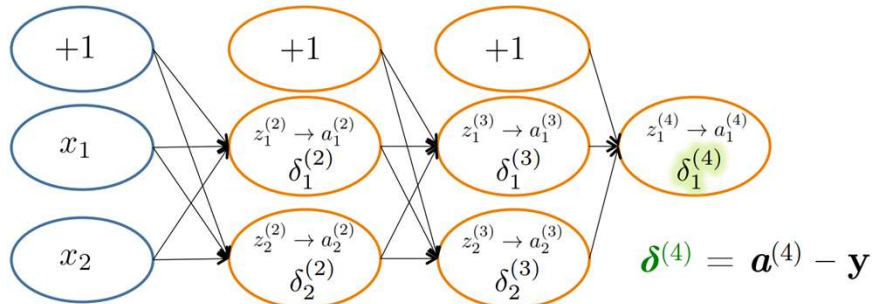
# Backpropagation Intuition



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

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = -\big[y_i \log h_\Theta(\mathbf{x}_i) + (1-y_i)\log(1 - h_\Theta(\mathbf{x}_i))\big]$

# Backpropagation Intuition



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

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

$\dfrac{\partial \boldsymbol{a}}{\partial \boldsymbol{z}} = g'(\boldsymbol{z}) = g(\boldsymbol{z})(1-g(\boldsymbol{z})) = \boldsymbol{a}(1-\boldsymbol{a})$

When $g$ is logistic sigmoid

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

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

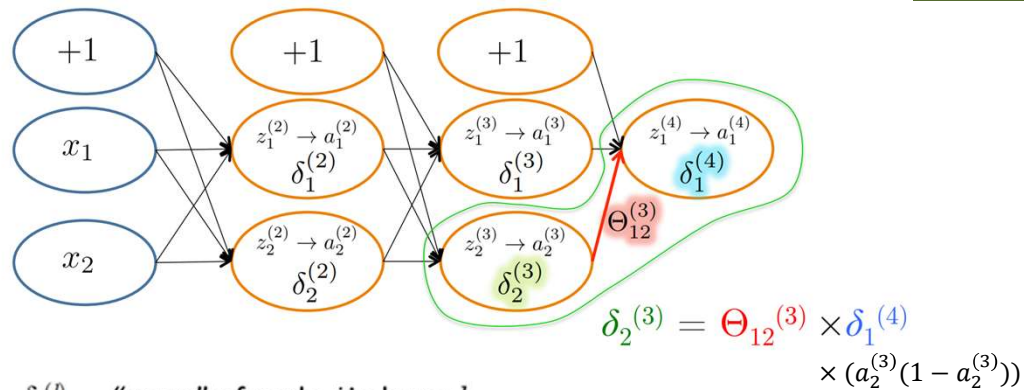where $\text{cost}(\mathbf{x}_i) = -\big[y_i \log h_\Theta(\mathbf{x}_i) + (1-y_i)\log(1 - h_\Theta(\mathbf{x}_i))\big]$

# Backpropagation Intuition

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

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = -\left[y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))\right]$
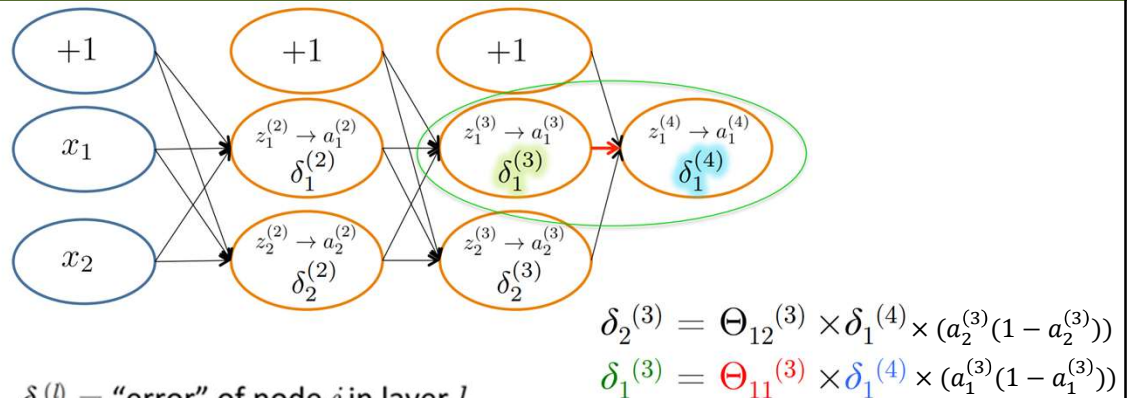
# Backpropagation Intuition

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

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = -\left[y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))\right]$

# Backpropagation Intuition

$$\delta_2^{(2)} = (\Theta_{12}^{(2)} \times \delta_1^{(3)} + \Theta_{22}^{(2)} \times \delta_2^{(3)}) \times (a_2^{(2)}(1 - a_2^{(2)}))$$
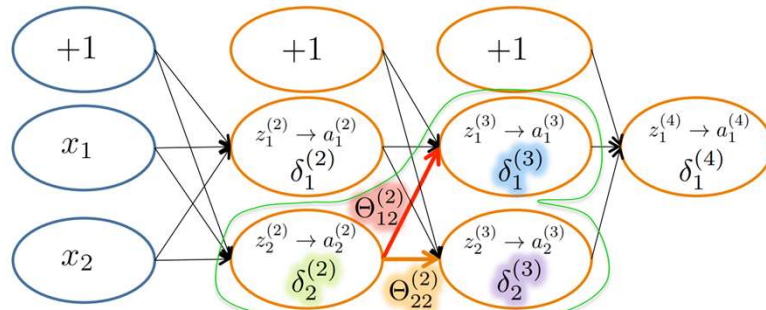
$\delta_j^{(l)}$ = "error" of node $j$ in layer $l$

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = -\big[y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))\big]$

---

# Backpropagation Intuition: Gradient Computation

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

(#layers $L$ = 4)

Element-wise product .*

## Backpropagation

- $\boldsymbol{\delta}^{(4)} = \boldsymbol{a}^{(4)} - \mathbf{y}$
- $\boldsymbol{\delta}^{(3)} = (\Theta^{(3)})^\top \boldsymbol{\delta}^{(4)} .* g'(\mathbf{z}^{(3)})$
- $\boldsymbol{\delta}^{(2)} = (\Theta^{(2)})^\top \boldsymbol{\delta}^{(3)} .* g'(\mathbf{z}^{(2)})$
- (No $\boldsymbol{\delta}^{(1)}$)

$g'(\mathbf{z}^{(3)}) = \mathbf{a}^{(3)} .* (1 - \mathbf{a}^{(3)})$

$g'(\mathbf{z}^{(2)}) = \mathbf{a}^{(2)} .* (1 - \mathbf{a}^{(2)})$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

(ignoring $\lambda$; if $\lambda = 0$)

## Training a Neural Network via Gradient Descent with Backprop

Given: training set $\{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)\}$
Initialize all $\Theta^{(l)}$ randomly (NOT to 0!)
Loop // each iteration is called an epoch
    Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$            (Used to accumulate gradient)
    For each training instance $(\mathbf{x}^{(s)}, y^{(s)})$:
        Set $\mathbf{a}^{(1)} = \mathbf{x}^{(s)}$
        Compute $\{\mathbf{a}^{(2)}, \ldots, \mathbf{a}^{(L)}\}$ via forward propagation
        Compute $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y^{(s)}$
        Compute errors $\{\boldsymbol{\delta}^{(L-1)}, \ldots, \boldsymbol{\delta}^{(2)}\}$
        Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$
    Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n}\Delta_{ij}^{(l)} + \lambda\Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n}\Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$
    Update weights via gradient step $\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \alpha D_{ij}^{(l)}$
Until weights converge or max #epochs is reached

Backpropagation

---

# Questions?