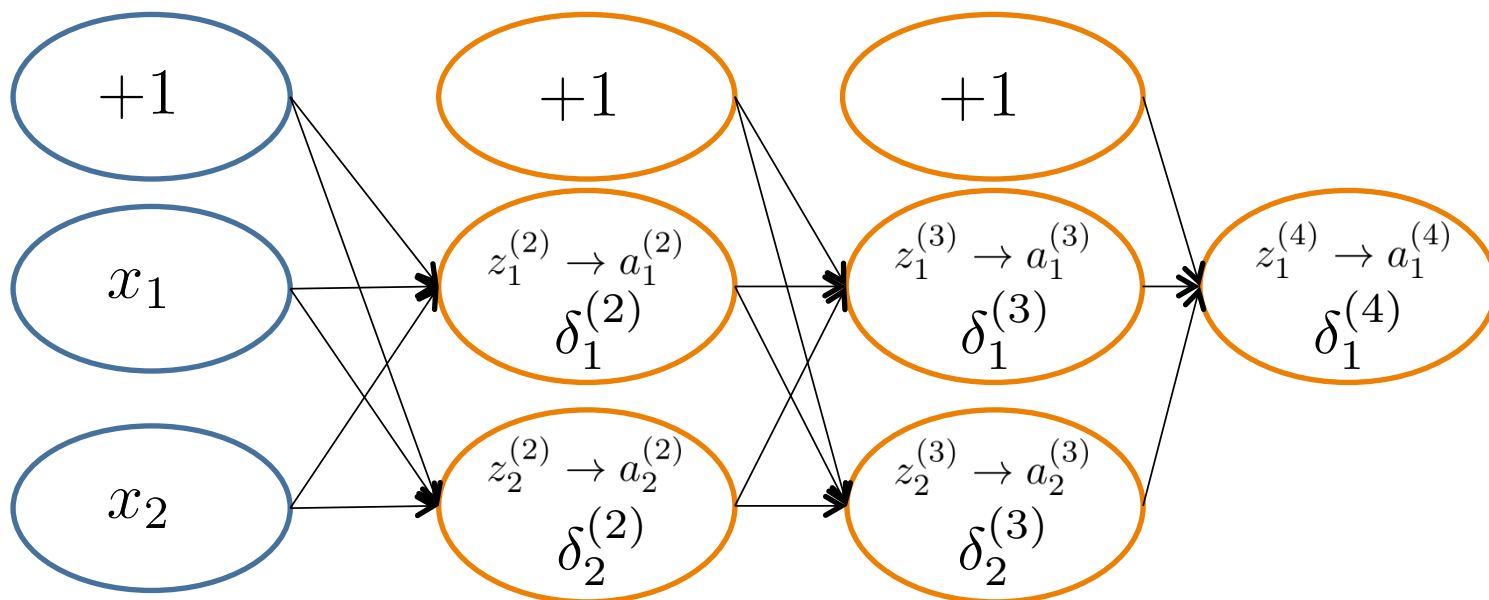# Backpropagation Intuition

- Each hidden node $j$ is "responsible" for some fraction of the error $\delta_j^{(l)}$ in each of the output nodes to which it connects

- $\delta_j^{(l)}$ is divided according to the strength of the connection between hidden node and the output node

- Then, the "blame" is propagated back to provide the error values for the hidden layer
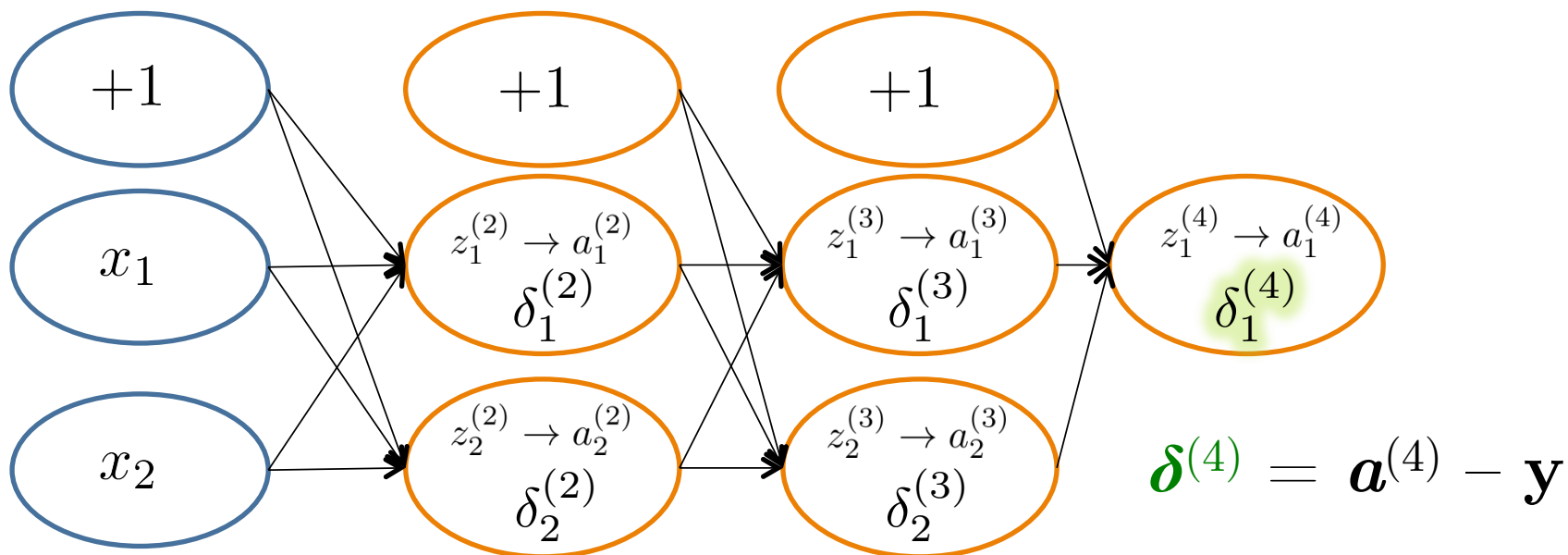
# Backpropagation Intuition



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

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

where $\mathrm{cost}(\mathbf{x}_i) = y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))$

# Backpropagation Intuition



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

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

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

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

Based on slide by Andrew Ng

# Backpropagation Intuition



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

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

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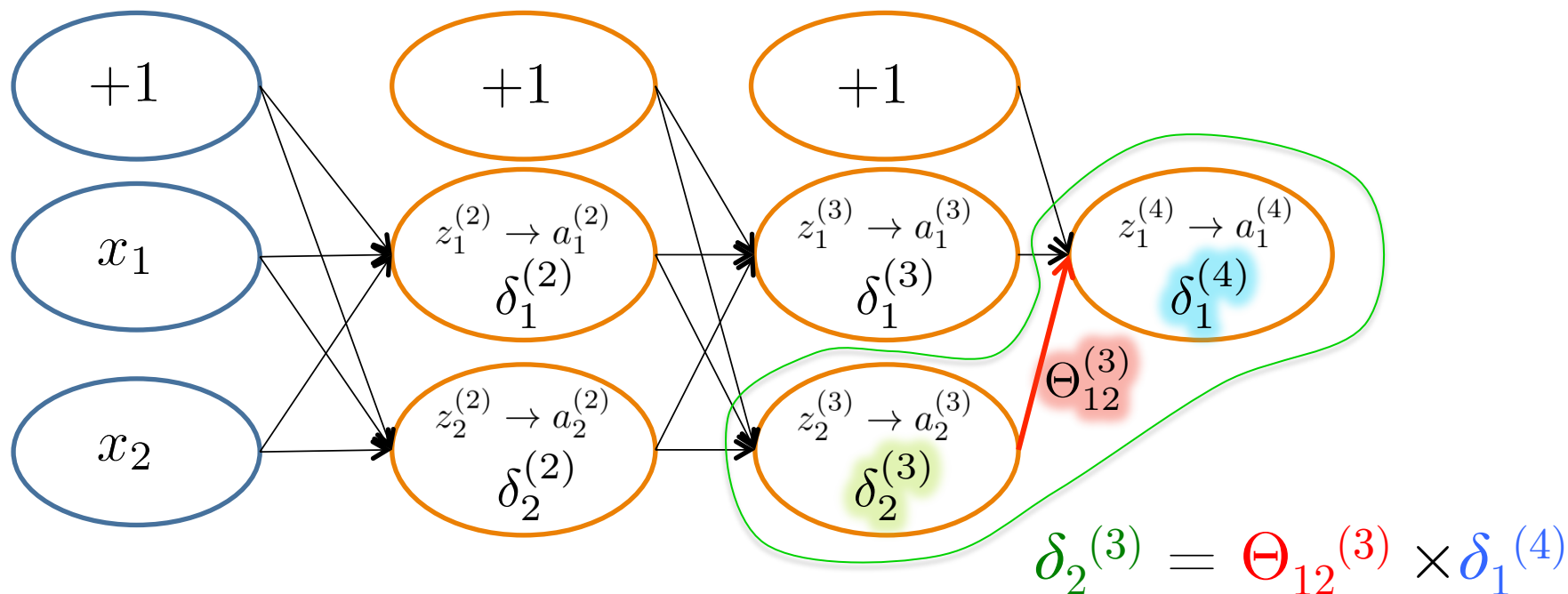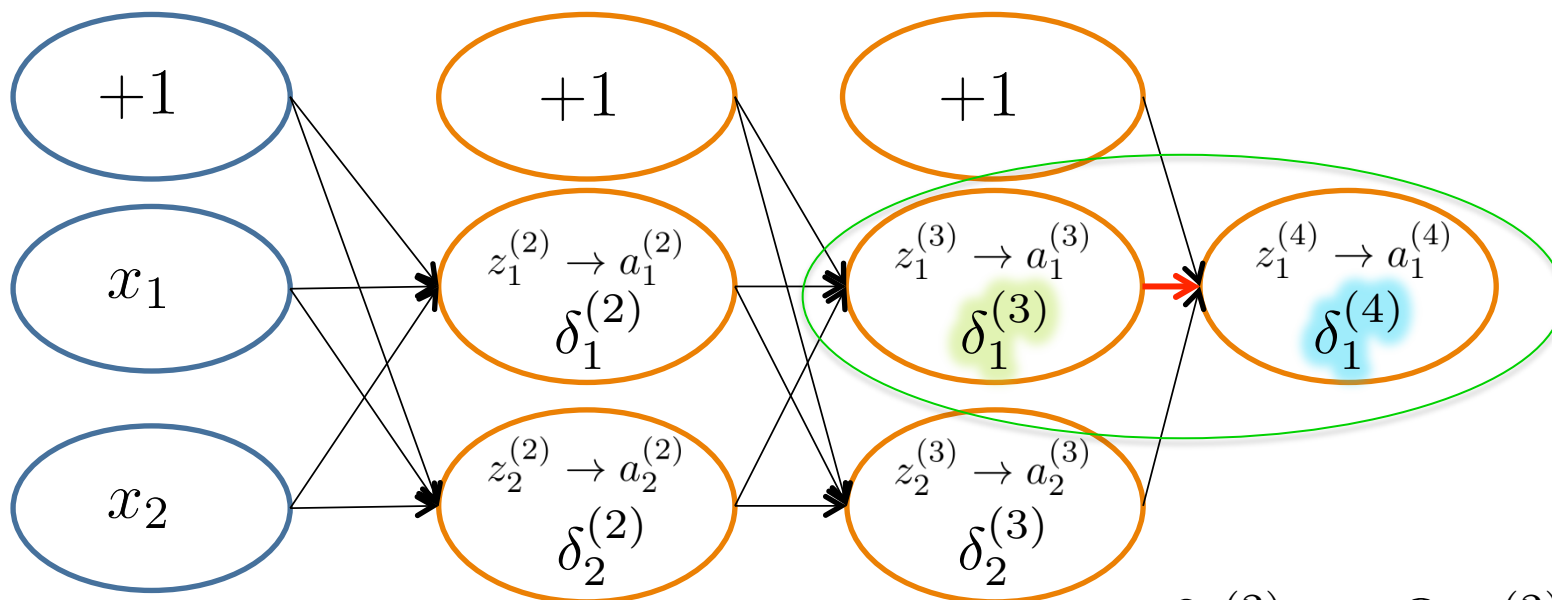
# Backpropagation Intuition



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

$$\delta_1{}^{(3)} = \Theta_{11}{}^{(3)} \times \delta_1{}^{(4)}$$

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

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

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

# Backpropagation Intuition


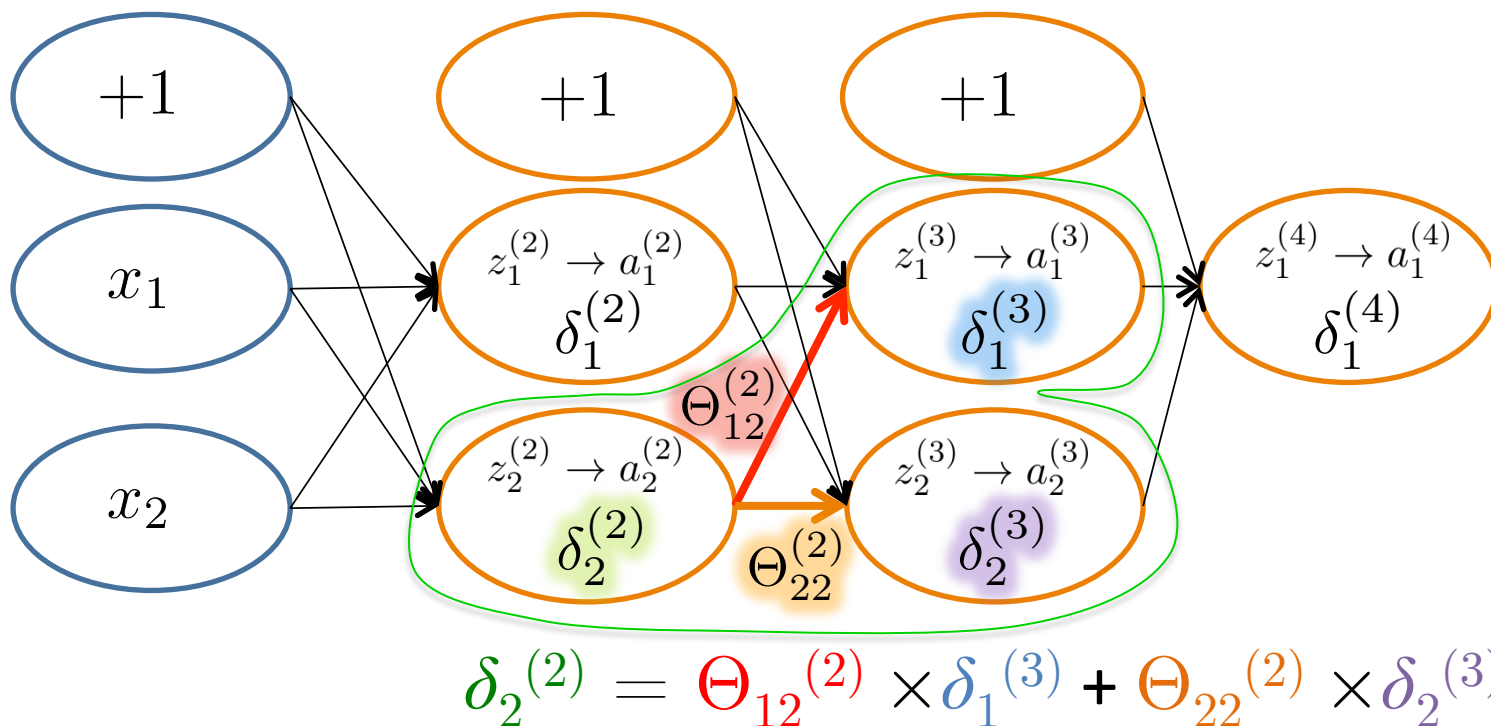
$$\delta_2^{(2)} = \Theta_{12}^{(2)} \times \delta_1^{(3)} + \Theta_{22}^{(2)} \times \delta_2^{(3)}$$

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

# Backpropagation: Gradient Computation

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

(#layers $L$ = 4)

## Backpropagation

Element-wise product .*

- $\boldsymbol{\delta}^{(4)} = \boldsymbol{a}^{(4)} - \mathbf{y}$
- $\boldsymbol{\delta}^{(3)} = (\Theta^{(3)})^\mathsf{T}\boldsymbol{\delta}^{(4)} \ .\!* \ g'(\mathbf{z}^{(3)})$
- $\boldsymbol{\delta}^{(2)} = (\Theta^{(2)})^\mathsf{T}\boldsymbol{\delta}^{(3)} \ .\!* \ g'(\mathbf{z}^{(2)})$
- $(\text{No } \boldsymbol{\delta}^{(1)})$

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

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

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

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

$\boldsymbol{\delta}^{(2)}$  $\boldsymbol{\delta}^{(3)}$  $\boldsymbol{\delta}^{(4)}$

# Backpropagation

Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$  (Used to accumulate gradient)

For each training instance $(\mathbf{x}_i, y_i)$:

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

Compute $\{\mathbf{a}^{(2)}, \ldots, \mathbf{a}^{(L)}\}$ via forward propagation

Compute $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y_i$

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 $D_{ij}^{(l)} = \begin{cases} \frac{1}{n}\Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n}\Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

$\boldsymbol{D}^{(l)}$ is the matrix of partial derivatives of $J(\Theta)$

Note:  Can vectorize $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ as $\boldsymbol{\Delta}^{(l)} = \boldsymbol{\Delta}^{(l)} + \boldsymbol{\delta}^{(l+1)} \mathbf{a}^{(l)\mathsf{T}}$

# 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 $(\mathbf{x}_i, y_i)$:

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

        Compute $\{\mathbf{a}^{(2)}, \ldots, \mathbf{a}^{(L)}\}$ via forward propagation

        Compute $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y_i$

        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**

# Backprop Issues

"Backprop is the cockroach of machine learning.  It's ugly, and annoying, but you just can't get rid of it."
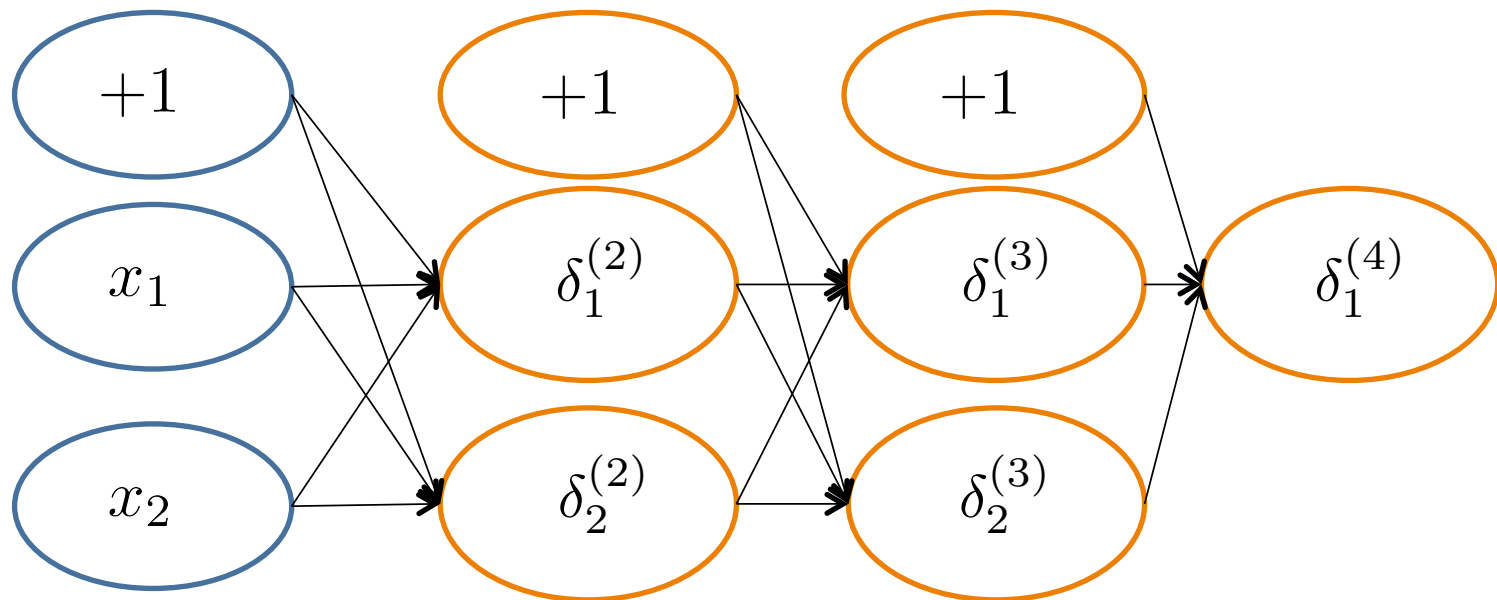            -Geoff Hinton


Problems:

- black box

- local minima

# Implementation Details

# Random Initialization

- Important to randomize initial weight matrices
- Can't have uniform initial weights, as in logistic regression
  - Otherwise, all updates will be identical & the net won't learn

# Implementation Steps

- Implement backprop to compute `DVec`
  - `DVec` is the unrolled $\{D^{(1)}, D^{(2)}, \dots\}$ matrices
- Implement numerical gradient checking to compute `gradApprox`
- Make sure `DVec` has similar values to `gradApprox`
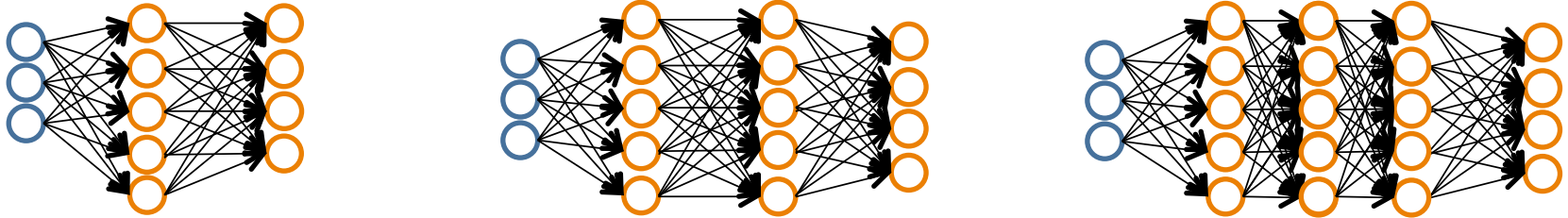- <u>Turn off gradient checking</u>. Using backprop code for learning.

**Important:** Be sure to disable your gradient checking code before training your classifier.

- If you run the numerical gradient computation on every iteration of gradient descent, your code will be <u>very</u> slow

# Putting It All Together

# Training a Neural Network

Pick a network architecture (connectivity pattern between nodes)



- # input units = # of features in dataset
- # output units = # classes

**Reasonable default:** 1 hidden layer

- or if >1 hidden layer, have same # hidden units in every layer (usually the more the better)

# Training a Neural Network

1. Randomly initialize weights

2. Implement forward propagation to get $h_\Theta(\mathbf{x}_i)$ for any instance $\mathbf{x}_i$

3. Implement code to compute cost function $J(\Theta)$

4. Implement backprop to compute partial derivatives $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$

5. Use gradient checking to compare $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ computed using backpropagation vs. the numerical gradient estimate.

   – Then, disable gradient checking code

6. Use gradient descent with backprop to fit the network