

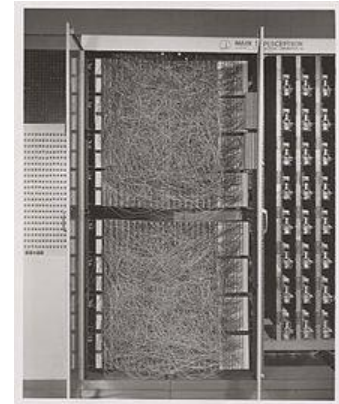
The background of the slide features a complex, glowing network of interconnected nodes and lines, resembling a neural network or a data mesh, set against a dark blue gradient.

Introduction to Neural Networks

Perceptron, XOR Problem, Multi-layer Perceptron (MLP),
Cost Functions, Activation functions and Output units

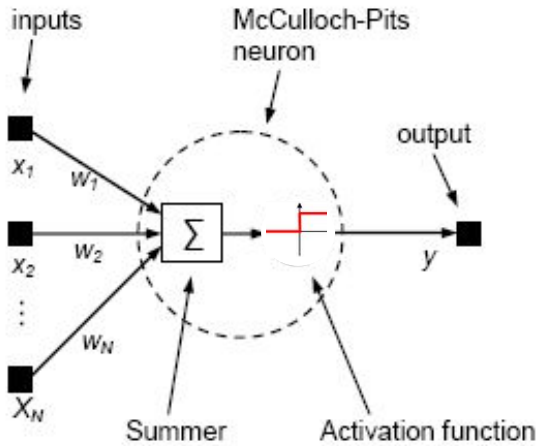
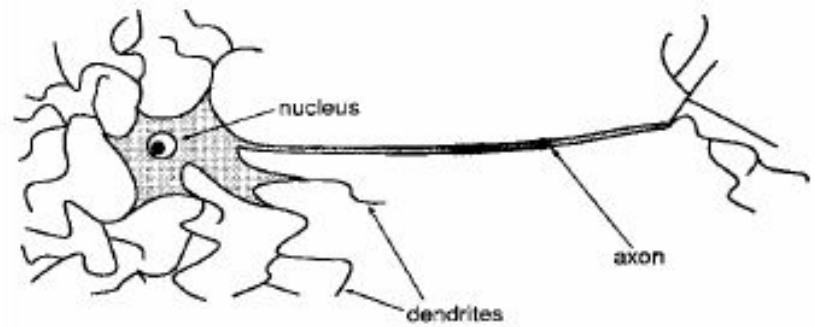
History

- 1932 - *The Integrative Action of the Nervous System*, **Sir Charles Scott Sherrington**
 - Nervous system - interconnection of individual entities (neuron)
- 1943 - *A Logical Calculus of Ideas Immanent in Nervous Activity*, **Warren McCulloch** and **Walter Pitts**
 - McCulloch and Pitt's model
- 1949 - *The Organization of Behavior*, **Donald Hebb**
 - Hebbian learning
- 1953 - *The Perceptron*, **Rosenblatt**
- 1969 - *Limitation of Perceptrons*, **Minsky** and **Papert**
- 1980's - *Connectionism*
 - *Error Back Propagation*, **Proposed simultaneously by Many**
- 20th century Deep learning
 - CNNs - LeNet, AlexNet, VGGNet, ResNet
 - Deep LSTMs,
- The **deep saga**..... what followed is discussed in the last class

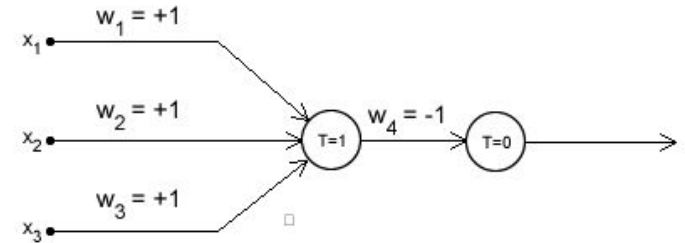
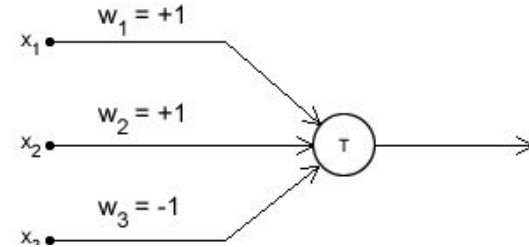


The Mark 1 Perceptron
By Rosenblatt

McCulloch - Pits model

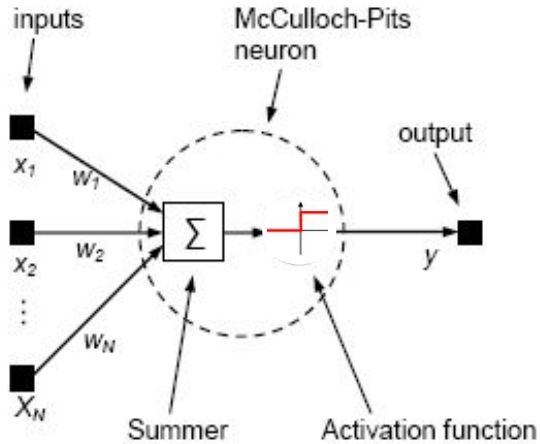


$$\sum_{i=1}^n w_i x_i > \mu$$

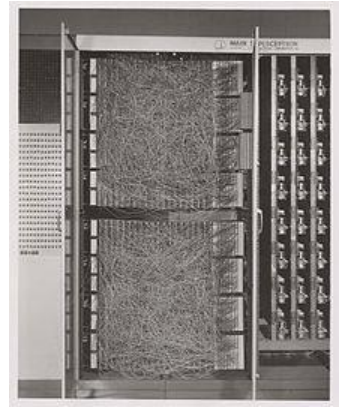


NOR gate

The Perceptron - Rosenblatt (1953)



$$\sum_{i=1}^n w_i x_i > \mu$$



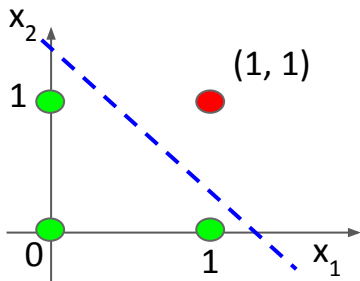
*Pic courtesy, wikipedia

The Mark 1 Perceptron
By Rosenblat
for digit recognition

Perceptron - geometrical interpretation

$$\sum_{i=1}^n w_i x_i > \mu \quad , \text{What does this inequality imply in 2D case? } \text{Half plane}$$

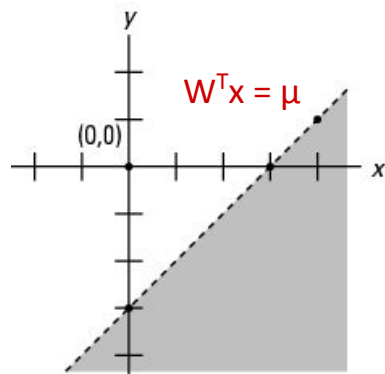
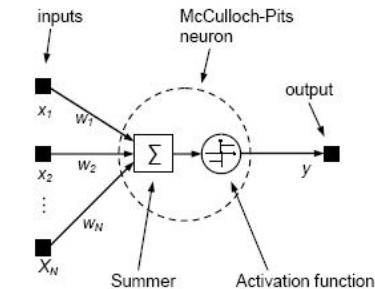
\mathbb{X}	AND
(0, 0)	0
(0, 1)	0
(1, 0)	0
(1, 1)	1



Solve for W, μ :

$$x_1 + x_2 > 1.5$$

$$w_1 = 1, w_2 = 1 \text{ and } \mu = 1.5$$



*Pic courtesy, [cliffsnotes](#)

Any function that is linearly separable can be computed by a perceptron

Perceptron Learning Algorithm (PLA)

How to learn the parameters $\{W\}$ given only the input, output pairs

$$w_i \leftarrow w_i + \Delta w_i$$

where

$$\Delta w_i = \eta(t - o)x_i$$

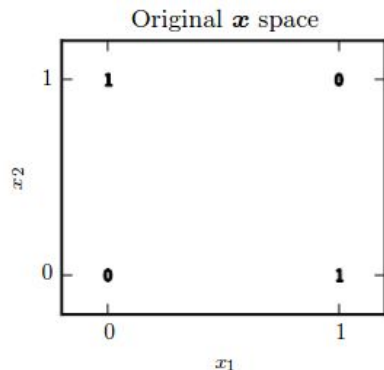
Where:

- $t = c(\vec{x})$ is target value
- o is perceptron output
- η is small constant (e.g., 0.1) called *learning rate*

Perceptron - Limitations

Goal: learn the XOR function (f^*)

\mathbf{x}	f^*
(0, 0)	0
(0, 1)	1
(1, 0)	1
(1, 1)	0



Task is adjust parameters θ , such that f is as close as to f^*

$$y = f(\mathbf{x}, \theta) \quad L = \sum_{\{\mathbf{x} \in \mathbf{X}\}} (f^*(\mathbf{x}) - f(\mathbf{x}, \theta))^2$$

Lets use our perceptron for f , $\theta = \{w, b\}$

$$f(\mathbf{x}; w, b) = w^T \mathbf{x} + b$$

Solve for $\{w, b\}$

$W = 0$, $b = 0.5$; output is 0.5 everywhere

Why this linear function can't model XOR?

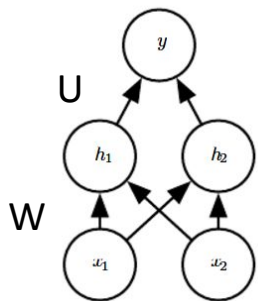
How to tackle this problem?

- Can we use more than one line?
- Yes, but **how**?

Perceptron - Limitations

How to tackle this problem?

- Add a hidden layer with two units



$$y = f^{(2)}(h; U, c)$$

$$y = f^{(2)}(f^{(1)}(x))$$

$$h = f^{(1)}(x; W, b)$$

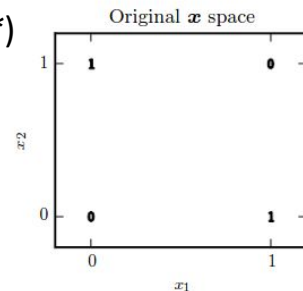
What should $f^{(1)}$ compute?

- If its linear again the composition still remains linear

$$f^{(2)}(h) = U^T h; \text{ since } h = Wx$$

$$y = U^T Wx = W'x$$

Goal: learn the XOR function (f^*)



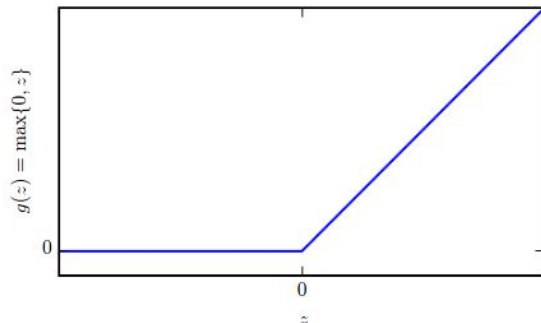
- $f^{(1)}$ should be nonlinear to extract useful features

$$h = f^{(1)}(x; W, b) = g(Wx+b)$$

- g is referred as **activation** function commonly

- We will use ReLU here

- ❑ Rectified Linear Unit (widely used)
- ❑ $g(z) = \max\{0, z\}$

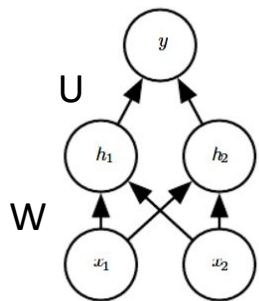


*Slide courtesy, Ian Goodfellow et al., deep learning book

Perceptron - Limitations

How to tackle this problem?

- Add a hidden layer with two units
- Use ReLU activation in 1st layer



$$y = U^T h + c; \quad y = U^T \underbrace{\max\{0, Wx+b\}}_{\text{ReLU}} + c$$

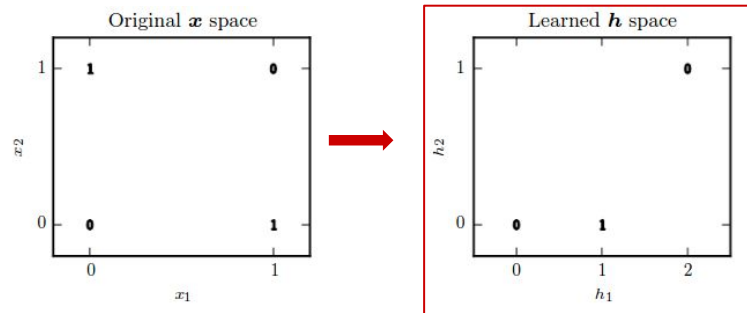
$$h = g(Wx+b)$$

Let,

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad U = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

$$c = 0$$

Goal: learn the XOR function (f^*)



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

$$WX = \begin{bmatrix} 0 & 1 & 1 & 2 \\ 0 & 1 & 1 & 2 \end{bmatrix}$$

$$WX + b = \begin{bmatrix} 0 & 1 & 1 & 2 \\ -1 & 0 & 0 & 1 \end{bmatrix}$$

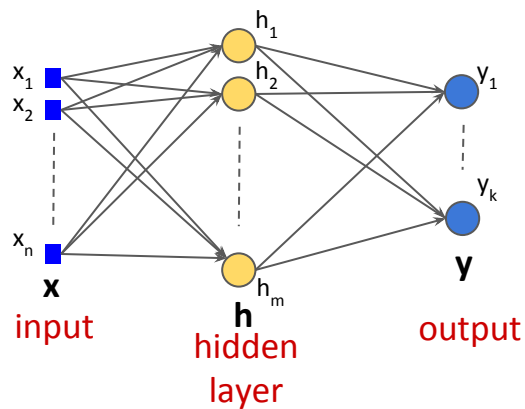
$$h = \begin{bmatrix} 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \text{Upon ReLU}$$

After layer2

$$[0 \ 1 \ 1 \ 0]$$

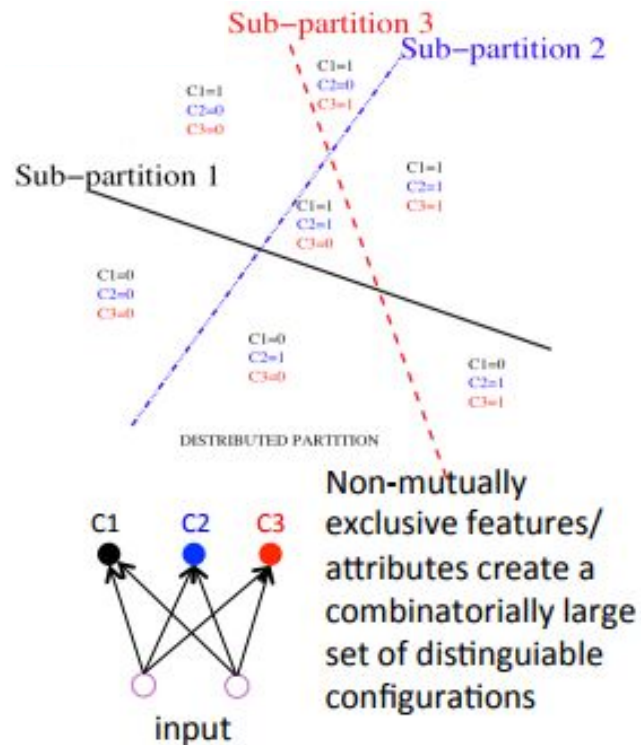
Multi-layer Perceptrons (MLP)

A typical feed forward neural network



$$\mathbf{h} = f(\mathbf{W}\mathbf{x} + \mathbf{b}_1); \quad \mathbf{y} = g(\mathbf{U}\mathbf{h} + \mathbf{b}_2)$$

With more hidden units network is more expressible



Feedforward Neural Networks - Cost functions

For regression,

$$J(\theta) = \frac{1}{2} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \|\mathbf{y} - f(\mathbf{x}; \theta)\|^2$$

$$\frac{1}{2} \sum_{\{x_i, y_i\}} \|y_i - f(x_i, \theta)\|^2$$

For classification,

- Typically outputs a probability vector $q(c = k | x) \forall k$
- How do you compare two distributions?
 - KL divergence, $KL(p \parallel q)$

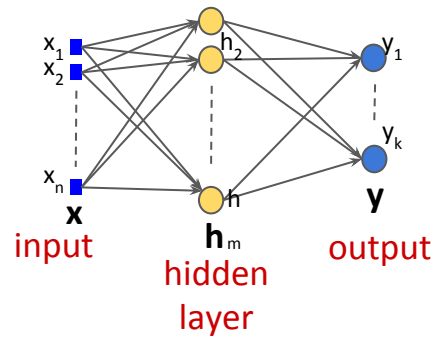
$$D_{KL}(p(x) \parallel q(x)) = \sum_{x \in X} p(x) \ln \frac{p(x)}{q(x)}$$

$$= \sum p(x) \ln p(x) - p(x) \ln q(x)$$

$$= \underbrace{-H(p)}_{\text{entropy}} + \underbrace{H(p, q)}_{\text{cross-entropy}}$$

$$J(\theta) = \sum_{x_i, y_i} H(p(x_i), q(x_i))$$

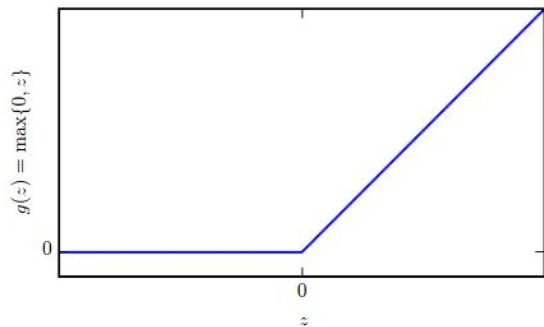
Activation functions



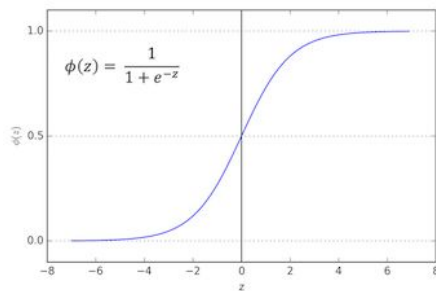
$h = g(Wx+b)$; Affine transformation followed by activation function, g

- Very important factor in learning features

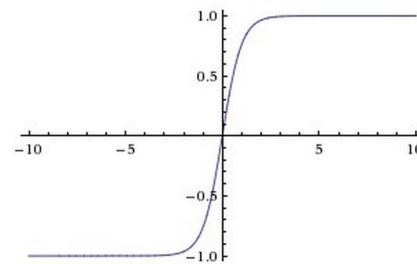
$$g(z) = \max\{0, z\}, \text{ ReLU}$$



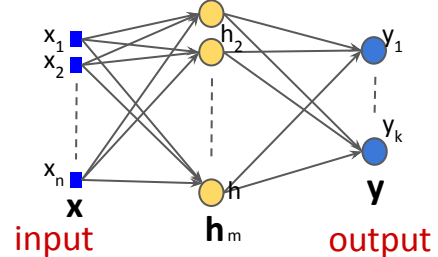
$$g(z) = \sigma(z), \text{ sigmoid}$$



$$\tanh(z) = 2\sigma(2z) - 1$$



Output units



- Linear units for real valued outputs
 - ❑ Activation function is left to be linear
 - ❑ Given features \mathbf{h} ,
$$\mathbf{y}' = \mathbf{W}\mathbf{h} + \mathbf{b}$$
 - ❑ Most commonly used with regression tasks

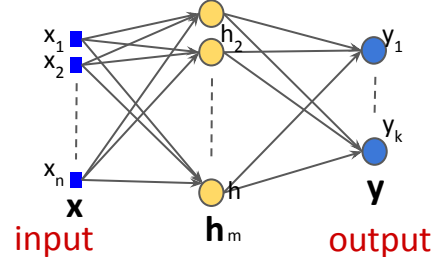
- Say you want to do binary classification
 - ❑ What kind of distribution describes output?
Bernouli
 - ❑ How to constrain the output - valid probability?
Can you use linear activation?

$$P(y = 1 \mid \mathbf{x}) = \max \left\{ 0, \min \left\{ 1, \mathbf{w}^\top \mathbf{h} + b \right\} \right\}.$$

- ❑ What is the problem? *Not amenable for gradient based learning*
- ❑ Instead, use sigmoid unit - output $\in [0,1]$

$$\hat{y} = \sigma(\mathbf{w}^\top \mathbf{h} + b)$$

Output units



- Now, say we want to do multi-class classification (K classes)

- ❑ Output should be K probabilities,
 $p_k = p(\text{class} = k | \mathbf{x}) \quad \forall k = 1 \text{ to } K$

- ❑ Can we use K sigmoid units?

Won't be sufficient, since probabilities are not constrained to sum to 1

$$\sum_k p_k = 1$$

- ❑ We will look at **softmax** unit for this

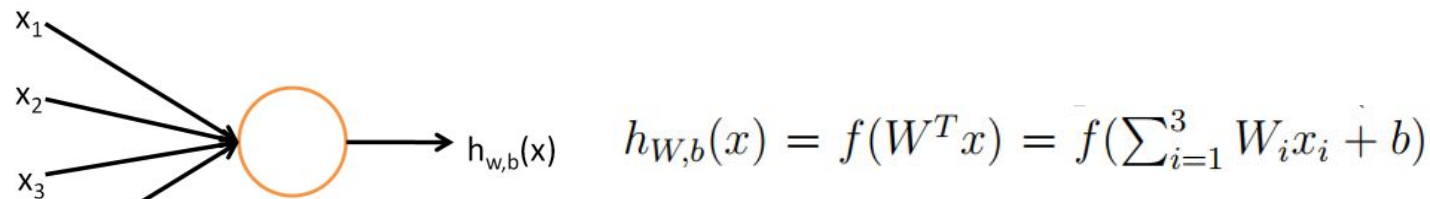
Idea is to convert a vector of real values to valid probabilities, **How?**

- ❑ Make all the elements positive
 - ❑ Normalize the values

- Let, $\mathbf{z} = [z_1, \dots, z_K]^T$; $\mathbf{z} = \mathbf{W}\mathbf{h} + \mathbf{b}$

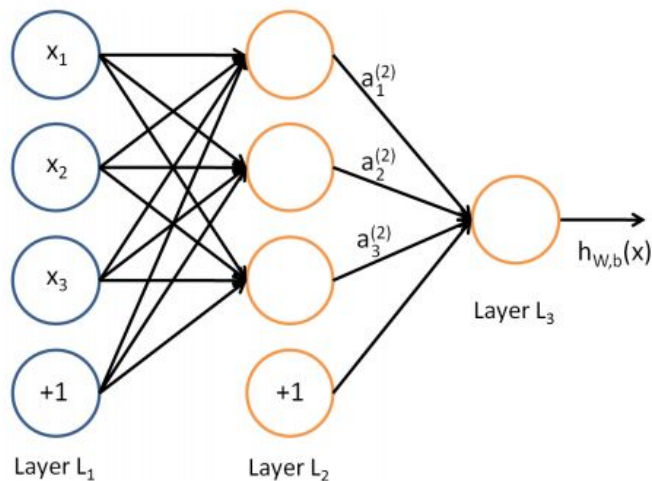
$$\text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}.$$

How to learn the parameters?



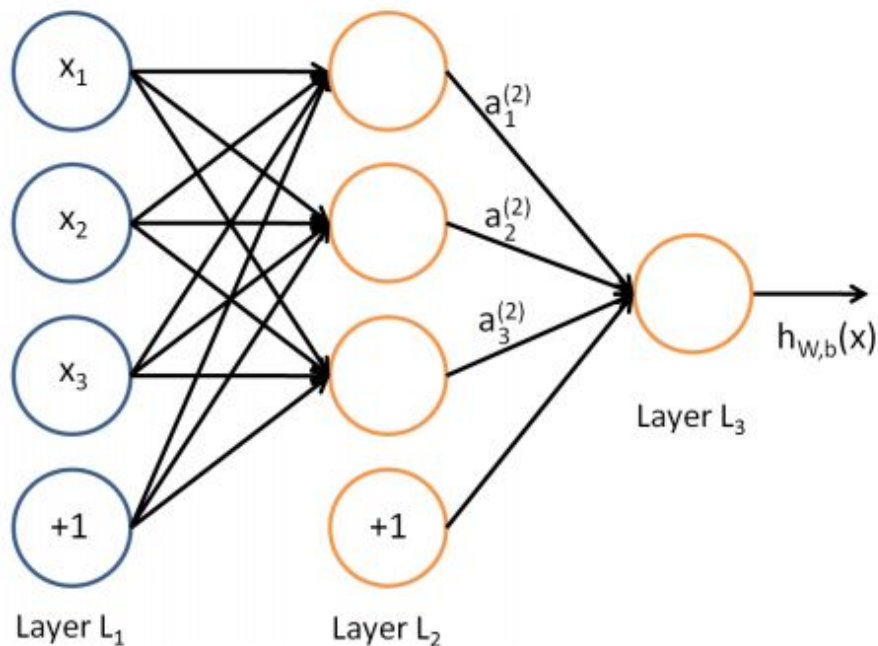
Bias

Typical NN with many such units



- One hidden layer
 - 3 neuron units
- One output

How to learn the parameters?



L_l – Layer l

$a_i^{(l)}$ – activation of unit i in layer l

$W_{ij}^{(l)}$ – Weight from j^{th} unit in l to i^{th} unit in $l+1$

$b_i^{(l)}$ – bias to unit i in layer $l+1$

Parameters:

$$(W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)})$$

$$W^{(1)} \in \mathbb{R}^{3 \times 3}, W^{(2)} \in \mathbb{R}^{1 \times 3}$$

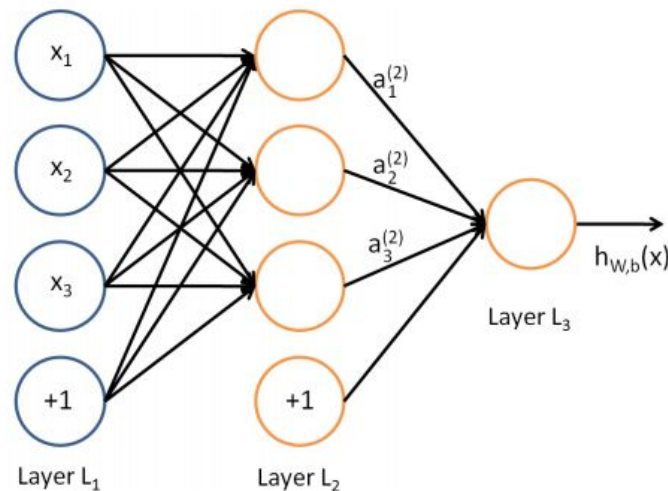
How to learn the parameters?

Layer 2,

$$a_1^{(2)} = f(W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)})$$

$$a_2^{(2)} = f(W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)})$$

$$a_3^{(2)} = f(W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)})$$



Layer 3,

$$h_{W,b}(x) = a_1^{(3)} = f(W_{11}^{(2)}a_1^{(2)} + W_{12}^{(2)}a_2^{(2)} + W_{13}^{(2)}a_3^{(2)} + b_1^{(2)})$$

Simplification

Let, $z_i^{(l)}$ denote weighted sum for the activation $a_i^{(l)}$

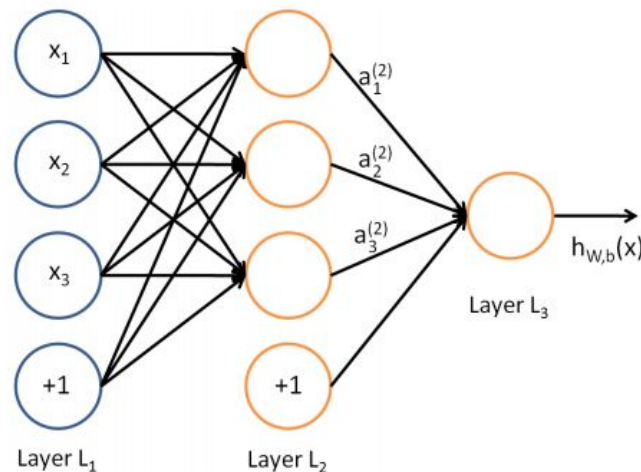
$$a_i^{(l)} = f(z_i^{(l)}) \quad f(.) \text{ applies the function point wise}$$

How to learn the parameters?

$$a_1^{(2)} = f(W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)})$$

$$a_2^{(2)} = f(W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)})$$

$$a_3^{(2)} = f(W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)})$$



$$h_{W,b}(x) = a_1^{(3)} = f(W_{11}^{(2)}a_1^{(2)} + W_{12}^{(2)}a_2^{(2)} + W_{13}^{(2)}a_3^{(2)} + b_1^{(2)})$$

Let, $z_i^{(l)}$ denote weighted sum for the activation $a_i^{(l)}$

$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$

$$h_{W,b}(x) = a^{(3)} = f(z^{(3)})$$

$$z^{(l+1)} = W^{(l)}a^{(l)} + b^{(l)}$$

$$a^{(l+1)} = f(z^{(l+1)})$$

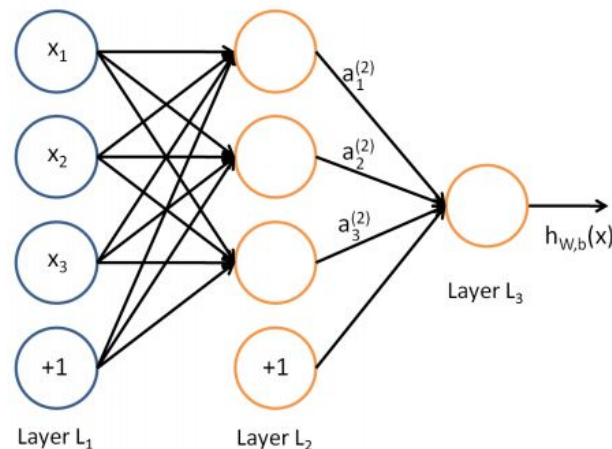
How to learn the parameters?

Given m training examples

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

Minimize:

$$J(W, b; x, y) = \frac{1}{2} \|h_{W,b}(x) - y\|^2$$



$$\begin{aligned} J(W, b) &= \left[\frac{1}{m} \sum_{i=1}^m J(W, b; x^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left(W_{ji}^{(l)} \right)^2 \quad (8) \\ &= \left[\frac{1}{m} \sum_{i=1}^m \left(\frac{1}{2} \|h_{W,b}(x^{(i)}) - y^{(i)}\|^2 \right) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left(W_{ji}^{(l)} \right)^2 \end{aligned}$$

How to learn the parameters?

Minimize:

$$J(W, b; x, y) = \frac{1}{2} \|h_{W,b}(x) - y\|^2$$

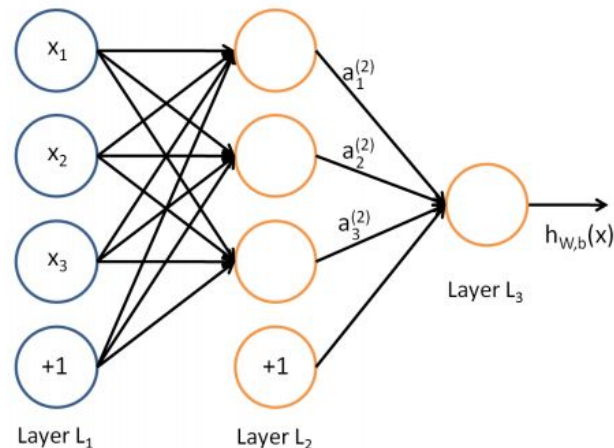
Gradient descent:

$$W_{ij}^{(l)} := W_{ij}^{(l)} - \alpha \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b)$$

$$b_i^{(l)} := b_i^{(l)} - \alpha \frac{\partial}{\partial b_i^{(l)}} J(W, b)$$

How to evaluate these
partial derivatives?

Error back-propagation



Back-propagation algorithm

Gradient descent:

$$W_{ij}^{(l)} := W_{ij}^{(l)} - \alpha \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b)$$

Idea:

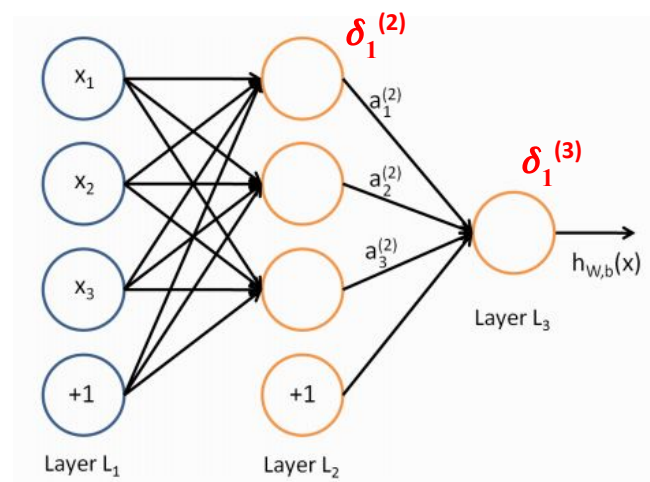
First, forward pass the data to calc. all responses

In backward pass, for each unit i in layer l calculate **error** term $\delta_i^{(l)}$ - measures how much unit i is responsible for output error

- For output unit in last layer (n_l), this is easy

$$\delta_i^{(n_l)} = \frac{\partial}{\partial z_i^{(n_l)}} \frac{1}{2} \|y - h_{W,b}(x)\|^2 = -(y_i - a_i^{(n_l)}) \cdot f'(z_i^{(n_l)})$$

- How to measure $\delta_i^{(l)}$ for hidden units?



Back-propagation algorithm

Gradient descent:

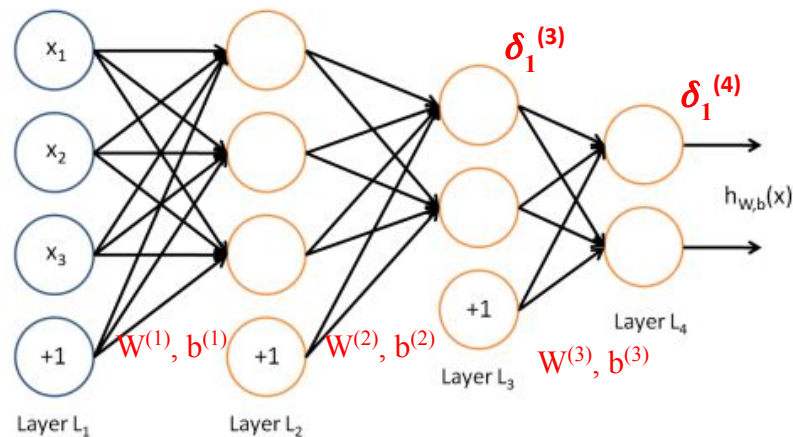
$$W_{ij}^{(l)} := W_{ij}^{(l)} - \alpha \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b)$$

$$J(W, b; x, y) = \frac{1}{2} \|h_{W,b}(x) - y\|^2$$

For last layer:

$$\frac{\partial J}{\partial W_{ij}^{(3)}} = \boxed{\frac{\partial J}{\partial z_i^4}} \boxed{\frac{\partial z_i^4}{\partial W_{ij}^{(3)}}}$$

$$\frac{\partial J}{\partial W_{ij}^{(l)}} = \delta_i^{(l+1)} a_j^{(l)} \quad \frac{\partial J}{\partial b_i^{(l)}} = \delta_i^{(l+1)}$$



$$h_{W,b}(x) = a^{(4)} = f(z^{(4)}); \quad z^{(4)} = W^{(3)}a^{(3)} + b^{(3)}$$

$$\boxed{\frac{\partial J}{\partial z_i^4} = -(y_i - a_i^{(4)}) \cdot f'(z_i^4)}$$

$\delta_i^{(4)}$ error term

$$\boxed{\frac{\partial z_i^4}{\partial W_{ij}^{(3)}} = a_j^{(3)}}$$

Back-propagation algorithm

Gradient descent:

$$W_{ij}^{(l)} := W_{ij}^{(l)} - \alpha \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b)$$

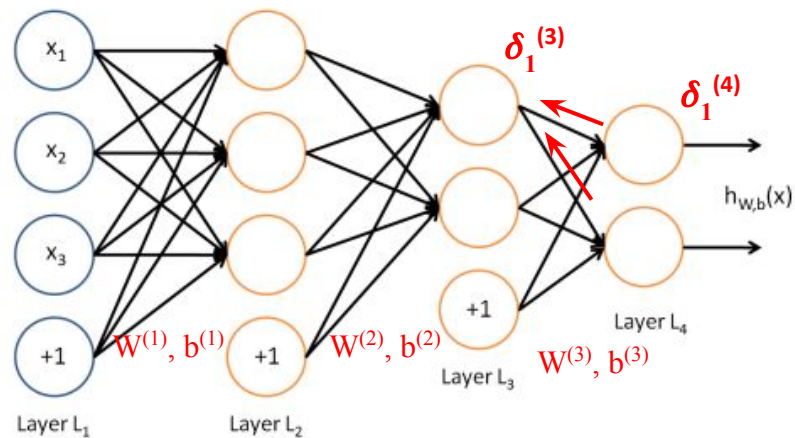
$$J(W, b; x, y) = \frac{1}{2} \|h_{W,b}(x) - y\|^2$$

For layers other than last:

$$\frac{\partial J}{\partial W_{ij}^{(2)}} = \frac{\partial J}{\partial z_i^{(3)}} \frac{\partial z_i^{(3)}}{\partial W_{ij}^{(2)}} a_j^{(2)}$$

$$\delta_i^{(l)} = \left(\sum_{j=1}^{s_{l+1}} W_{ji}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$$

$$\frac{\partial J}{\partial W_{ij}^{(l)}} = \delta_i^{(l+1)} a_j^{(l)} \quad \frac{\partial J}{\partial b_i^{(l)}} = \delta_i^{(l+1)}$$



$$h_{W,b}(x) = a^{(4)} = f(z^{(4)}); \quad z^{(4)} = W^{(3)} a^{(3)} + b^{(3)}$$

$$a^{(3)} = f(z^{(3)}); \quad z^{(3)} = W^{(2)} a^{(2)} + b^{(2)}$$

$\delta_i^{(3)}$
error term

$$\begin{aligned} \frac{\partial J}{\partial z_i^{(3)}} &= \frac{\partial J}{\partial a_i^{(3)}} \frac{\partial a_i^{(3)}}{\partial z_i^{(3)}} \\ &= \left(\sum_j \frac{\partial J}{\partial z_j^{(4)}} \frac{\partial z_j^{(4)}}{\partial a_i^{(3)}} \right) f'(z_i^{(3)}) \end{aligned}$$

$\delta_j^{(4)}$
Layer - (l+1) $W_{ji}^{(3)}$

Back-propagation algorithm

1. Perform a feedforward pass
 - Computing activations L_1 , L_2 and so on ...

2. For each output unit i in layer L_4 (output layer), set

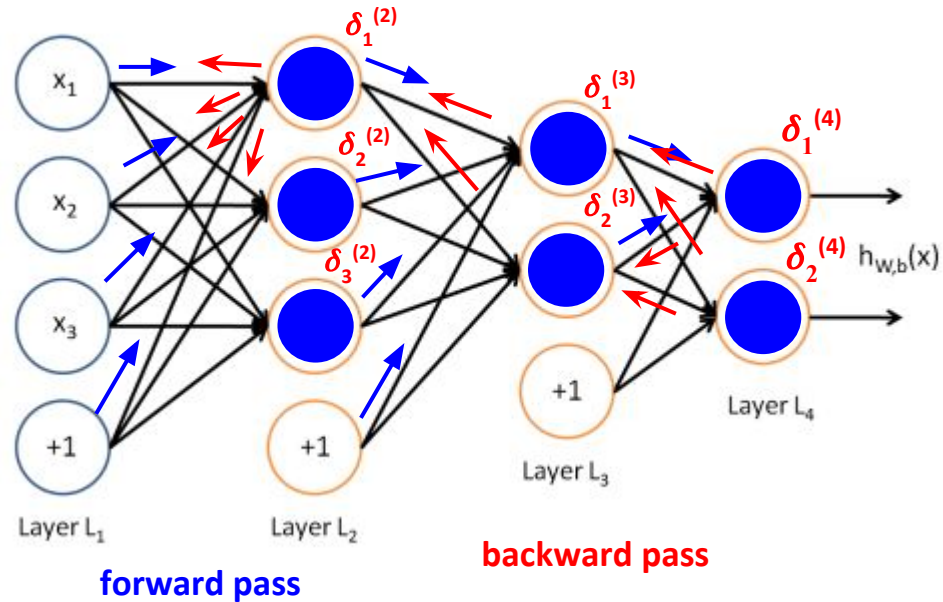
$$\delta_i^{(n_l)} = \frac{\partial}{\partial z_i^{(n_l)}} \frac{1}{2} \|y - h_{W,b}(x)\|^2 = -(y_i - a_i^{(n_l)}) \cdot f'(z_i^{(n_l)})$$

3. Starting from last but one layer to 2nd layer;
 $l = n_l - 1, n_l - 2, \dots, 2$
 - For each node i in layer l , set

$$\delta_i^{(l)} = \left(\sum_{j=1}^{s_{l+1}} W_{ji}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$$

4. Compute the desired partial derivatives, as:

$$\frac{\partial}{\partial W_{ij}^{(l)}} J(W, b; x, y) = a_j^{(l)} \delta_i^{(l+1)} \quad \frac{\partial}{\partial b_i^{(l)}} J(W, b; x, y) = \delta_i^{(l+1)}.$$



Back-propagation algorithm

Gradient descent:

$$W_{ij}^{(l)} := W_{ij}^{(l)} - \alpha \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b)$$

$$J(W, b; x, y) = \frac{1}{2} \|h_{W,b}(x) - y\|^2$$

Partial derivatives:

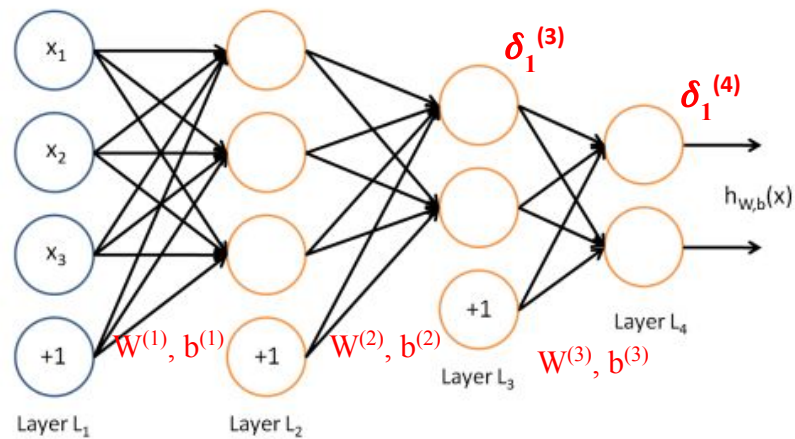
$$\delta_i^{(l)} = \left(\sum_{j=1}^{s_{l+1}} W_{ji}^{(l)} \delta_j^{(l+1)} \right) f'(z_i^{(l)})$$

$$\frac{\partial J}{\partial W_{ij}^{(l)}} = \delta_i^{(l+1)} a_j^{(l)} \quad \frac{\partial J}{\partial b_i^{(l)}} = \delta_i^{(l+1)}$$

Matrix notation:

$$\delta^{(l)} = ((W^{(l)})^T \delta^{(l+1)}) \bullet f'(z^{(l)})$$

$$\frac{\partial J}{\partial W^{(l)}} = \delta^{(l+1)} (a^{(l)})^T \quad \frac{\partial J}{\partial b^{(l)}} = \delta^{(l+1)}$$



$$h_{W,b}(x) = a^{(4)} = f(z^{(4)}); \quad z^{(4)} = W^{(3)}a^{(3)} + b^{(3)}$$

Back-propagation algorithm

1. Perform a feedforward pass
 - Computing activations L_1 , L_2 and so on ...

2. For each output unit i in layer L_4 (output layer), set

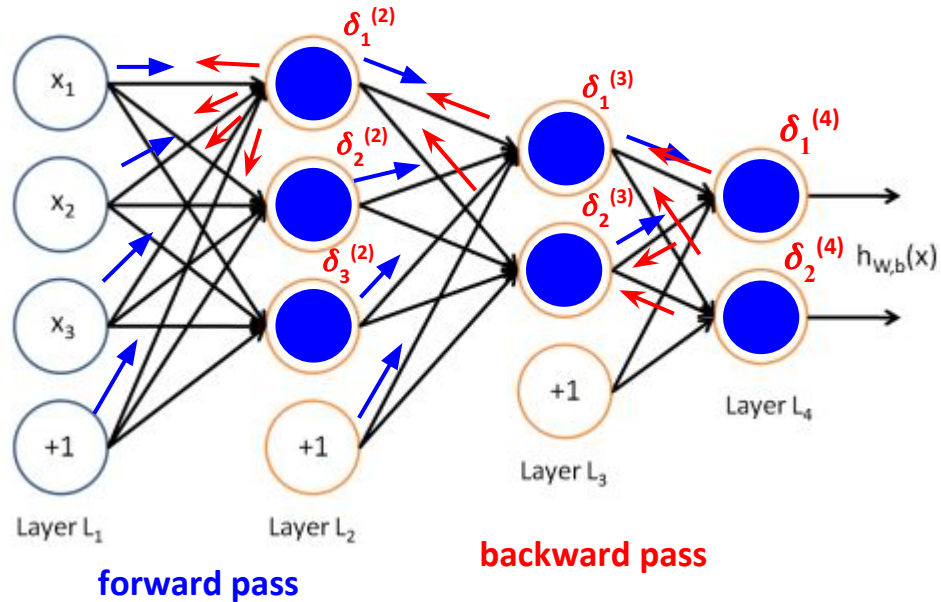
$$\delta^{(n_i)} = -(y - a^{(n_i)}) \bullet f'(z^{(n)})$$

3. Starting from last but one layer to 2nd layer;
 $l = n_l - 1, n_l - 2, \dots, 2$

$$\delta^{(l)} = ((W^{(l+1)})^T \delta^{(l+1)}) \bullet f'(z^{(l)})$$

4. Compute the desired partial derivatives, as:

$$\begin{aligned}\nabla_{W^{(l)}} J(W, b; x, y) &= \delta^{(l+1)} (a^{(l)})^T, \\ \nabla_{b^{(l)}} J(W, b; x, y) &= \delta^{(l+1)}.\end{aligned}$$



END