

# **DASC7606– 3A**

# **Deep Learning**

## **Training Neural Networks**

**Dr Bethany Chan**  
**Professor Francis Chin**  
**2024**

# Revision - More on Logistic Regression

- Metrics for accuracy of predictions:
  - Precision (accurate), Recall (positive), Specificity (reject)
  - F1 (Harmonic mean)
  - Receiver Operating Characteristic ROC (AUC)

- Binary classification (sigmoid  $\frac{1}{1+e^{-s}}$ )  
with binary cross-entropy cost function

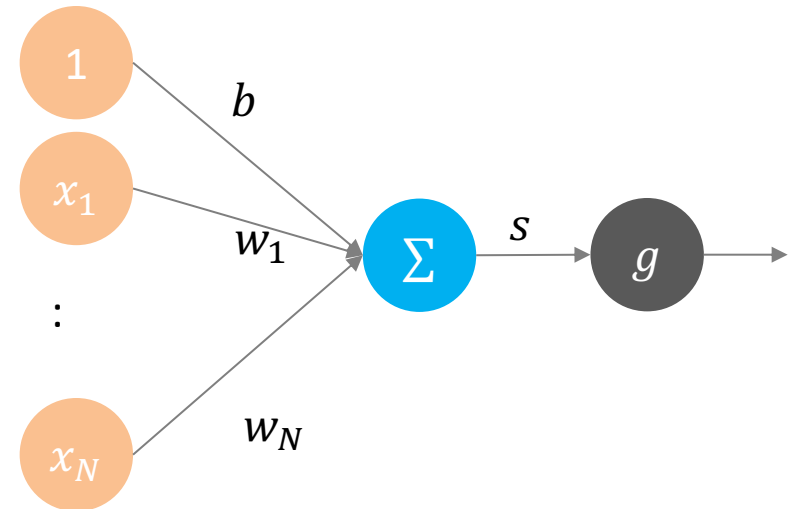
$$J(\boldsymbol{\theta}) = -\frac{1}{M} \sum_{i=1}^M [y_i \log h(\mathbf{x}_i) + (1 - y_i) \log(1 - h(\mathbf{x}_i))]$$

- Multi-classification -  $K$  classes (softmax  $\frac{e^{s_k}}{\sum e^{s_i}}$ )  
with categorical cross-entropy cost function

$$J(\boldsymbol{\theta}) = -\frac{1}{M} \sum_{i=1}^M \sum_K [y_i \log \mathbf{h}(\mathbf{x}_i)] = -\frac{1}{M} \sum_{i=1}^M \sum_{k=1}^K [y_{i,k} \log h_k(\mathbf{x}_i)]$$

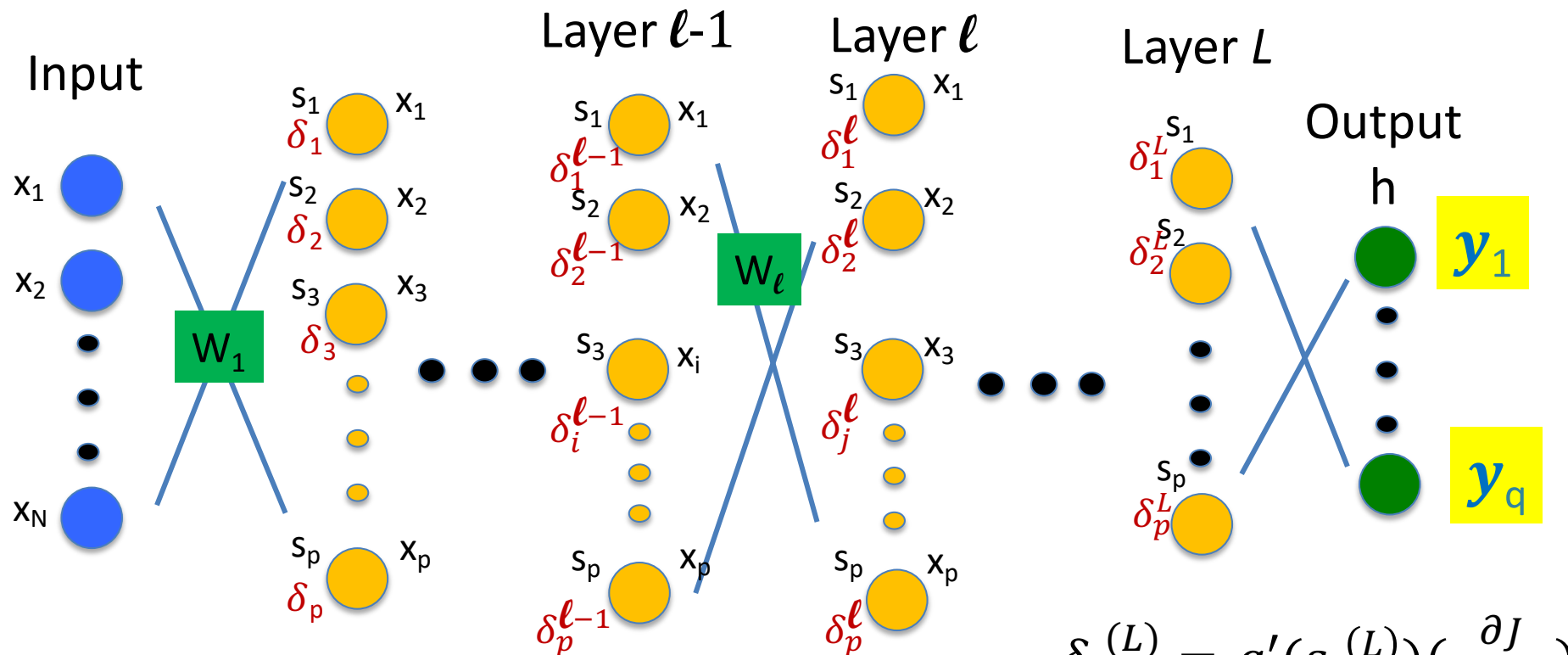
# Revision: Perceptron and Neural Network

- Logistic Regression  
≈ **Perceptron with Sigmoid**  
as the activation function
- Perceptron with **nonlinear**  
activation functions  $g$   
(Sigmoid, TanH, ReLu, Step)  
as the building block of  
neural networks
- **Any logic functions,  
classification and  
Regression problems  
can be approximated**  
by the Neural Network  
with sufficient complexity,  
i.e., number of layers and  
number of neurons



Perceptron

# Revision: Feedforward and Backpropagation

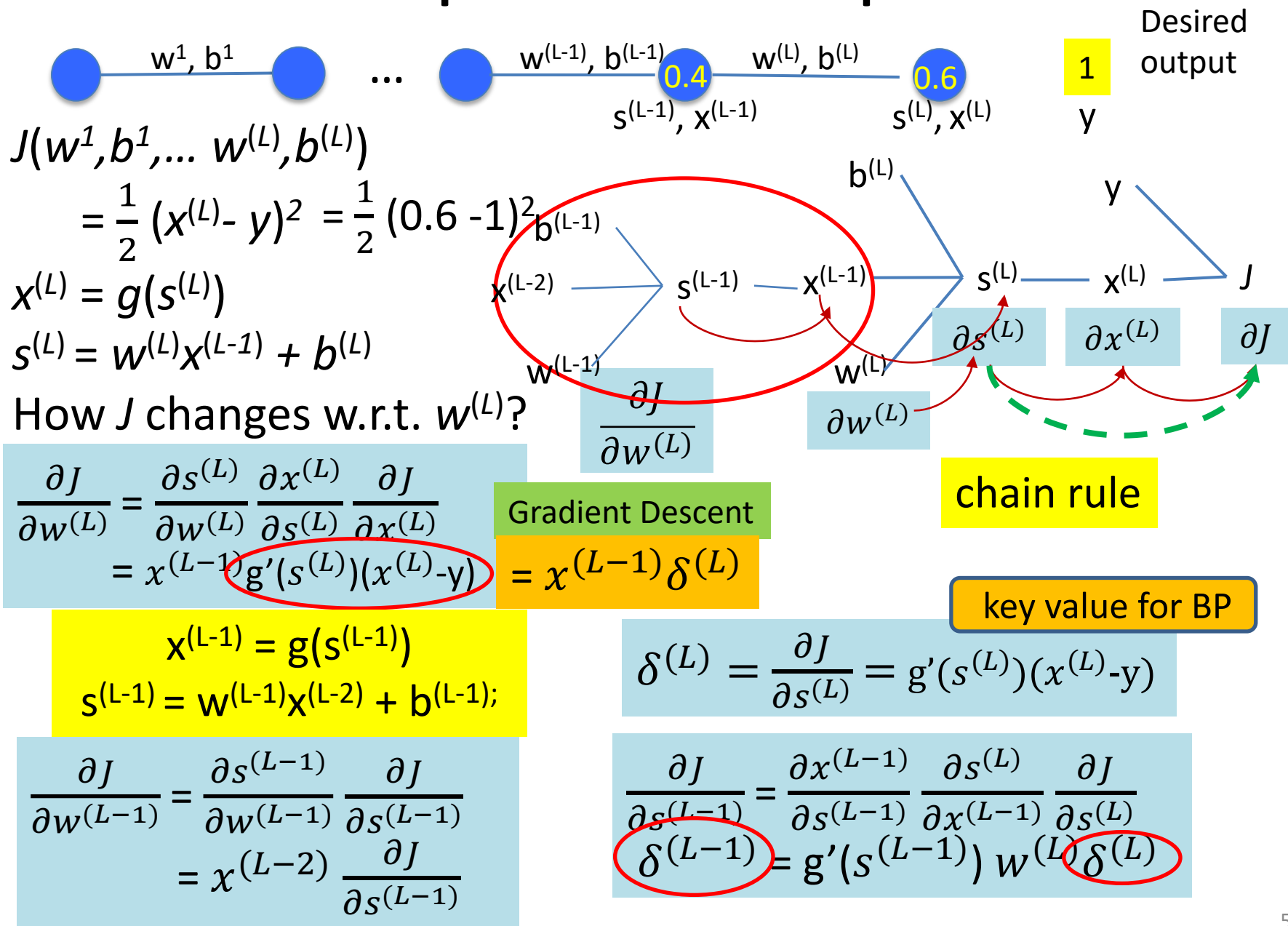


$$s_j^{(l)} = \sum (w_{ij}^{(l)} x_i^{(l-1)} + b_i) \quad x_j^{(l)} = g(s_j^{(l)})$$

$$\delta_i^{(l-1)} = g'(s_i^{(l-1)}) (\sum w_{ij}^{(l)} \delta_j^{(l)}); \quad l = L, L-1, \dots, \quad \text{where } \delta_i^{(l-1)} = \frac{\partial J}{\partial s_i^{(l-1)}}$$

$$\delta_i^{(L)} = g'(s_i^{(L)}) \left( \frac{\partial J}{\partial h(x)} \right)$$

# A simple BP example

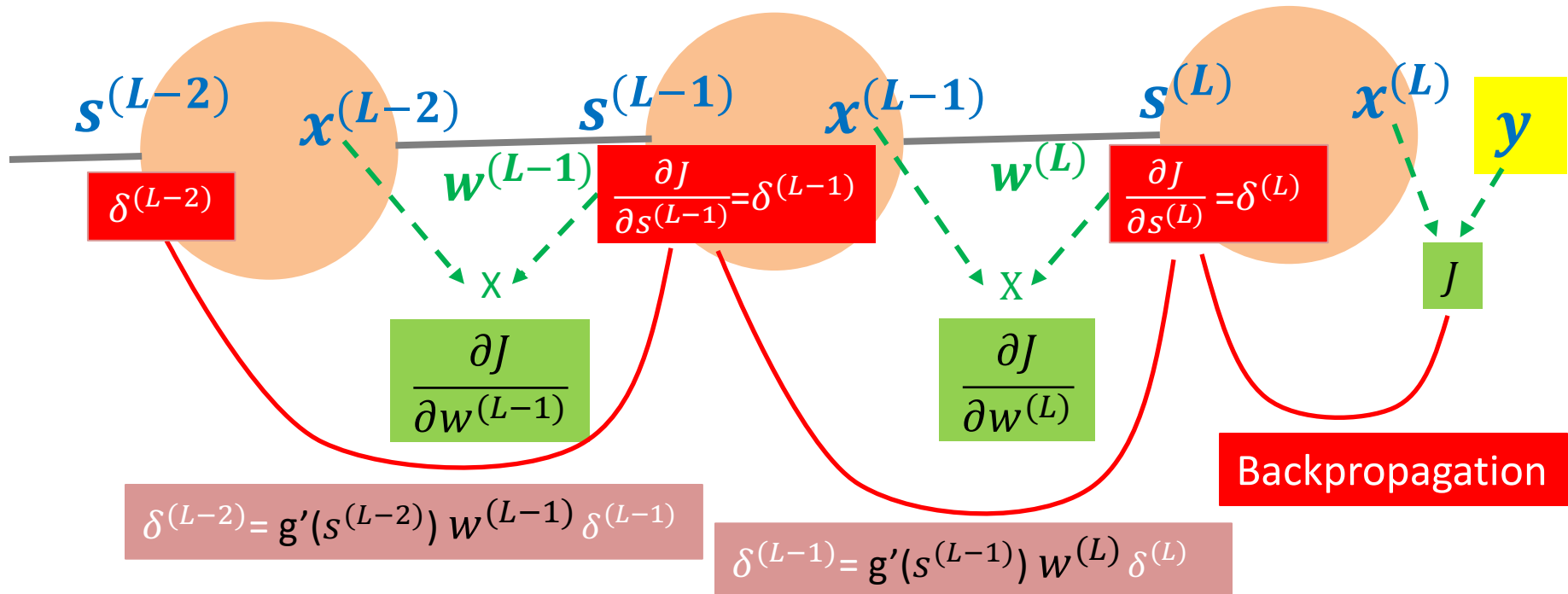


# How $w$ 's are adjusted?

Layer  $L-2$

Layer  $L-1$

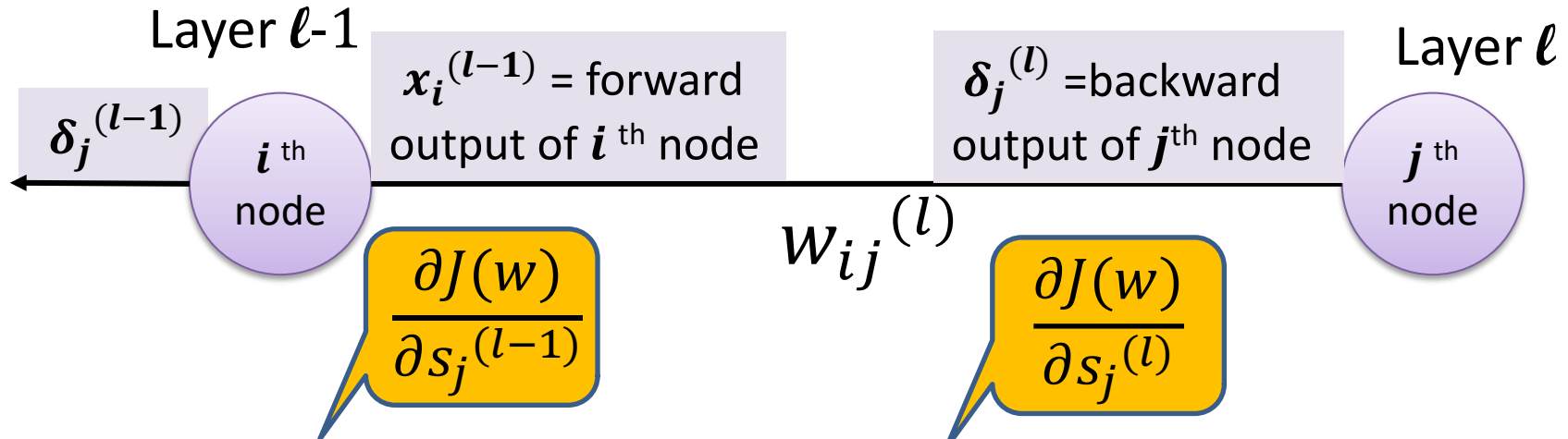
Layer  $L$



Backpropagate  $\frac{\partial J}{\partial s^{(l)}} = \delta^{(l)}$

$\frac{\partial J}{\partial w^{(l)}} = x^{(l-1)} \delta^{(l)}$  based on output  $x^{(l-1)}$  and BP value  $\frac{\partial J}{\partial s^{(l)}} = \delta^{(l)}$

# Backpropagation



$$\delta_i^{(\ell-1)} = g'(s_i^{(\ell-1)}) (\sum w_{ij}^{(\ell)} \delta_j^{(\ell)}); l = L, L-1, L-2 \dots$$

As  $J(w)$  depends on  $\{h(x), y\}$  and  $h(x) = g(s)$

$$\text{where } \delta_i^{(L)} = \frac{\partial J(w)}{\partial s_i^{(L)}} = \frac{\partial h(x)}{\partial s_i^{(L)}} \frac{\partial J}{\partial h(x)} = g'(s_i^{(L)}) \left( \frac{\partial J}{\partial h(x)} \right)$$

$$w_{ij}^{(l)} \rightarrow w_{ij}^{(l)} - \alpha x_i^{(l-1)} \delta_j^{(l)}$$

need to prove

$$\frac{\partial J(w)}{\partial w_{ij}^{(l)}} = x_i^{(l-1)} \delta_j^{(l)} \quad \text{and} \quad \delta_j^{(l)} = \frac{\partial J(w)}{\partial s_j^{(l)}}$$

# Training NN with backpropagation

- Backpropagation
- Explanation through examples
- **Showing gradient  $\frac{\partial J(W)}{\partial w_{ij}^{(l)}}$  is indeed equal to  $x_i^{(l-1)} \delta_j^{(l)}$**
- BP with different activation and cost functions
- Coding backpropagation



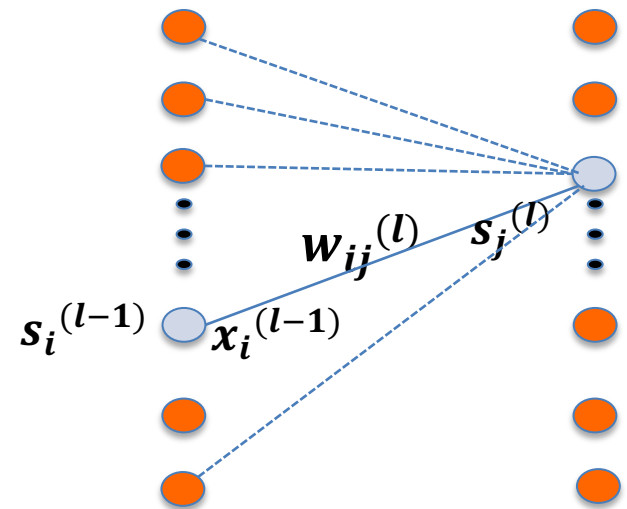
GOAL: Want to show:  $\frac{\partial J(\mathbf{w})}{\partial w_{ij}^{(l)}} = \delta_j^{(l)} \times x_i^{(l-1)}$

By induction (next few slides)

$$\frac{\partial J(\mathbf{w})}{\partial w_{ij}^{(l)}} = \frac{\partial J(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}} = \frac{\partial J(\mathbf{w})}{\partial s_j^{(l)}} \times x_i^{(l-1)}$$

$$s_j^{(l)} = \left( \sum_i w_{ij}^{(l)} x_i^{(l-1)} \right) + b$$

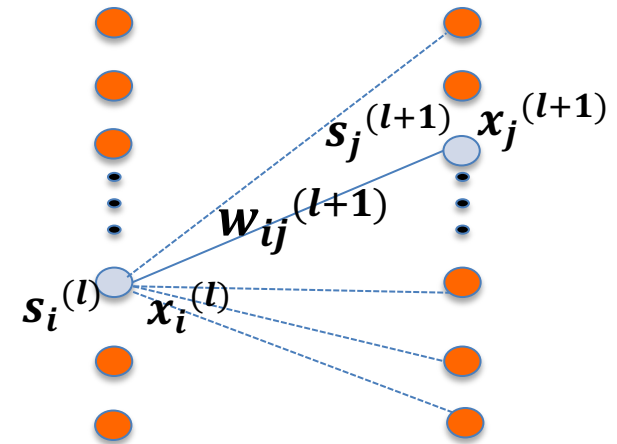
$$\Rightarrow \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}} = x_i^{(l-1)}$$



GOAL now becomes: want to show that  $\delta_j^{(l)} = \frac{\partial J(\mathbf{w})}{\partial s_j^{(l)}}$

NEW GOAL : want to show that  $\delta_j^{(l)} = \frac{\partial J(\mathbf{w})}{\partial s_j^{(l)}}$

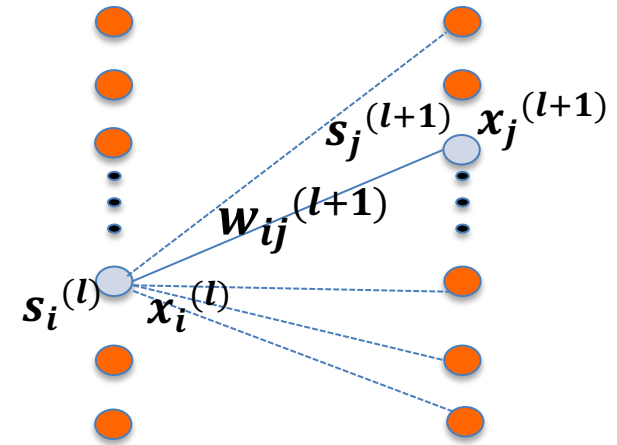
$$\begin{aligned} \frac{\partial J(\mathbf{w})}{\partial s_i^{(l)}} &= \sum_j \frac{\partial J(\mathbf{w})}{\partial s_j^{(l+1)}} \times \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}} \\ &= \sum_j \delta_j^{(l+1)} \times \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}} \end{aligned}$$



By Induction (Backward)  
with  $\delta_j^{(k)} = \frac{\partial J(\mathbf{w})}{\partial s_j^{(k)}}$   
for  $k = L, L-1, \dots, l+1$

NEW GOAL : want to show that  $\delta_j^{(l)} = \frac{\partial J(\mathbf{w})}{\partial s_j^{(l)}}$

$$\begin{aligned} \frac{\partial J(\mathbf{w})}{\partial s_i^{(l)}} &= \sum_j \frac{\partial J(\mathbf{w})}{\partial s_j^{(l+1)}} \times \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}} \\ &= \sum_j \delta_j^{(l+1)} \times \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}} \end{aligned}$$



$$\begin{aligned} s_j^{(l+1)} &= (\sum_i w_{ij}^{(l+1)} x_i^{(l)}) + b \\ \Rightarrow \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} &= w_{ij}^{(l+1)} \end{aligned}$$

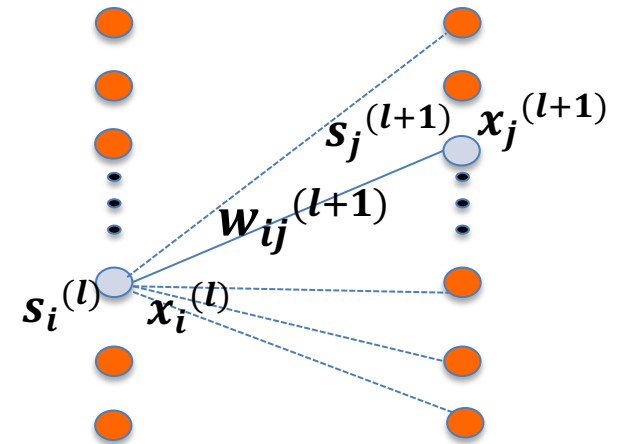
NEW GOAL : want to show that  $\delta_j^{(l)} = \frac{\partial J(\mathbf{w})}{\partial s_j^{(l)}}$

$$\frac{\partial J(\mathbf{w})}{\partial s_i^{(l)}} = \sum_j \frac{\partial J(\mathbf{w})}{\partial s_j^{(l+1)}} \times \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}}$$

$$= \sum_j \delta_j^{(l+1)} \times \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}}$$

$$= \sum_j \delta_j^{(l+1)} \times w_{ij}^{(l+1)} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}}$$

$$= \sum_j \delta_j^{(l+1)} \times w_{ij}^{(l+1)} \times g'(s_i^{(l)}) \quad [\text{since } x_i^{(l)} = g(s_i^{(l)})]$$



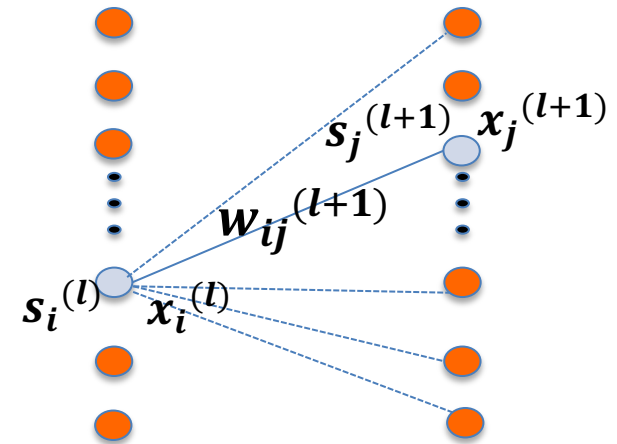
NEW GOAL : want to show that  $\delta_j^{(l)} = \frac{\partial J(\mathbf{w})}{\partial s_j^{(l)}}$

$$\frac{\partial J(\mathbf{w})}{\partial s_i^{(l)}} = \sum_j \frac{\partial J(\mathbf{w})}{\partial s_j^{(l+1)}} \times \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}}$$

$$= \sum_j \delta_j^{(l+1)} \times \frac{\partial s_j^{(l+1)}}{\partial x_i^{(l)}} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}}$$

$$= \sum_j \delta_j^{(l+1)} \times w_{ij}^{(l+1)} \times \frac{\partial x_i^{(l)}}{\partial s_i^{(l)}}$$

$$= \sum_j \delta_j^{(l+1)} \times w_{ij}^{(l+1)} \times g'(s_i^{(l)}) = \delta_i^{(l)}$$

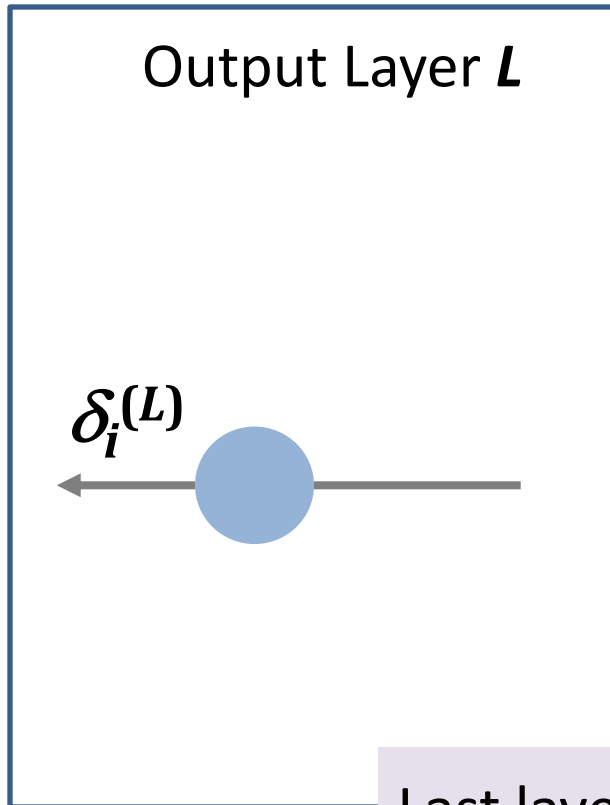


Recall:  $\delta_i^{(l)} = g'(s_i^{(l)}) (\sum_j w_{ij}^{(l+1)} \delta_j^{(l+1)})$

# Training NN with backpropagation

- Backpropagation
- Explanation through examples
- Showing gradient  $\frac{\partial J(\mathbf{W})}{\partial w_{ij}^{(l)}}$  is indeed equal to  $x_i^{(l-1)} \delta_j^{(l)}$
- **BP with different activation and cost functions**
- Coding backpropagation

# For the special case of output layer



For earlier layers:

$$\delta_i^{(l)} = g'(s_i^{(l)}) (\sum w_{ij}^{(l+1)} \delta_j^{(l+1)})$$

For output layer  $L$ :

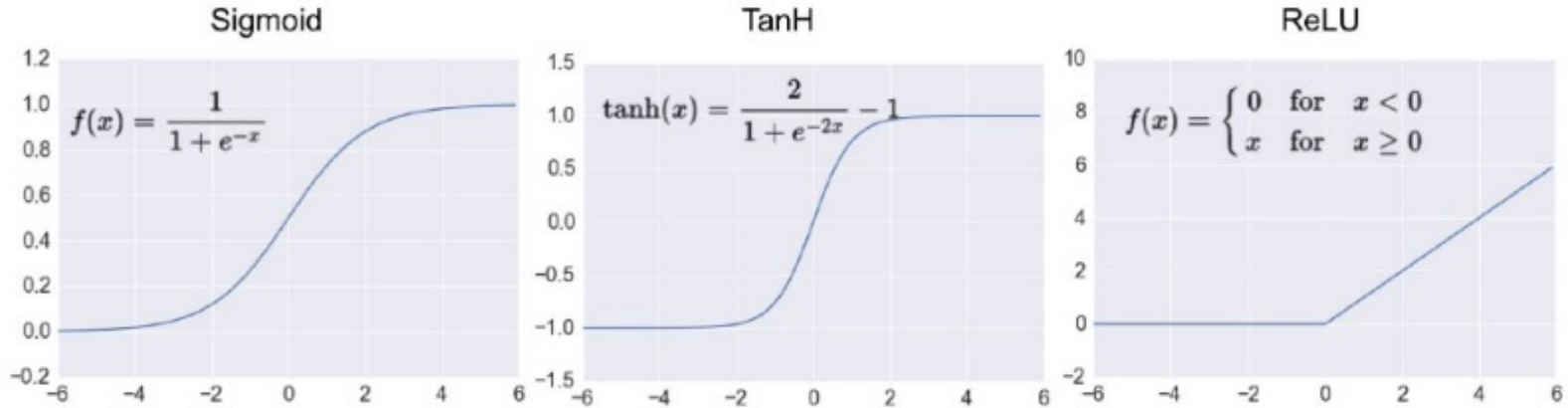
$$\delta_i^{(L)} = g'(s_i^{(L)}) \left( \frac{\partial J}{\partial h(x)} \right)$$

Note:  $J(h(x), y)$ ,  $h(x) = h(g(s))$  at output layer

Last layer  $\delta_i^{(L)}$  depends on (the derivative of) :

- activation function  $g$  in the last layer
- cost function  $J$

# $g$ : Common Activation Functions



Sigmoid:  $g(s) = \sigma(s) = \frac{1}{1 + e^{-s}}$

$$\rightarrow g'(s) = \frac{dg(s)}{ds} = \frac{e^{-s}}{(1 + e^{-s})^2}$$

$$= \left( \frac{1 + e^{-s} - 1}{1 + e^{-s}} \right) \left( \frac{1}{1 + e^{-s}} \right)$$

$$= (1 - g(s))g(s)$$



# ***g* : Common Activation Functions**

Tanh: 
$$g(s) = \frac{2}{1 + e^{-2s}} - 1 = \frac{1 - e^{-2s}}{1 + e^{-2s}}$$
$$= \frac{1 - e^{-s}/e^s}{1 + e^{-s}/e^s} = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$

Quotient rule  
$$\left(\frac{u}{v}\right)' = \frac{u'v - uv'}{v^2}$$

$$g'(s) = \frac{dg(s)}{ds} = \frac{(e^s + e^{-s})(e^s + e^{-s}) - (e^s - e^{-s})(e^s - e^{-s})}{(e^s + e^{-s})^2}$$
$$= 1 - \frac{(e^s - e^{-s})^2}{(e^s + e^{-s})^2} = 1 - g(s)^2$$

ReLU: 
$$g(s) = \begin{cases} 0 & \text{if } s < 0 \\ s & \text{if } s \geq 0 \end{cases} \quad \rightarrow \quad g'(s) = \begin{cases} 0 & \text{if } s < 0 \\ 1 & \text{if } s \geq 0 \end{cases}$$

# Output Layer

	Size	Loss Function $J$	Activation Function $g$
Regression	1	Mean square error	$g$

Mean square error:  $J = \frac{1}{2}(h(x) - y)^2$

Binary cross entropy loss:  $J = -[y \log h(x) + (1 - y) \log(1 - h(x))]$

Cross entropy loss:  $J = -\sum y_k \log h(x)$

# Output Layer $\delta^{(L)}$ - Mean Square Error

	Size	Loss Function $J$	Activation Function $g$	$\delta^{(L)}$
Regression	1	Mean square error	$g$	$g'(s)(h(x) - y)$
Binary Classification	1	Binary cross entropy	Sigmoid	$h(x) - y$
Multi Classification	K	Cross entropy	Softmax	$\mathbf{h}(\mathbf{x}) - \mathbf{y}$

$$J = \frac{1}{2} (h(x) - y)^2$$

For output layer  $L$  (with one output):

$$\begin{aligned}\delta^{(L)} &= g'(s^{(L)}) \left( \frac{\partial J}{\partial h(x)} \right) \\ &= g'(s^{(L)}) (h(x) - y)\end{aligned}$$

# Output Layer $\delta^{(L)}$ - Binary Classification

	Size	Loss Function $J$	Activation Function $g$	$\delta^{(L)}$
Regression	1	Mean square error	$g$	$g'(s)(h(x) - y)$
Binary Classification	1	Binary cross entropy	Sigmoid	$h(x) - y$
Multi Classification	K	Cross entropy	Softmax	$\mathbf{h}(x) - \mathbf{y}$

Note that  $h(x) = g(s) = \sigma(s)$

$$J = -[y \log h(x) + (1 - y) \log(1 - h(x))]$$

For output layer  $L$  (with one output):

$$\begin{aligned}\delta^{(L)} &= g'(s^{(L)}) \left( \frac{\partial J}{\partial h(x)} \right) \\ &= \sigma'(s^{(L)}) \left( -\frac{y}{h(x)} + \frac{1-y}{1-h(x)} \right) \\ &= (1 - \sigma(s))\sigma(s) \left( -\frac{y-h(x)}{h(x)(1-h(x))} \right) = h(x) - y\end{aligned}$$

# Output Layer $\delta^{(L)}$ - Multi Classification

	Size	Loss Function $J$	Activation Function $g$	$\delta^{(L)}$
Regression	1	Mean square error	$g$	$g'(s)(h(x) - y)$
Binary Classification	1	Binary cross entropy	Sigmoid	$h(x) - y$
Multi Classification	K	Cross entropy	Softmax	$h(x) - y$

Cross entropy loss:  $J = -\sum y_k \log p_k$

Note that  $\mathbf{h}(\mathbf{x})$ ,  $\mathbf{y}$  and  $\delta^{(L)}$  are matrices

$$\mathbf{y} = (y_1, \dots, y_K)$$

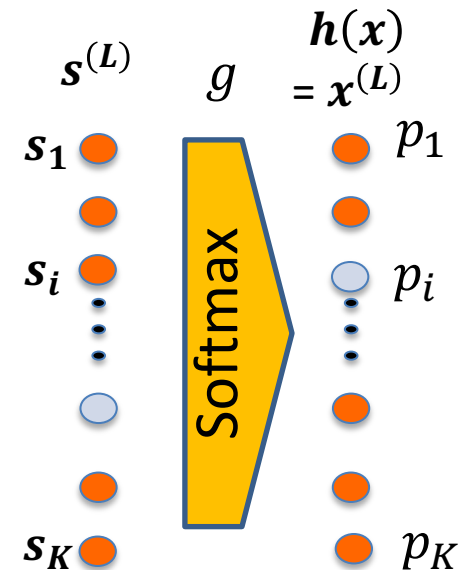
$$\mathbf{h}(\mathbf{x}) = \mathbf{x}^{(L)} = (p_1, \dots, p_K)$$

$$p_i = g(s_i) \text{ (note that } g \text{ is softmax)}$$

$$\delta^{(L)} = g'(\mathbf{s}^{(L)}) \left( \frac{\partial J}{\partial \mathbf{h}(\mathbf{x})} \right)$$

$$\frac{\partial J}{\partial \mathbf{h}(\mathbf{x})} = \left( \frac{\partial J}{\partial p_1}, \frac{\partial J}{\partial p_2}, \dots, \frac{\partial J}{\partial p_K} \right)$$

$$g'(\mathbf{s}^{(L)}) = \left( \frac{\partial p_i}{\partial s_j} \right), \text{ for } 1 \leq i, j \leq K \text{ is a matrix}$$



# Output Layer $\delta^{(L)}$ - Multi Classification

As  $\delta^{(L)} = \mathbf{g}'(\mathbf{s}^{(L)}) \left( \frac{\partial J}{\partial \mathbf{h}(\mathbf{x})} \right)$ ; let  $s_i = s_i^{(L)}$  and  $p_i = x_i^{(L)}$

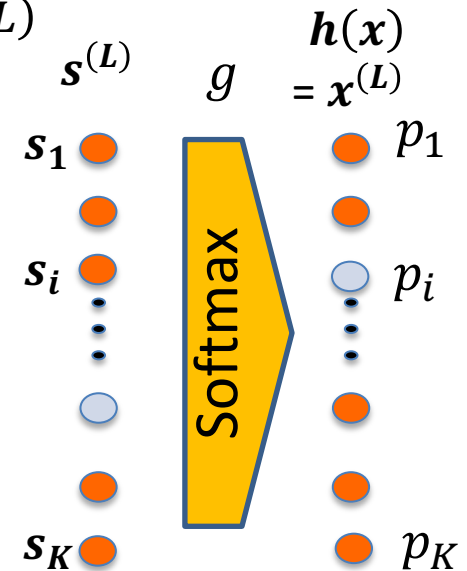
$\mathbf{h}(\mathbf{x}) = (p_1, \dots, p_K)$  where  $g(s_i) = p_i = \frac{e^{s_i}}{\sum e^{s_k}}$

$\mathbf{g}'(\mathbf{s}^{(L)}) = \left( \frac{\partial p_i}{\partial s_j} \right)$ , for  $1 \leq i, j \leq K$

If  $i = j$ ;  $\frac{\partial p_i}{\partial s_i} = (e^{s_i} \sum e^{s_k} - e^{s_i} e^{s_i}) / (\sum e^{s_k})^2 = p_i(1-p_i)$

If  $i \neq j$ ;  $\frac{\partial p_i}{\partial s_j} = -(e^{s_j} e^{s_i}) / (\sum e^{s_k})^2 = -p_j p_i$

$$\mathbf{g}'(\mathbf{s}^{(L)}) = \begin{pmatrix} p_1(1-p_1) & -p_1p_2 & -p_1p_3 & \dots & -p_1p_K \\ -p_2p_1 & p_2(1-p_2) & -p_2p_3 & \dots & -p_2p_K \\ \dots & \dots & \dots & \dots & \dots \\ -p_Kp_1 & -p_Kp_2 & -p_Kp_3 & \dots & p_K(1-p_K) \end{pmatrix} \quad K \times K \text{ matrix}$$



# Output Layer $\delta^{(L)}$ - Multi Classification

	Size	Loss Function $J$	Activation Function $g$	$\delta^{(L)}$
Regression	1	Mean square error	$g$	$g'(s)(h(x) - y)$
Binary Classification	1	Binary cross entropy	Sigmoid	$h(x) - y$
Multi Classification	K	Cross entropy	Softmax	$\mathbf{h(x) - y}$

Cross entropy loss:  $J = -\sum y_k \log p_k$

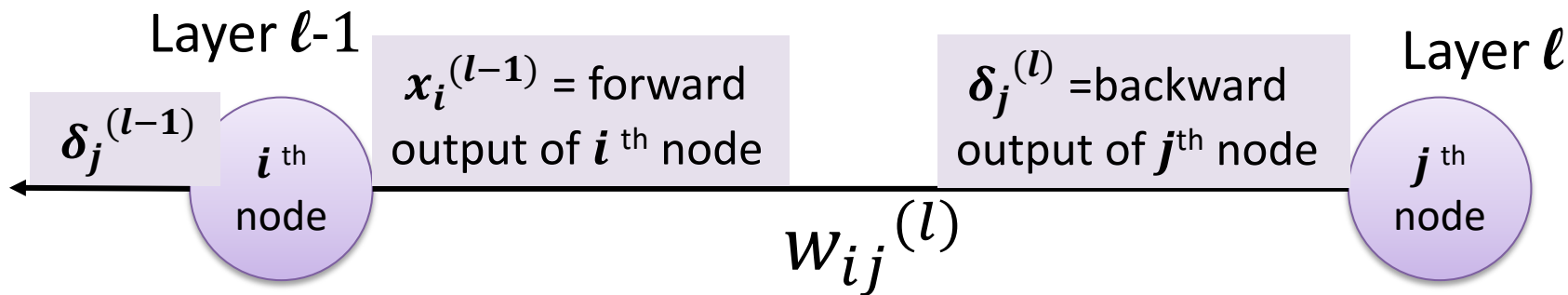
$$\frac{\partial J}{\partial \mathbf{h}(x)} = \left( \frac{\partial J}{\partial p_1}, \frac{\partial J}{\partial p_2}, \dots, \frac{\partial J}{\partial p_C} \right) \text{ where } \frac{\partial J}{\partial p_i} = -\frac{y_i}{p_i}$$

$$\delta^{(L)} = \mathbf{g}'(\mathbf{s}^{(L)}) \left( \frac{\partial J}{\partial \mathbf{h}(x)} \right) = \begin{bmatrix} p_1(1-p_1) & -p_1p_2 & -p_1p_3 & \dots & -p_1p_K \\ -p_2p_1 & p_2(1-p_2) & -p_2p_3 & \dots & -p_2p_K \\ \dots & \dots & \dots & \dots & \dots \\ -p_Kp_1 & -p_Kp_2 & -p_Kp_3 & \dots & p_K(1-p_K) \end{bmatrix} \begin{bmatrix} -y_1/p_1 \\ -y_2/p_2 \\ \dots \\ -y_K/p_K \end{bmatrix}$$

$$= (-y_1 + p_1 \sum y_k, -y_2 + p_2 \sum y_k, \dots, -y_K + p_K \sum y_k)$$

$$= (-y_1 + p_1, -y_2 + p_2, \dots, -y_K + p_K) \text{ as } \sum y_k = 1$$

$$= \mathbf{h(x) - y}$$



- $w_{ij}^{(l)} \rightarrow w_{ij}^{(l)} - \alpha x_i^{(l-1)} \delta_j^{(l)}$  where  $\frac{\partial J(\mathbf{w})}{\partial w_{ij}^{(l)}} = \delta_j^{(l)} \times x_i^{(l-1)}$
- $\delta_i^{(l-1)} = g'(s_i^{(l-1)}) (\sum w_{ij}^{(l)} \delta_j^{(l)})$ ;  $l = L, L-1, L-2, \dots$

<b>sigmoid</b>	$g(s) = \frac{1}{1+e^{-s}}$	$g'(s) = (1 - g(s))g(s)$
<b>Tanh</b>	$g(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$	$g'(s) = 1 - g(s)^2$
<b>Relu</b>	$g(s) = \begin{cases} 0 & \text{if } s < 0 \\ s & \text{if } s \geq 0 \end{cases}$	$g'(s) = \begin{cases} 0 & \text{if } s < 0 \\ 1 & \text{if } s \geq 0 \end{cases}$

## BP Summary

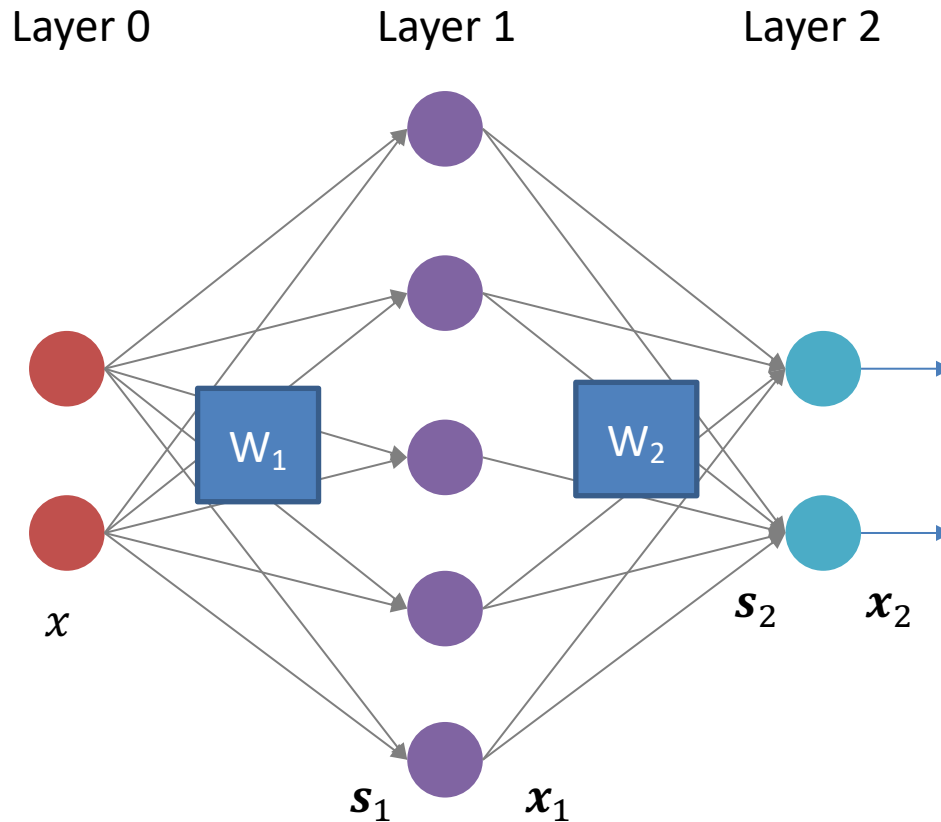
$g$	$\delta^{(L)} = g'(s)(h(x) - y)$
<b>Sigmoid</b>	$\delta^{(L)} = h(x) - y$
<b>Softmax</b>	$\delta^{(L)} = \mathbf{h}(x) - \mathbf{y}$



# Training NN with backpropagation

- Backpropagation
- Explanation through examples
- Showing gradient  $\frac{\partial J(\mathbf{W})}{\partial w_{ij}^{(l)}}$  is indeed equal to  $x_i^{(l-1)} \delta_j^{(l)}$
- BP with different activation and cost functions
- **Coding backpropagation**

# NN Classifier with one hidden layer



Forward pass equations:

$$s_1 = xW_1 + b_1$$

$$x_1 = \tanh(s_1)$$

$$s_2 = x_1W_2 + b_2$$

$$h = x_2 = \text{softmax}(s_2)$$

Dimensions of matrices:

$$x: 1 \times 2$$

$$W_1: 2 \times 5$$

$$b_1: 1 \times 5$$

$$s_1, x_1: 1 \times 5$$

$$W_2: 5 \times 2$$

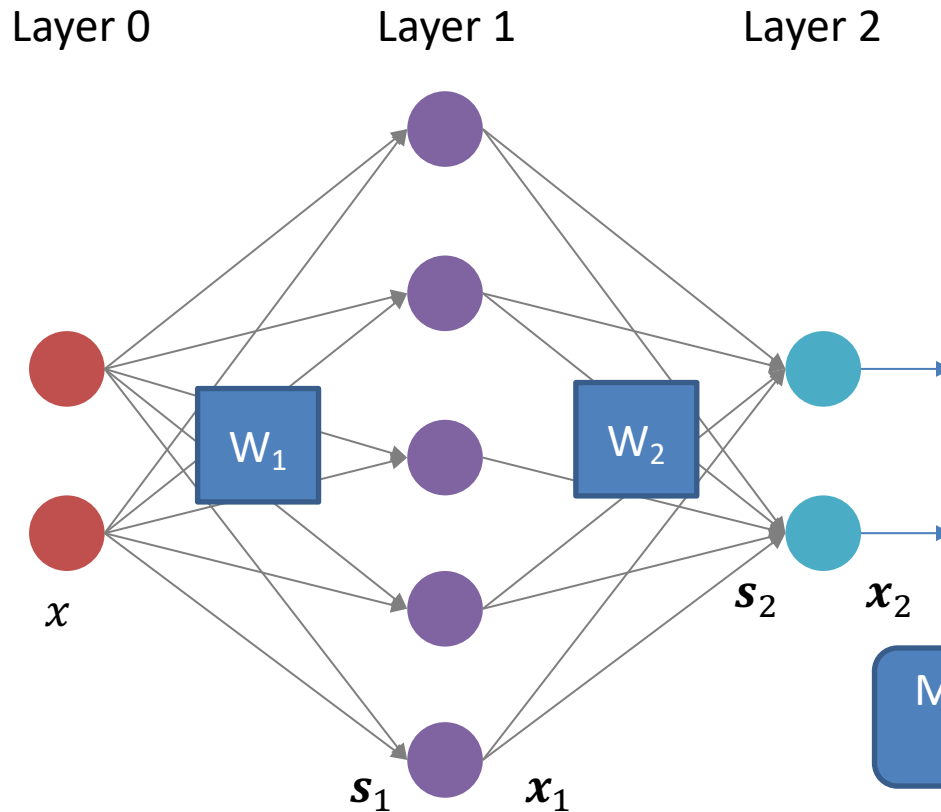
$$b_2: 1 \times 2$$

$$s_2, x_2: 1 \times 2$$

Let  $x_i$  is the forward output of layer  $i$  (here  $x_1, x_2$ )

Let  $\delta_i$  is the backward output of layer  $i$  (here  $\delta_1, \delta_2$ )

# NN Classifier with one hidden layer



Forward pass equations:

$$s_1 = xW_1 + b_1$$

$$x_1 = \tanh(s_1)$$

$$s_2 = x_1W_2 + b_2$$

$$h = x_2 = \text{softmax}(s_2)$$

Dimensions of matrices:

$$x: \mathbf{M} \times 2$$

$$W_1: 2 \times 5$$

$$b_1: 1 \times 5$$

$$s_1, x_1: \mathbf{M} \times 5$$

$$W_2: 5 \times 2$$

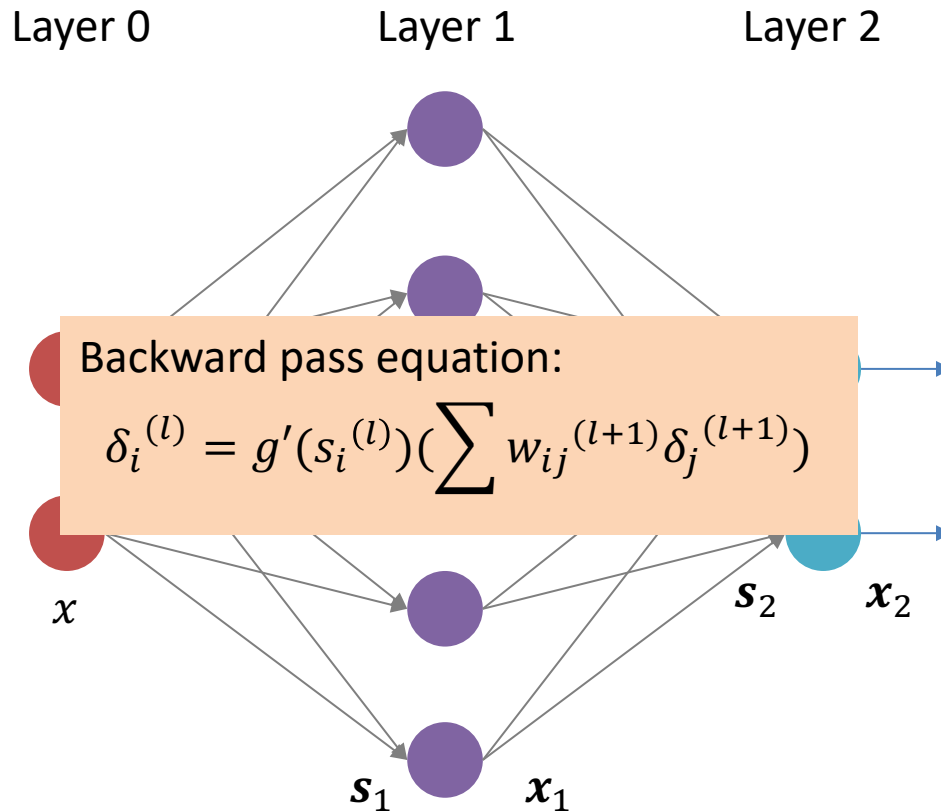
$$b_2: 1 \times 2$$

$$s_2, x_2: \mathbf{M} \times 2$$

$M$  = the number  
of input data

Let  $x_i$  is the forward output of layer  $i$  (here  $x_1, x_2$ )  
Let  $\delta_i$  is the backward output of layer  $i$  (here  $\delta_1, \delta_2$ )

# NN Classifier with one hidden layer



Forward pass equations:

$$s_1 = xW_1 + b_1$$

$$x_1 = \tanh(s_1)$$

$$s_2 = x_1W_2 + b_2$$

$$h = x_2 = \text{softmax}(s_2)$$

Backward pass equations:

$$\delta_2 = h - y$$

$$\begin{aligned} \delta_1 &= g'(s_1) \delta_2 W_2^T \\ &= (1 - \tanh^2 s_1) \delta_2 W_2^T \end{aligned}$$

Dimensions of matrices:

$$\delta_2, x : M \times 2$$

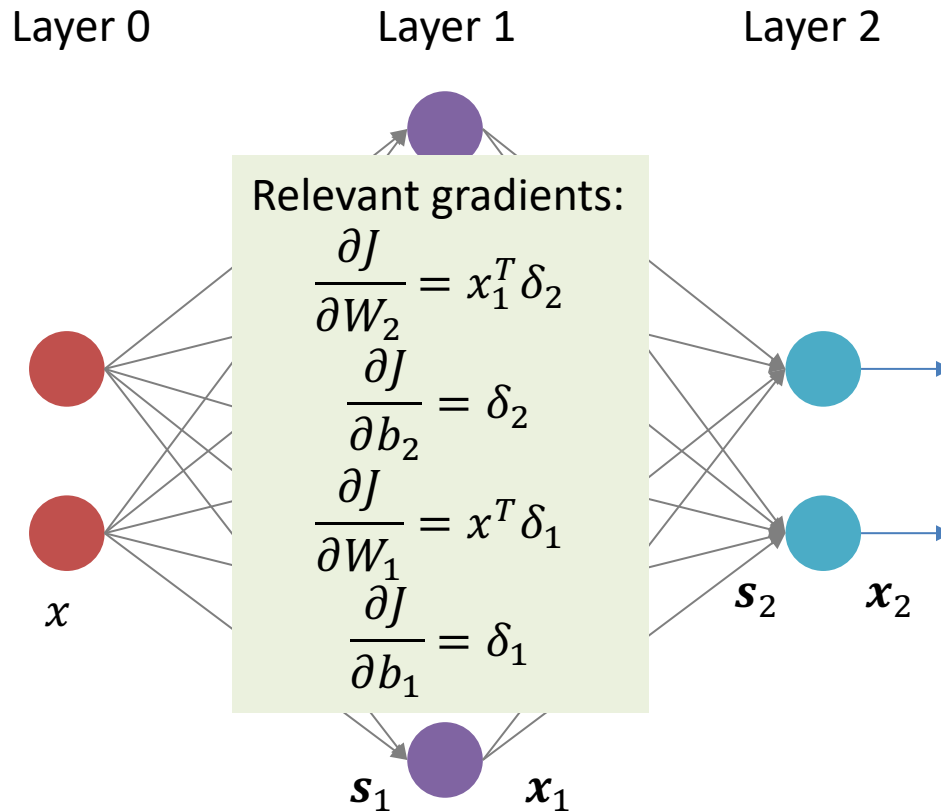
$$W_2 : 5 \times 2$$

$$\delta_1, s_1 : M \times 5$$

Let  $x_i$  is the forward output of layer  $i$  (here  $x_1, x_2$ )

Let  $\delta_i$  is the backward output of layer  $i$  (here  $\delta_1, \delta_2$ )

# NN Classifier with one hidden layer



Forward pass equations:

$$s_1 = xW_1 + b_1$$

$$x_1 = \tanh(s_1)$$

$$s_2 = x_1W_2 + b_2$$

$$h = x_2 = \text{softmax}(s_2)$$

Backward pass equations:

$$\delta_2 = h - y$$

$$\delta_1 = (1 - \tanh^2 s_1) \delta_2 W_2^T$$

Adjust weights:

$$W_2 = W_2 - \alpha x_1^T \delta_2$$

$$b_2 = b_2 - \alpha \delta_2$$

$$W_1 = W_1 - \alpha x^T \delta_1$$

$$b_1 = b_1 - \alpha \delta_1$$

Let  $x_i$  is the forward output of layer  $i$  (here  $x_1, x_2$ )

Let  $\delta_i$  is the backward output of layer  $i$  (here  $\delta_1, \delta_2$ )

# Framework for training Neural Network

1. Initialize all weights  $w_{ij}^{(l)}$  at random
2. Repeat until convergence:
3. Pick a random example to feed into layer 0
4. *Forward*: Compute all  $x_j^{(l)}$
5. *Backward*: Compute all  $\delta_j^{(l)}$
6. Update weights:  $w_{ij}^{(l)} \rightarrow w_{ij}^{(l)} - \alpha x_i^{(l-1)} \delta_j^{(l)}$
7. Return the final weights  $w_{ij}^{(l)}$

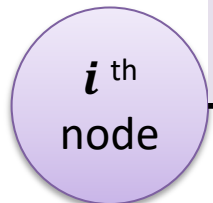
$$\frac{\partial J(\mathbf{W})}{\partial w_{ij}^{(l)}}$$



$$x_i^{(l-1)} \delta_j^{(l)}$$

$$\frac{\partial J}{\partial s_j^{(l)}}$$

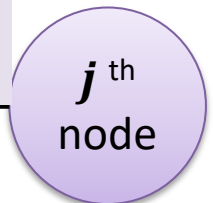
Layer  $\ell-1$



$x_i^{(l-1)}$  = forward  
output of  $i^{\text{th}}$  node

$\delta_j^{(l)}$  = backward  
output of  $j^{\text{th}}$  node

Layer  $\ell$



$$w_{ij}^{(l)}$$

# Techniques to Improve Training

- **Choice of activation function**
- Training by mini-batches
- Finding the right learning rate
- Overfitting and techniques to avoid it
  1. Regularization
  2. Dropout (Minimal effort backpropagation)
  3. Early stopping

# Key equations for backpropagation: $\delta_i^{(l)}$

For output layer  $L$ :  $\delta_i^{(L)} = g'(s_i^{(L)}) \left( \frac{\partial J}{\partial h(x)} \right)$

For earlier layers for  $l = L-1, \dots, 1$ :

$$\delta_i^{(l)} = g'(s_i^{(l)}) \left( \sum w_{ij}^{(l+1)} \delta_j^{(l+1)} \right)$$

$$= g'(s_i^{(l)}) \left( \sum w_{ij}^{(l+1)} g'(s_j^{(l+1)}) \left( \sum w_{jk}^{(l+2)} \delta_k^{(l+2)} \right) \right)$$



# Key equations for backpropagation: $\delta_i^{(l)}$

For output layer  $L$ :  $\delta_i^{(L)} = g'(s_i^{(L)}) \left( \frac{\partial J}{\partial h(x)} \right)$

For earlier layers for  $l = L-1, \dots, 1$ :

$$\delta_i^{(l)} = g'(s_i^{(l)}) \left( \sum_j w_{ij}^{(l+1)} \delta_j^{(l+1)} \right)$$

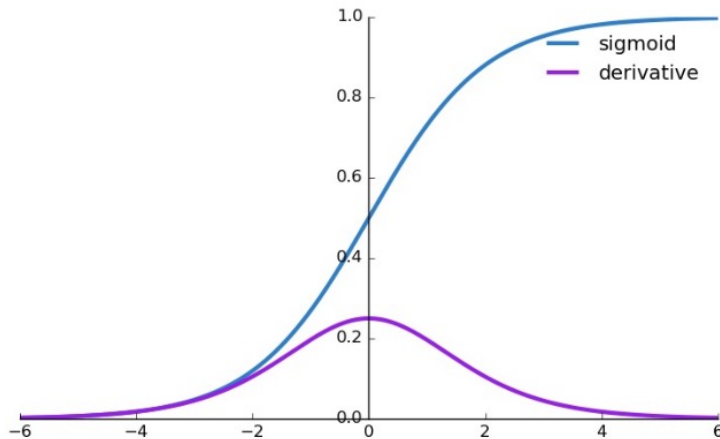
$$= g'(s_i^{(l)}) \left( \sum_j w_{ij}^{(l+1)} g'(s_j^{(l+1)}) \left( \sum_k w_{jk}^{(l+2)} \delta_k^{(l+2)} \right) \right)$$

$$= g'(\ ) \dots g'(\ ) \dots g'(\ ) \dots \quad [up\ to\ L\ occurrences]$$

# Gradient with Sigmoid

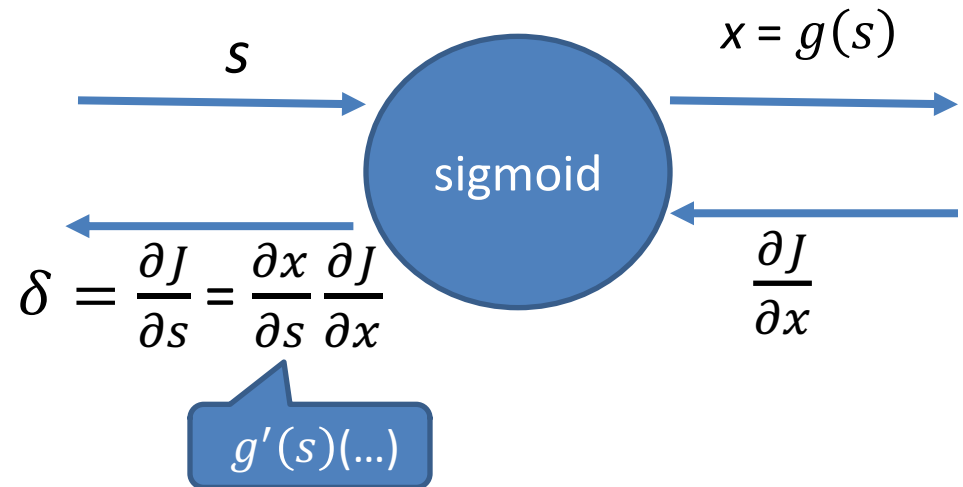
$$g(s) = \sigma(s) = \frac{1}{1 + e^{-s}}$$

$$\rightarrow g'(s) = (1 - g(s))g(s)$$



When  $g'(s)$  is very small,  
 $\frac{\partial J}{\partial s}$  is very small

=> tiny steps in gradient descent!



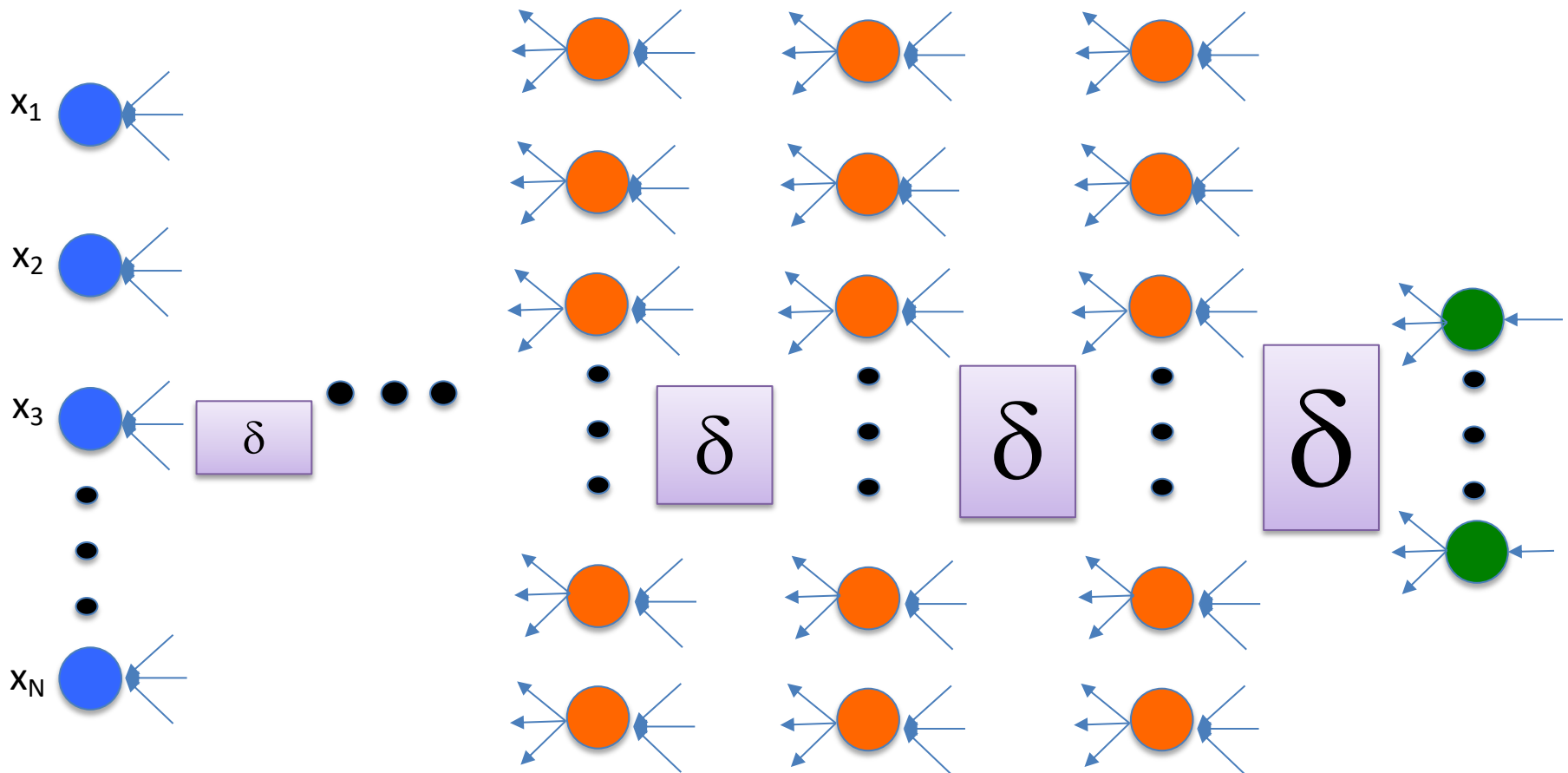
- When  $s$  is very small or very large,  $g'(s)$  is very small.
- Even when  $s$  isn't,  $g'(s) \leq 0.25$ .

# Vanishing Gradient Problem

$\delta$ 's get smaller and smaller as you go back because

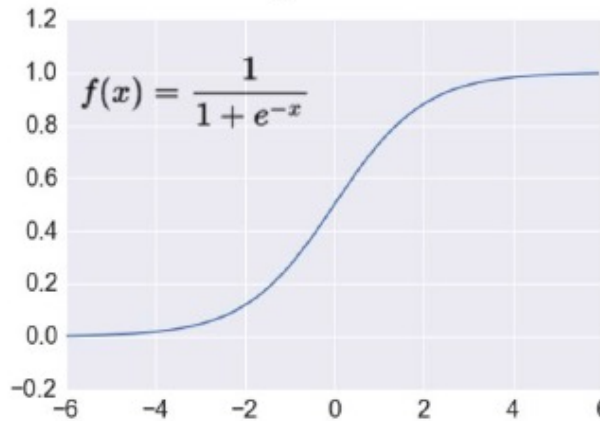
$$\delta = g'() \dots g'() \dots g'() \dots \text{ and } g'() \leq 0.25$$

$\Rightarrow$  smaller changes to weights as you go back

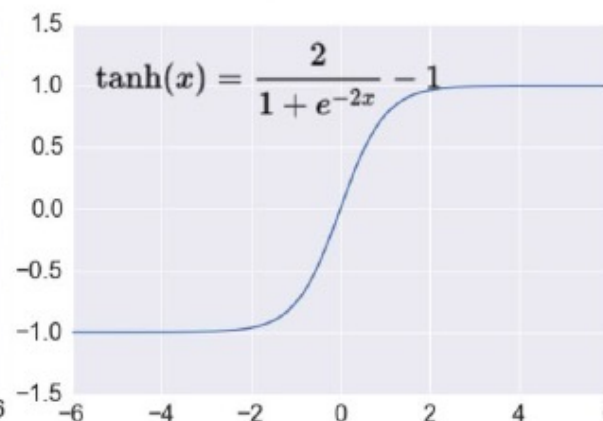


# What about other activation functions?

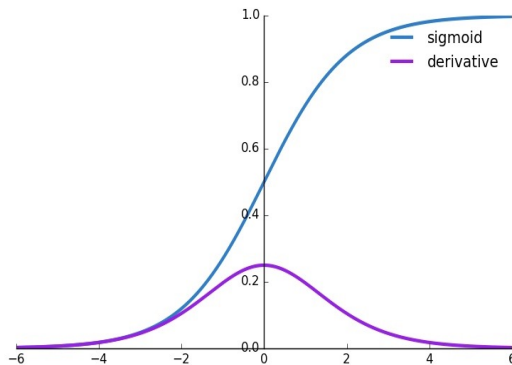
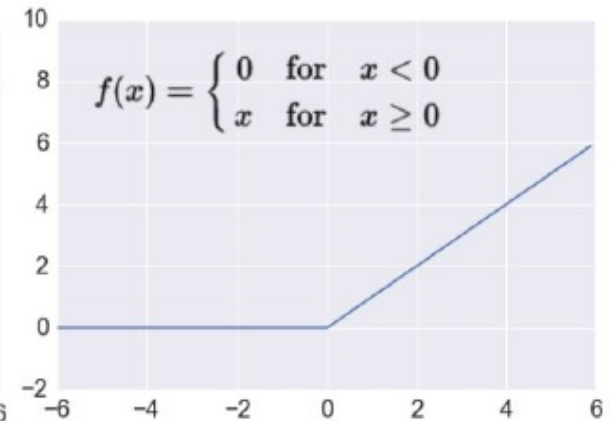
Sigmoid



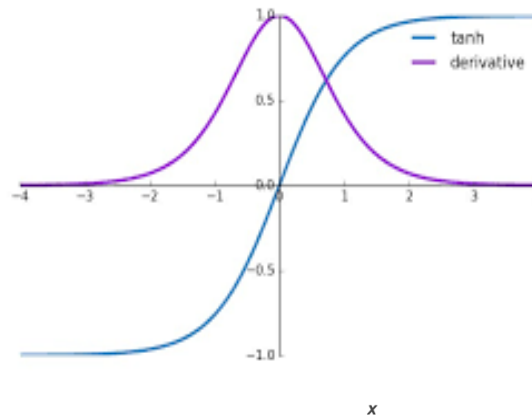
TanH



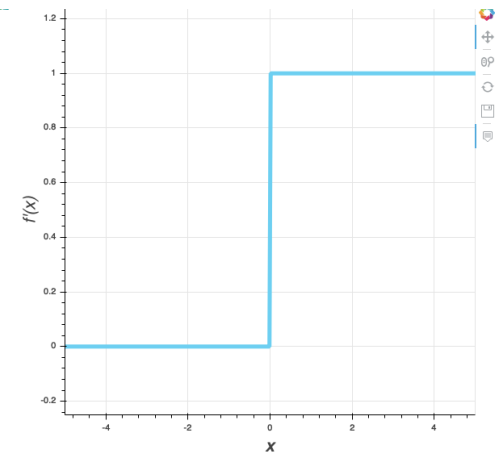
ReLU



$$g'(s) = (1 - g(s))g(s)$$

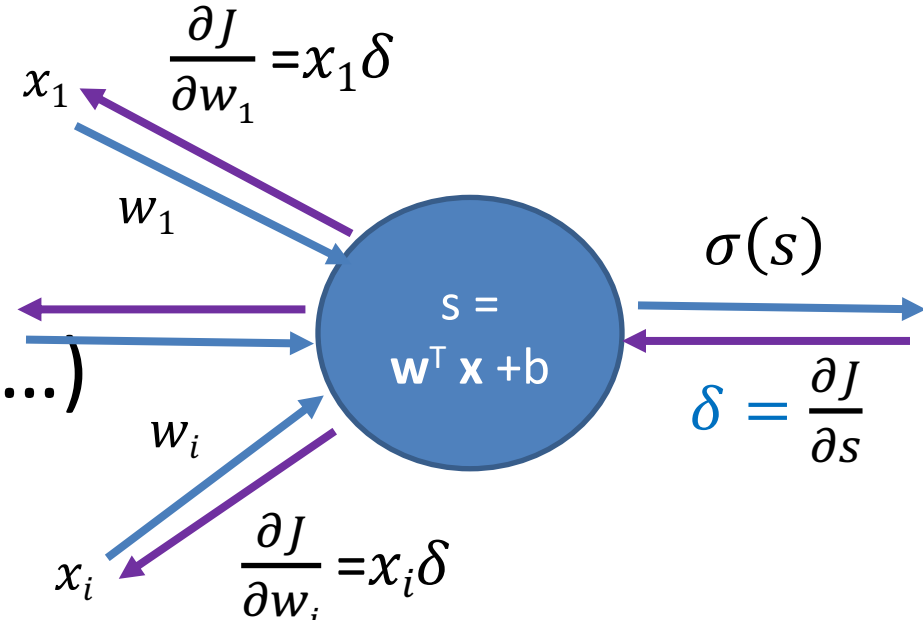


$$g'(s) = 1 - g(s)^2$$



$$g'(s) = \begin{cases} 0 & \text{if } s < 0 \\ 1 & \text{if } s \geq 0 \end{cases}$$

# Not Zero-centered

$$\begin{aligned}
 \frac{\partial J}{\partial \mathbf{w}} &= \left( \frac{\partial J}{\partial w_1}, \dots, \frac{\partial J}{\partial w_i}, \dots \right) \\
 &= \left( \frac{\partial s}{\partial w_1} \frac{\partial J}{\partial s}, \dots, \frac{\partial s}{\partial w_i} \frac{\partial J}{\partial s}, \dots \right) \\
 &= (x_1 \delta, \dots, x_i \delta, \dots)
 \end{aligned}$$


The diagram illustrates a neuron model. A central blue circle represents the neuron, labeled with the equation  $s = \mathbf{w}^T \mathbf{x} + b$ . 
   
Forward pass (blue arrows): Inputs  $x_1$  and  $x_i$  are multiplied by weights  $w_1$  and  $w_i$  respectively, and then summed with a bias  $b$  to produce the net input  $s$ . The output of the neuron is  $\sigma(s)$ .
   
Backward pass (purple arrows): The error  $\delta = \frac{\partial J}{\partial s}$  is propagated back from the output. This error is then multiplied by the input values  $x_1$  and  $x_i$  to determine the contribution of each weight to the total error gradient.

What happens when we have only positive values  $x = \sigma(s)$  (i.e., sigmoid output)? What about  $\frac{\partial J}{\partial \mathbf{w}}$ ?

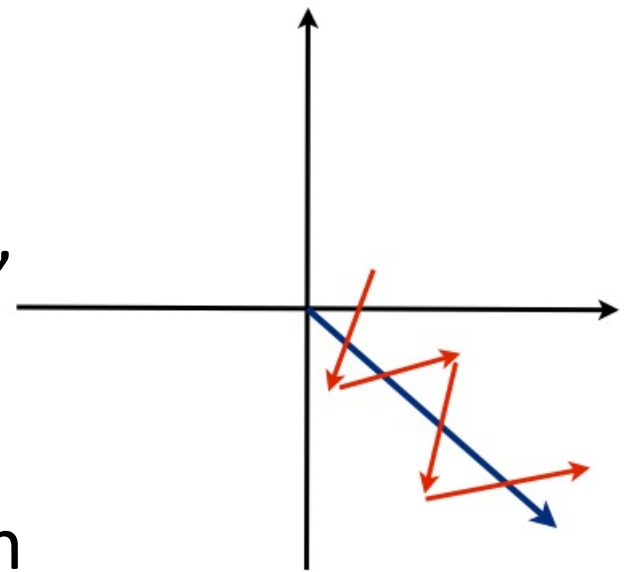
$(x_1 \delta, \dots, x_i \delta, \dots)$  will be either all positive or negative depending on  $\delta$ , i.e., all in the same direction.

# Zig-Zag Path for Gradient Descent

Suppose the blue vector is the optimal direction for the gradient (assume a simple case with 2 weights). Since the weights either all increase or decrease (same direction), then we get a zig-zag path, not good for gradient descent.

Similarly, non zero-centered activation functions are not ideal, they produce more positive or more negative output values.

Thus, we want activation function to be **zero-centered**.



# Sigmoid Activation Function

Historically popular because of nice interpretation of probability by squeezing numbers in  $[0, 1]$  for classification.

## Shortcomings:

- Vanishing gradient
- Sigmoid outputs are not zero-centered (mean  $> 0$ )
- $\exp()$  is expensive to compute

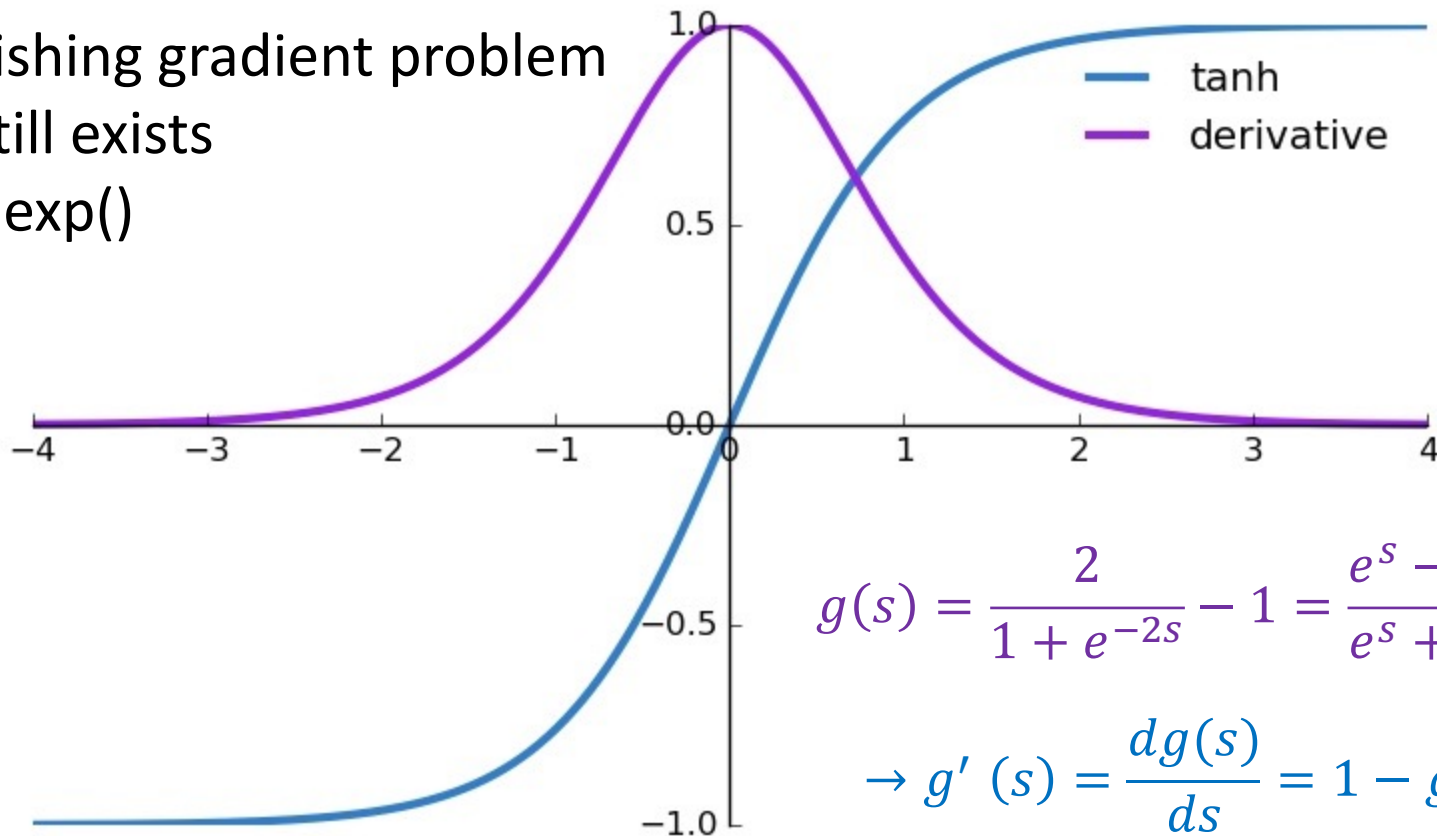
# Tanh and its derivative

Squashes numbers in  $[-1, 1]$

Zero centered

Vanishing gradient problem  
still exists

Still exp()

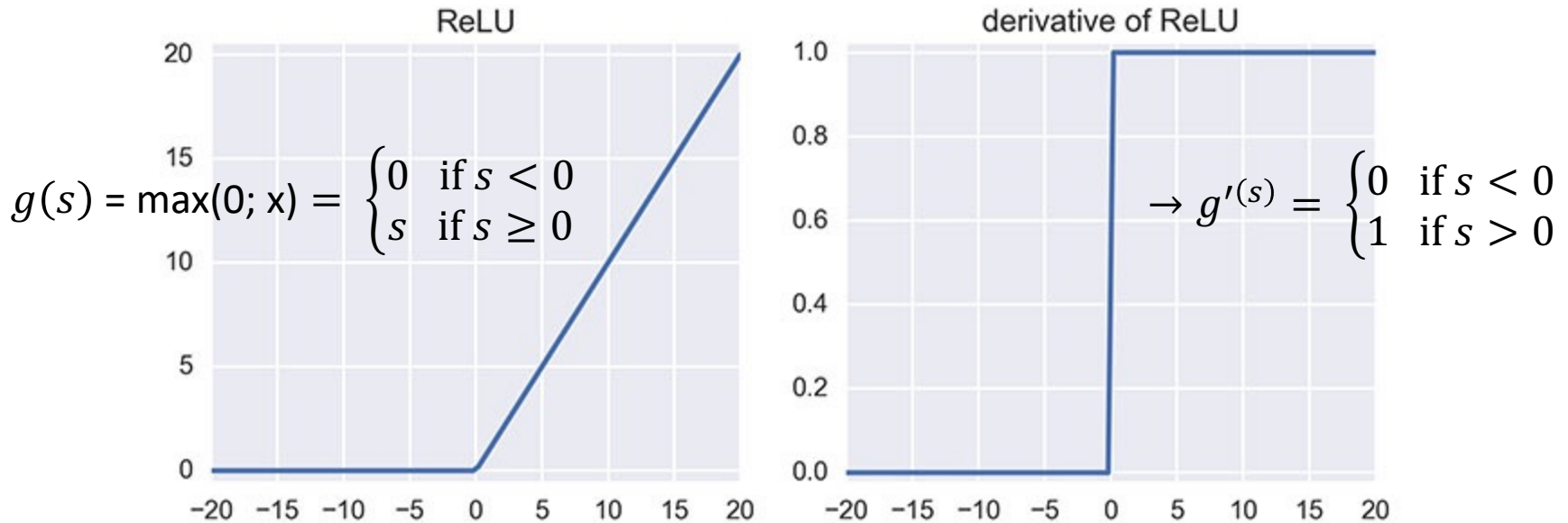


$$g(s) = \frac{2}{1 + e^{-2s}} - 1 = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$

$$\rightarrow g'(s) = \frac{dg(s)}{ds} = 1 - g(s)^2$$



# ReLU (Rectified Linear Unit)

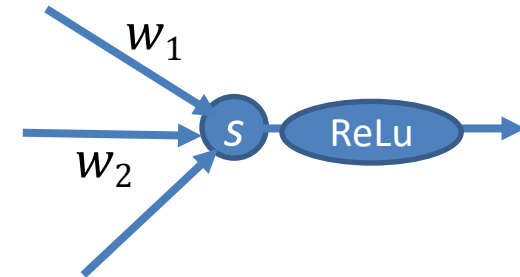


- No vanishing gradient problem when positive, but still exists when negative (gradient = 0 at 0)
- Converges much faster than sigmoid/tanh (a factor of 6)
- Computationally efficient
- Non-zero centered
- Widely used in CNN (since Alexnet) – to be learned later
- Seems to model better (neurons has the dropout effect)

# ReLU-like Activation Functions

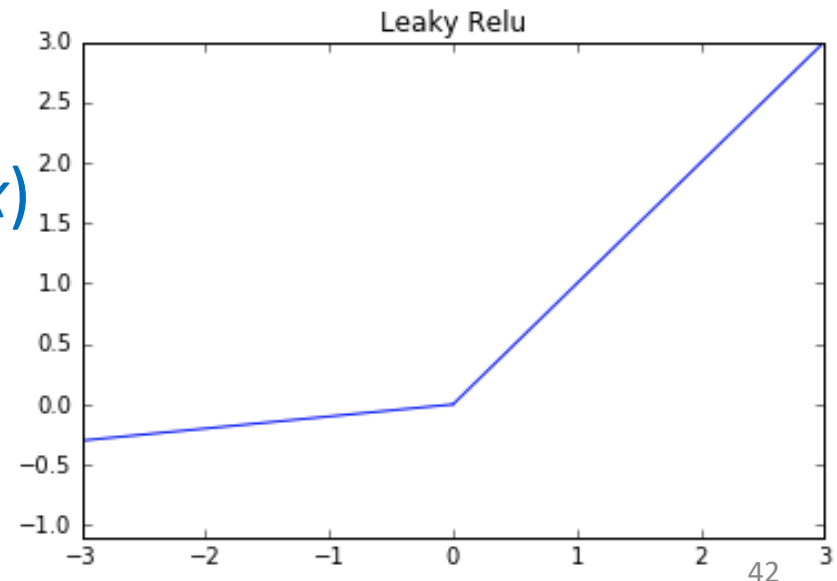
Problem for ReLU:

- Gradient is 0 for negative  $s$
- Weights will not get adjusted
- **Dying ReLU problem:** neurons will stop (die) simply because gradient is 0, nothing changes ( $s$  remains  $-ve$  and gradient remains 0).
- As much as 40% of your network can be “dead” (i.e., neurons never activate across the entire training dataset)



Leaky Relu -  $f(x) = \max(0.01x, x)$

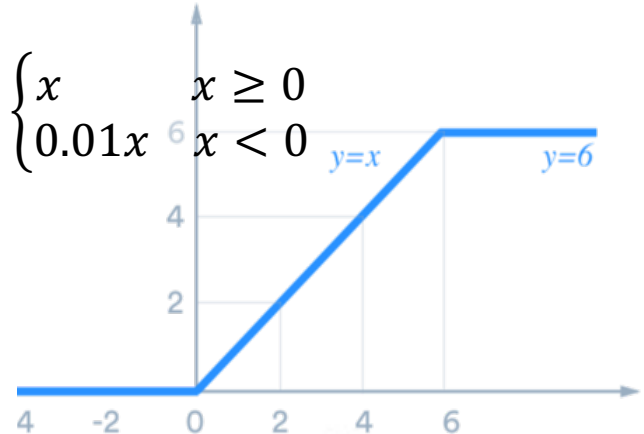
- All benefits of ReLU
- closer to zero-centered outputs
- nonzero gradient when negative



# ReLU-like Activation Functions

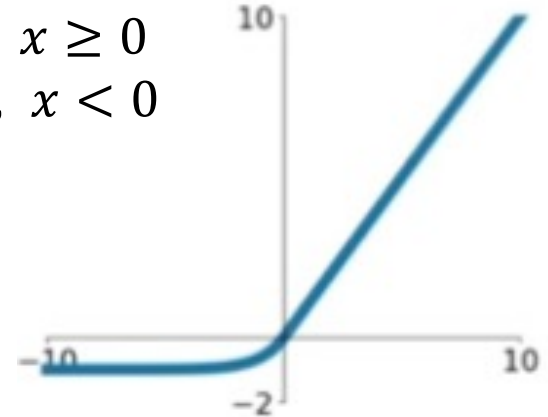
## Variations in ReLU

- **Leaky ReLU** -  $f(x) = \max(0.01x, x)$   $f(x) = \begin{cases} x & x \geq 0 \\ 0.01x & x < 0 \end{cases}$
- **Parametric ReLU** -  $f(x) = \max(ax, x)$
- **ReLU6** -  $f(x) = \min(\max(x, 0), 6)$
- clipping the gradient at some value  $n$



**Exponential Linear Units (ELU):**  $f(x) = \begin{cases} x, & x \geq 0 \\ \alpha(e^x - 1), & x < 0 \end{cases}$

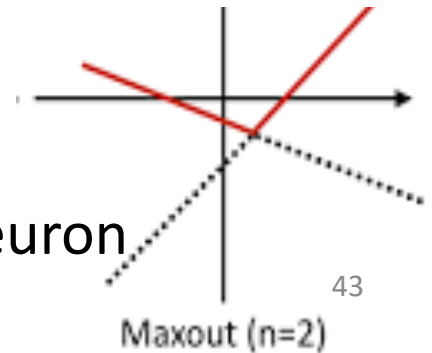
- Similar to leaky ReLU
- Differentiable at 0 when  $\alpha=1$
- Drawback: Higher computation for  $\exp()$



## Maxout "Neuron"

$\text{Max}(w_1x + b_1, w_2x + b_2)$

- Generalizes ReLU and Leaky ReLU
- Drawback: doubles number of parameters per neuron



# Issue with Symmetry

## Network Initialization?

### Neural Network with

- same activation function
- same no. of neurons/layer
- same initial weights and bias

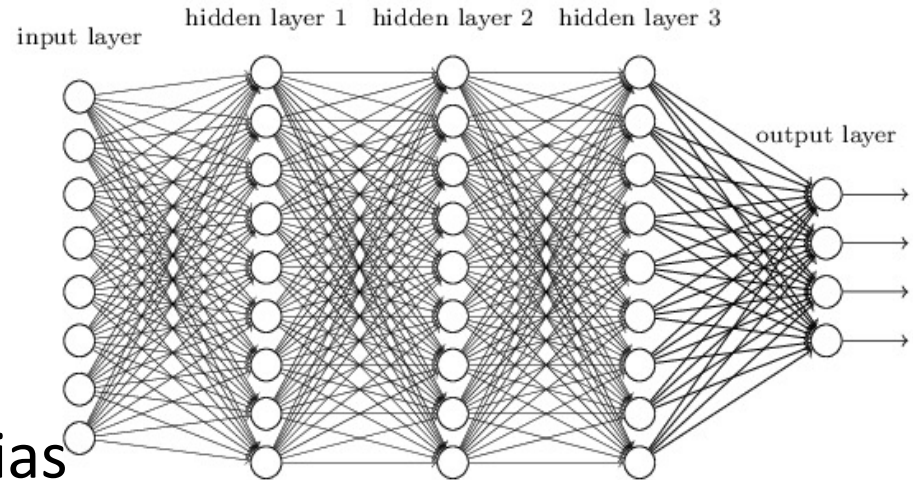
### Problem?

- Same activation values in each layer
- Same gradient at every neuron in each layer

Limited number of directions for gradient descent

Solution: weight initialization –

sample from the normal distribution  $N(0, 10^{-2})$   
with bias=0



# Techniques

- Choice of activation function
- **Training by mini-batches**
- Finding the right learning rate
- Overfitting and techniques to avoid
  1. regularization
  2. Dropout (Minimal effort backpropagation)
  3. Early stopping
- Error Analysis

# Gradient Descent (GD)

- Minimize  $J(\theta) = \frac{1}{M} \sum_{i=1}^M e(h(\mathbf{x}_i), y_i)$
- By taking iterative steps along  $-\nabla J(\theta)$  :  
$$\Delta\theta = -\alpha \nabla J(\theta)$$
- $\nabla J$  computed based on all  $M$  examples

# Stochastic Gradient Descent (SGD)

- Minimize  $J(\theta) = \frac{1}{M} \sum_{i=1}^M e(h(\mathbf{x}_i), y_i)$
- By taking iterative steps along  $-\nabla J(\theta)$  :  
$$\Delta\theta = -\alpha \nabla J(\theta)$$
- $\nabla J$  computed based on all  $M$  examples
- Stochastic GD: randomly pick one  $(\mathbf{x}_i, y_i)$
- Expected value of negative gradient of  $(\mathbf{x}_i, y_i)$  :  
$$E[-\nabla e(h(\mathbf{x}_i), y_i)] = \frac{1}{M} \sum_{i=1}^M -\nabla e(h(\mathbf{x}_i), y_i) = -\nabla J$$

# Stochastic Gradient Descent (SGD)

- Minimize  $J(\theta) = \frac{1}{M} \sum_{i=1}^M e(h(\mathbf{x}_i), y_i)$
- By taking iterative steps along  $-\nabla J(\theta)$  :  
$$\Delta\theta = -\alpha \nabla J(\theta)$$
- $\nabla J$  computed based on all  $M$  examples
- Stochastic GD : randomly pick one  $(\mathbf{x}_i, y_i)$
- Expected value of negative gradient of  $(\mathbf{x}_i, y_i)$  :  
$$E[-\nabla e(h(\mathbf{x}_i), y_i)] = \frac{1}{M} \sum_{i=1}^M -\nabla e(h(\mathbf{x}_i), y_i) = -\nabla J$$

Alternatively: SGD can pick few samples.



# Mini-batch SGD

- Minimize  $J(\theta) = \frac{1}{B} \sum_{i=1}^B e(h(x_i), y_i)$
- By taking iterative steps along  $-\nabla J(\theta)$  :  
$$\Delta\theta = -\alpha \nabla J(\theta)$$
- $\nabla J$  computed based on ~~all~~  $B$  examples

# SGD Algorithm

1. Initialize all weights  $w_{ij}^{(l)}$  at random
2. Repeat until convergence:
3. Pick a **mini-batch** of (can be one) examples to feed into layer 0
4. **Forward**: Compute all  $x_j^{(l)}$
5. Compute error
6. **Backward**: Compute all  $\delta_j^{(l)}$
7. Update weights:  $w_{ij}^{(l)} \rightarrow w_{ij}^{(l)} - \alpha x_i^{(l-1)} \delta_j^{(l)}$
8. Return the final weights  $w_{ij}^{(l)}$

# Mini-batches in SGD

- Instead of updating weights based on single random example as input, use a **batch** of random examples
- Advantages:
  - More accurate estimation of gradient
  - Smoother convergence
  - Allows for larger learning rates
  - Mini-batches lead to fast training!
  - Can parallelize computation + achieve significant speed increases on GPU's

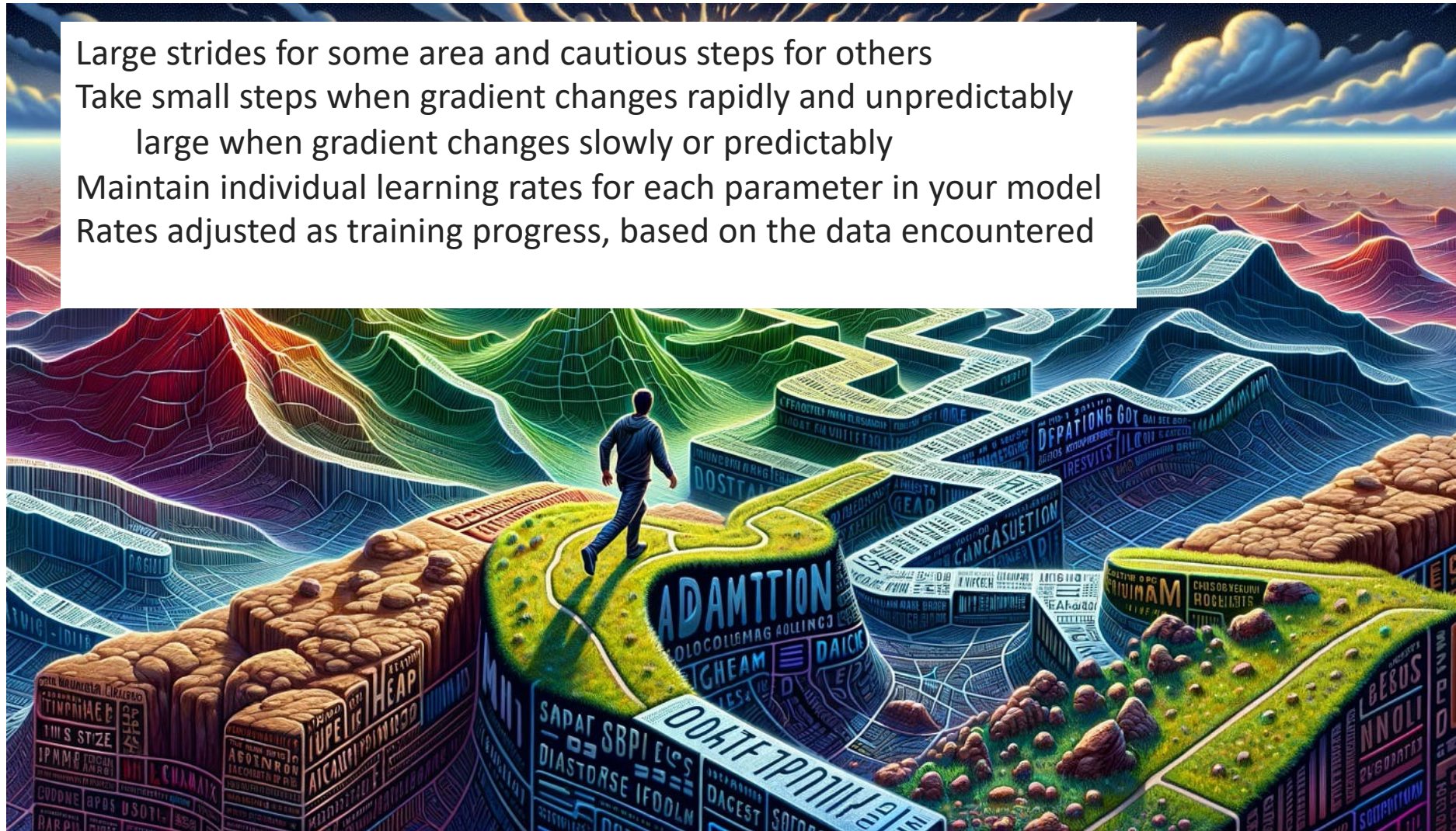
# Techniques

- Choice of activation function
- Training by mini-batches
- **Finding the right learning rate**
- Overfitting and techniques to avoid
  1. regularization
  2. Dropout (Minimal effort backpropagation)
  3. Early stopping

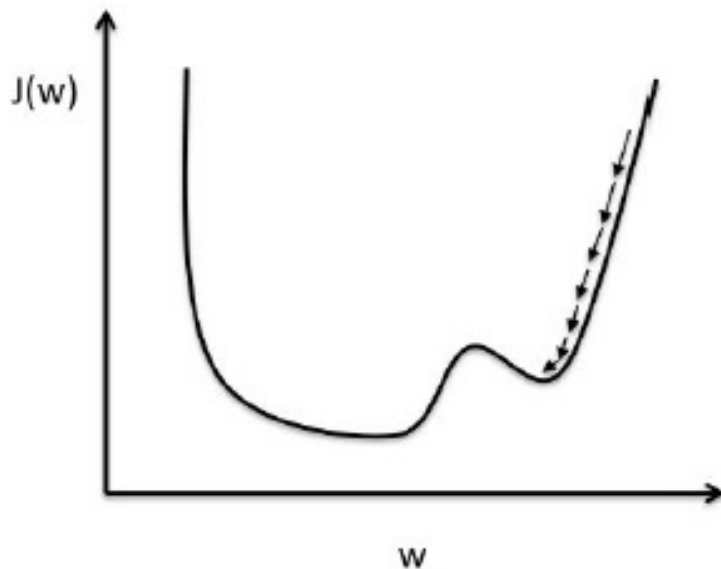
# Navigating a Complex Terrain

<https://towardsdatascience.com/the-math-behind-adam-optimizer-c41407efe59b>

Large strides for some area and cautious steps for others  
Take small steps when gradient changes rapidly and unpredictably  
large when gradient changes slowly or predictably  
Maintain individual learning rates for each parameter in your model  
Rates adjusted as training progress, based on the data encountered

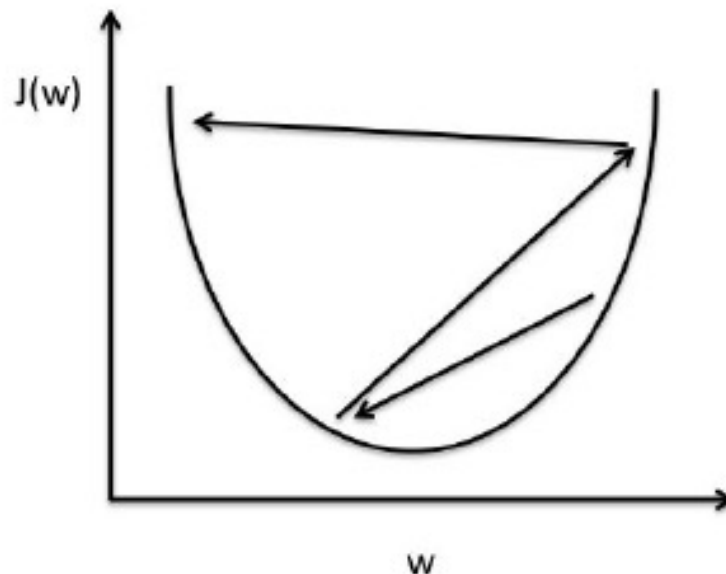


# Problem: Learning Rate $\alpha$



## Small learning rate:

- Many iterations till convergence
- Trapped in local minimum



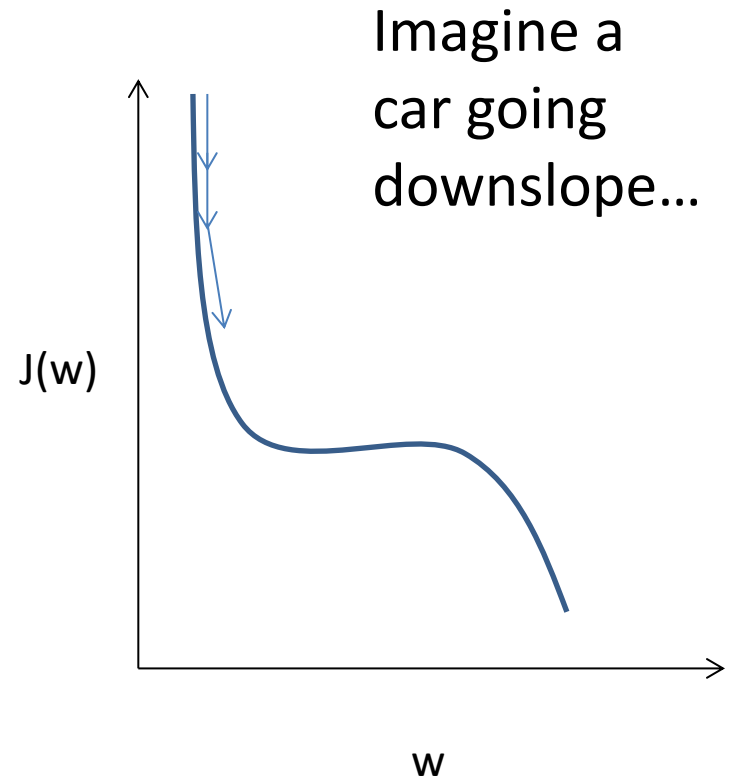
## Large learning rate:

- Overshooting
- No convergence



# Momentum

- Trapped in a local minimum:  
What if you get stuck in a flattish region where gradient is small?
- Let past steps affect current step so that you **keep on moving in the direction you were heading**
- Weights updated by a function of past steps and current gradient



# Momentum

- Weights updated by a function of past steps and current gradient
- Introduce new variable to capture past steps:

$v_t$  = velocity at time  $t$

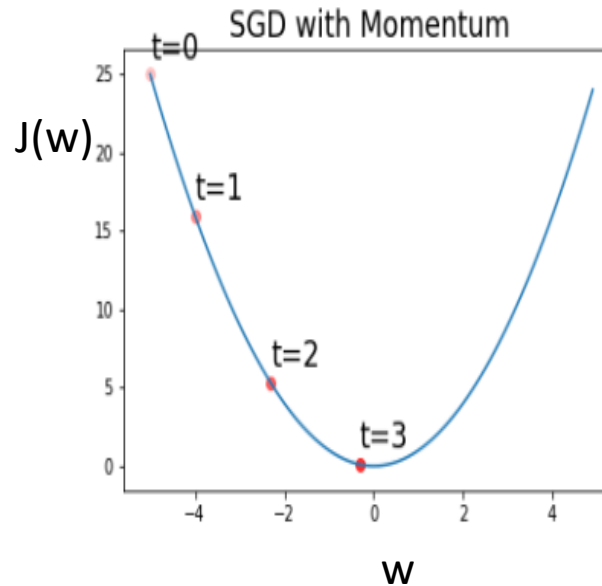
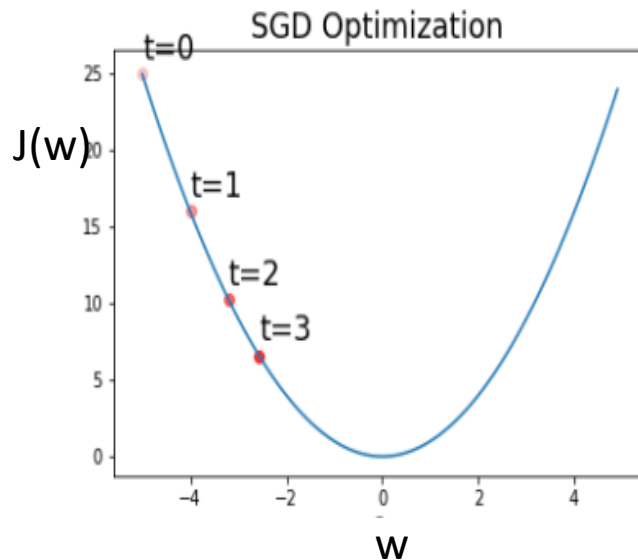
- Update of weights:  $w_t \leftarrow w_{t-1} + v_t$
- Update of velocity:

$$v_t \leftarrow \mu v_{t-1} - \alpha \nabla J(w_{t-1})$$

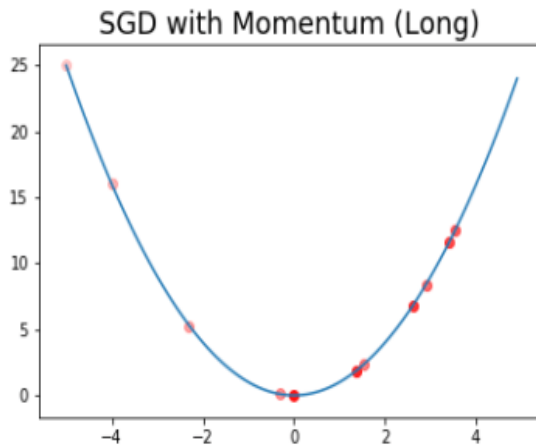
$$[\text{with } v_1 = -\alpha \nabla J(w_0)]$$

- Momentum step:  $\mu v_{t-1}$  [ $\mu$  is typically set to 0.9]
- Gradient step:  $-\alpha \nabla J(w_{t-1})$





Momentum  
gets to minimum  
faster

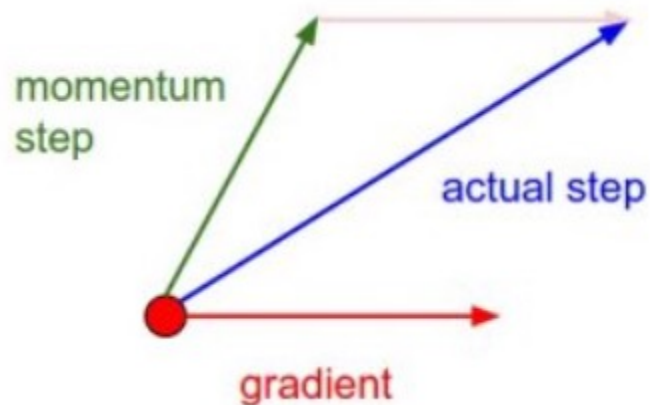


But overshoots  
the minimum  
if run too long

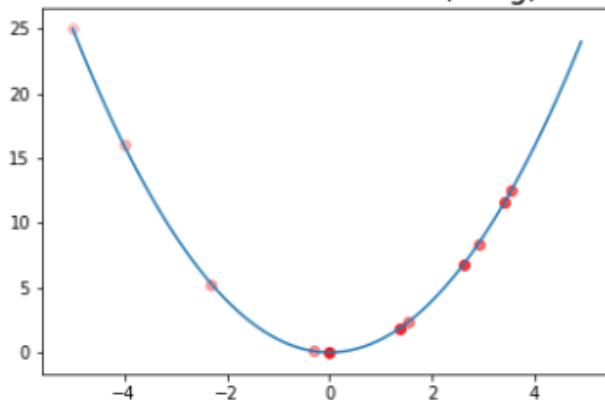
Wouldn't it be  
great if we could:  
detect the  
overshooting, i.e.  
change in the  
direction of the  
gradient earlier?

# Momentum vs. Nesterov Momentum

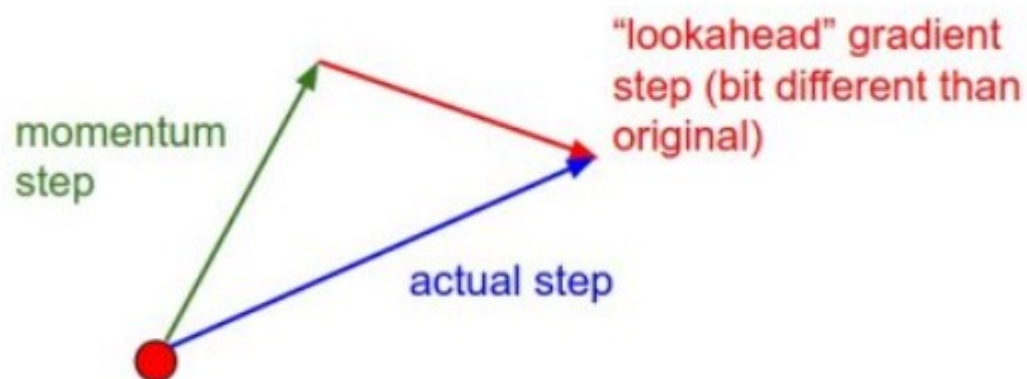
- $v_t \leftarrow \mu v_{t-1} - \alpha \nabla J(w_{t-1})$
- $w_t \leftarrow w_{t-1} + v_t$



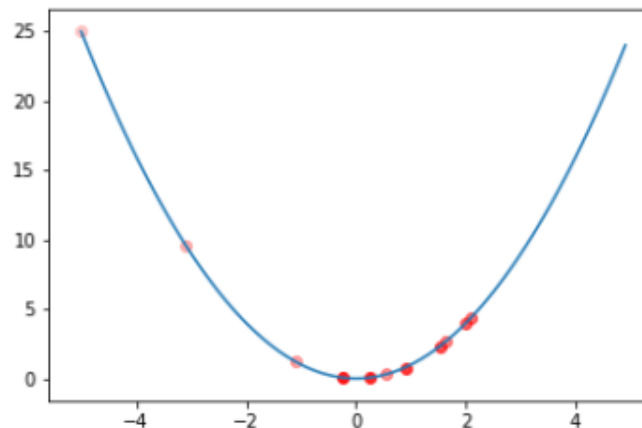
SGD with Momentum (Long)



- $v_t \leftarrow \mu v_{t-1} - \alpha \nabla J(w_{t-1} + \mu v_{t-1})$
- $w_t \leftarrow w_{t-1} + v_t$



Nesterov's Accelerated Gradient Method



## SGD

```
CLASS torch.optim.SGD(params, lr=<required parameter>, momentum=0, dampening=0,  
    weight_decay=0, nesterov=False, *, maximize=False, foreach=None, differentiable
```

Example

```
>>> optimizer = torch.optim.SGD(model.parameters(), lr=0.1, momentum=0.9)  
>>> optimizer.zero_grad()  
>>> loss_fn(model(input), target).backward()  
>>> optimizer.step()
```

Implements stochastic gradient descent (optionally with momentum).

---

**input** :  $\gamma$  (lr),  $\theta_0$  (params),  $f(\theta)$  (objective),  $\lambda$  (weight decay),  
 $\mu$  (momentum),  $\tau$  (dampening), *nesterov*, *maximize*

---

**for**  $t = 1$  **to** ... **do**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

**if**  $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

**if**  $\mu \neq 0$

**if**  $t > 1$

$\mathbf{b}_t \leftarrow \mu \mathbf{b}_{t-1} + (1 - \tau) g_t$

**else**

$\mathbf{b}_t \leftarrow g_t$

**if** *nesterov*

$g_t \leftarrow g_t + \mu \mathbf{b}_t$

**else**

$g_t \leftarrow \mathbf{b}_t$

**if** *maximize*

$\theta_t \leftarrow \theta_{t-1} + \gamma g_t$

**else**

$\theta_t \leftarrow \theta_{t-1} - \gamma g_t$

---

**return**  $\theta_t$

---

### • NOTE

The implementation of SGD with Momentum/Nesterov subtly differs from Sutskever et. al. and implementations in some other frameworks.

Considering the specific case of Momentum, the update can be written as

$$v_{t+1} = \mu * v_t + g_{t+1},$$

$$p_{t+1} = p_t - lr * v_{t+1},$$

where  $p$ ,  $g$ ,  $v$  and  $\mu$  denote the parameters, gradient, velocity, and momentum respectively.

This is in contrast to Sutskever et. al. and other frameworks which employ an update of the form

$$v_{t+1} = \mu * v_t + lr * g_{t+1},$$

$$p_{t+1} = p_t - v_{t+1}.$$

The Nesterov version is analogously modified.

- $v_t \leftarrow \mu v_{t-1} - \alpha \nabla J(w_{t-1})$
- $w_t \leftarrow w_{t-1} + v_t$

# AdaGrad (separate learning rate for each w)

- **Adaptive Gradient**
- Idea:

Historical gradients for w	Learning rate for w
High	Low
Low	High

- Introduce variable **G** to capture historical gradient information (second moment), with the following update rule :

$$G_t \leftarrow G_{t-1} + (\nabla J(w_{t-1}))^2$$

- Update weights:

$$w_t \leftarrow w_{t-1} - \frac{\alpha}{\sqrt{G_{t-1} + \epsilon}} \nabla J(w_{t-1})$$

$\epsilon$  (epsilon) is a small scalar (e.g.,  $10^{-8}$ ) added to prevent division by zero and maintain numerical stability.

- Learning rate diminishes quickly as G gets bigger and bigger

# AdaGrad vs. RMSProp

## AdaGrad

$$G \leftarrow G + (\nabla J(w))^2$$

G gets bigger and bigger

$$w \leftarrow w - \frac{\alpha}{\sqrt{G+\epsilon}} \nabla J(w)$$

coefficient gets smaller and smaller

## RMSProp

$$G \leftarrow \beta G + (1-\beta) (\nabla J(w))^2$$

$\beta$  is the decay rate (typically 0.9, 0.99,..) preventing G from getting too big

$$w \leftarrow w - \frac{\alpha}{\sqrt{G+\epsilon}} \nabla J(w)$$

Same update rule for w but problem goes away

# RMSProp vs. Adam

## RMSProp

$$G \leftarrow \beta G + (1-\beta) (\nabla J(w))^2$$

G does not get too big

$$w \leftarrow w - \frac{\alpha}{\sqrt{G+\epsilon}} \nabla J(w)$$

Same update rule as AdaGrad

## Adam

$$v \leftarrow \beta v + (1-\beta) (\nabla J(w))^2$$

Same idea as RMSProp (except called v instead of G)

$$m \leftarrow \gamma m + (1-\gamma) (\nabla J(w))$$

New first order variable m introduced to “smooth” gradient

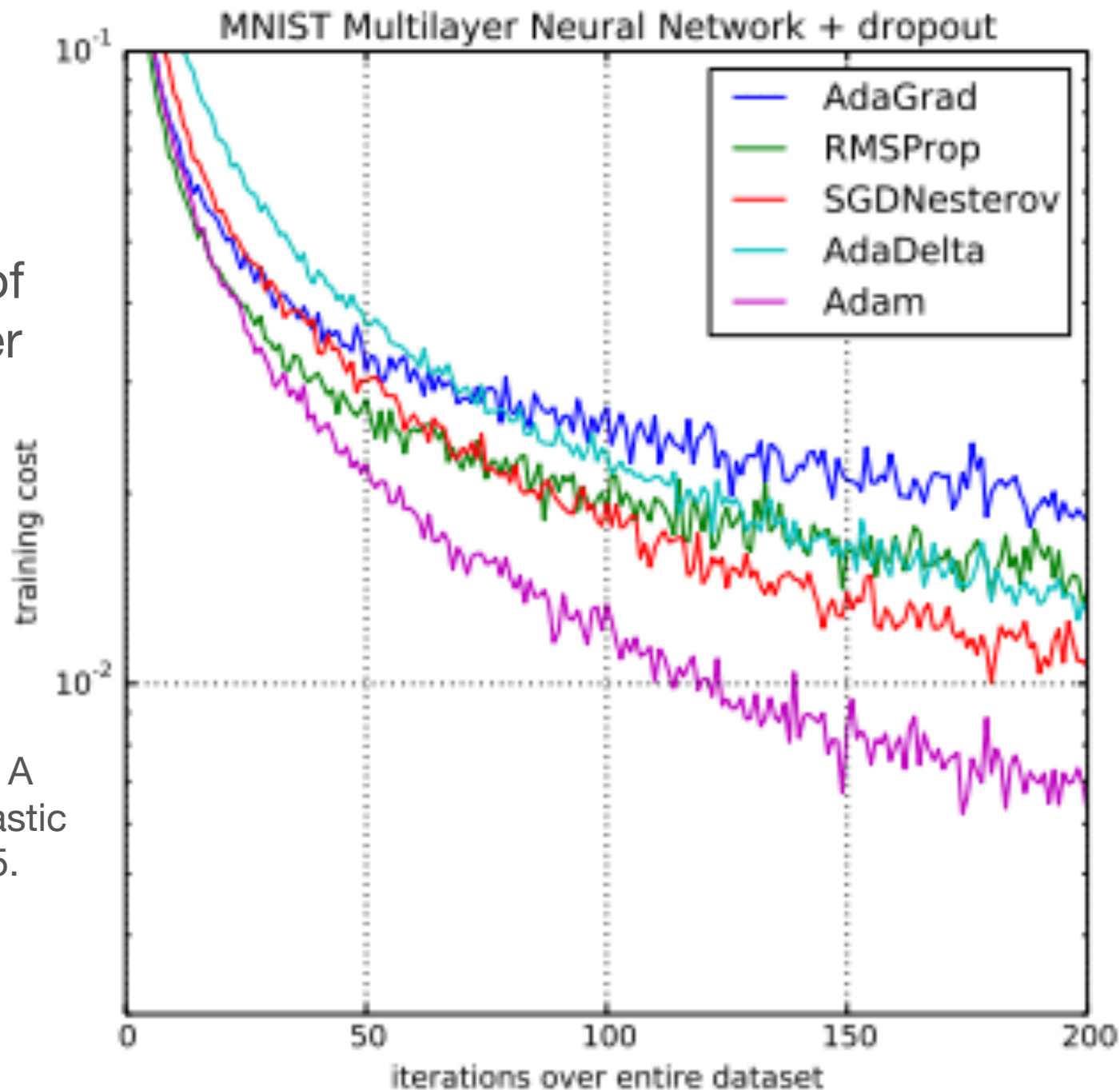
$$w \leftarrow w - \frac{\alpha}{\sqrt{v+\epsilon}} m$$

Both v and m used to update w

Parameter values:  $\beta=0.999$ ,  $\gamma=0.9$ ,  $\alpha=0.001$  ( $10^{-5}$ , 0.3),  $\epsilon=10^{-8}$

# Comparison of Adam to Other Optimization Algorithms Training a Multilayer Perceptron

Taken from Adam: A Method for Stochastic Optimization, 2015.



## ADAM

```
CLASS torch.optim.Adam(params, lr=0.001, betas=(0.9, 0.999), eps=1e-08,  
weight_decay=0, amsgrad=False, *, foreach=None, maximize=False,  
capturable=False, differentiable=False, fused=False) [SOURCE]
```

Implements Adam algorithm.

---

**input** :  $\gamma$  (lr),  $\beta_1, \beta_2$  (betas),  $\theta_0$  (params),  $f(\theta)$  (objective)  
 $\lambda$  (weight decay), *amsgrad*, *maximize*

**initialize** :  $m_0 \leftarrow 0$  ( first moment),  $v_0 \leftarrow 0$  (second moment),  $\widehat{v}_0^{max} \leftarrow 0$

---

**for**  $t = 1$  **to** ... **do**

**if** *maximize* :

$g_t \leftarrow -\nabla_{\theta} f_t(\theta_{t-1})$

**else**

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$

**if**  $\lambda \neq 0$

$g_t \leftarrow g_t + \lambda \theta_{t-1}$

$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$

$v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

$\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t)$

$\widehat{v}_t \leftarrow v_t / (1 - \beta_2^t)$

**if** *amsgrad*

$\widehat{v}_t^{max} \leftarrow \max(\widehat{v}_t^{max}, \widehat{v}_t)$

$\theta_t \leftarrow \theta_{t-1} - \gamma \widehat{m}_t / (\sqrt{\widehat{v}_t^{max}} + \epsilon)$

**else**

$\theta_t \leftarrow \theta_{t-1} - \gamma \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$

---

**return**  $\theta_t$

---

## Adam

$$v \leftarrow \beta v + (1-\beta) (\nabla J(w))^2$$

$$m \leftarrow \gamma m + (1-\gamma) (\nabla J(w))$$

$$w \leftarrow w - \frac{\alpha}{\sqrt{v} + \epsilon} m$$



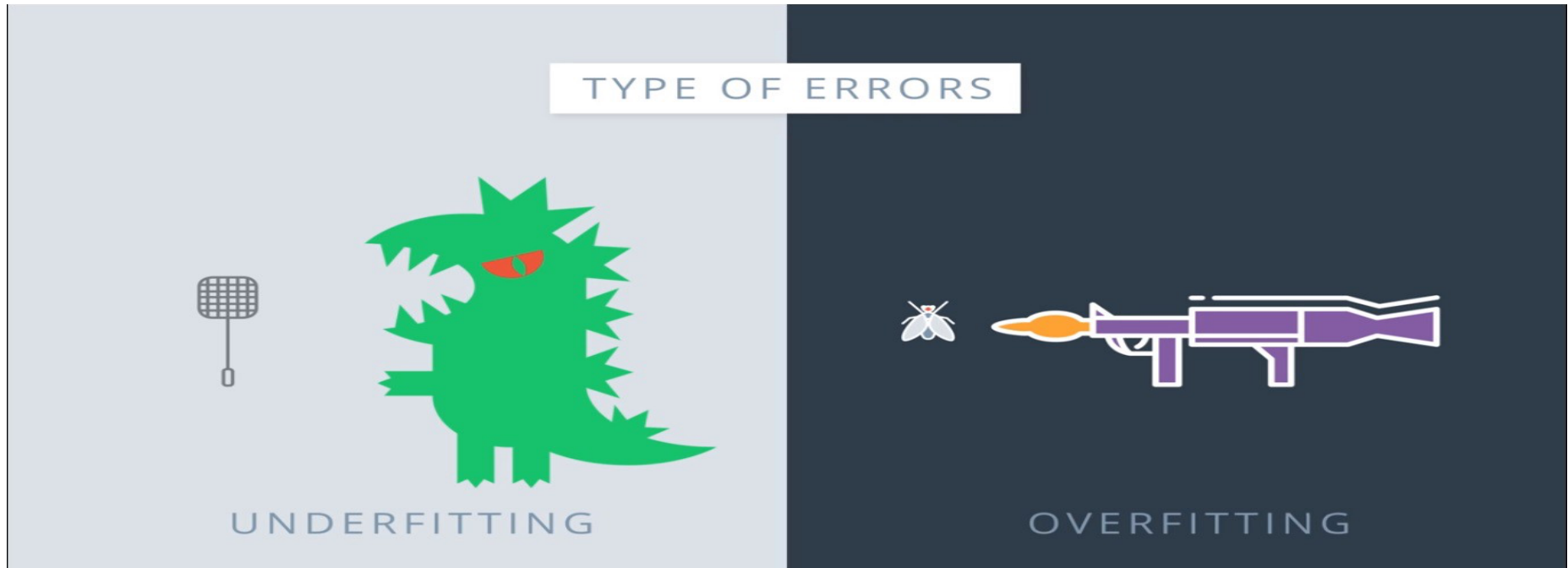
# Techniques

- Choice of activation function
- Training by mini-batches
- Finding the right learning rate
- **Overfitting and techniques to avoid**
  1. regularization
  2. Dropout (Minimal effort backpropagation)
  3. Early stopping

# Which NN is suitable for my problem?

## # of Neurons and Layers?

### Mistakes in Life

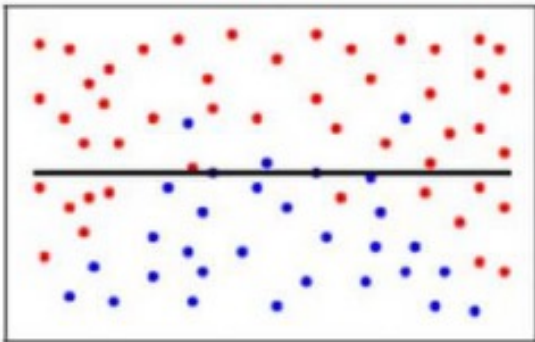


Kill Godzilla  
using a flyswatter

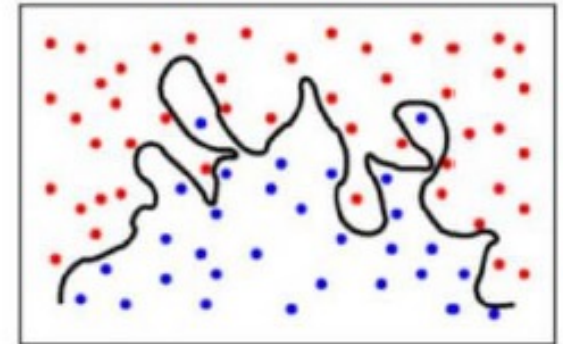
Kill a fly  
using a machine gun

# Problem of Overfitting

Underfitting



Overfitting



NN with a few neurons

NN with many layers  
and neurons

Difficult to find the right model

Approach:

Start with an overly complicated model

Then apply certain techniques to prevent overfitting

(Analogy: buying clothing without knowing the size)

# Overfitting is a familiar concept



*Memorization vs Generalization*

# Solution 1: **Weight Regularization**

Cost function  $J$  includes another term known as the regularization term

$$\text{Cost function} = \text{Loss (say, binary cross entropy)} \\ + \text{Regularization term}$$

- This regularization term is to suppress the values of weight matrices because a neural network with smaller weights leads to simpler models (reduce overfitting).

# Solution 1: **Weight Regularization**

- Observation: Large weights  $\rightarrow$  overfitting
- Add the size of the weights to loss function

$$\min \frac{1}{M} \sum_i J(h_{\theta}(\mathbf{x}_i), y_i) + \lambda \sum_j |\theta_j|$$

L1 regularization  
(Lasso Regression)

$$\min \frac{1}{M} \sum_i J(h_{\theta}(\mathbf{x}_i), y_i) + \lambda \sum_j \theta_j^2$$

L2 regularization  
(Ridge Regression)

# Solution 1: **Weight Regularization**

- Observation: Large weights  $\rightarrow$  overfitting
- Add the size of the weights to loss function

$$\min \frac{1}{M} \sum_i J(h_{\theta}(\mathbf{x}_i), y_i) + \lambda \sum_j |\theta_j|$$

L1 regularization  
(Lasso Regression)

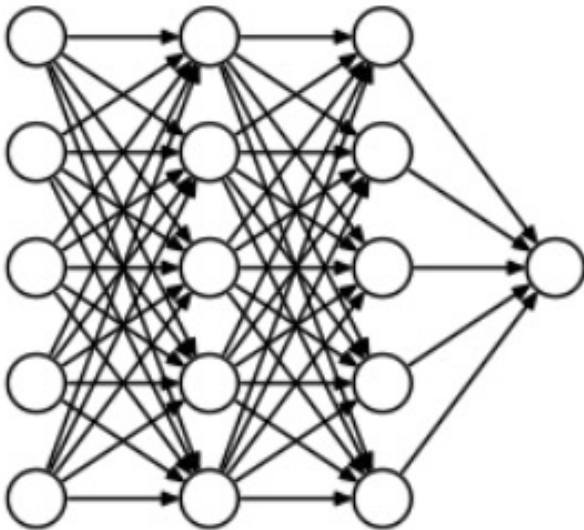
Sparse output (lots of zeros)  $\rightarrow$  feature selection

$$\min \frac{1}{M} \sum_i J(h_{\theta}(\mathbf{x}_i), y_i) + \lambda \sum_j \theta_j^2$$

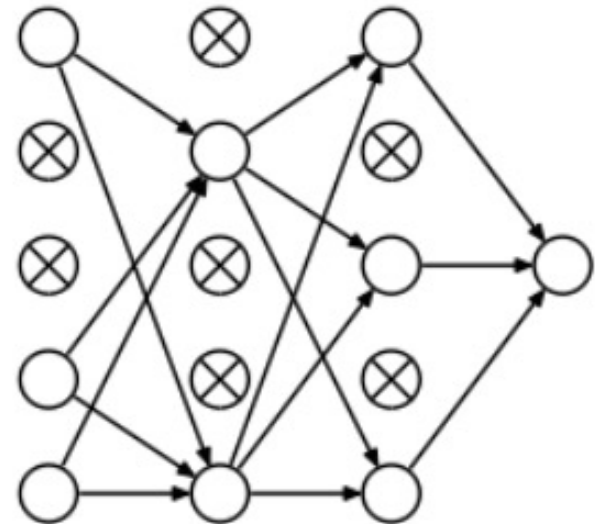
L2 regularization  
(Ridge Regression)

Computationally efficient (analytical solutions)

# Solution 2: Dropout



Standard Neural Network



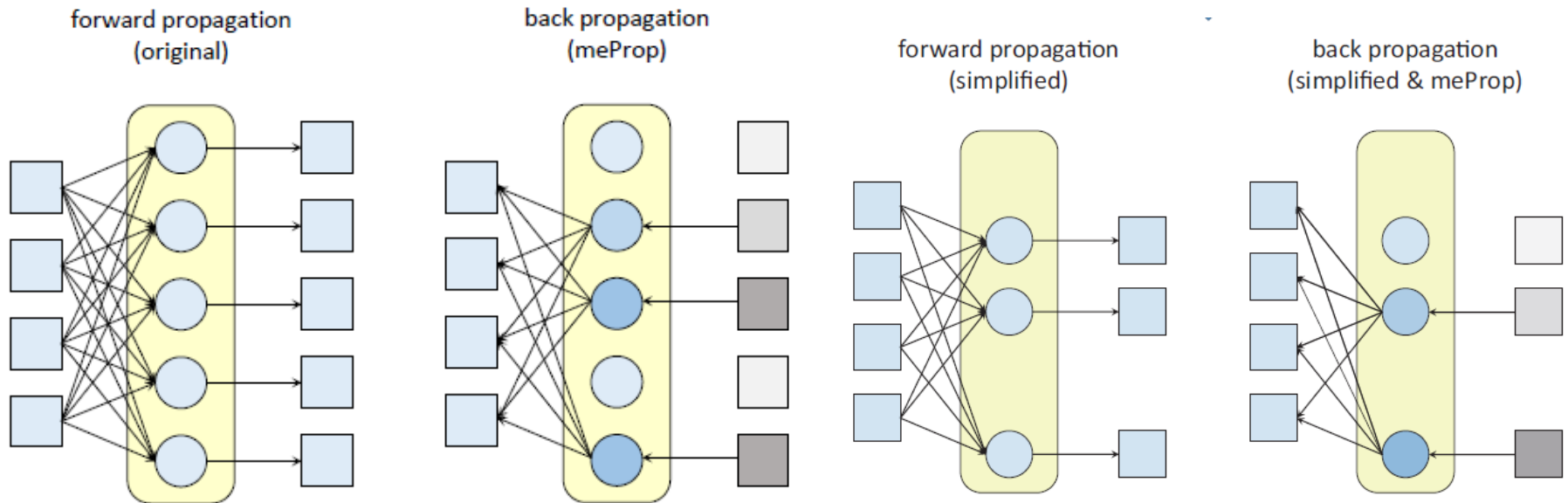
**After applying dropout –  
Randomly select 50% to  
drop out, prevents reliance  
on any one neuron**



# Minimal Effort Back Propagation Method

- Work of Xu SUN et. al. of Peking University (Nov 2017)
- Motivation 1: Reduce cost of costly backprop
- Only a **small subset** of full gradient is computed to update weights (only  $k$  rows or columns of weight matrix modified)
- Most of the time only **need to update < 5%** of the weights at each back propagation pass
- **Accuracy improved** rather than degraded
- **Motivation 2: Network size reduction**
- Size of network could be reduced by up to **9x**
- **Accuracy** of simplified model **improved** rather than degraded

# Only **keep** essential features



**Model is trained for several iterations. Keep track of the “activeness” of path – indicated by shade of neurons.**

**Model simplification – inactive neurons are eliminated, repeat as needed**

## Solution 3: **Early Stopping**

- The choice of the number of training epochs is important for training neural networks
- Too many epochs can lead to overfitting, whereas too few may result in an underfit model.
- Early stopping is a method that allows you to stop training once the model performance stops improving (how to detect this?).

# How did we learn in school?



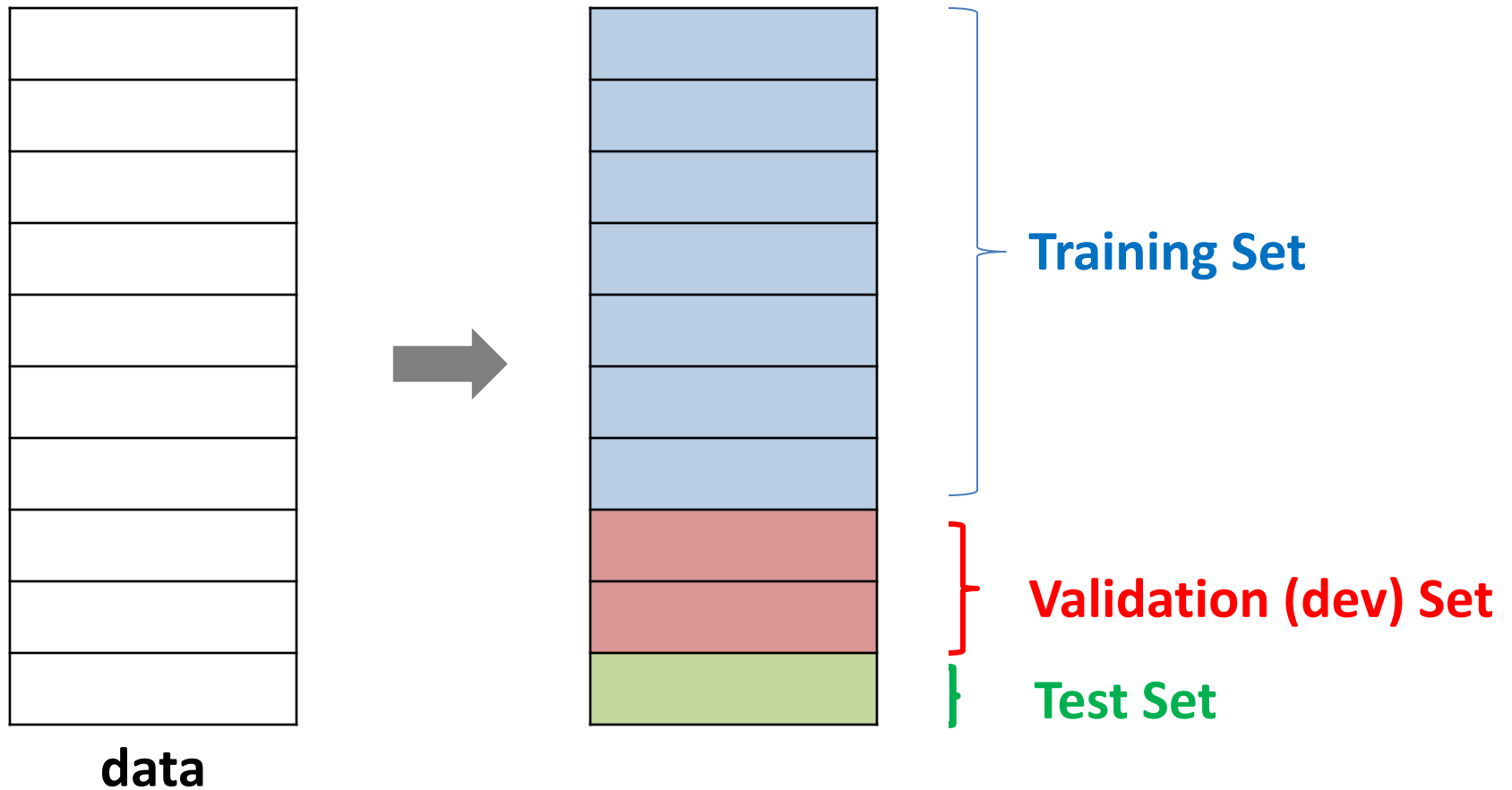
Asked **Questions**  
given **Answers**

Questions & Answers

- Practice exercises  
(**training data**)
- Mock exams (**validation**)
- Final exams (**testing**)

**Supervised Learning**

# Divide up the data



Standard split ratios: (50,25,25) or (60,20,20)

# Iterative Process

1. Train the model with the **training** set.
2. Evaluate the performance on the **training** set and **validation** set.
3. Repeat Step 1 & 2 until performance (accuracy) stops improving. [EARLY STOPPING]
4. If performance is unsatisfactory, improve model and repeat Steps 1-4.
5. Evaluate performance on **test** set.



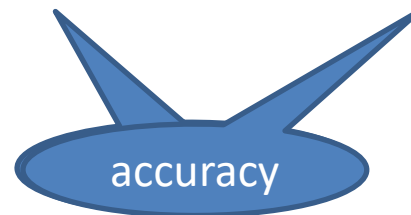
# Early Stopping by Learning Curves

- Don't give the network time to over-fit

...

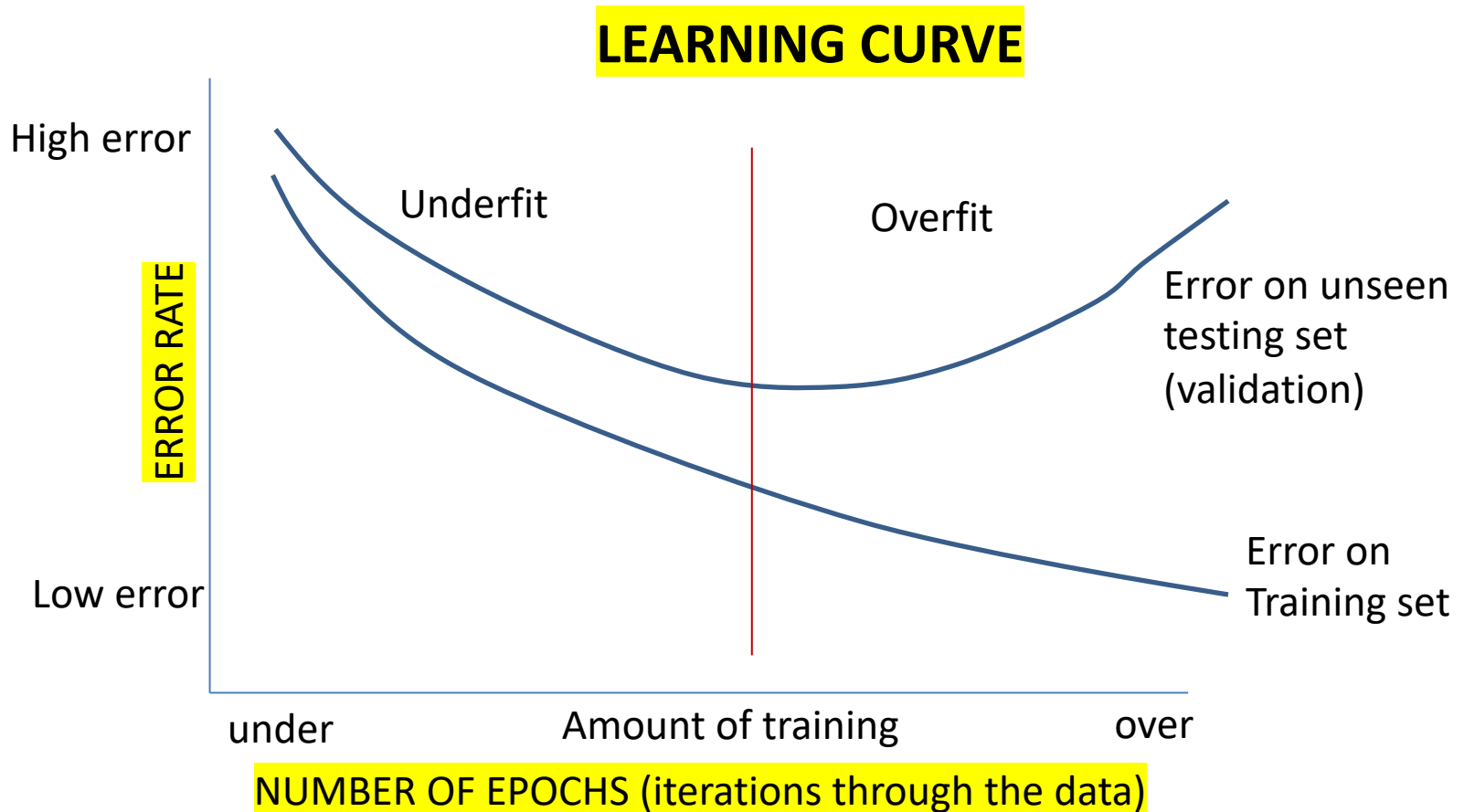
- Epoch 15: Train: 85% Validation: 80%
- Epoch 16: Train: 87% Validation: 82%
- Epoch 17: Train: 90% Validation: 85%
- Epoch 18: Train: 95% Validation: 83%
- Epoch 19: Train: 97% Validation: 78%
- Epoch 20: Train: 98% Validation: 75%

**STOP!!!**  
when  
accuracy  
stops  
improving



# Phenomenon in training

More training is equivalent to better fitting with a higher degree of polynomial





# Two Types of Errors - (Bias and Variance)

Given data set of  $(x, y)$ , where  
 $y = f(x) + \varepsilon$ ;  $\varepsilon$  is the noise

$f$  is the unknown true function

Let  $h(x)$  = trained function of fixed complexity (i.e, linear, quad, cubic..)

MSE error =  $(f(x) + \varepsilon - h(x))^2$

Let  $\bar{h} = E\{h\}$  of different training sets

$$\text{Bias}^2 = (f - \bar{h})^2$$

**Variance** =  $E[(h - \bar{h})^2]$  = error which  $h(x)$  differs from  $\bar{h}$

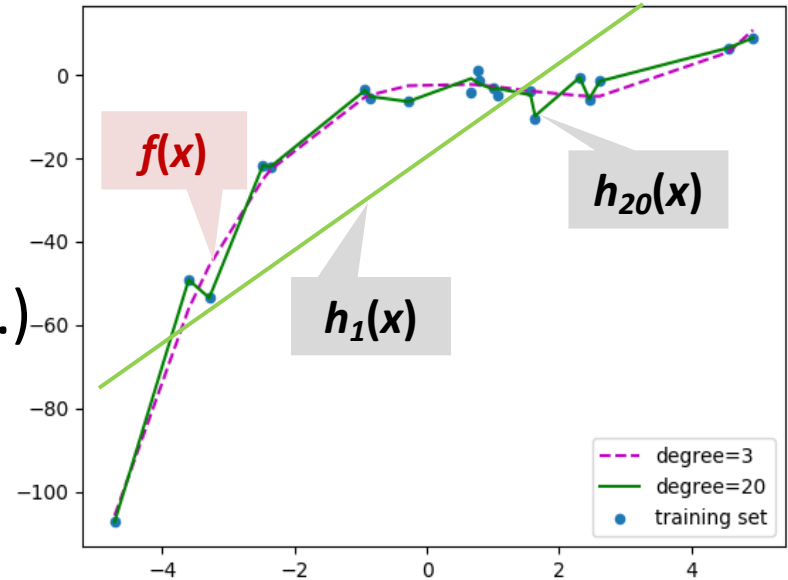
**Expected prediction error (MSE) = Bias<sup>2</sup> + Variance**

$$\text{MSE} = E \{ (f(x) + \varepsilon - h(x))^2 \} = E((f - \bar{h})^2) + E[(h - \bar{h})^2]$$

When  $h$  is of low complexity; **Bias<sup>2</sup>** is high due to model's simplistic assumptions

When  $h$  is of high complexity; **Variance** is high

- $\bar{h}$  approximates  $f$
- $h(x)$  fits the noise and deviates from  $f$



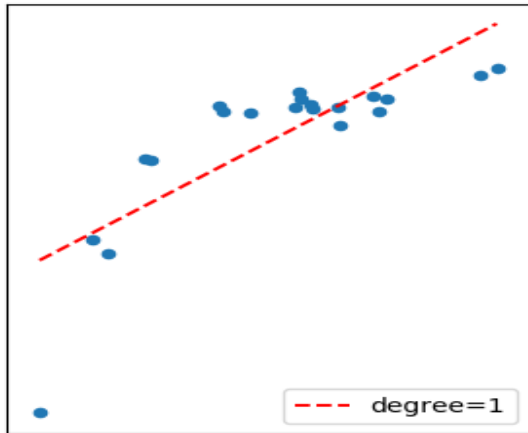
**Variance** - between expected model  $\bar{h}$  and predicted model  $h$

**Bias<sup>2</sup>** - between expected model  $\bar{h}$  and  $f$

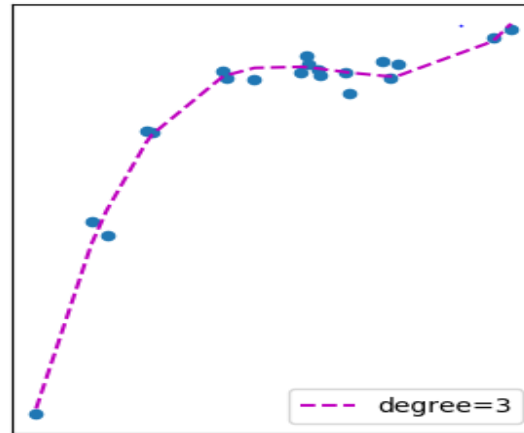
# Bias vs Variance Trade-off

Expected prediction error (MSE) =  $\text{Bias}^2 + \text{Variance}$

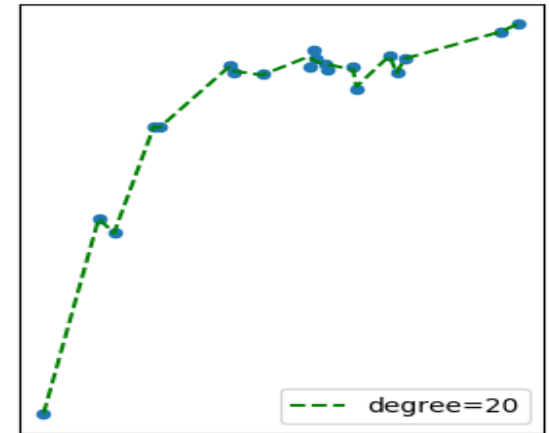
- A **high bias** means that the model is unable to capture the patterns in the data, and results in **under-fitting**.
- **High variance** means the model passes through most of the data points, and results in **over-fitting** the data.
- In order to achieve a good model that performs well both on the training and unseen data, a **trade-off** is made.



Underfit  
High Bias, Low Variance

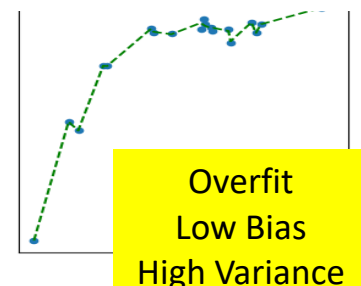
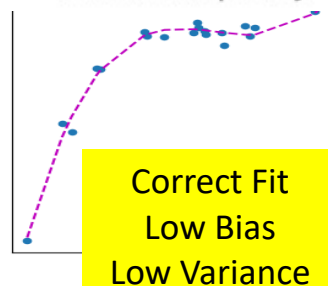
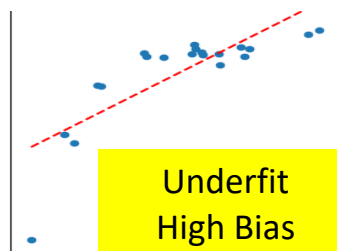
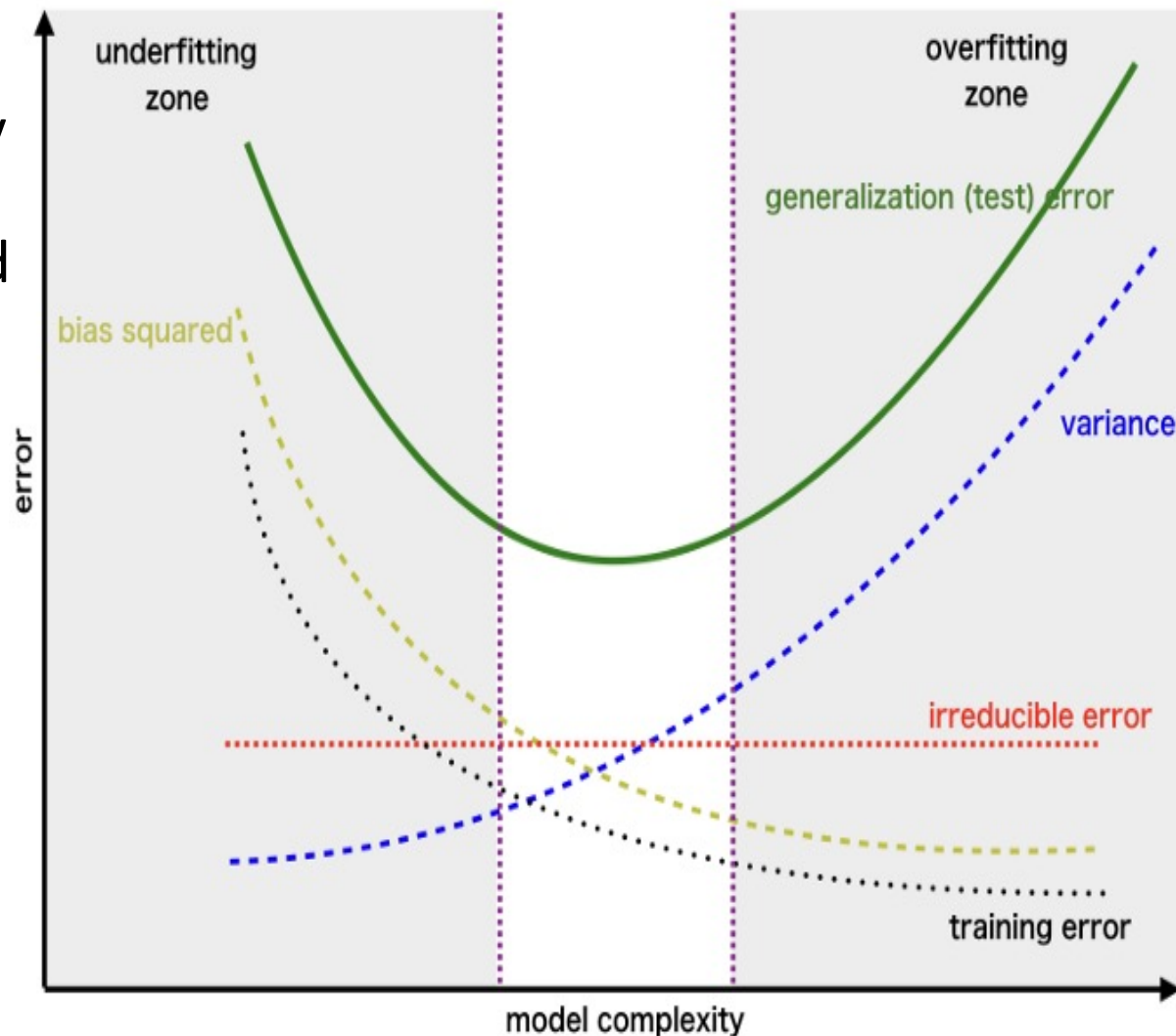


Correct Fit  
Low Bias, Low Variance

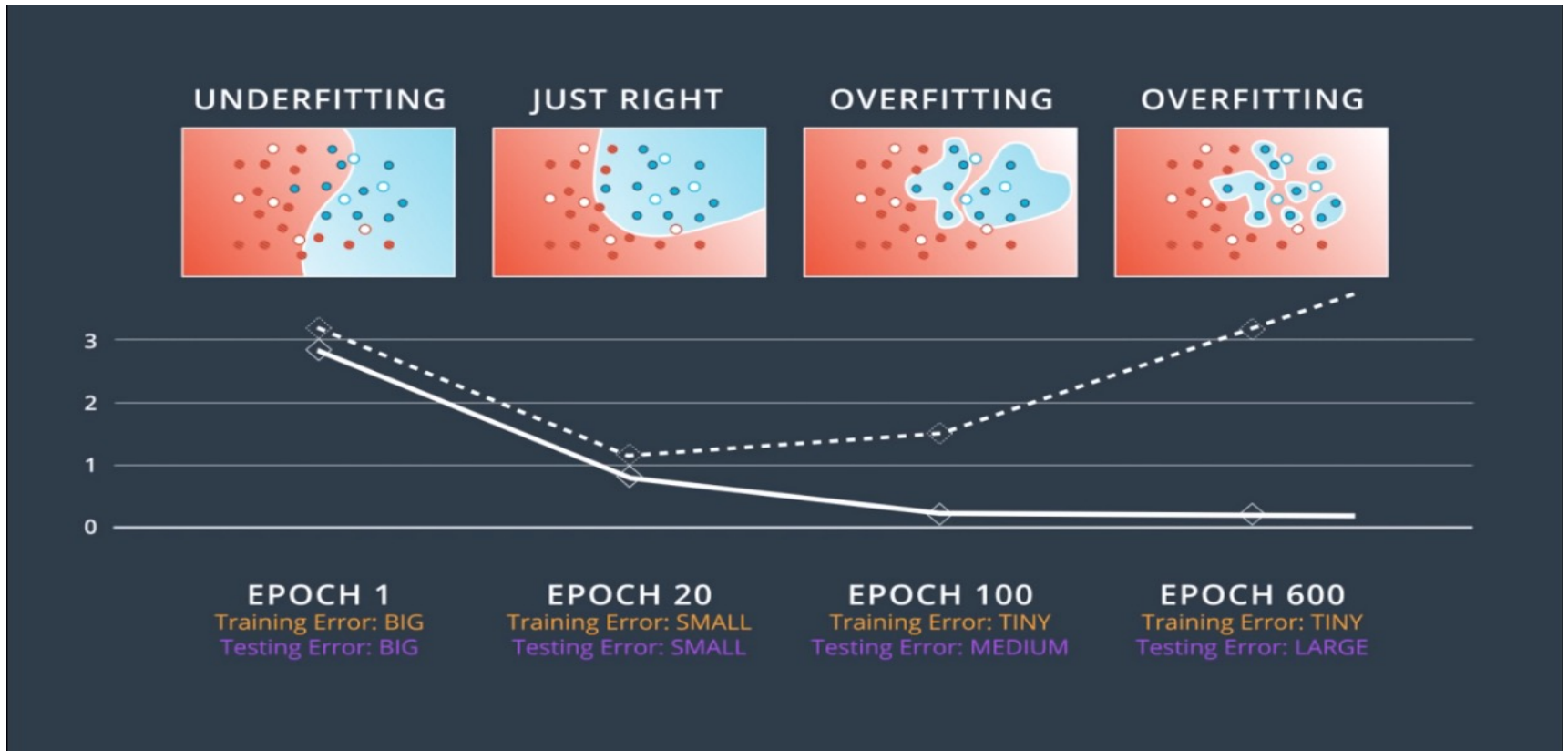


Overfit  
Low Bias, High Variance

- As model complexity increases, “bias” decreases and “variance” increases and vice-versa.
- The combined test error is U-shape.
- A machine learning model should have **low variance and low bias**



- Start with a complicated network architecture
- Evaluate these models by plotting the error in the training and testing set with respect to each epoch



- Testing set error decreases as the model generalizes well until it gets to a ***minimum point — the Goldilocks spot*** (stops under-fitting and starts over-fitting)