

Neural Network

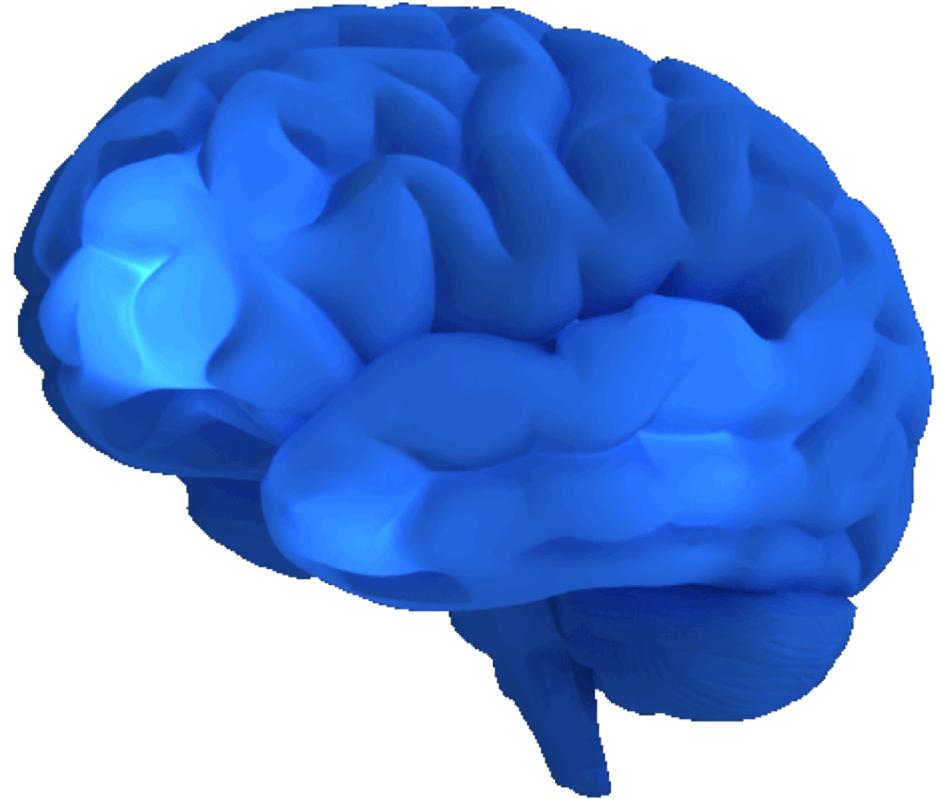
Anahita Zarei, Ph.D.

Overview

- Reading: Section 1.1.2 and e-chapter 7 from Learning from Data

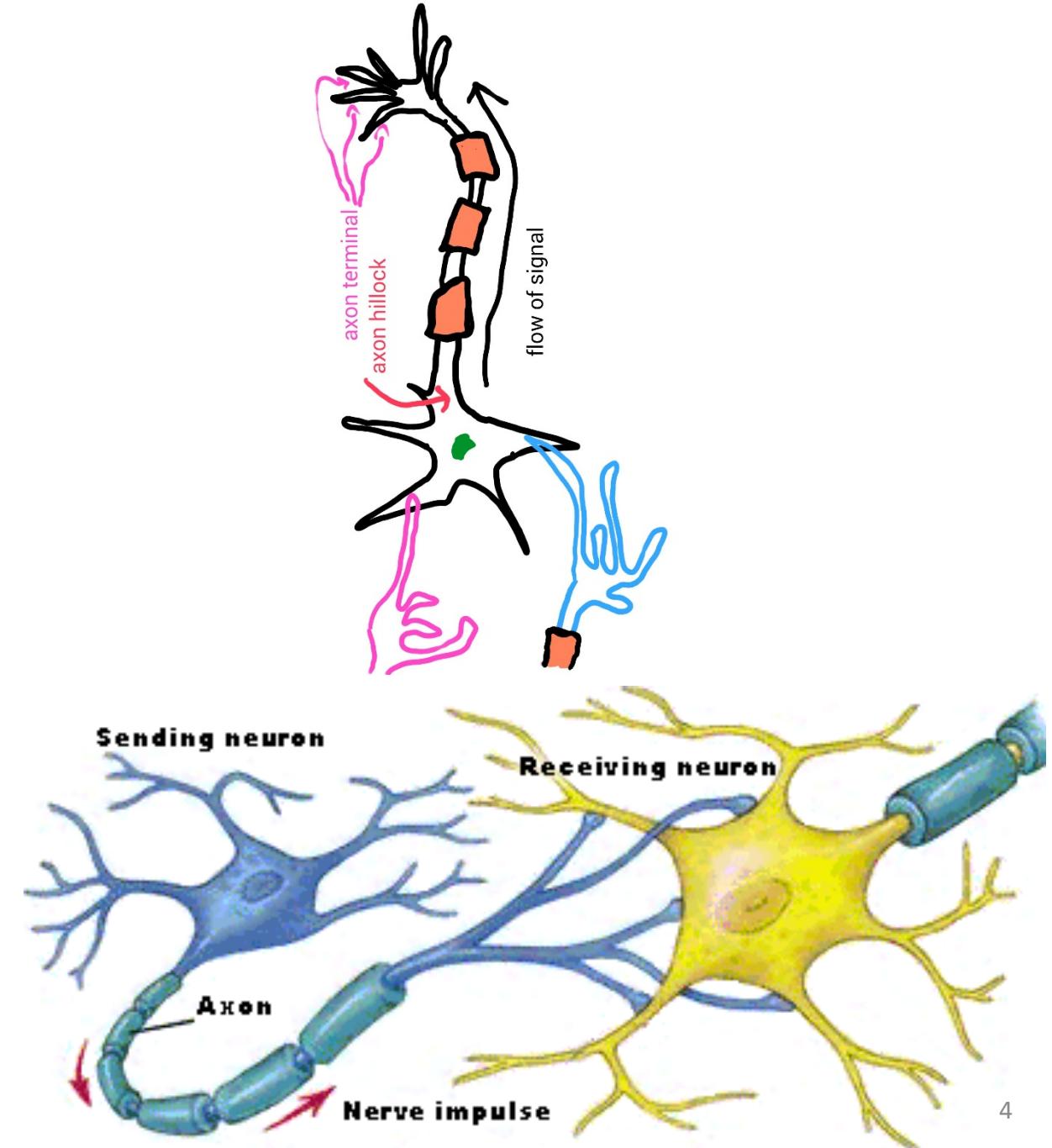
Why Neural Network?

- It's hip!
- It can learn complex functions from examples based on generalized training data in supervised learning.
- It can learn by extracting patterns and finding the underlying structure of the data in unsupervised learning.

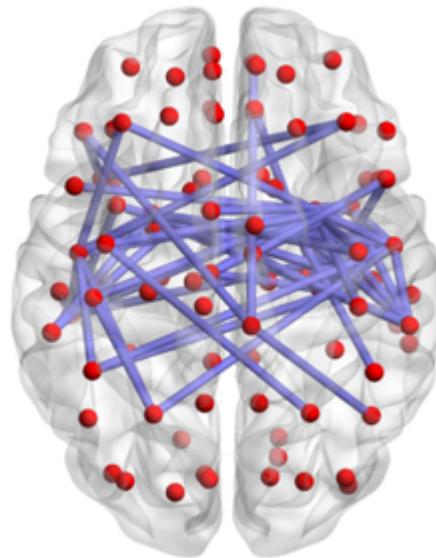
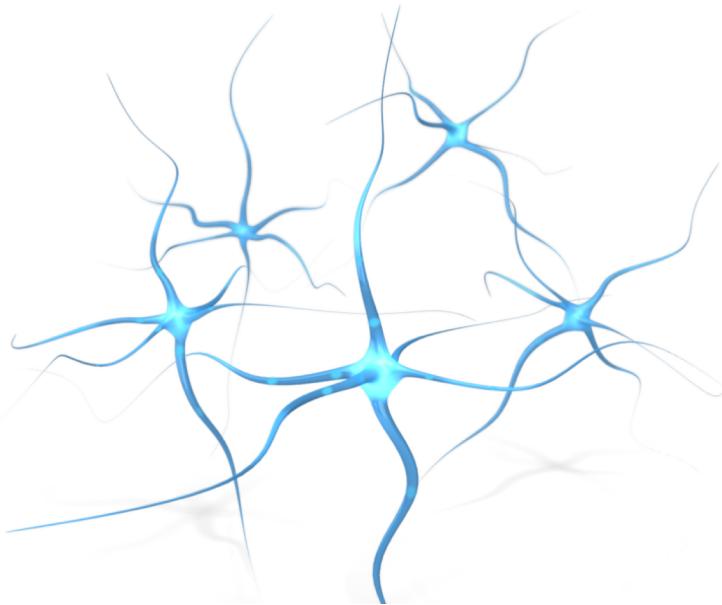


Human Neural Network

- Neuron is a fundamental cellular unit of the brain's nervous system.
- The biological neuron consists of four main parts:
 - Dendrites: resemble roots of a tree and act as input channels that receive impulses through synapses of other neurons.
 - Cell Body: processes (integrates) the signals received by the dendrites. If the combined input signal is strong enough the neuron "fires".
 - Axon: resembles tree trunk. It conducts electrical impulses and transmit information to neighboring neurons.
 - Synapse: are gaps between neurons where neurons communicate with another. This junction is filled with neurotransmitter fluid.

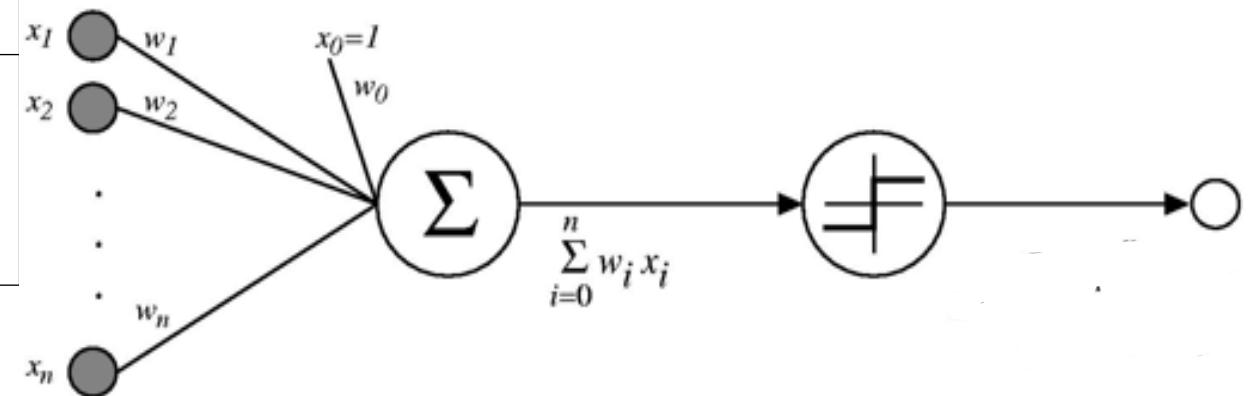
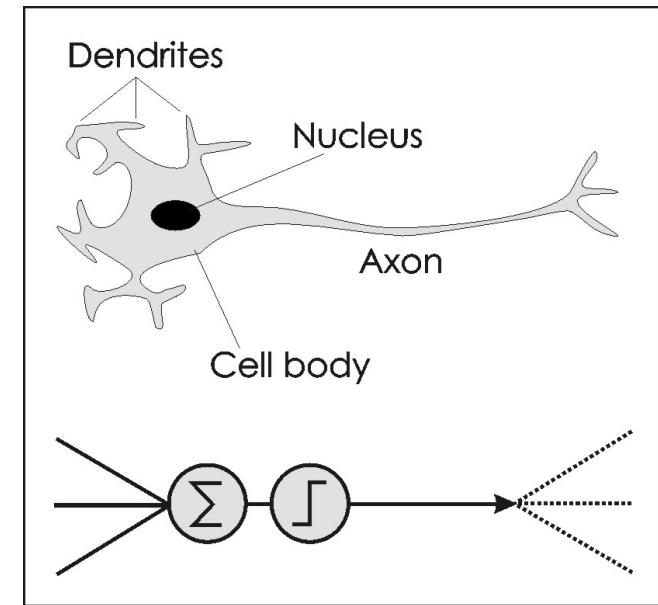


Human Neural Network



Artificial Neuron VS. Biological Neuron

Artificial Neuron	Biological Neuron
Lines that connect the input features to the summation processing element.	Dendrites
Processing element that has two parts: summation and the nonlinearity that decides if there's an action potential or not.	Cell Body
The output of a neuron that is used by other neurons	Axon



Single Neuron: Perceptron

Example: Approve or Deny Credit

This guy is my hero!



Single Neuron: Perceptron

Example: Approve or Deny Credit

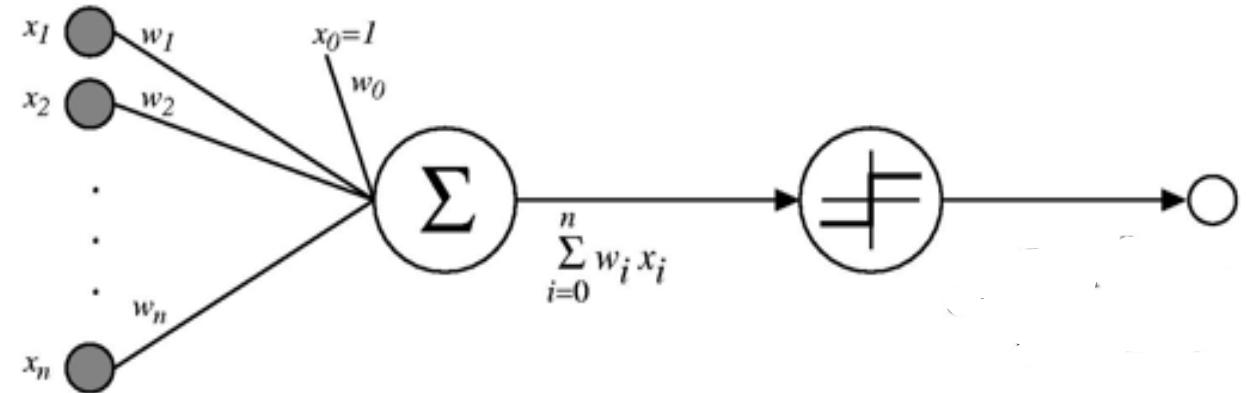
For input $\mathbf{x} = (x_1, \dots, x_d)$

Approve credit if $\sum_{i=1}^d w_i x_i > \text{threshold}$

Deny credit if $\sum_{i=1}^d w_i x_i < \text{threshold}$

$$h(\mathbf{x}) = \text{sign} \left(\left(\sum_{i=1}^d w_i x_i \right) - \text{threshold} \right)$$

$$h(\mathbf{x}) = \text{sign} \left(\left(\sum_{i=1}^d w_i x_i \right) + w_0 \right)$$

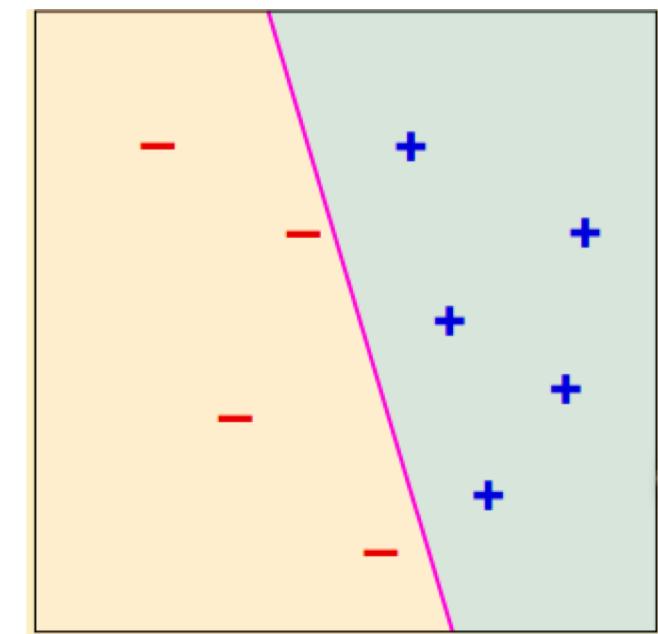
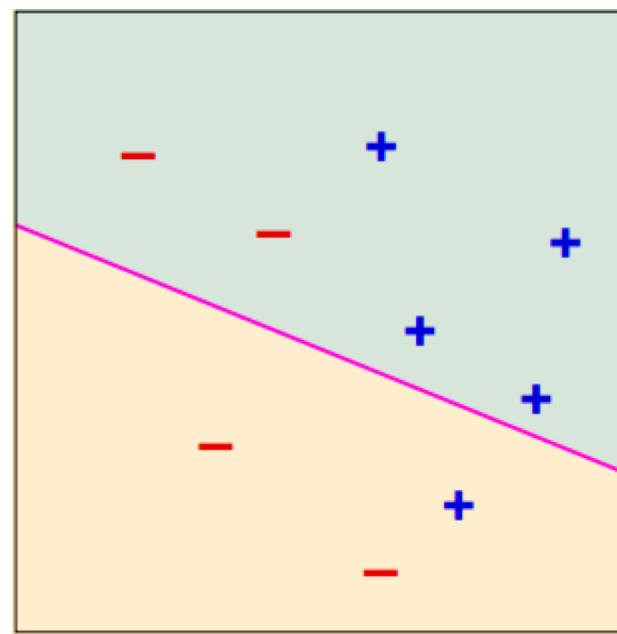


$$h(\mathbf{x}) = \text{sign} \left(\sum_{i=0}^n w_i x_i \right)$$

Perceptron Decision Boundary

If vector \vec{X} is a row vector and vector \vec{w} is a column vector then $h(\vec{X})$ can be expressed in terms of the dot products of these two vectors as follows:

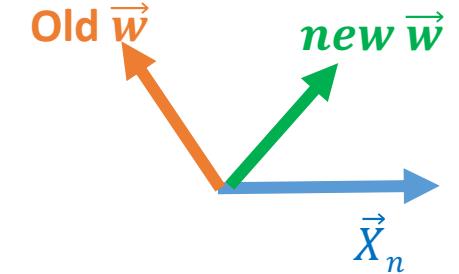
$$h(\vec{X}) = \text{sign}(\vec{X}\vec{w})$$



Perceptron Learning Algorithm (PLA)

- PLA Pseudo Code:
 - a) Choose a misclassified point (\vec{X}_n, y_n) from the following training set: $(\vec{X}_1, y_1), \dots, (\vec{X}_N, y_N)$ where vector $\vec{X} = (x_1, \dots, x_d)$
 - b) Update the weight vector with the following rule
$$\text{new } \vec{w} = \text{old } \vec{w} + y_n \vec{X}_n$$
 - c) Repeat the above process until all points are correctly classified.

Rational Behind PLA Algorithm



- Case 1: If the weight and X vector obtuse angle then the dot product will give you a negative value.

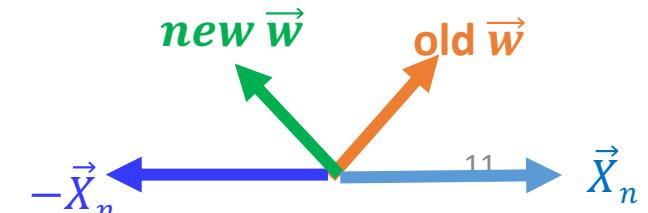
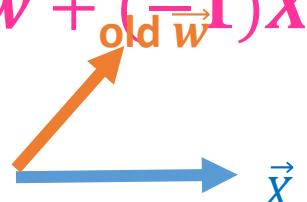
- $y_n = +1 \quad h(\vec{X}_n) = -1$

$$\begin{aligned} \text{new } \vec{w} &= \text{old } \vec{w} + y_n \vec{X}_n \\ \text{new } \vec{w} &= \text{old } \vec{w} + (+1) \vec{X}_n \end{aligned}$$

- Case 2: If the weight and X vector acute angle then the dot product will give you a positive value.

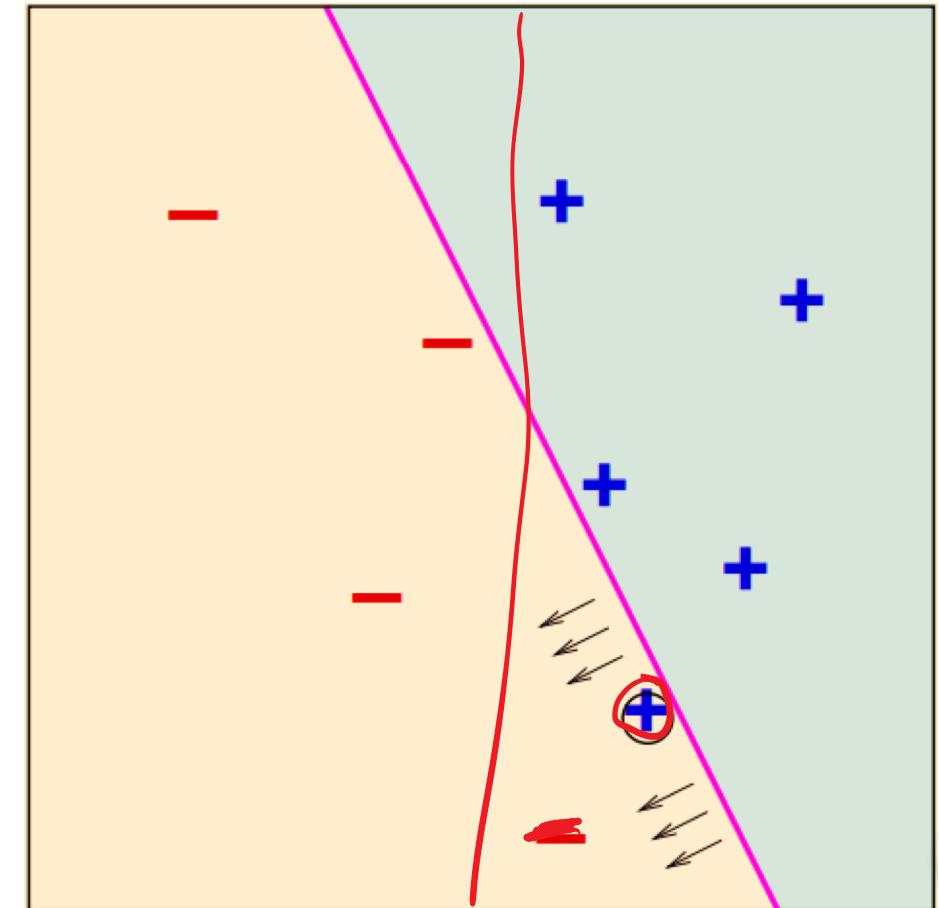
- $y_n = -1 \quad h(\vec{X}_n) = +1$

$$\begin{aligned} \text{new } \vec{w} &= \text{old } \vec{w} + y_n \vec{X}_n \\ \text{new } \vec{w} &= \text{old } \vec{w} + (-1) \vec{X}_n \end{aligned}$$



PLA- cont.

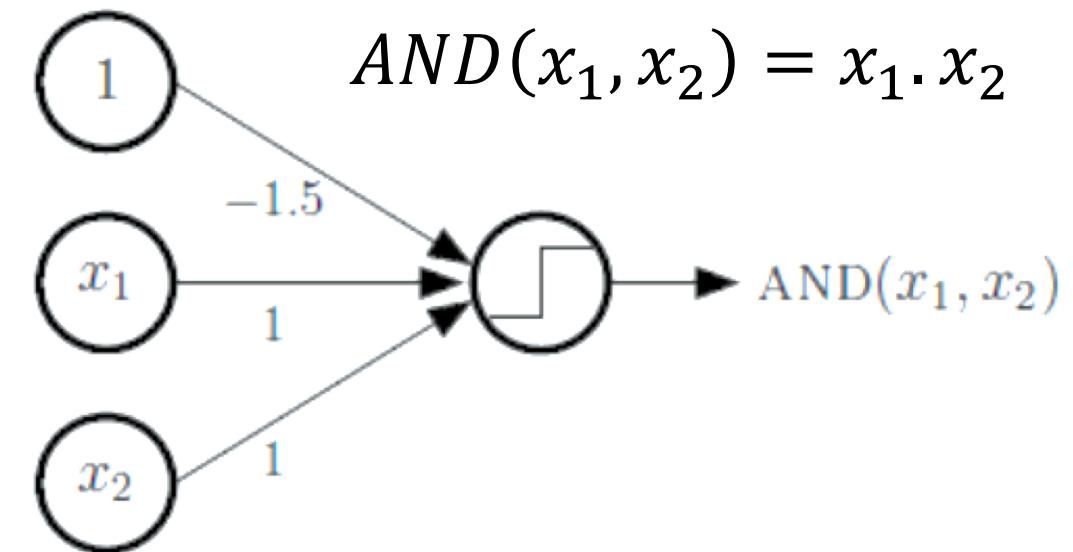
- The update rule moves the boundary in the direction of classifying point X correctly, as showed in the figure.
- PLA considers only one training example at a time. In this process, it may misclassify some of the previously correctly classified points.
- However, it's proved that there's a guarantee that PLA converges to the correct boundary decision.
- Does this mean that this hypothesis will also be successful in classifying new data points that were not in the training set?



Example: Perceptron Implementation of an AND operator

- Recall the boundary equation:
- If $w_1x_1 + w_2x_2 > \text{Threshold}$, System fires ($y=1$)
- If $w_1x_1 + w_2x_2 < \text{Threshold}$, System doesn't fire ($y=-1$)
- If $w_1 = w_2 = 1$, what threshold value implements the AND operator?
 - $-1 - 1 < \text{Threshold} \Rightarrow -2 < \text{Threshold}$
 - $1 + 1 < \text{Threshold} \Rightarrow 0 < \text{Threshold}$
 - $1 + 1 > \text{Threshold} \Rightarrow 2 > \text{Threshold}$
- Any value between 0 and 2 would work for Threshold!
- For example if Threshold = 1.5 then $w_0 = -1.5$
- $\text{AND}(x_1, x_2) = \text{sign}(x_1 + x_2 - 1.5)$.

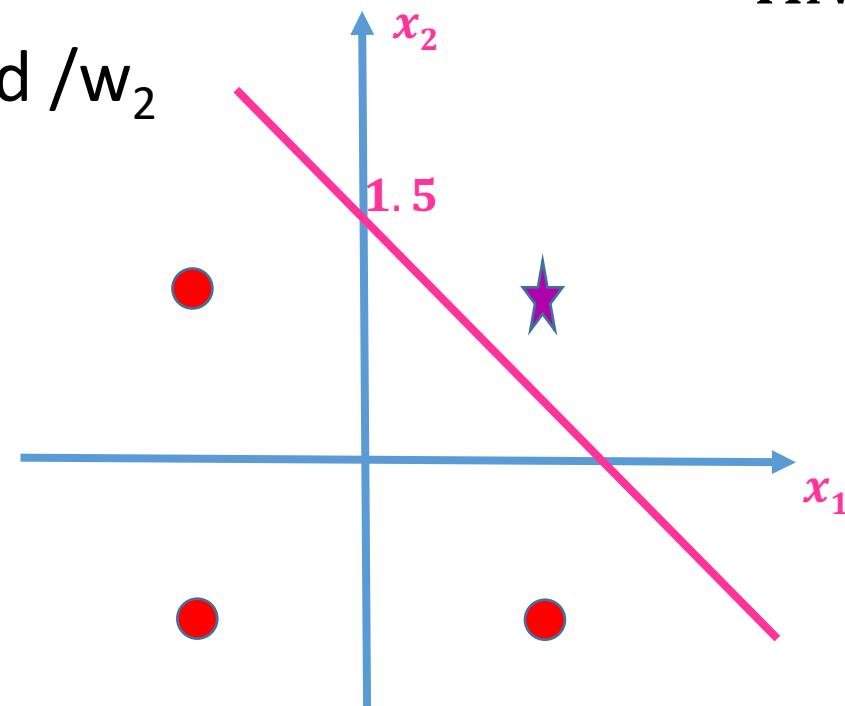
x_1	x_2	$y = \text{AND}(x_1, x_2) = x_1 \cdot x_2$
false (-1)	false (-1)	false (-1)
false (-1)	true (1)	false (-1)
true (1)	false (-1)	false (-1)
true (1)	true (1)	true (1)



Example – cont.

Boundary Decision

- $w_1 x_1 + w_2 x_2 = \text{Threshold}$, Dividing Line
- $x_2 = \frac{-w_1}{w_2} x_1 + \frac{\text{Threshold}}{w_2}$
- Slope = $-w_1/w_2$
- y-intercept = Threshold / w_2
- $x_2 = -x_1 + 1.5$



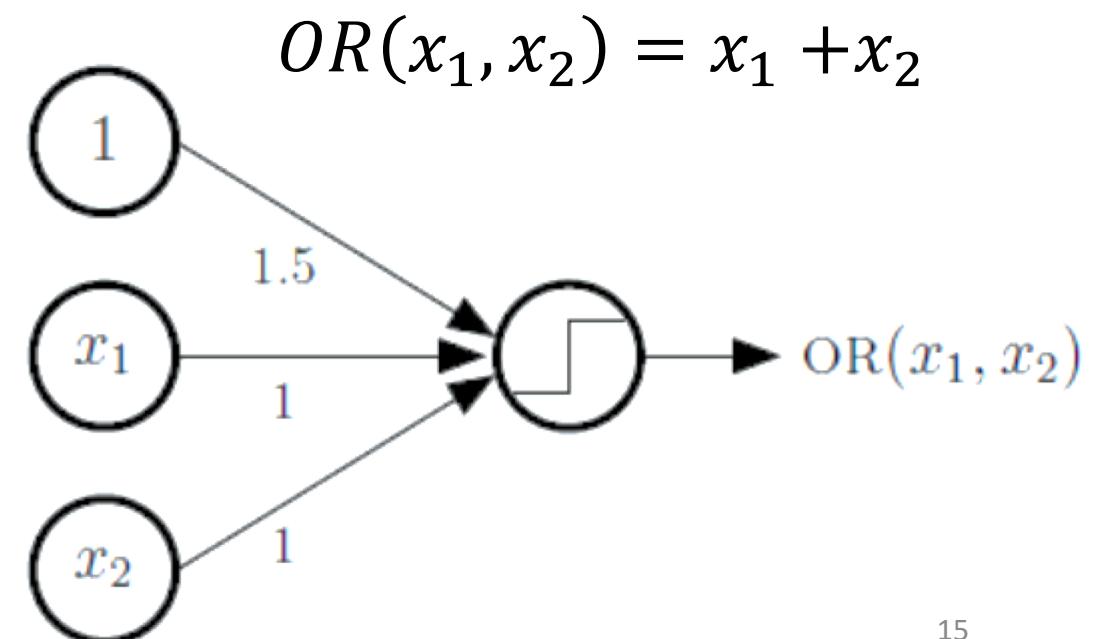
x_1	x_2	$y = \text{AND}(x_1, x_2)$
false (-1)	false (-1)	false (-1)
false (-1)	true (1)	false (-1)
true (1)	false (-1)	false (-1)
true (1)	true (1)	true (1)

$$\text{AND}(x_1, x_2) = x_1 \cdot x_2$$

Example: Perceptron Implementation of an OR operator

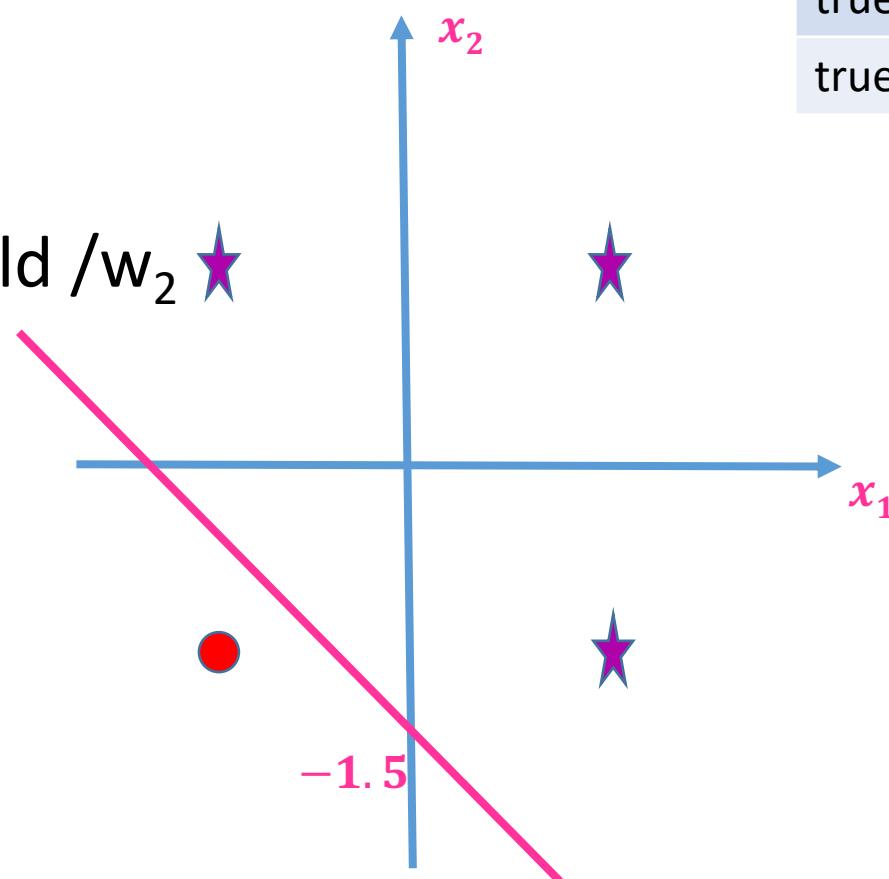
- Recall the boundary equation:
- $w_1x_1 + w_2x_2 > \text{Threshold}$ System fires ($y=1$)
- $w_1x_1 + w_2x_2 < \text{Threshold}$ System doesn't fire ($y=-1$)
- If $w_1 = w_2 = 1$, what threshold value implements the OR operator?
 - $-1 < \text{Threshold} \Rightarrow -2 < \text{Threshold}$
 - $1 > \text{Threshold} \Rightarrow 0 > \text{Threshold}$
 - $1+1 > \text{Threshold} \Rightarrow 2 > \text{Threshold}$
- Any value between 0 and -2 would work for Threshold!
- For example if Threshold=-1.5 then $w_0 = 1.5$
- $\text{OR}(x_1, x_2) = \text{sign}(x_1 + x_2 + 1.5)$.

x_1	x_2	$y = \text{OR}(x_1, x_2) = x_1 + x_2$
false (-1)	false (-1)	false (-1)
false (-1)	true (1)	true (1)
true (1)	false (-1)	true (1)
true (1)	true (1)	true (1)



Example – cont. Boundary Decision

- $w_1x_1 + w_2x_2 = \text{Threshold, Dividing Line}$
- $x_2 = \frac{-w_1}{w_2}x_1 + \frac{T}{w_2}$
- Slope = $-w_1/w_2$
- y-intercept = Threshold / w_2
- $x_2 = -x_1 - 1.5$

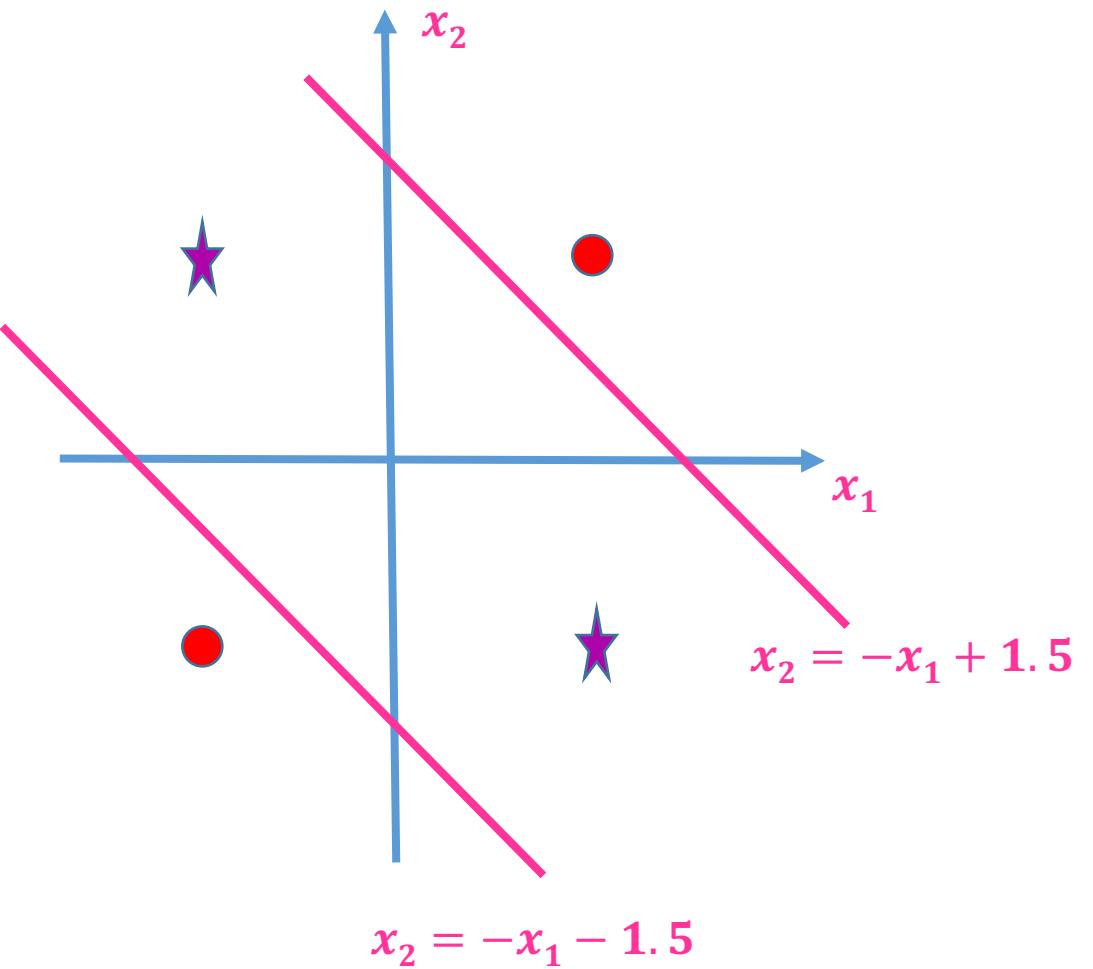


x_1	x_2	$y = OR(x_1, x_2)$
false (-1)	false (-1)	false (-1)
false (-1)	true (1)	true (1)
true (1)	false (-1)	true (1)
true (1)	true (1)	true (1)

$$OR(x_1, x_2) = x_1 + x_2$$

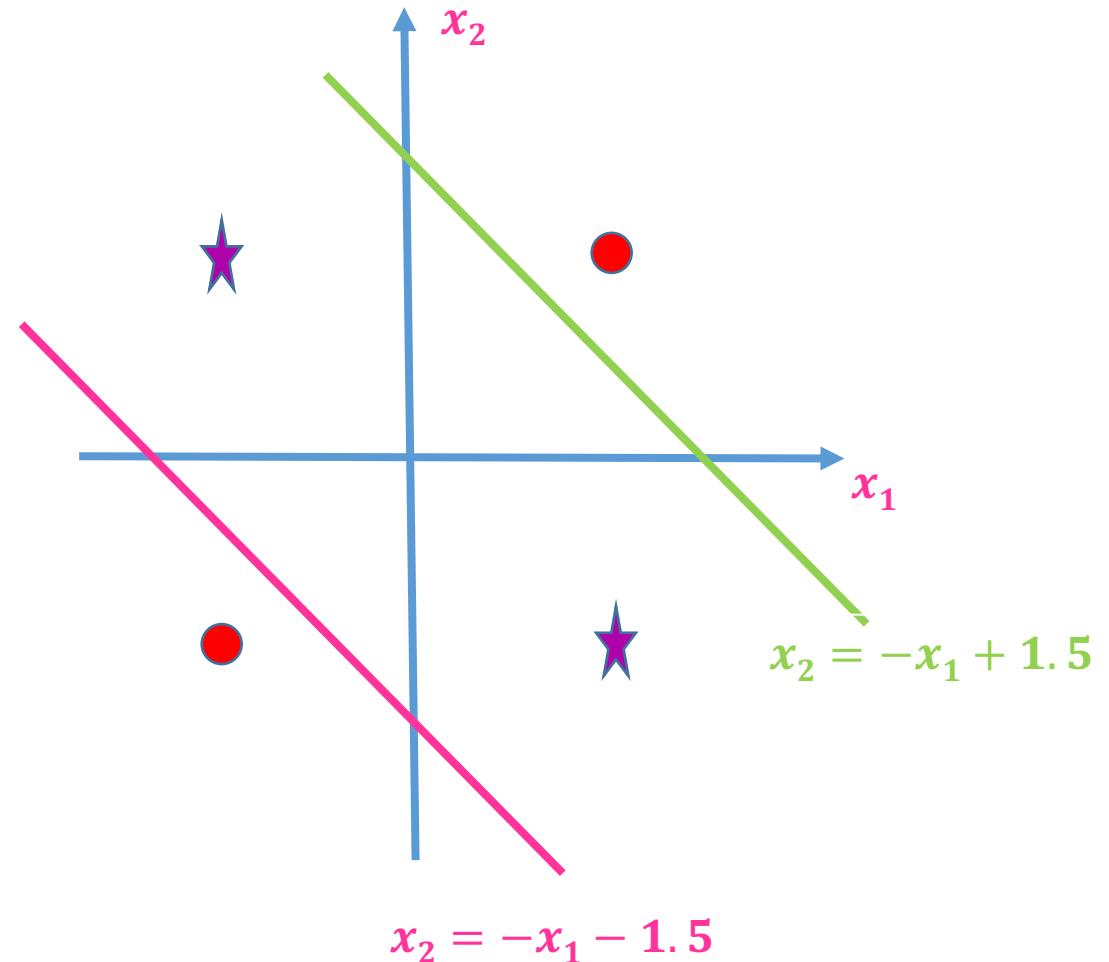
Example: Perceptron Implementation of an XOR operator?

x_1	x_2	$y = \text{XOR}(x_1, x_2)$
false (-1)	false (-1)	false (-1)
false (-1)	true (1)	true (1)
true (1)	false (-1)	true (1)
true (1)	true (1)	false (-1)



Example – cont.

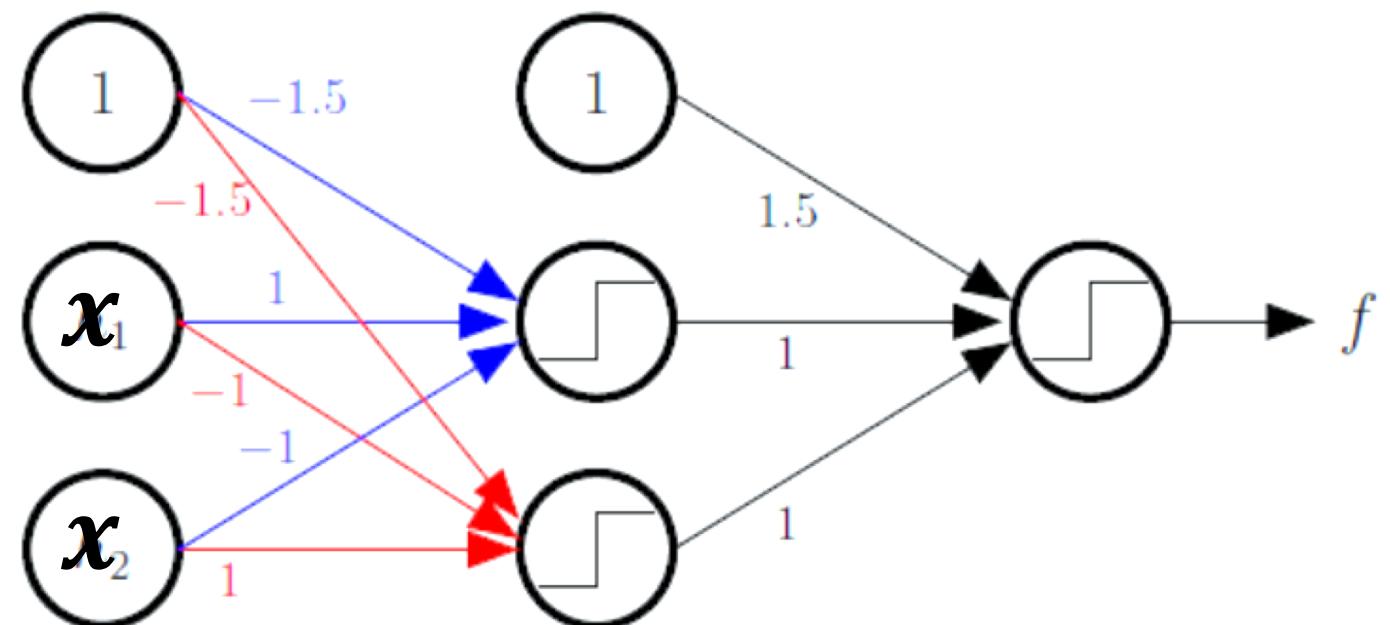
- These 2 lines will classify the points correctly.
- E.g. let $(x_1, x_2) = (-1, +1)$
- $x_1 + x_2 > -1.5 \Rightarrow (-1, +1)$ is above this line.
- $x_1 + x_2 < +1.5 \Rightarrow (-1, +1)$ is below this line.
- Intersection of the two regions will be the area between the 2 lines.



Combining Perceptrons to Achieve XOR

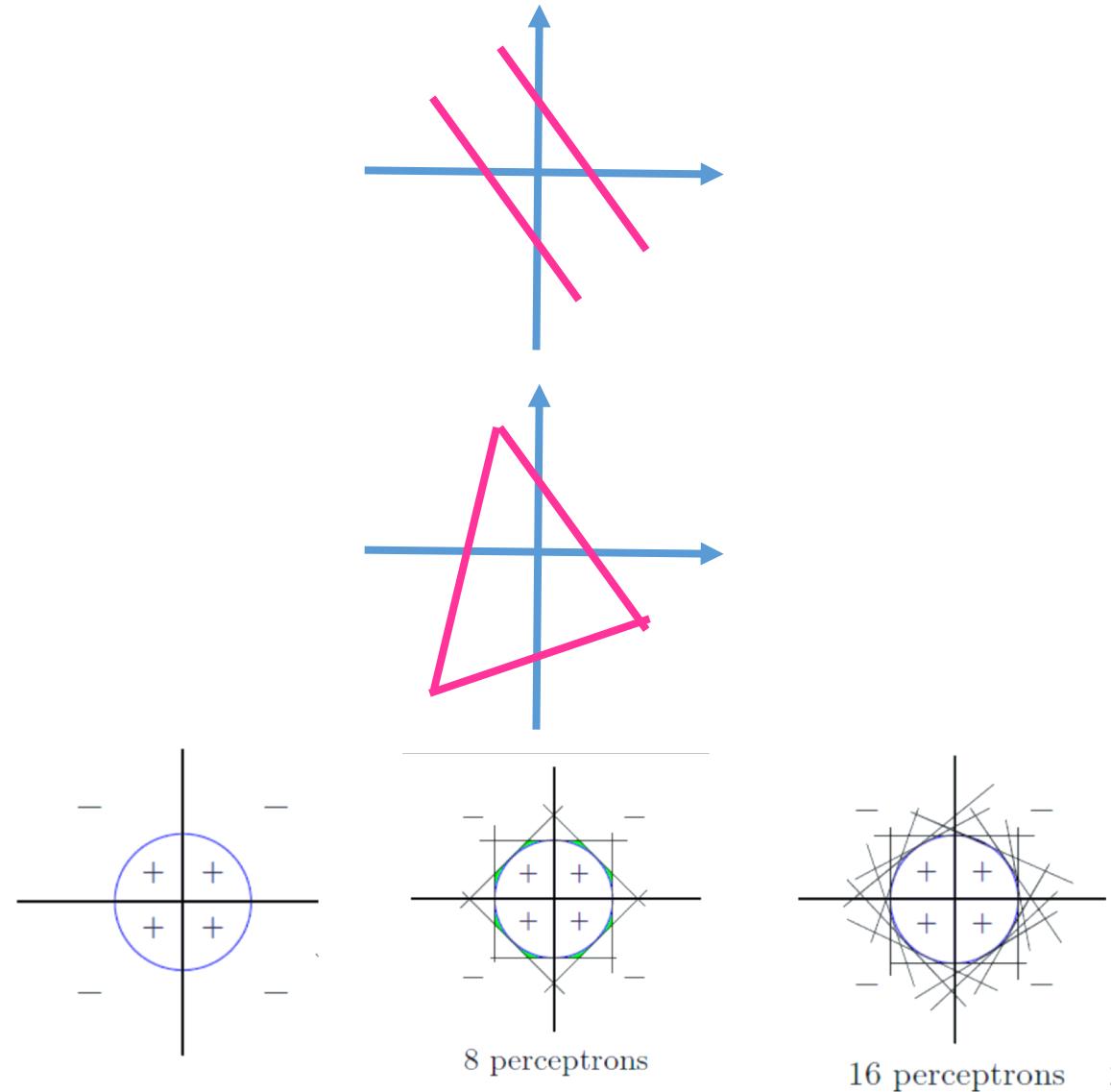
x_1	x_2	$y = \text{XOR}(x_1, x_2)$
false (-1)	false (-1)	false (-1)
false (-1)	true (1)	true (1)
true (1)	false (-1)	true (1)
true (1)	true (1)	false (-1)

$$XOR(x_1, x_2) = x_1 \overline{x_2} + \overline{x_1} x_2$$



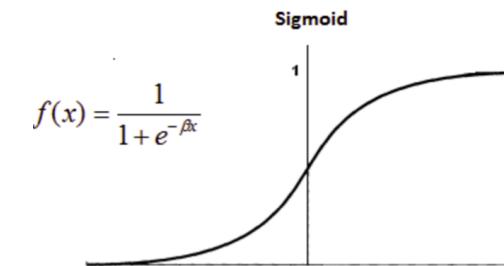
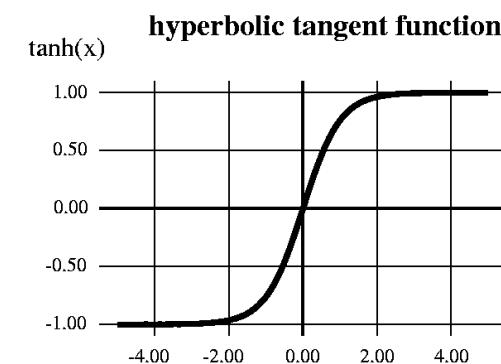
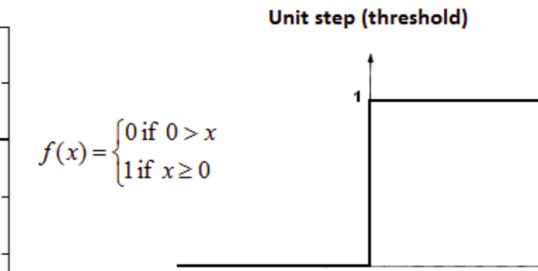
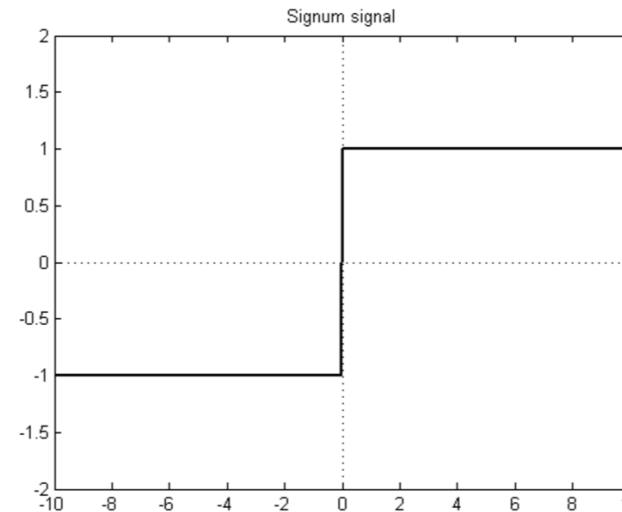
What can Multiple Layer Perceptron (MLP) learn?

- Use of 2 processing elements (excluding bias) in the hidden layer provides double division of the plane.
- Use of 3 hidden neurons (excluding bias) in the middle layer subdivides the plane with three lines, producing a triangular closed region.
- Use of additional neurons allows us to generate virtually any type of enclosed area from polygon to an approximation of a circle.



Soft Threshold Activation Function

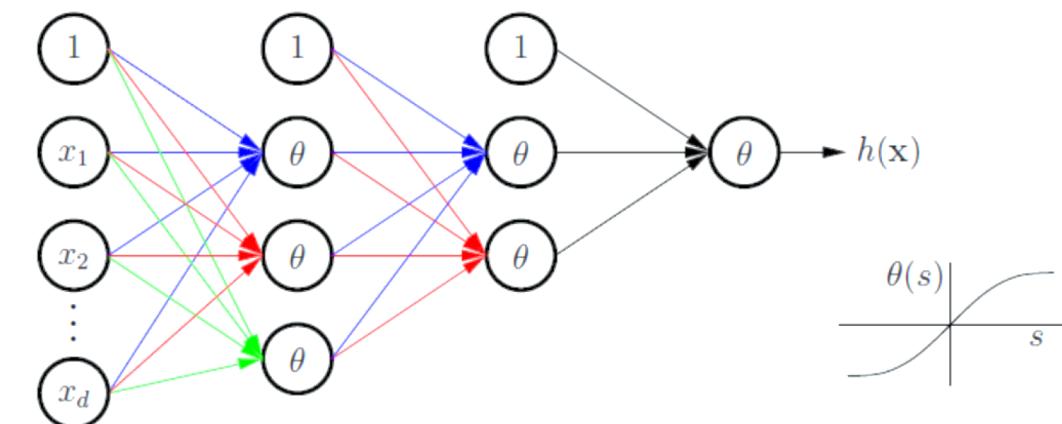
- Using hard threshold activation functions such as step, or sign are ok for a single perceptron, however it'll be a difficult combinatorial optimization problem.
- A smooth, differentiable approximation to step or sign will allow us to use analytic methods such as GD rather than purely combinatorial methods, to find the optimal weights.
- We therefore approximate, sign function by using the tanh function. (step can be approximated by sigmoid)
- Note that tanh is in between linear and the hard threshold: nearly linear for $x \approx 0$ and nearly ± 1 for $|x|$ large.



$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

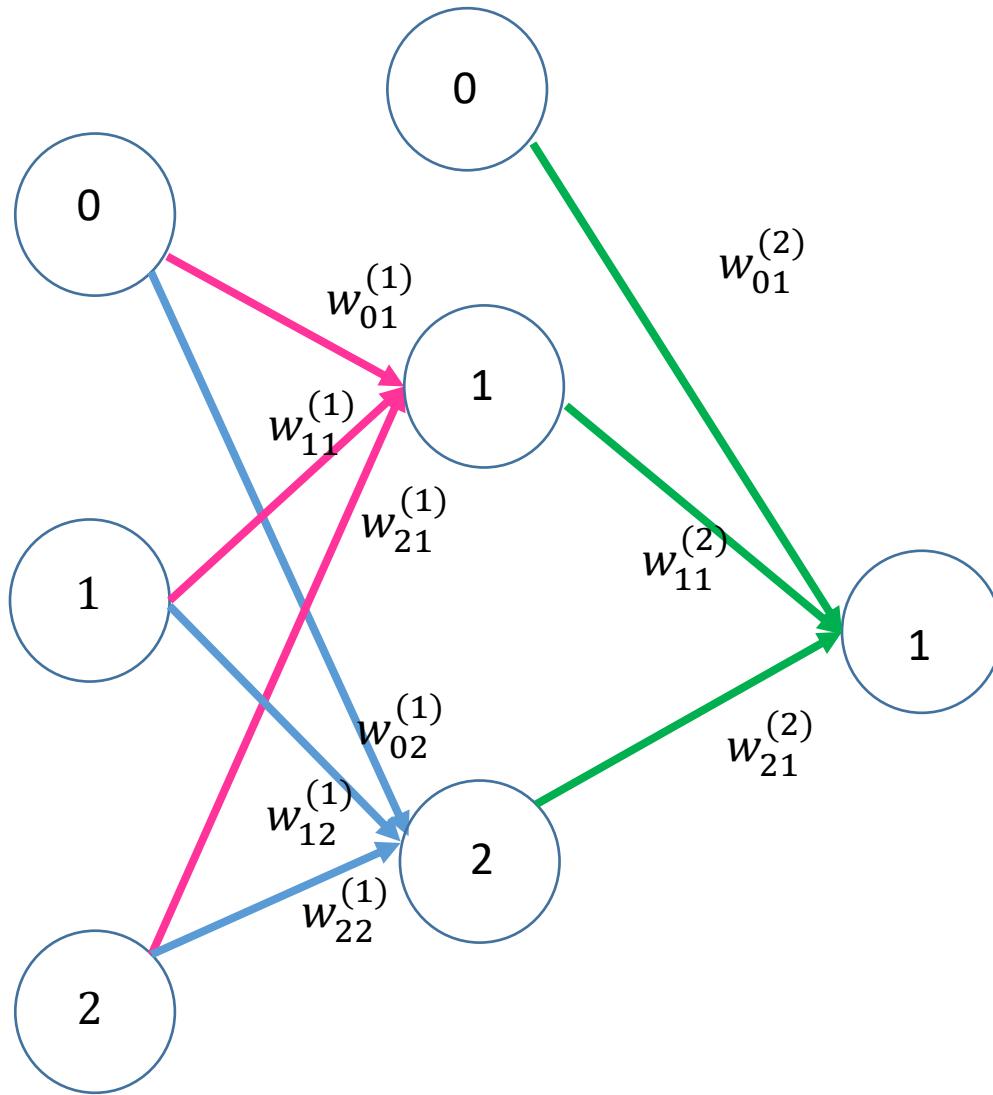
Neural Network

- The neurons are organized into a sequence of layers:
 - Input Layer: It presents data to the network. (It's not a computing layer, because the nodes have no input weights and activation function).
 - Hidden Layer: It's the intermediate layer and has no connections to the outside world.
 - Output Layer: It's the last layer that presents the output response to a given input.
- A ***feedforward*** neural network is fully connected ; there are no lateral connections between neurons in a given layer and none back to previous layers.



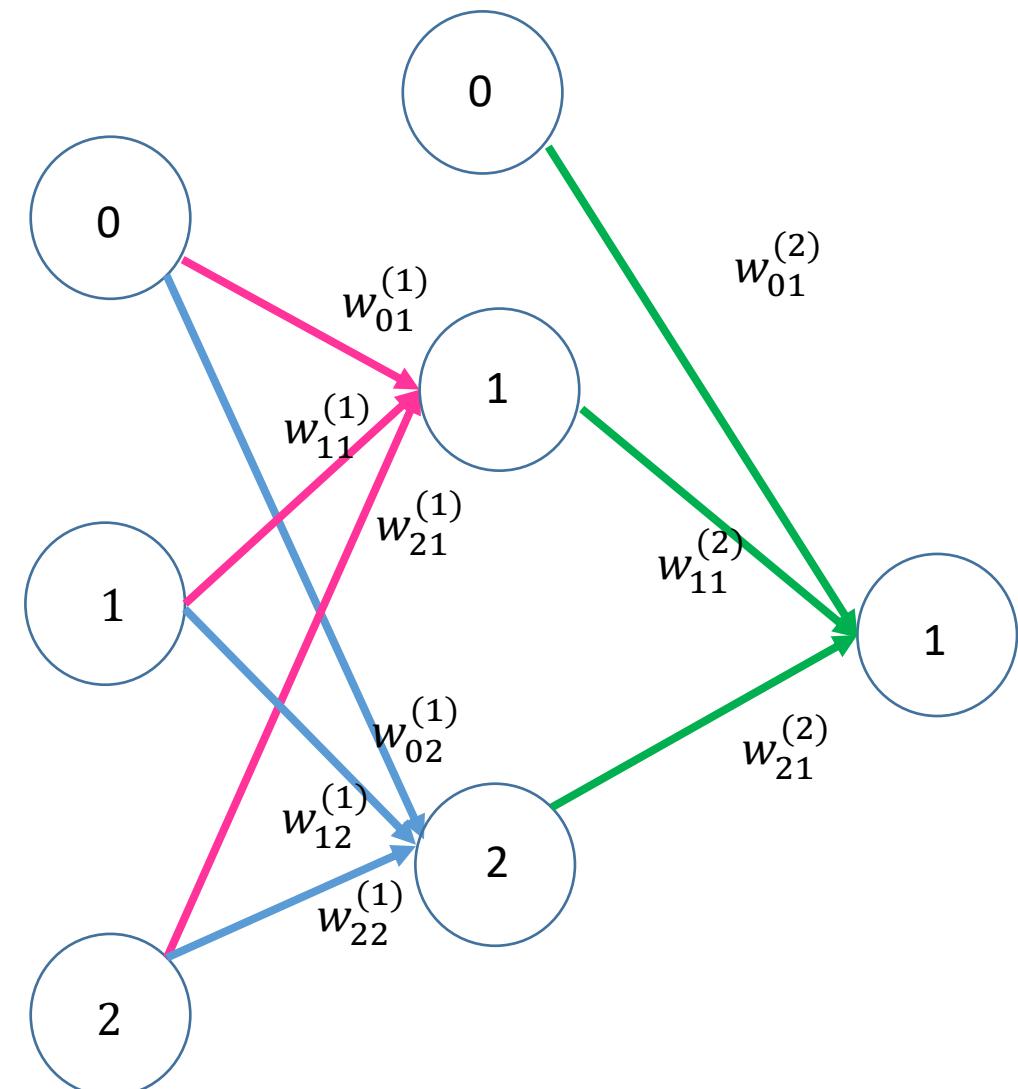
Weights

- Each of the connections between neurons has an adjustable weight.
- Each is designated by symbol $w_{ij}^{(l)}$ where i , j , and l denote starting neuron, ending neuron, and layer number, respectively.



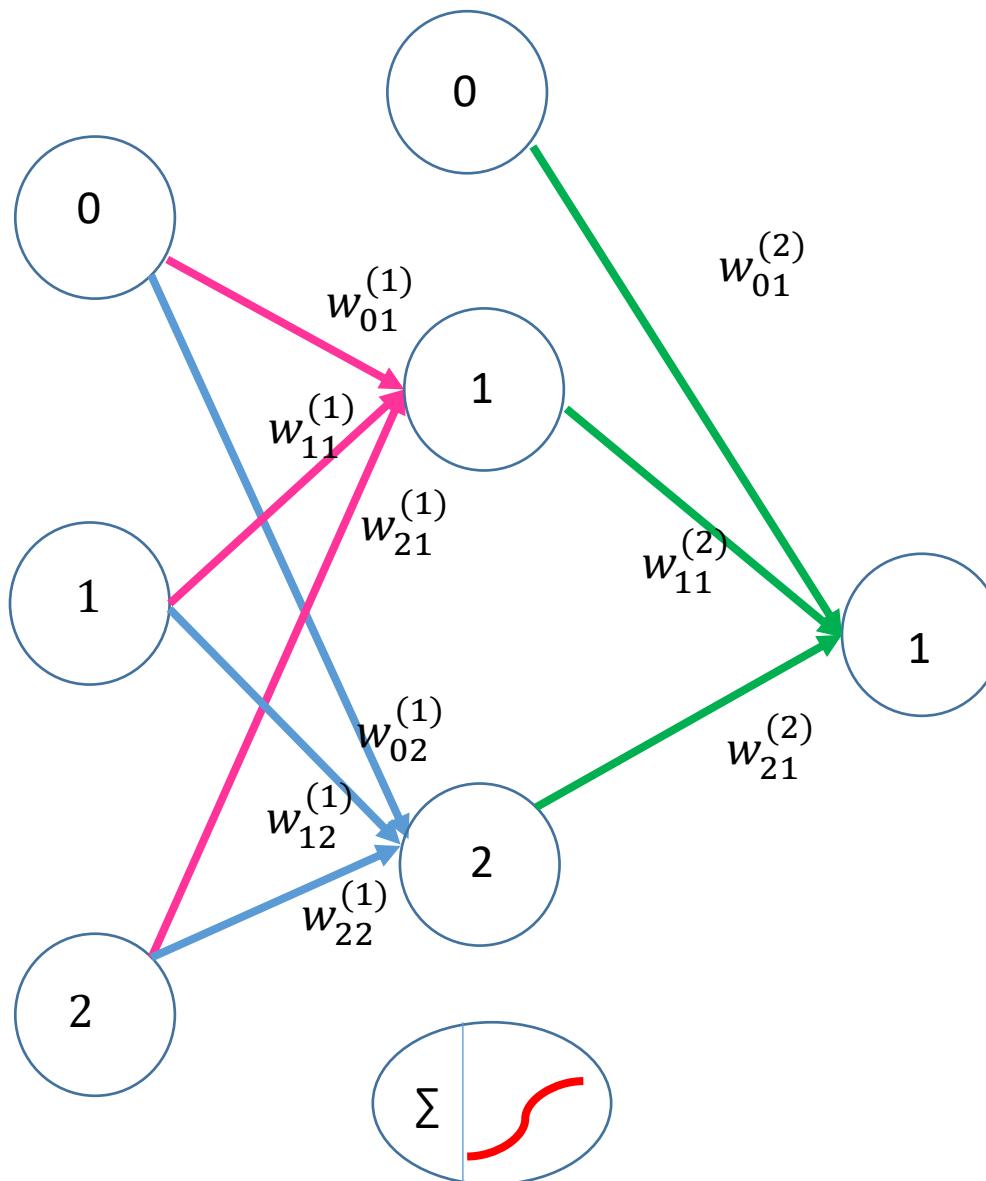
Forward Pass

- When a signal \vec{x}_1 is applied to neuron 1 in the input layer, the output goes to the hidden layer, passing through weights $w_{10}^{(1)}$, and $w_{11}^{(1)}$.
- The input signal \vec{x}_2 behaves in the similar manner, sending signals to hidden nodes through the appropriate weights.



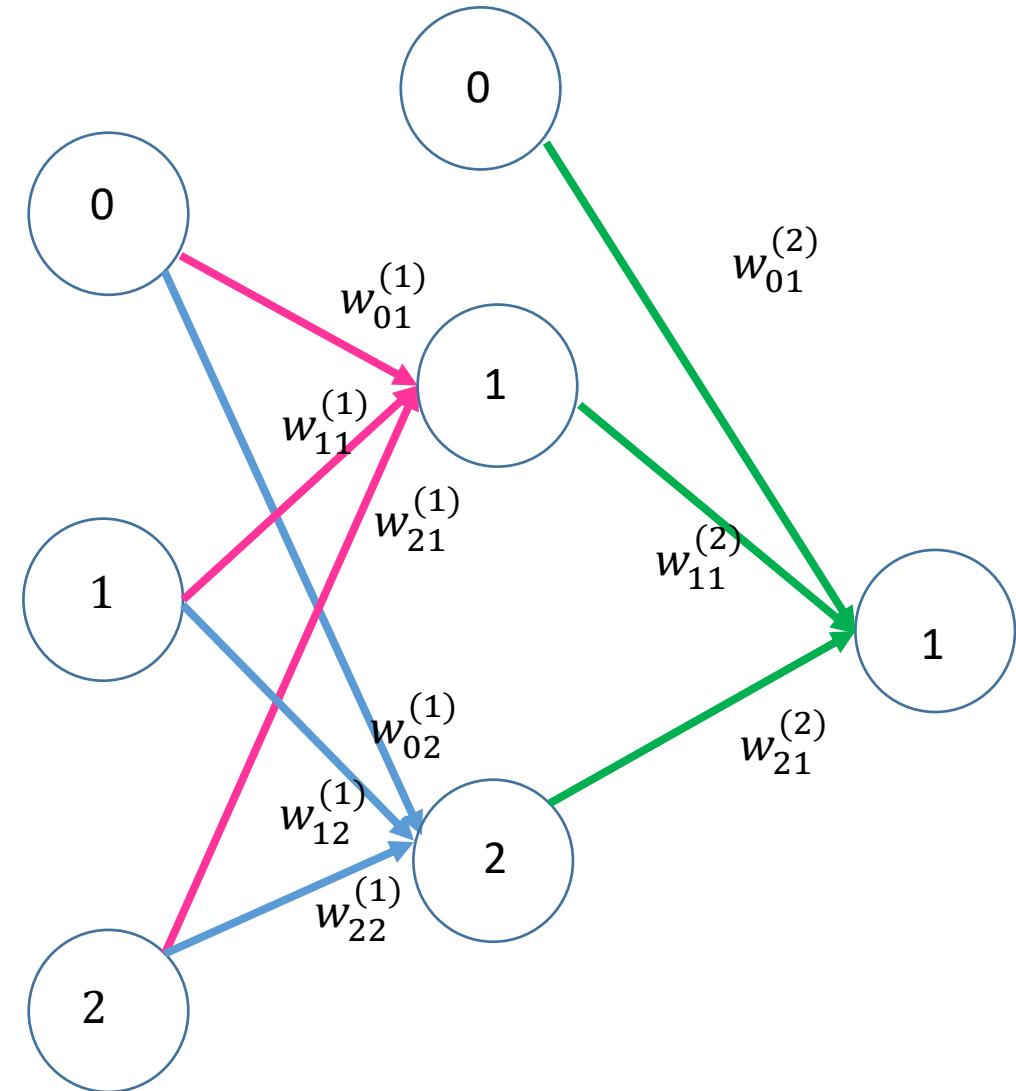
Forward Pass – cont.

- Now consider neuron 1 in the hidden layer:
- It has 3 inputs from the input layer that have been modified by the connection weights $w_{01}^{(1)}, w_{11}^{(1)}, w_{21}^{(1)}$
- The first part of this neuron simply sums up these two weighted inputs.
- The summation then passes through the second part (activation function).
- The output of this activation function is then sent to the output neuron.
- Neuron 2 in the hidden layer behaves in a similar manner.



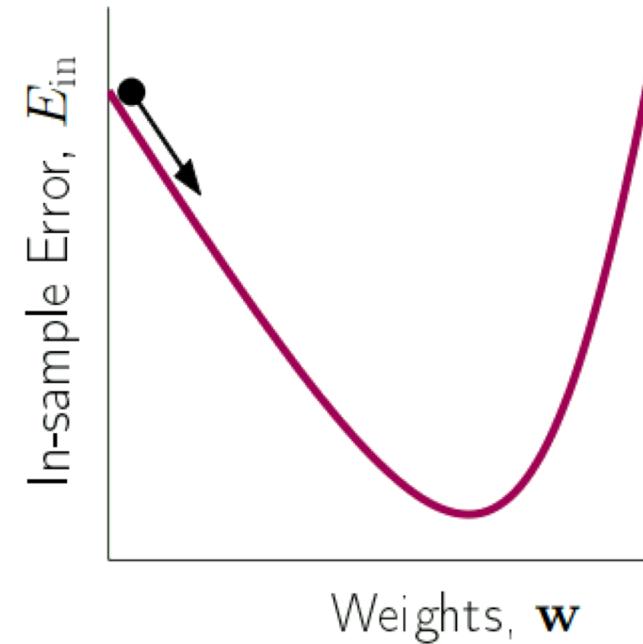
Forward Pass – cont.

- Output neuron collects the weighted inputs from hidden neurons, sum them, and pass the sums through the activation functions to produce the output.
- We use linear activation in the output neuron for regression, and a sigmoidal for classification problems.



How to learn the weights?

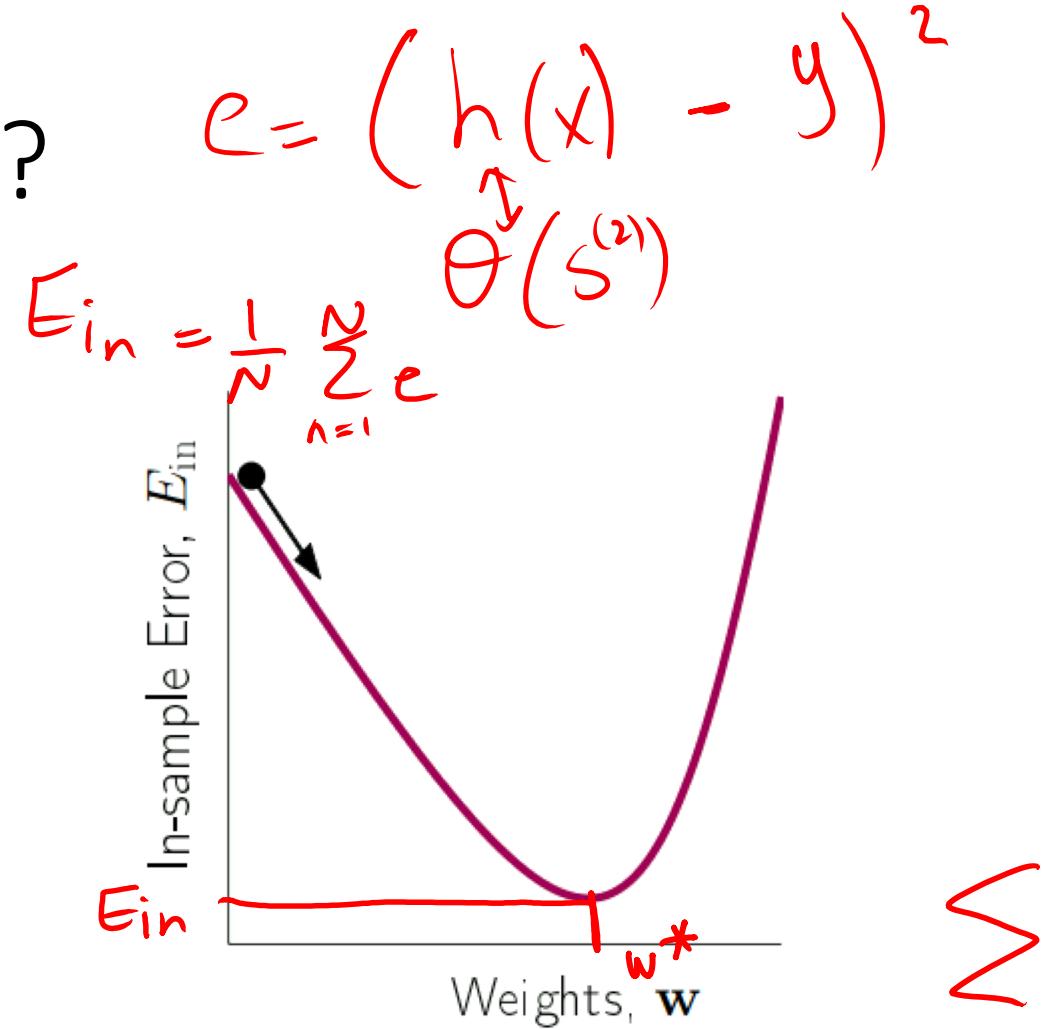
- The objective is to find a set of weights that minimizes the error.
- Recall the update rule in GD:
$$\text{new } \vec{w} = \text{old } \vec{w} - \eta \nabla E(\vec{w})$$
- Therefore we need to find the derivative of error with respect to all w's.



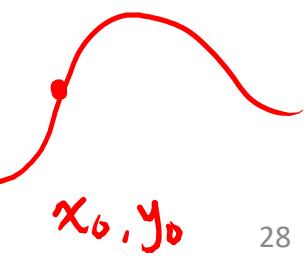
How to learn the weights?

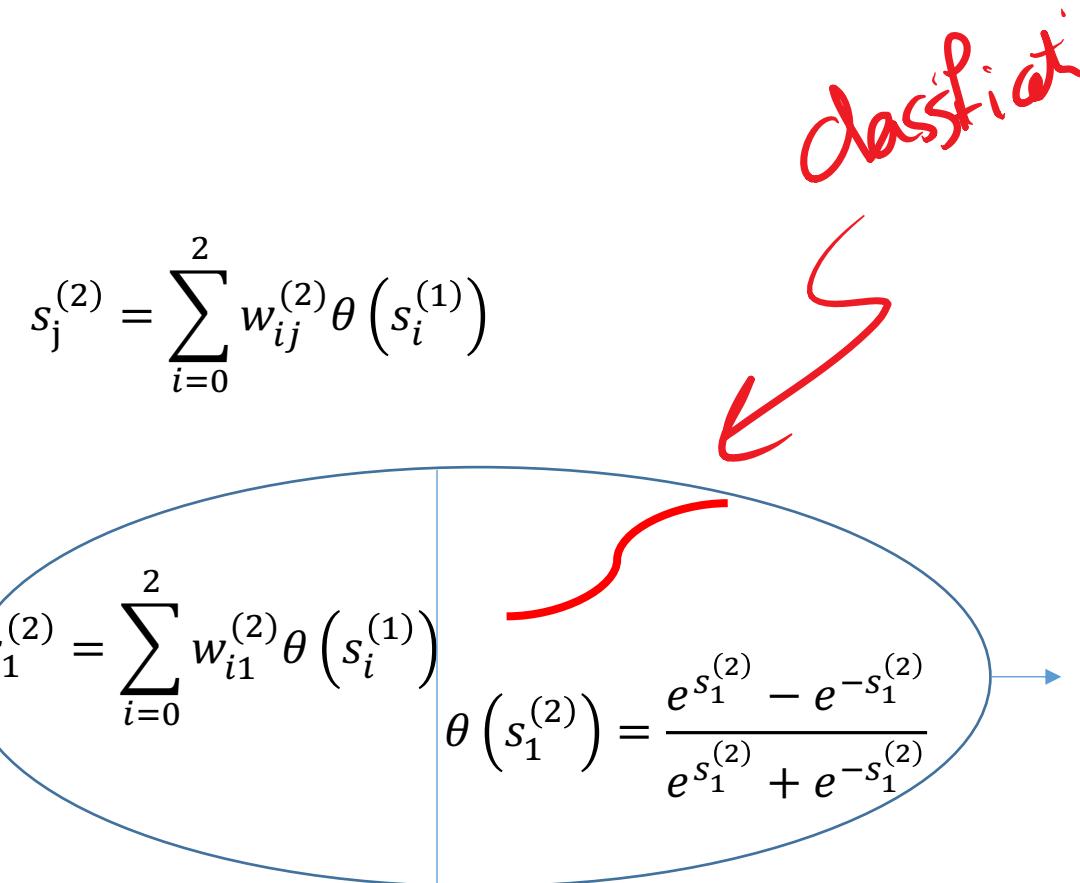
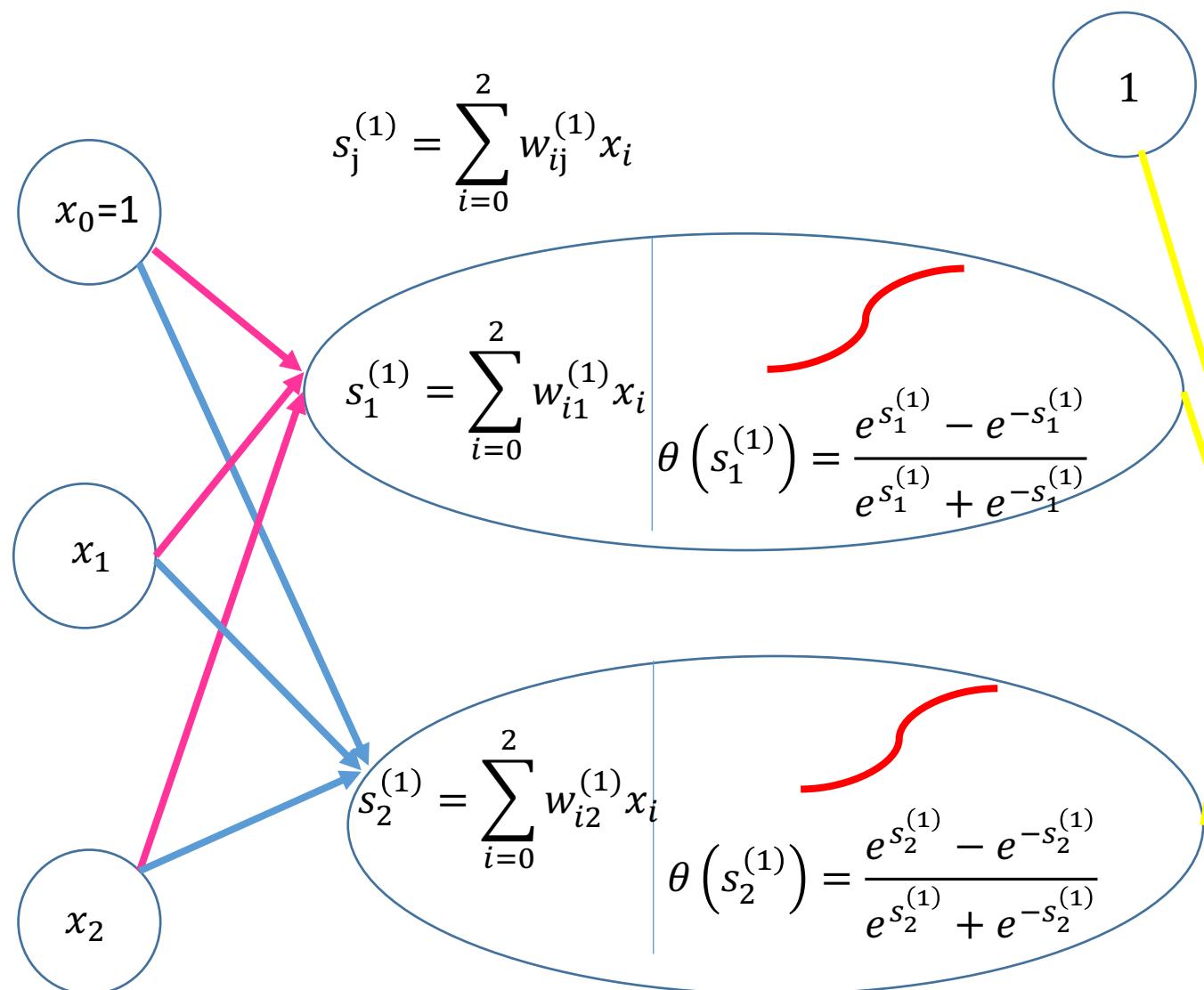
- The objective is to find a set of weights that minimizes the error.
- Recall the update rule in GD:
 $\text{new } \vec{w} = \text{old } \vec{w} - \eta \nabla E(\vec{w})$
- Therefore we need to find the derivative of error with respect to all w's.

$$\nabla E = \begin{bmatrix} \frac{\partial E}{\partial w_0} \\ \frac{\partial E}{\partial w_1} \end{bmatrix}$$



$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$$





dass ist

y = x

reg

Red annotations on the right side of the diagram include the handwritten text "dass ist" at the top, and at the bottom, a red line drawn through the sigmoid curves with the label "y = x" and the word "reg" written next to it.

Exercise

Based on this figure determine the following quantities:

$$\frac{\partial s_1^{(1)}}{\partial w_{21}^{(1)}}, \frac{\partial s_j^{(1)}}{\partial w_{ij}^{(1)}}, \frac{\partial s_1^{(2)}}{\partial w_{21}^{(2)}}, \frac{\partial s_j^{(2)}}{\partial w_{ij}^{(2)}}, \frac{d\theta(s)}{d(s)}$$

Solution

$$s_1^{(1)} = \sum_{i=0}^2 w_{i1}^{(1)} x_i = w_{01}^{(1)}(1) + w_{11}^{(1)}x_1 + w_{21}^{(1)}x_2$$

$$\frac{\partial s_1^{(1)}}{\partial w_{21}^{(1)}} = x_2$$

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

$$s_1^{(2)} = w_{01}^{(2)}(1) + w_{11}^{(2)}\theta(s_1^{(1)}) + w_{21}^{(2)}\theta(s_2^{(1)})$$

$$\frac{\partial s_1^{(2)}}{\partial w_{21}^{(2)}} = \theta(s_2^{(1)}) \Rightarrow \frac{\partial s_j^{(2)}}{\partial w_{ij}^{(2)}} = \theta(s_i^{(1)})$$

$$\frac{d\theta(s)}{d(s)} = 1 - \theta(s)^2$$

Proof: $\theta'(s) = 1 - \theta^2(s)$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Let $u = e^x - e^{-x}$ and $v = e^x + e^{-x}$

$$\frac{u'}{v} - \frac{(e^x - e^{-x})(e^x + e^{-x})}{v^2} = \frac{(e^x - e^{-x})(e^x + e^{-x})}{v'}$$

$$\frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} = 1 - \left(\frac{e^x - e^{-x}}{e^x + e^{-x}}\right)^2 = 1 - \theta^2(x)$$

Example (Forward Pass)

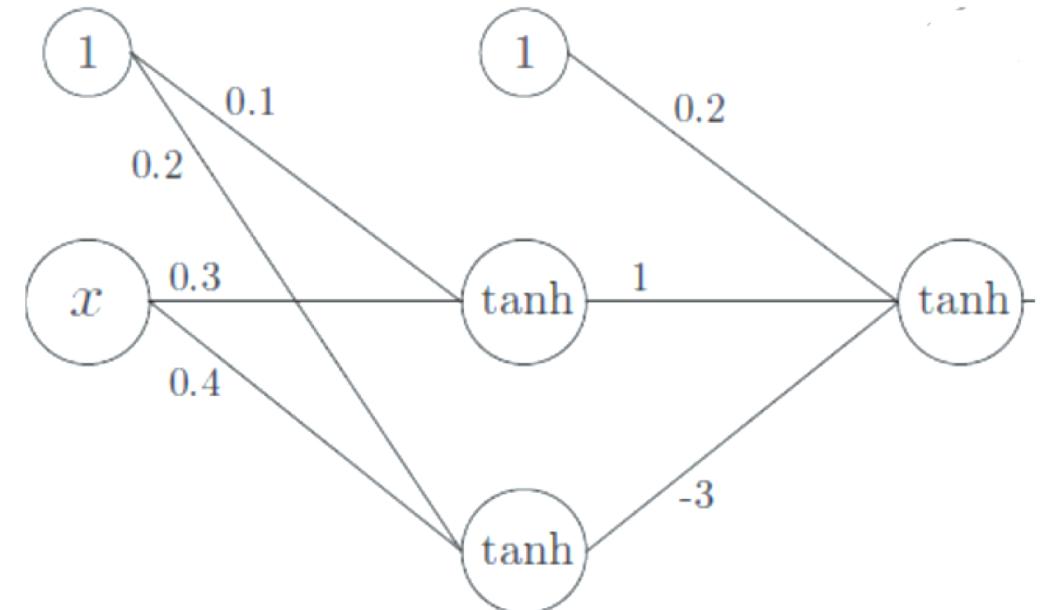
Assume the following weight values are provided. Draw a network that corresponds to these weight matrices and determine the error if the data point is $(x=2, y=1)$.



$$W^{(1)} = \begin{bmatrix} 0.1 & 0.2 \\ 0.3 & 0.4 \end{bmatrix}; \quad W^{(2)} = \begin{bmatrix} 0.2 \\ 1 \\ -3 \end{bmatrix}$$

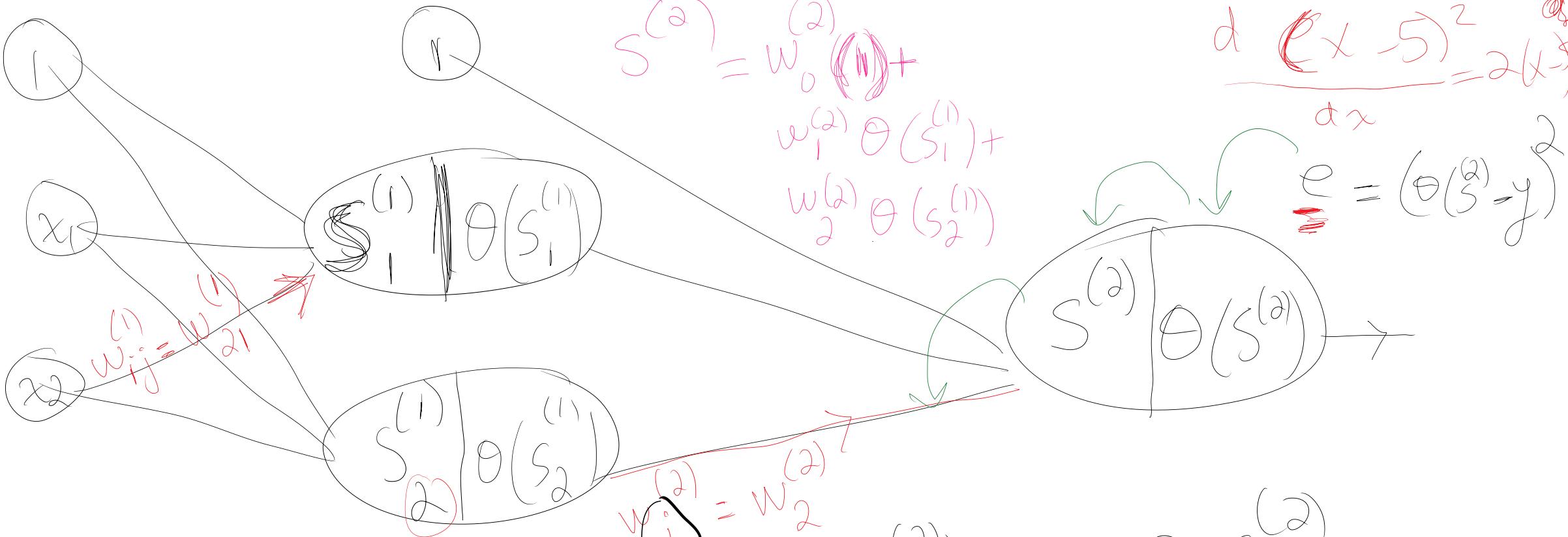
Solution

- $\vec{x}^T = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$
- $s^1 = w^T x^T = \begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0.7 \\ 1 \end{bmatrix}$
- $\theta(s^1) = \tanh \begin{bmatrix} 0.7 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.60 \\ 0.76 \end{bmatrix}$
- $s^2 = [0.2 \quad 1 \quad -3] \begin{bmatrix} 1 \\ 0.6 \\ 0.76 \end{bmatrix} = -1.48$
- $e = (\tanh(-1.48) - 1)^2 = 0.81$



Backpropagation

Finding the gradient for layer 2



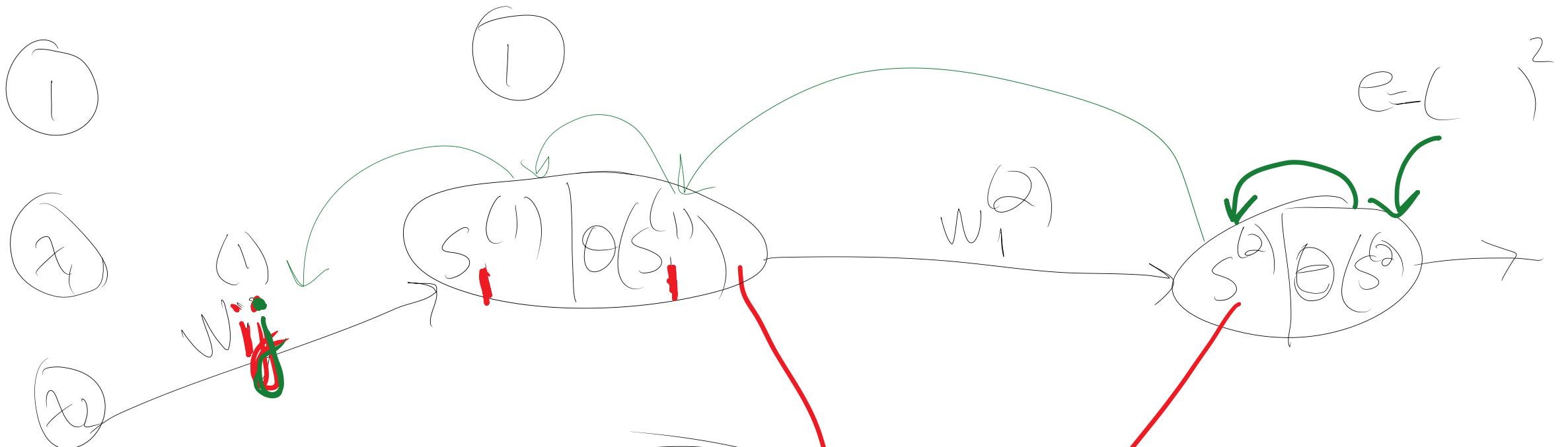
$$\frac{\partial e}{\partial w_i^{(2)}} = \frac{\partial e}{\partial \theta(S^{(2)})} \times \frac{\partial \theta(S^{(2)})}{\partial S^{(2)}} \times \frac{\partial S^{(2)}}{\partial w_i}$$

2(\theta(S^{(2)}) - y)
 ↓ 1 - \theta^2(S^{(2)})
 see for proof next slide

if $i \neq 0 \Rightarrow i = 0 \Rightarrow$ 35

Backpropagation

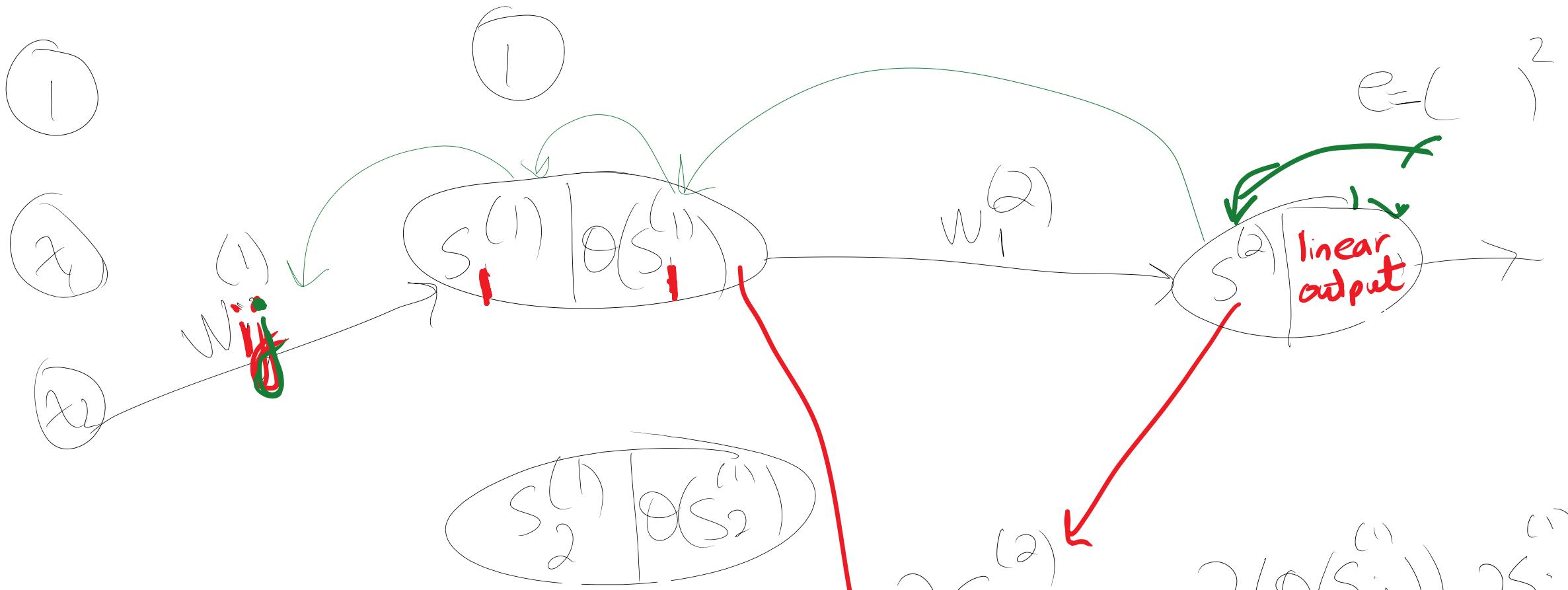
Finding the gradient for layer 1



$$\frac{\partial e}{\partial w_{ij}^{(1)}} = \frac{\partial e}{\partial \theta(S^{(2)})} \times \frac{\partial \theta(S^{(2)})}{\partial S^{(2)}} \times \frac{\partial S^{(2)}}{\partial \theta(S_j^{(1)})} \times \frac{\partial (\theta(S_j))}{\partial S_j^{(1)}} \times \frac{\partial S_j^{(1)}}{\partial w_{ij}^{(1)}}$$

2(\theta(S_j^{(2)}) - y) 1 - \theta^2(S^{(2)}) w_j^{(2)} 1 - \theta^2(S_j^{(1)}) x_i

37



$$\frac{\partial e}{\partial w_{ij}^{(1)}} = \frac{\partial e}{\partial (s_j^{(2)})} \times \frac{\partial s_j^{(2)}}{\partial \theta(s_j^{(1)})} \times \frac{\partial \theta(s_j^{(1)})}{\partial s_j^{(1)}} \times \frac{\partial s_j^{(1)}}{\partial w_{ij}^{(1)}}$$

Below the equation, the terms are highlighted with brackets:

- $\frac{\partial e}{\partial (s_j^{(2)})}$ is highlighted with a red bracket under the term $2(s_j^{(2)} - y)$.
- $\frac{\partial s_j^{(2)}}{\partial \theta(s_j^{(1)})}$ is highlighted with a green bracket under the term $w_{j(2)}$.
- $\frac{\partial \theta(s_j^{(1)})}{\partial s_j^{(1)}}$ is highlighted with a purple bracket under the term $1 - \theta^2(s_j^{(1)})$.
- $\frac{\partial s_j^{(1)}}{\partial w_{ij}^{(1)}}$ is highlighted with a yellow bracket under the term x_i .

38

Backpropagation Algorithm (SGD)

- Initialize the weights at random
- Loop until it's time to stop (e.g. reach max iterations)
 - for every training point
 - Do all forward calculations (s , θ)
 - Do all backward calculations (partial derivatives)
 - Update W 's from the rule:
$$\text{new } \vec{w} = \text{old } \vec{w} - \eta \nabla e(\vec{w})$$
- End for
- End Loop

Backpropagation Algorithm

Batch Gradient Descent

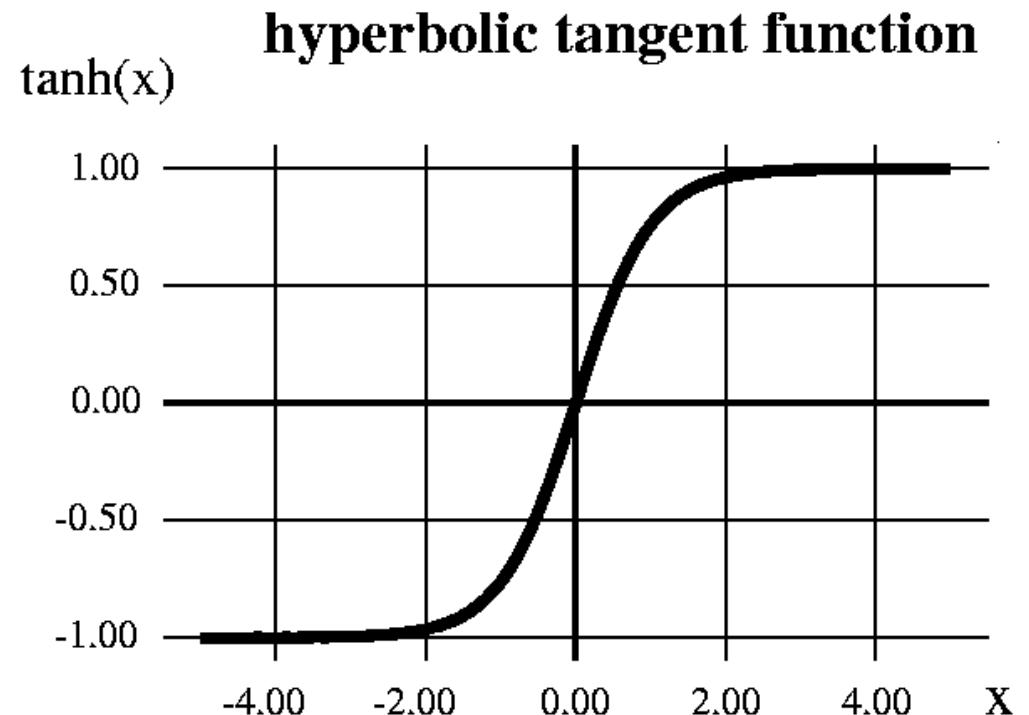
- Initialize the weights at random
- Loop until it's time to stop (e.g. reach max iterations)
- for every training point
 - Do all forward calculations (s , theta)
 - Do all backward calculations (partial derivatives)
 - Accumulate $\nabla e_n(\vec{w})$
- End for
- Update W's from the rule:
$$\text{new } \vec{w} = \text{old } \vec{w} - \eta \nabla E(\vec{w})$$
- End Loop

Stochastic Gradient Descent

- Initialize the weights at random
- Loop until it's time to stop (e.g. reach max iterations)
- for every training point
 - Do all forward calculations (s , theta)
 - Do all backward calculations (partial derivatives)
 - Update W's from the rule:
$$\text{new } \vec{w} = \text{old } \vec{w} - \eta \nabla e_n(\vec{w})$$
- End for
- End Loop

Initialization of Weights in GD

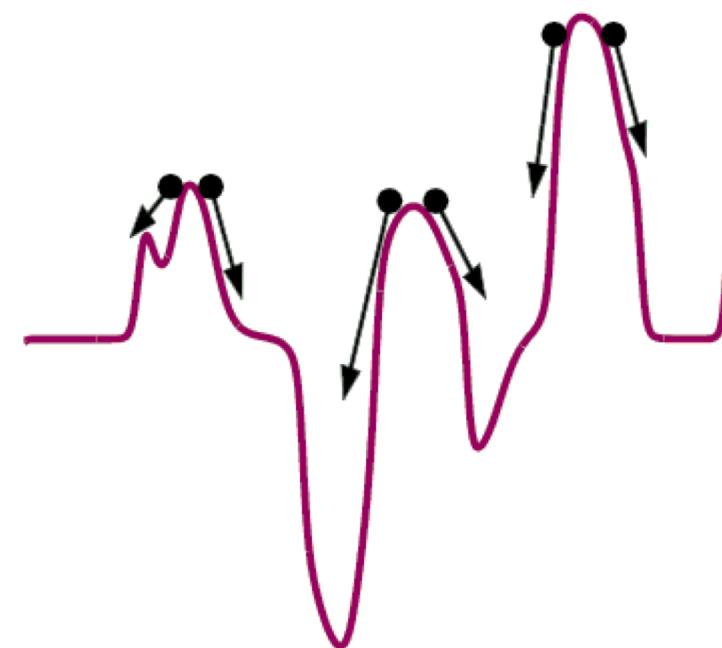
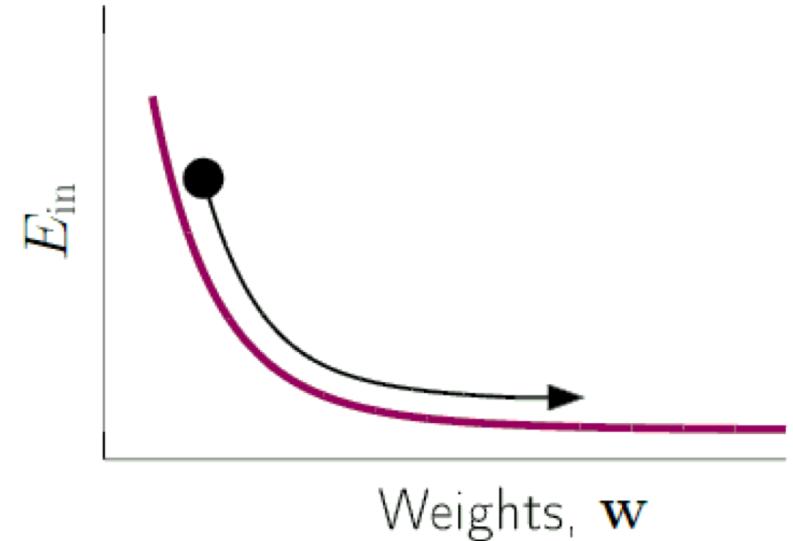
- Choosing the initial weights and stopping criteria can be tricky, as compared with logistic regression, because the in-sample error is not convex anymore.
- In NN, it's the best to initialize the weights to small random values. (why?)
- When $\tanh(w^t x) \approx 0$ algorithm has the flexibility to move the weights easily to fit the data.
- In practice choosing weights from a normal distribution with zero mean and small variance typically works well.



Termination of GD



- Termination is a non-trivial decision in optimization.
- One approach is based on the gradient being zero at a minimum, therefore we could stop once the gradient of error drops below a threshold.
- However, this may happen too early when you've hit a local minima.
- To overcome this problem we stop when the error falls below a threshold.
- However, that threshold may never be reached.
- To overcome this we can stop when an upper bound on the number of iterations is reached.
- In practice a combination of stopping criteria is used.

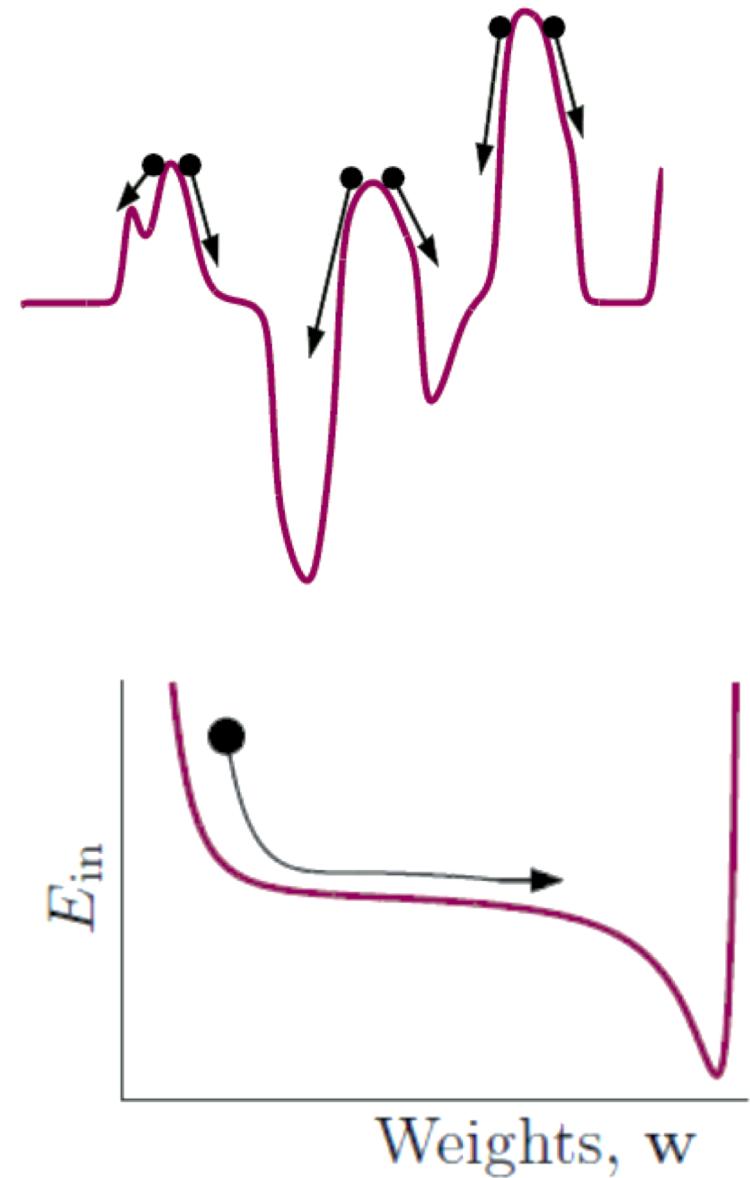


Why does SGD converge?

- Recall that in batch GD the gradient is calculated on the whole data set before a weight is updated. i.e. $\text{new } \vec{w} = \text{old } \vec{w} - \eta \nabla E(\vec{w})$ where $\nabla E(\vec{w}) = \frac{1}{N} \sum_{n=1}^N \nabla e_n(\vec{w})$.
- In contrast in SGD, we choose a point uniformly at random from the training set and update the weight based on this individual gradient. i.e. $\text{new } \vec{w} = \text{old } \vec{w} - \eta \nabla e_n(\vec{w})$.
- We can prove that “**on average**” the minimization proceeds in the right direction (but a bit wiggly).
- A similar situation was seen in PLA where minimizing the error on one data point could have interfered with the error on the rest of data set, but this effect cancels out on average.

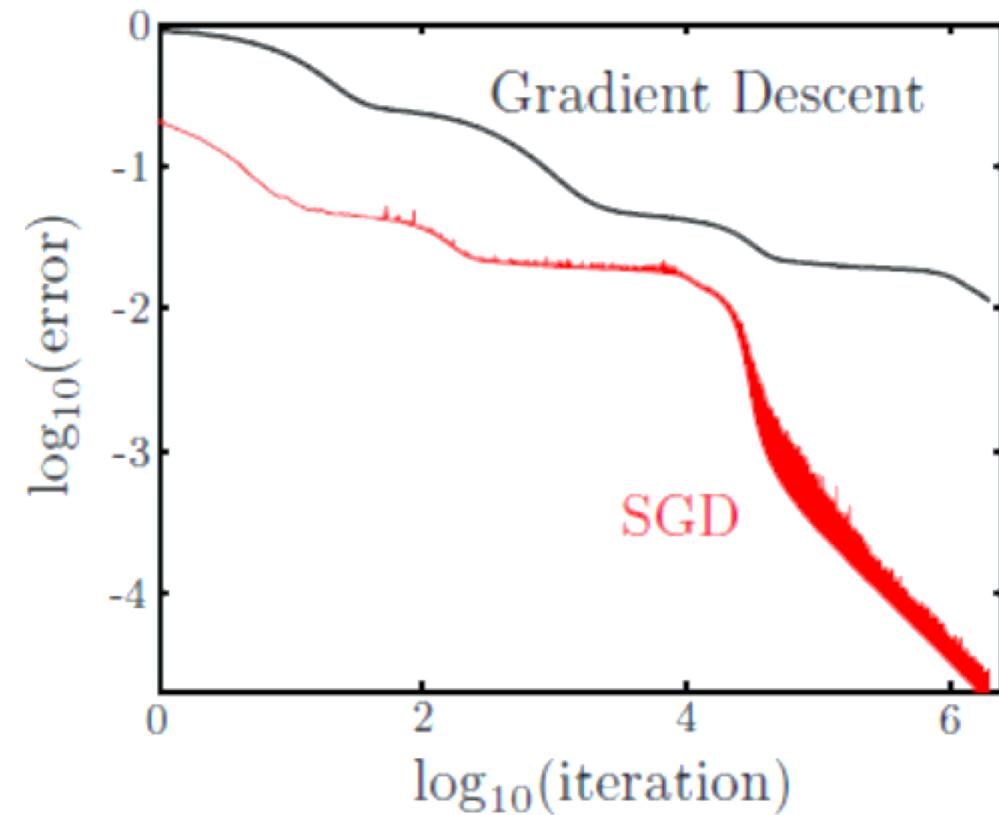
Benefits of SGD over Batch GD

- Speed: Instead of calculating all N gradients before making an update, you will make an update after evaluating every gradient.
- Randomness: This is an appealing feature in optimization problems.
- Because you're not going in a deterministic direction (unlike the batch mode), there's a chance that you skip a shallow local minima and avoid premature termination.
- In practice, run the algorithm several times every time from a different starting point and then take the minimum of those minimums to get to the best local minima.



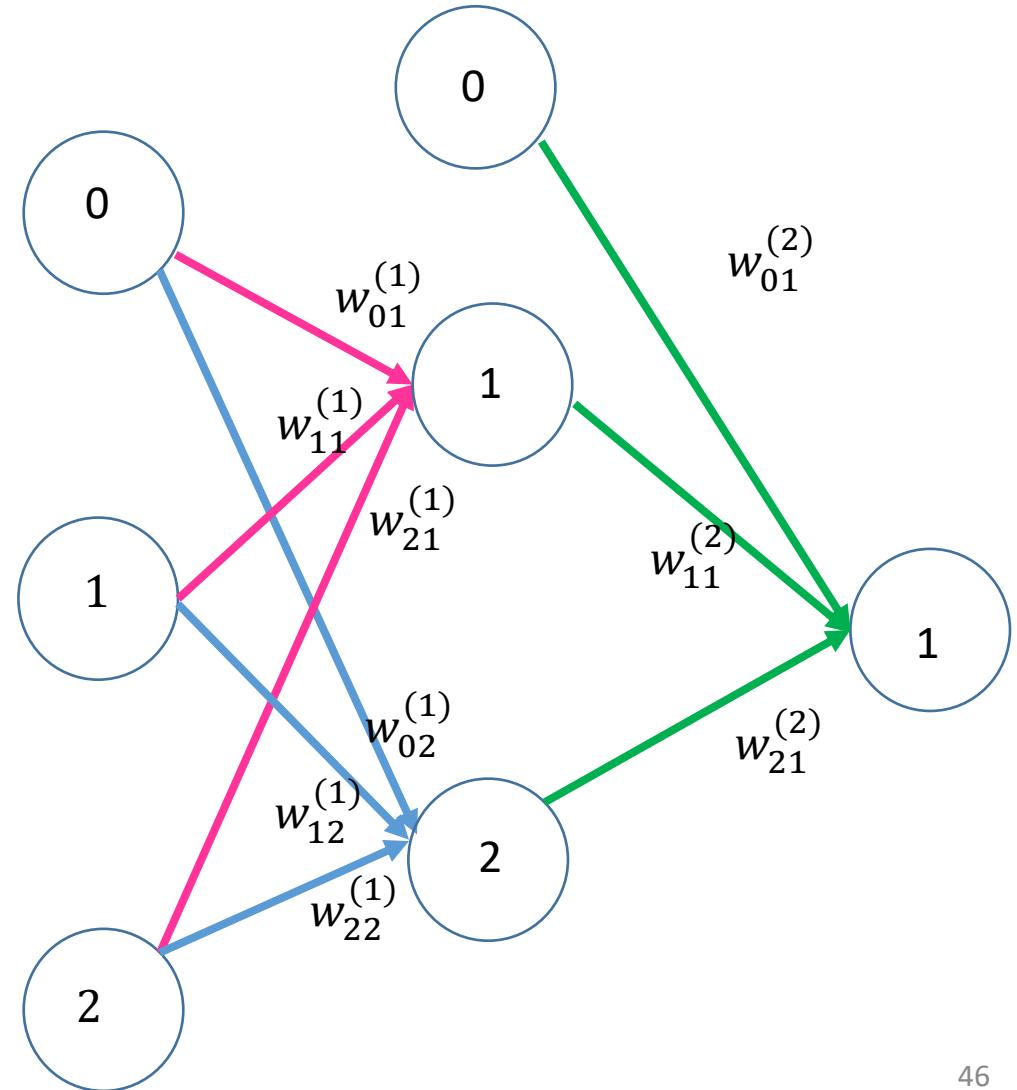
Example

- 500 training examples from the digits data
- 1 hidden layer
- 5 hidden units
- learning rate $\eta = 0.01$.
- The SGD curve is erratic because the weights are updated based error on a single data point in each iteration.
- One method to overcome the wiggly behavior is to decrease the learning rate as the minimization proceeds.



Rule of Thumb for Choosing Network Parameters

- Learning rate $\eta = 0.1$ is a good starting point
- Approximately the number of data points should be at least 10 times the number of parameters you need to estimate.
- i.e. $N > 10w$



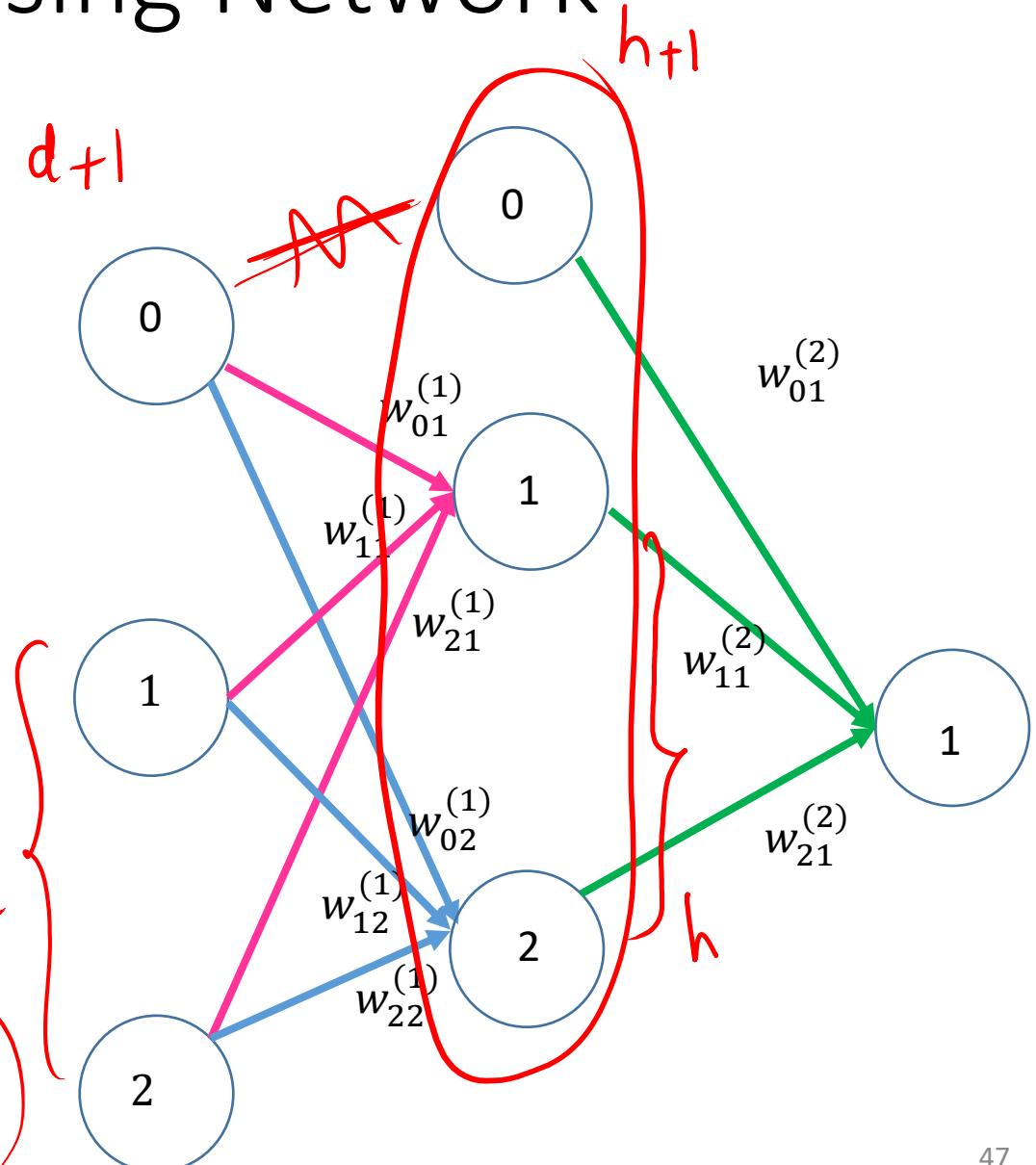
Rule of Thumb for Choosing Network Parameters

- Learning rate $\eta = 0.1$ is a good starting point
- Approximately the number of data points should be at least 10 times the number of parameters you need to estimate.
- i.e. $N > 10w$

$$(d+1)h + h + 1 \leq \frac{N}{10}$$

$$h(d+a) \leq \frac{N-1}{10}$$

$$h \leq \frac{\frac{N-1}{10}}{d+a}$$



References

- Learning from Data by Abu-Mustafa, Ismail, and Lin