

Topic 6 Neural Networks

<http://deeplearning.stanford.edu/tutorial/> and then go
to MultiLayerNeuralNetworks
<http://neuralnetworksanddeeplearning.com/>

INTRODUCTION TO MACHINE LEARNING
PROF. LINDA SELLIE

Some of these slides are from Prof. Rangan

Overview

Used for both regression and classification

Neural networks *extend* logistic regression and linear regression

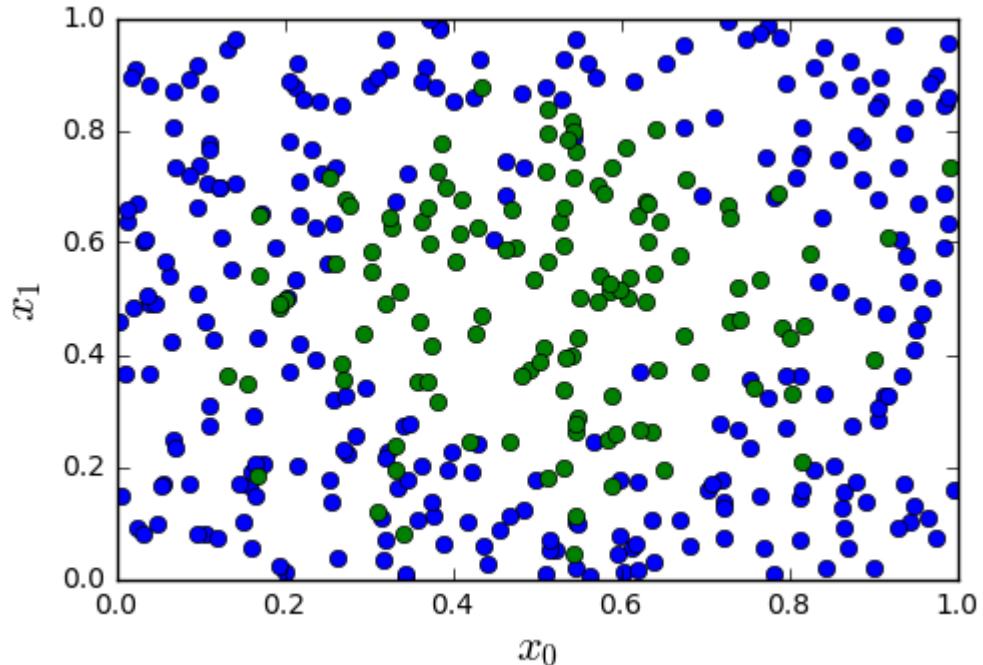
Neural networks are universal approximators (it is possible to approximate any bounded function)

Outline



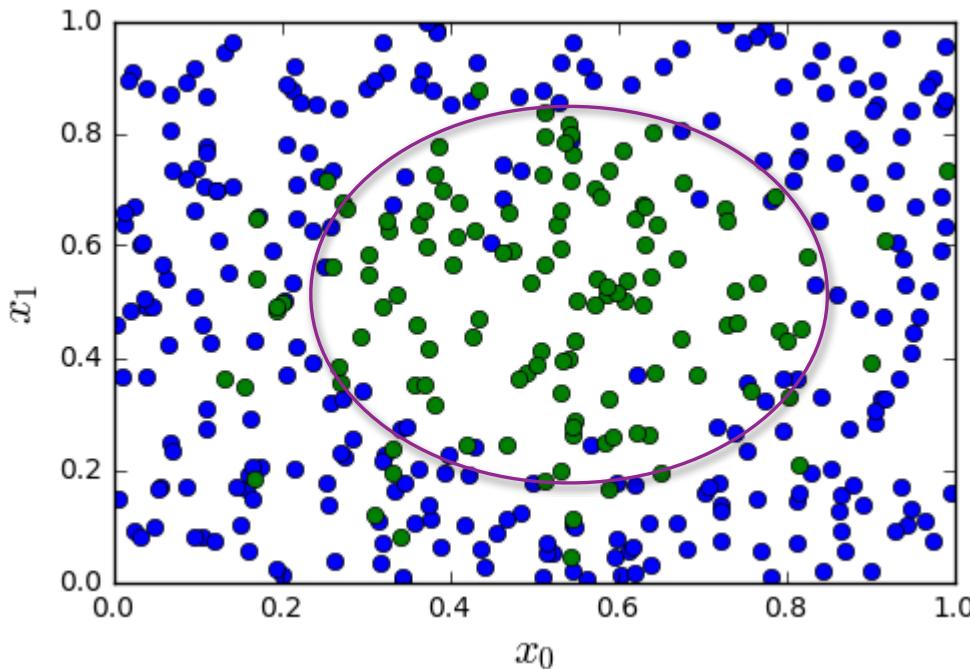
- ❑ Motivation Introduction to neurons
- ❑ Nonlinear classifiers from linear features
- ❑ Neural networks notation
- ❑ Pseudocode for prediction
- ❑ Training a neural network
- ❑ Implementing gradient descent for neural networks
 - Vectorization
 - Pseudocode
- ❑ Preprocessing
- ❑ Initialization
- ❑ Activations

Most Datasets are not Linearly Separable



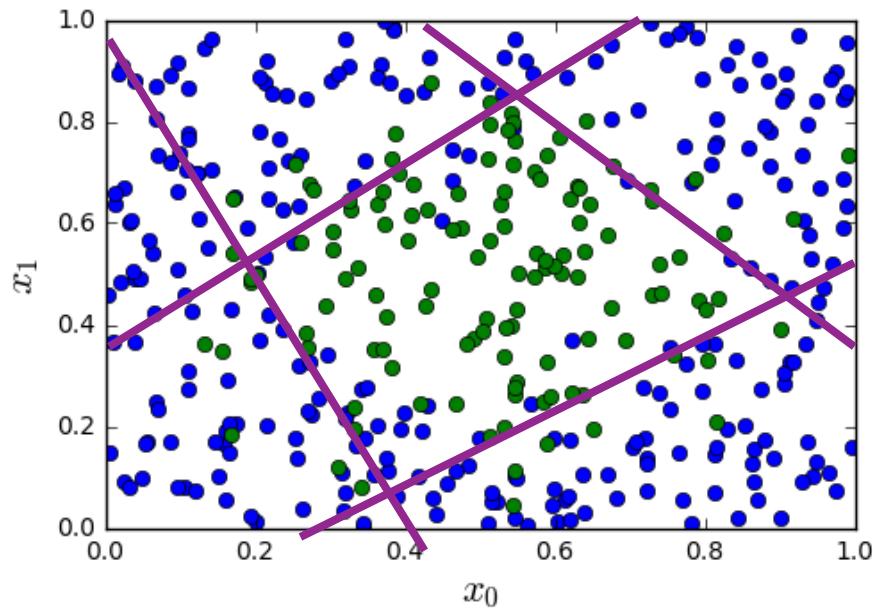
- ❑ Consider simple synthetic data
 - See figure to the left
 - 2D features
 - Binary class label
- ❑ Not separated linearly
- ❑ We can use Logistic Regression with non-linear features

From Linear to Nonlinear



❑ Idea: Build nonlinear region from linear decisions

From Linear to Nonlinear



❑ Idea: Build nonlinear region from linear decisions

How can we learn the right
feature transformations
(aka *functions* to
transform our features)?

THAT IS THE MAIN TOPIC OF THIS LECTURE

How can we learn the right feature transformations (aka functions to transform our features):

This is what neural
networks do! We input the
original features and the
network learns different
function of the features

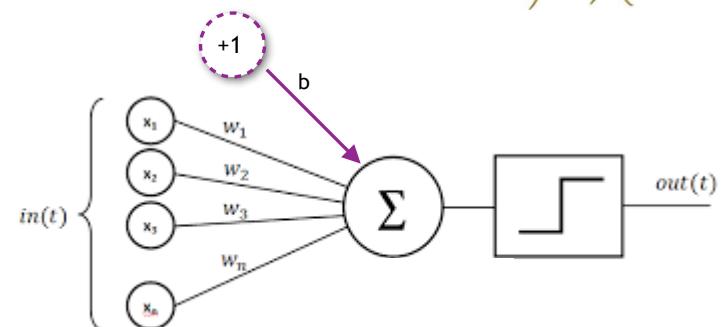
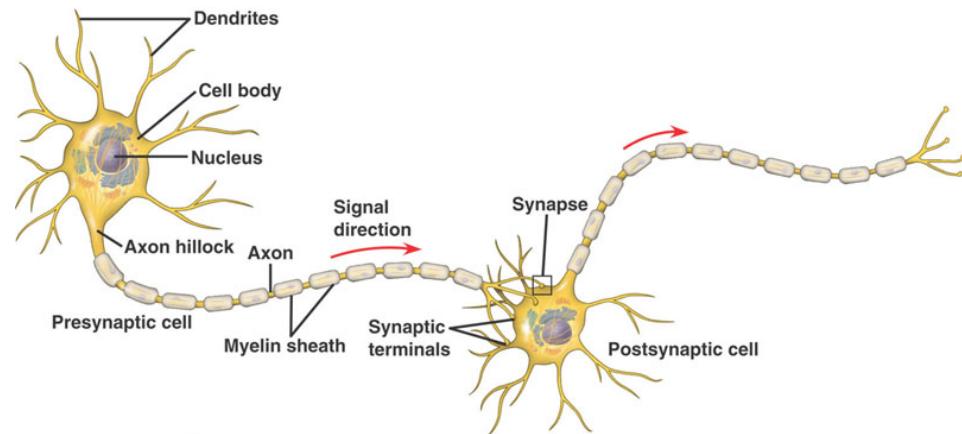
THAT IS THE MAIN TOPIC OF THIS LECTURE

Recent Resurgence and Developments

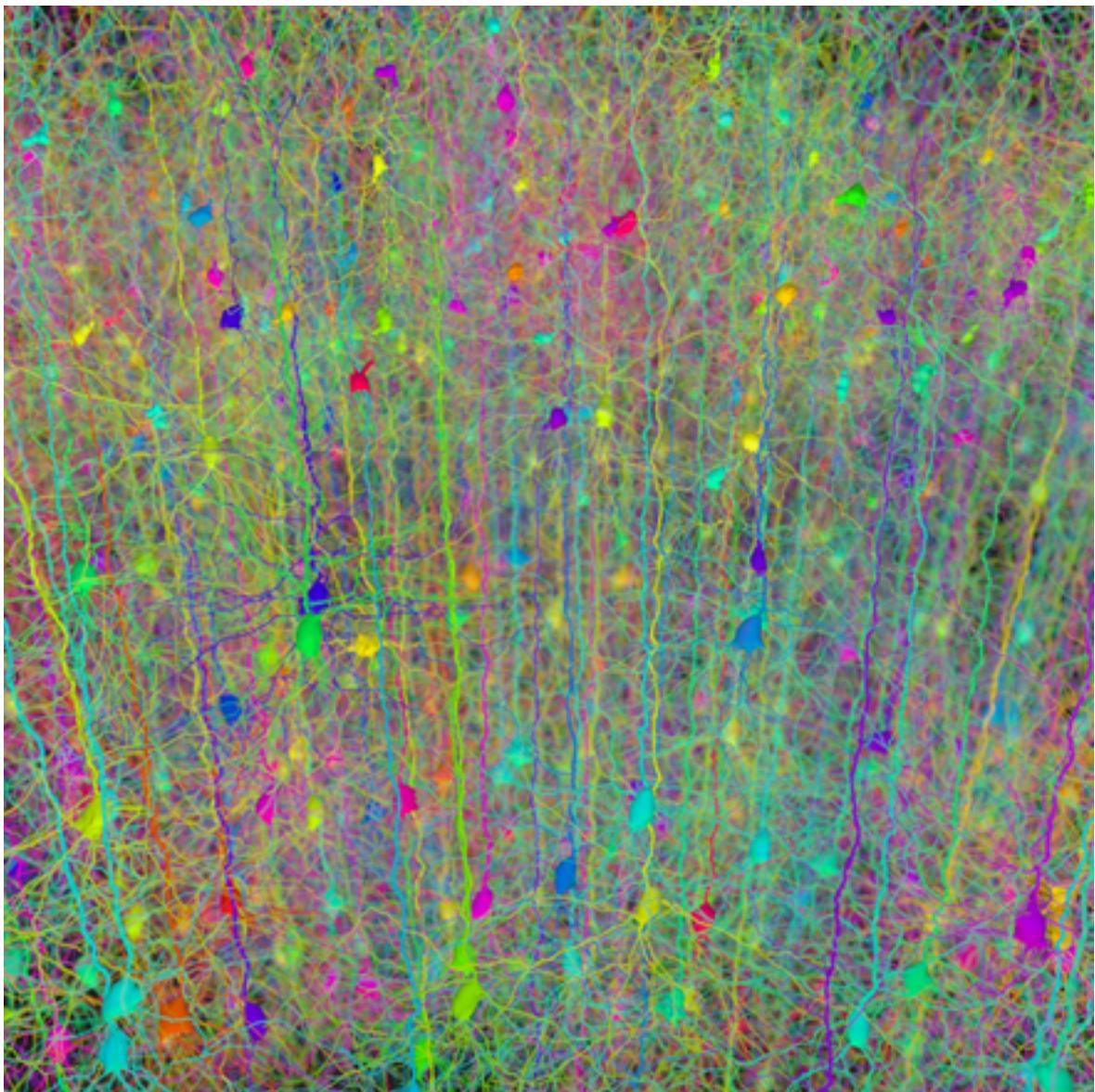
- ❑ In 1980's neural networks were used.
- ❑ However interest in them declined due to computation power and not enough data
- ❑ Advances in hardware (GPU and high performance computation, etc) and the massive amounts of data being collected have removed has significantly increased the use of neural networks since 2012
- ❑ The past 5-7 years have seen algorithmic innovations. Many of these are used in self-driving cars, facial recognition, speech recognition, etc.
- ❑ Recently, deep learning has had the most impact to expand what is possible to learn. (Deep learning is just neural networks with many "layers". Later in today's class we will discuss what a "layer" is)

Inspiration from Biology

Abstract away biological construct into a mathematical construct



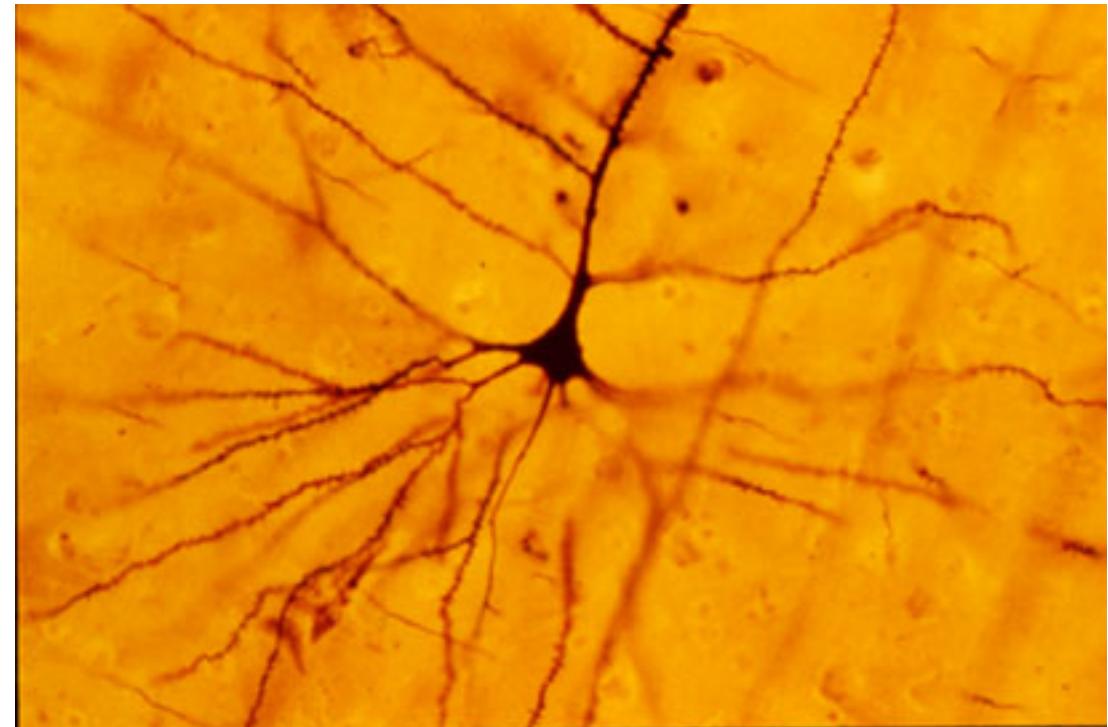
- ❑ Simple model of neurons
 - Dendrites: Input currents from other neurons
 - Soma: Cell body, accumulation of charge
 - Axon: Outputs to other neurons
 - Synapse: Junction between neurons
- ❑ Operation:
 - Take weighted sum of input current
 - Outputs when sum reaches a threshold
- ❑ Each neuron is like one unit in neural network
- ❑ No one knows how the brain really works - but just like people were inspired by birds to build airplanes, our neurons do not work the same but are inspired by the neurons in our brains



Computer simulation of the branching architecture of the dendrites of pyramidal neurons.^[6]

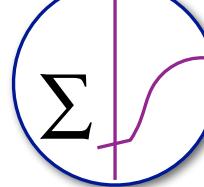
https://en.wikipedia.org/wiki/Neural_network

"**Pyramidal cells**, or **pyramidal neurons**, are a type of multipolar neuron found in areas of the brain including the cerebral cortex, the hippocampus, and the amygdala"



https://en.wikipedia.org/wiki/Pyramidal_cell

Perceptron/Neuron



Each node has a left and right side. The left side is a weighted linear sum, the right size is a non-linear function

“A **perceptron** is a simple model of a biological neuron in an artificial neural network. ... The **perceptron algorithm** was designed to classify visual inputs, categorizing subjects into one of two types and separating groups with a line.“

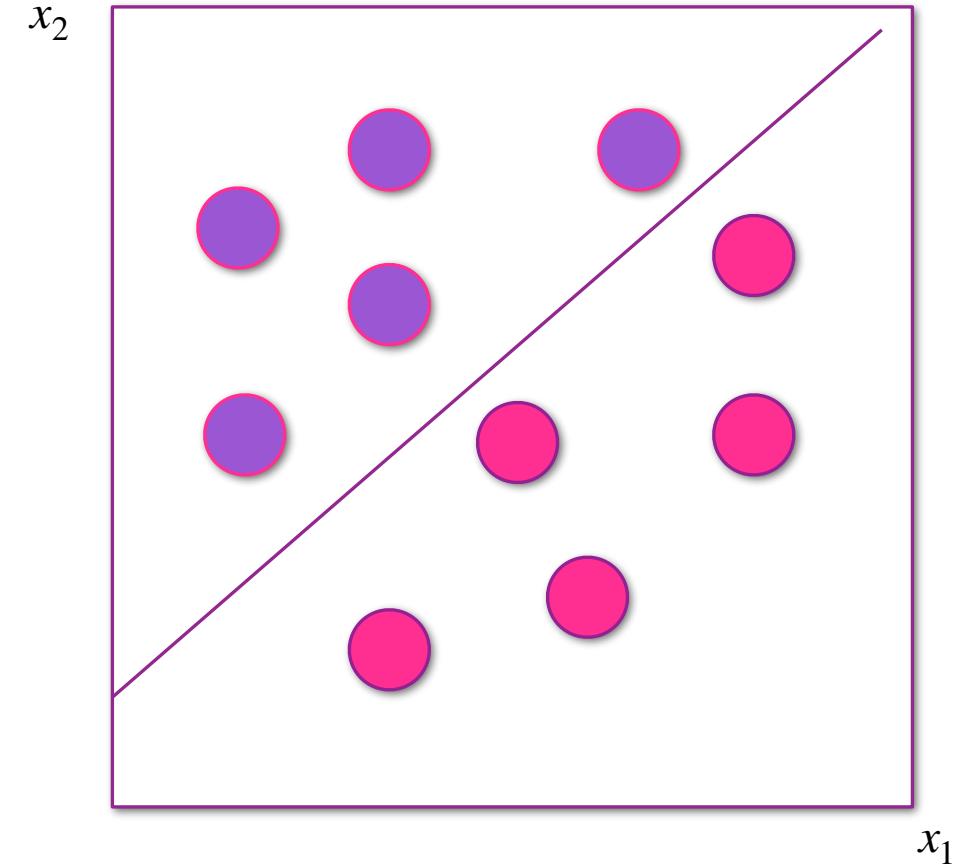
What kind of decision boundary can a single “neuron” create?

The diagram shows a perceptron model. It has two input nodes labeled x_1 and x_2 , a bias node labeled $+1$, and a single output node. The output node contains a blue Greek letter sigma (Σ) and an S-shaped curve. Arrows point from x_1 to the output node via weight $W_{11}^{(1)}$, from x_2 to the output node via weight $W_{12}^{(1)}$, and from the bias node $+1$ to the output node via weight b . The output of the perceptron is given by the equation:

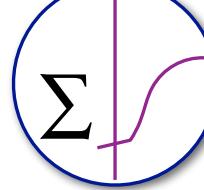
$$h_{W,b}(x) = f\left(\sum_{i=1}^2 W_{1i}^{(1)} x_i + b\right)$$

where $f(z) = f(W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 + b)$.

$$z = [W_{11}^{(1)}, W_{12}^{(1)}] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b$$



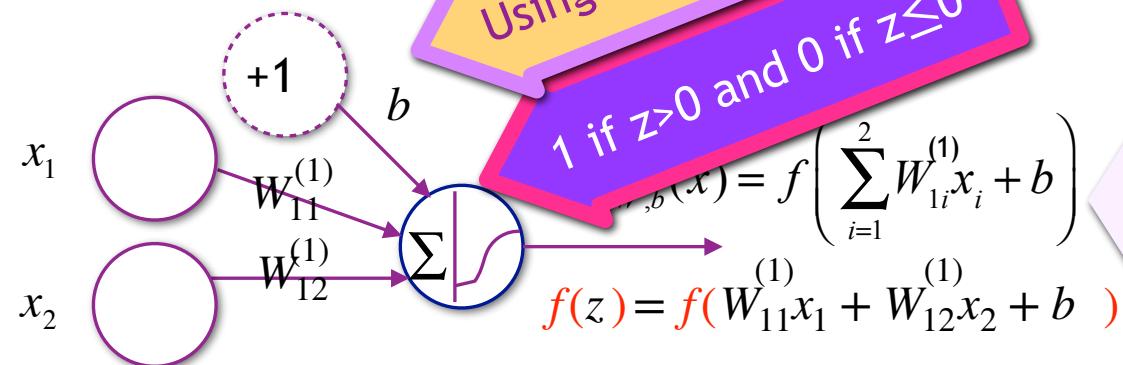
Perceptron/Neuron



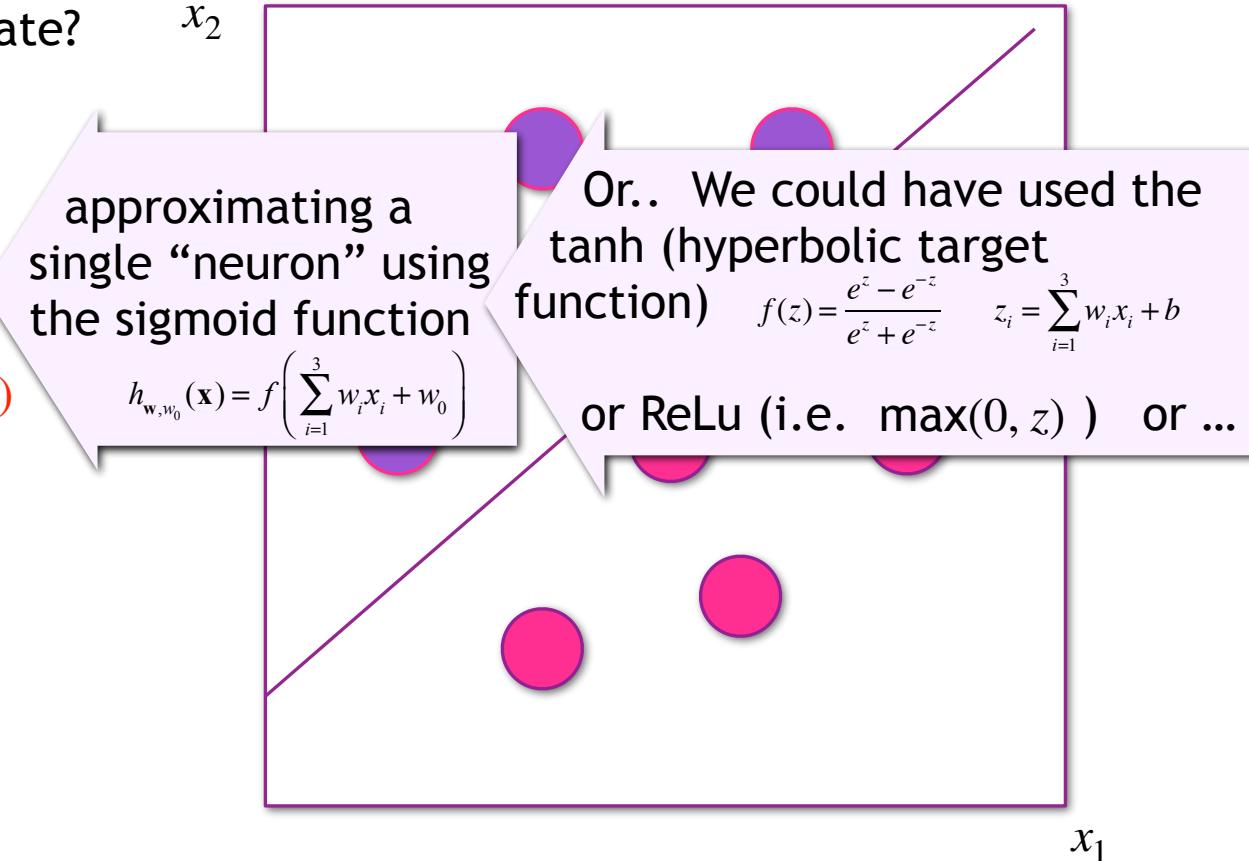
For this lecture will use the sigmoid function

"A **perceptron** is a simple model of a biological neuron that can be used to classify visual inputs, categorizing subjects into two types and separating groups with a line."

What kind of decision boundary?



" create?

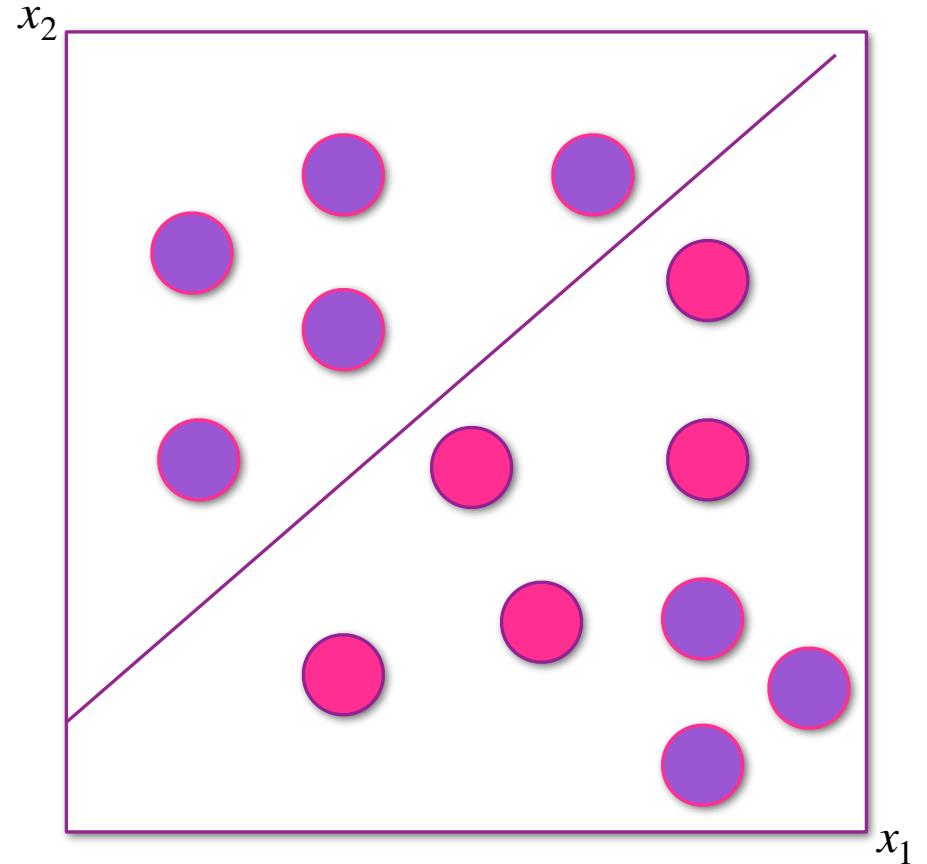
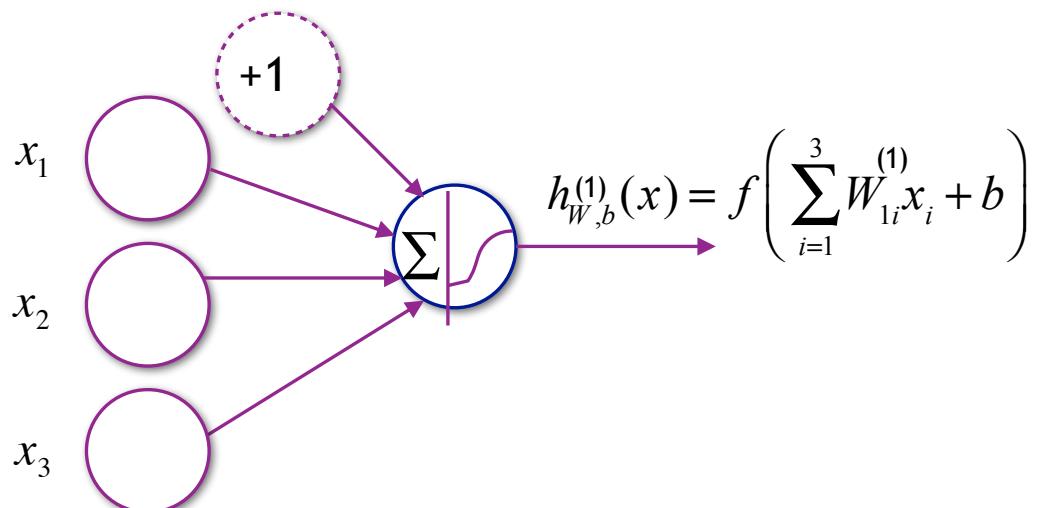


Outline

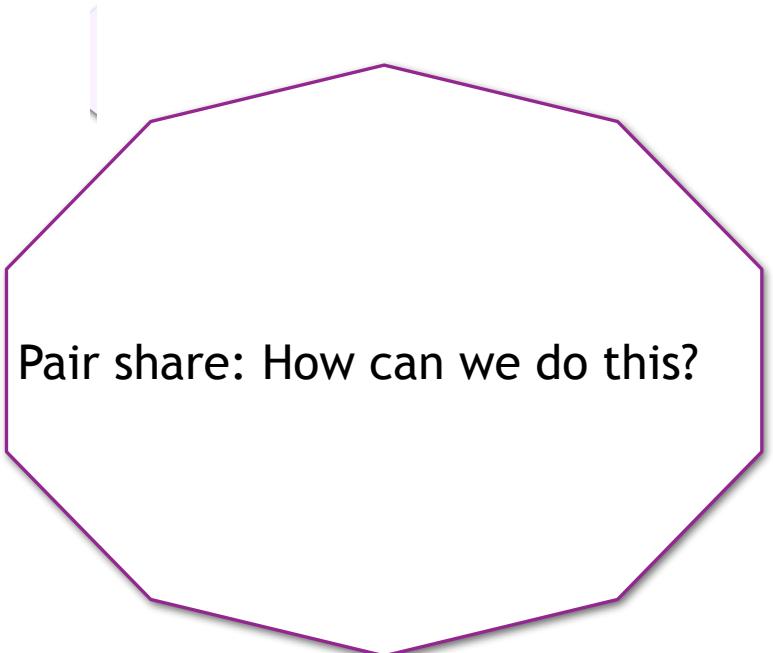
- 
- ❑ Introduction to neurons
 - ❑ Nonlinear classifiers from linear features
 - ❑ Neural networks notation
 - ❑ Pseudocode for prediction
 - ❑ Training a neural network
 - ❑ Implementing gradient descent for neural networks
 - Vectorization
 - Pseudocode
 - ❑ Preprocessing
 - ❑ Initialization
 - ❑ Activations

A more complicated decision boundary?

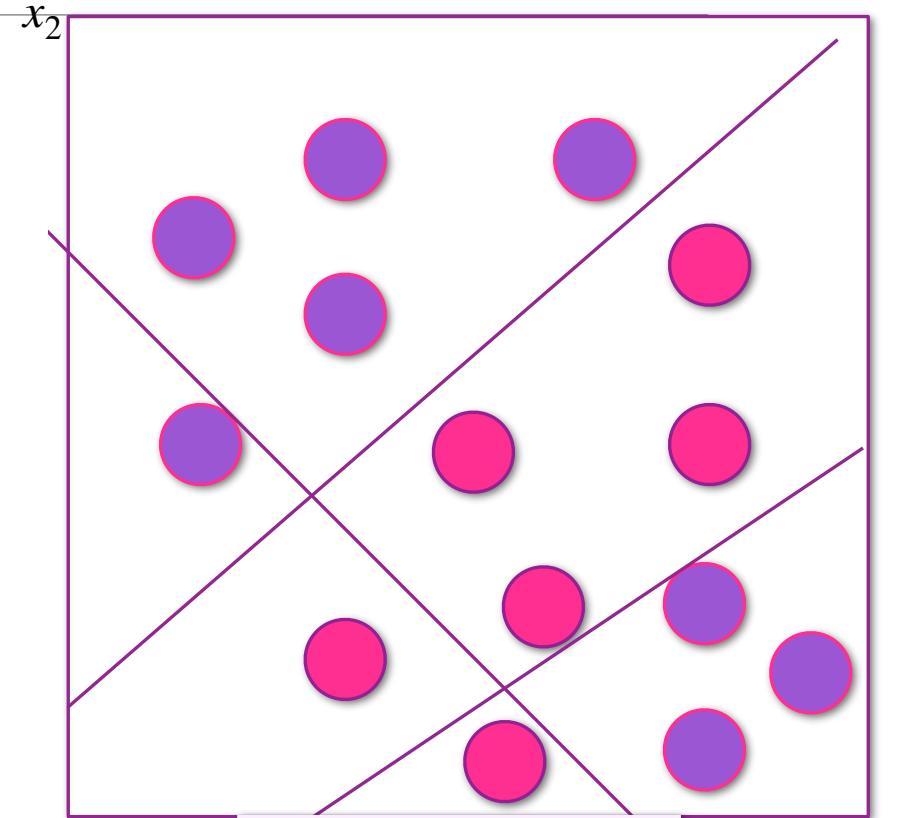
How can we get a more complicated decision boundary?



How could we learn a more complicated decision boundary?



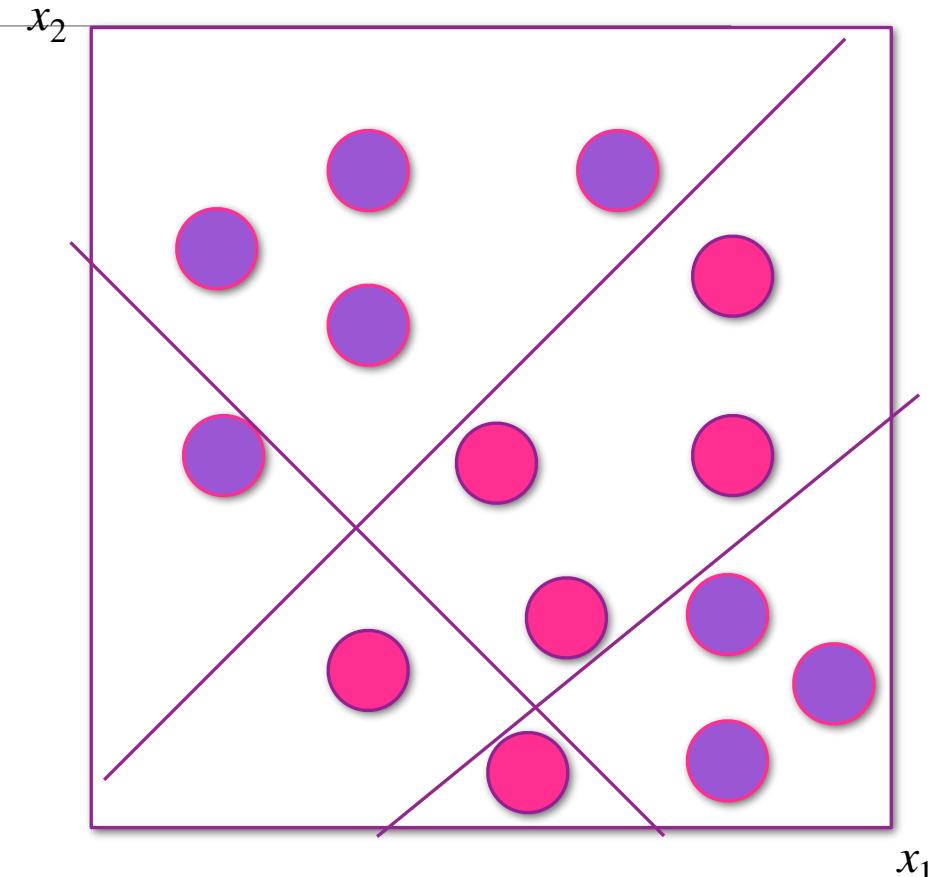
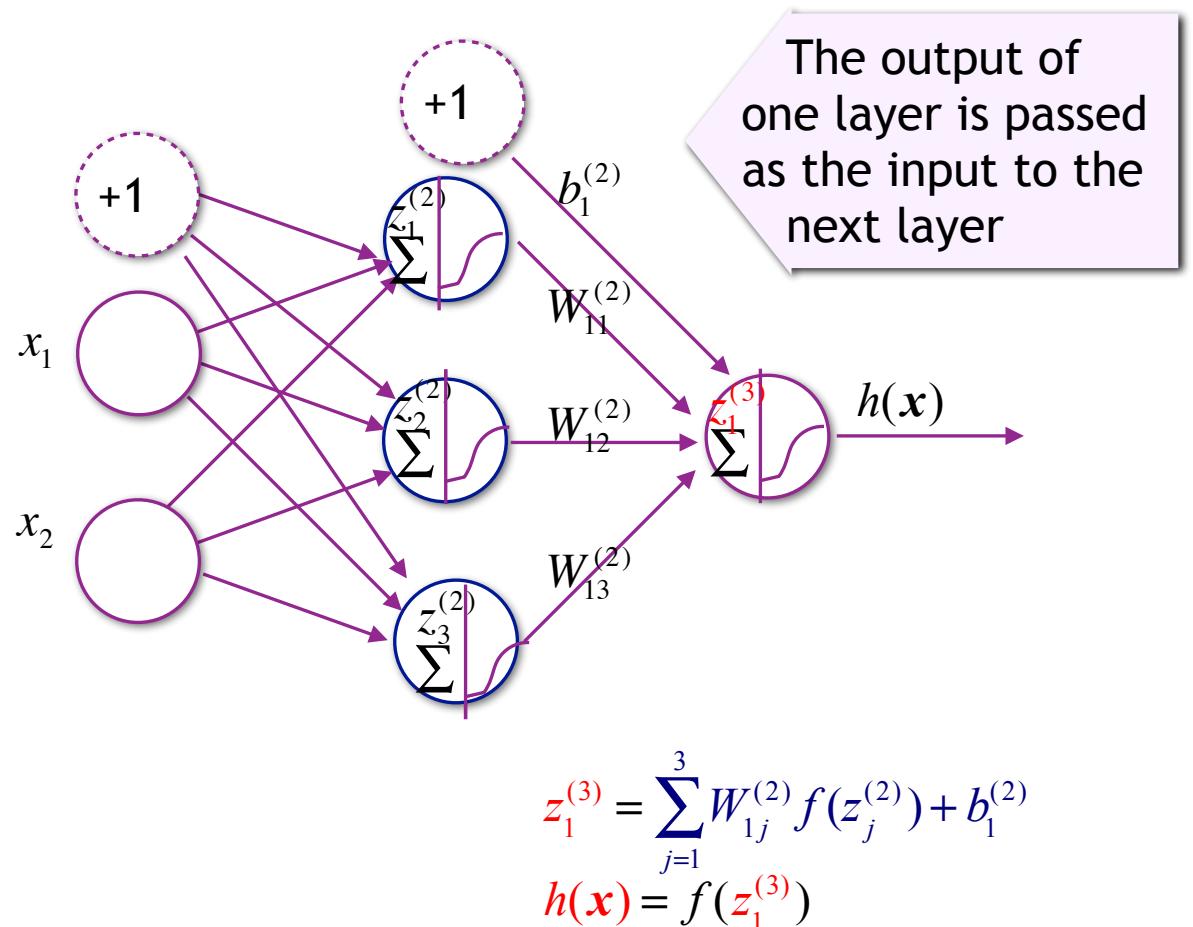
It would be easy if our features were the 3 hyperplanes...



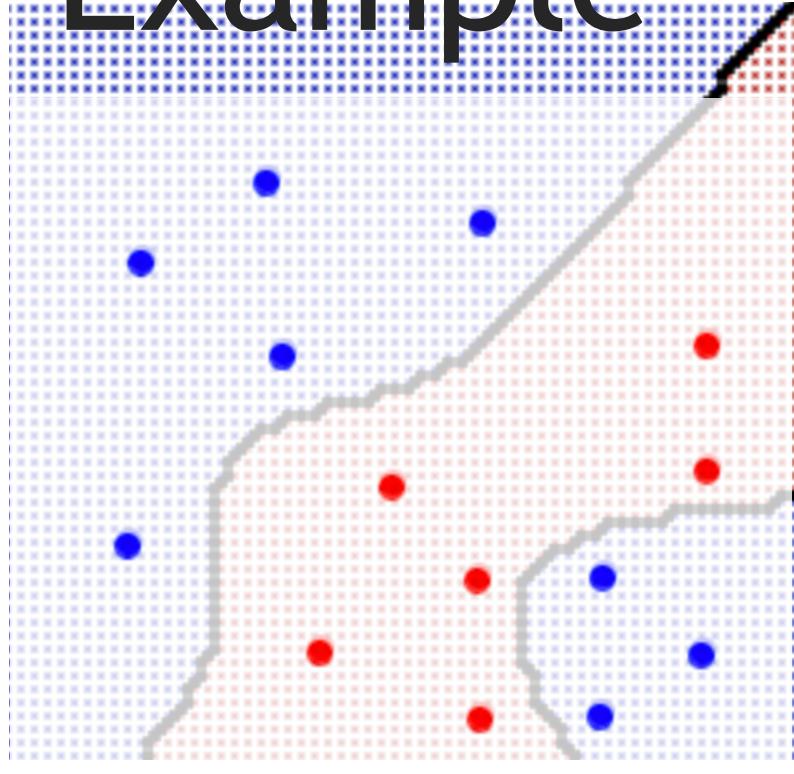
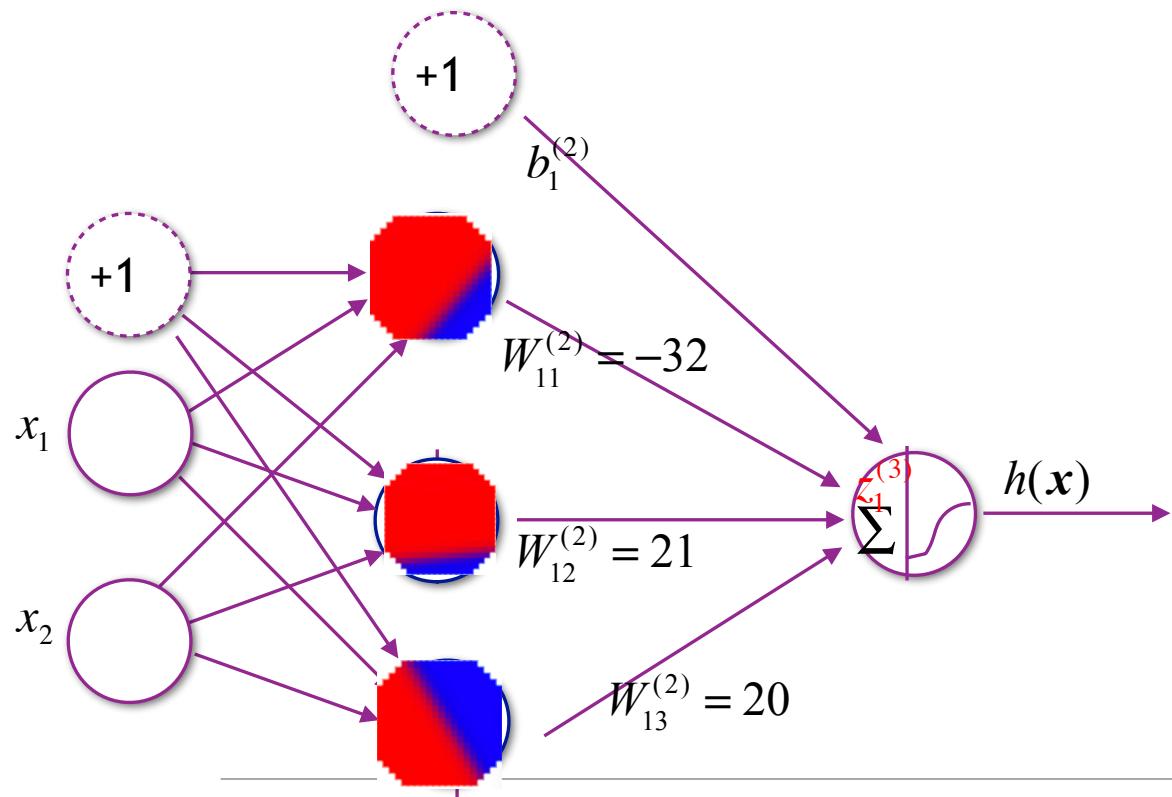
If we could use these as features we can create a more complex decision boundary.

Feature Construction!

Creating a more useful set of features that allow for a linear decision boundary



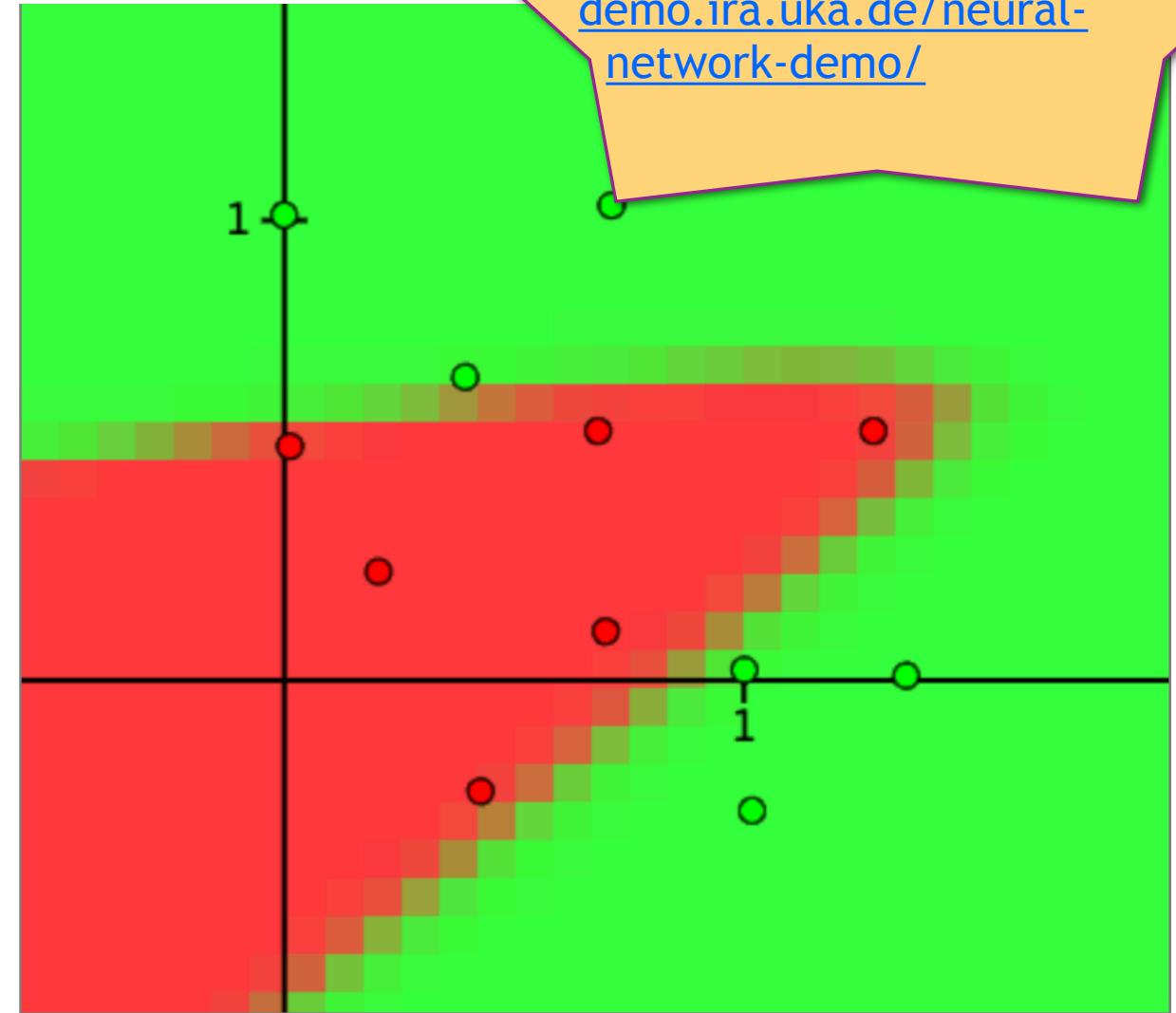
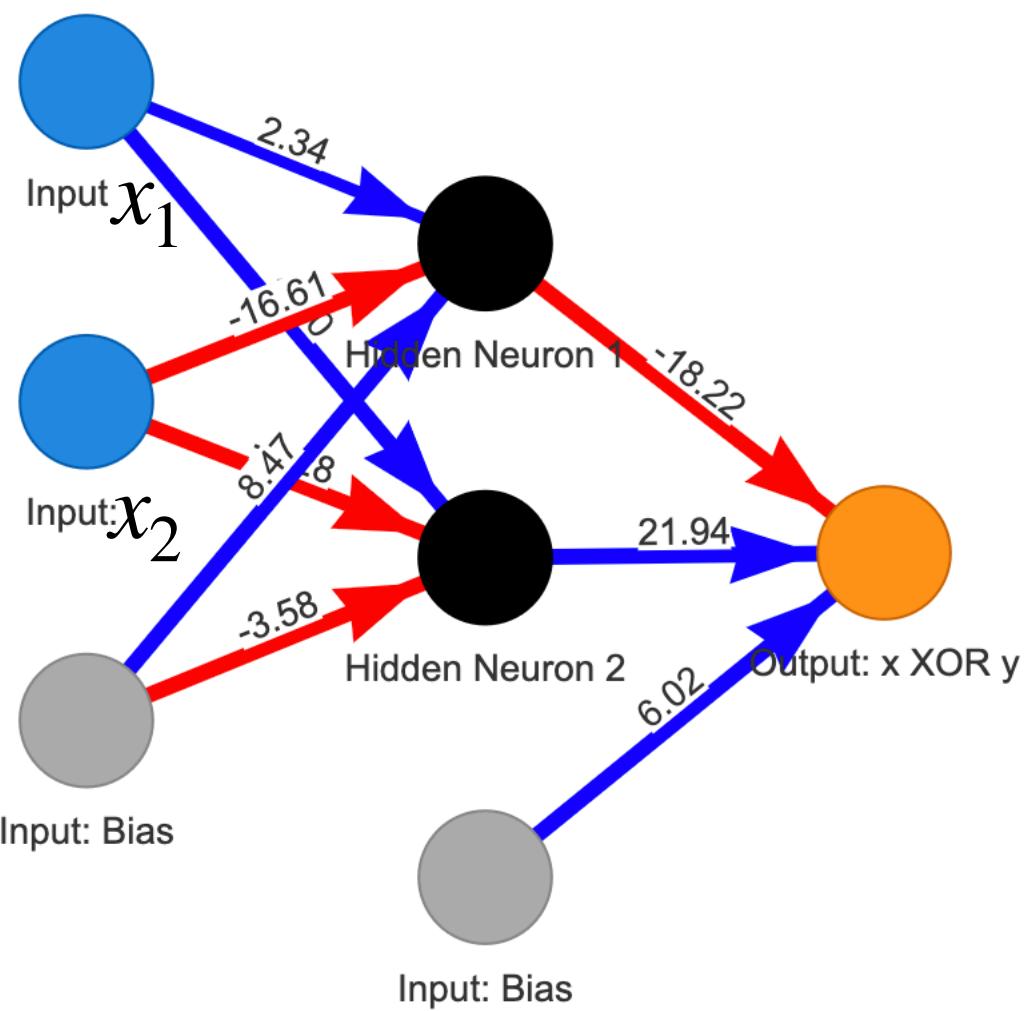
Example



Example generated from

<http://www.ccom.ucsd.edu/~cdeotte/programs/neuralnetwork.html>

Prediction using a neural network





$x_1 = 1.00$



$x_2 = 0.02$



Bias (1)



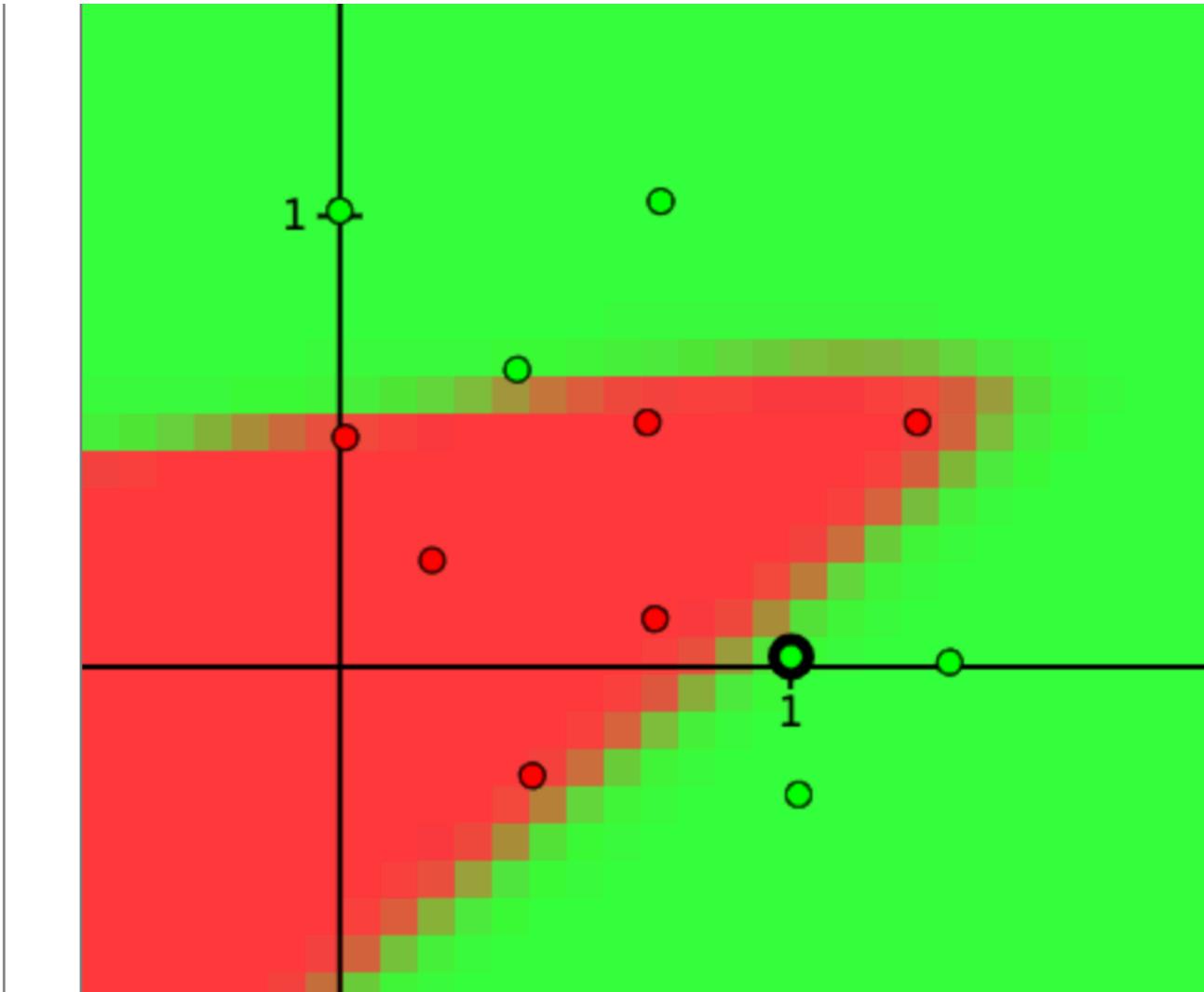
0

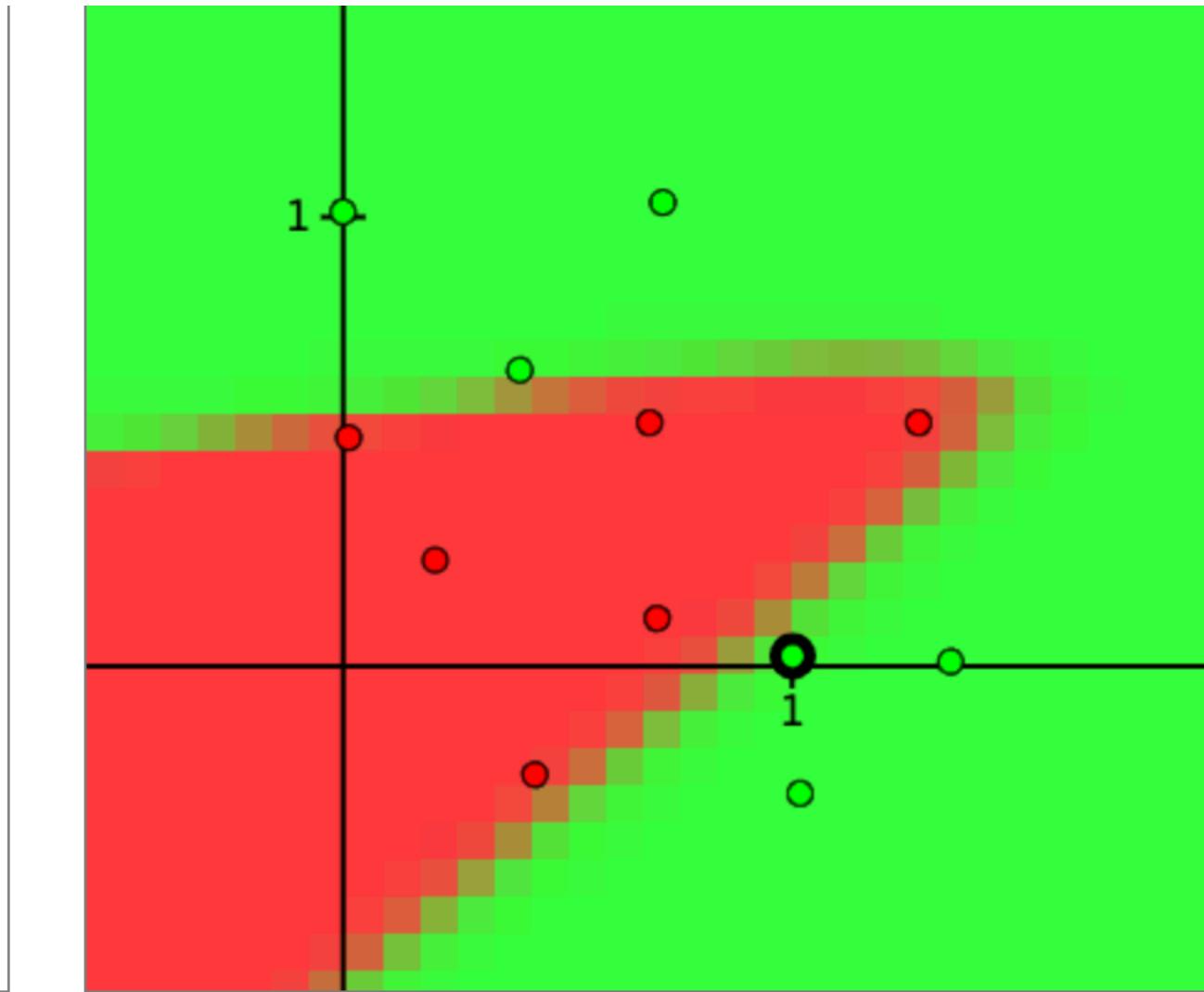
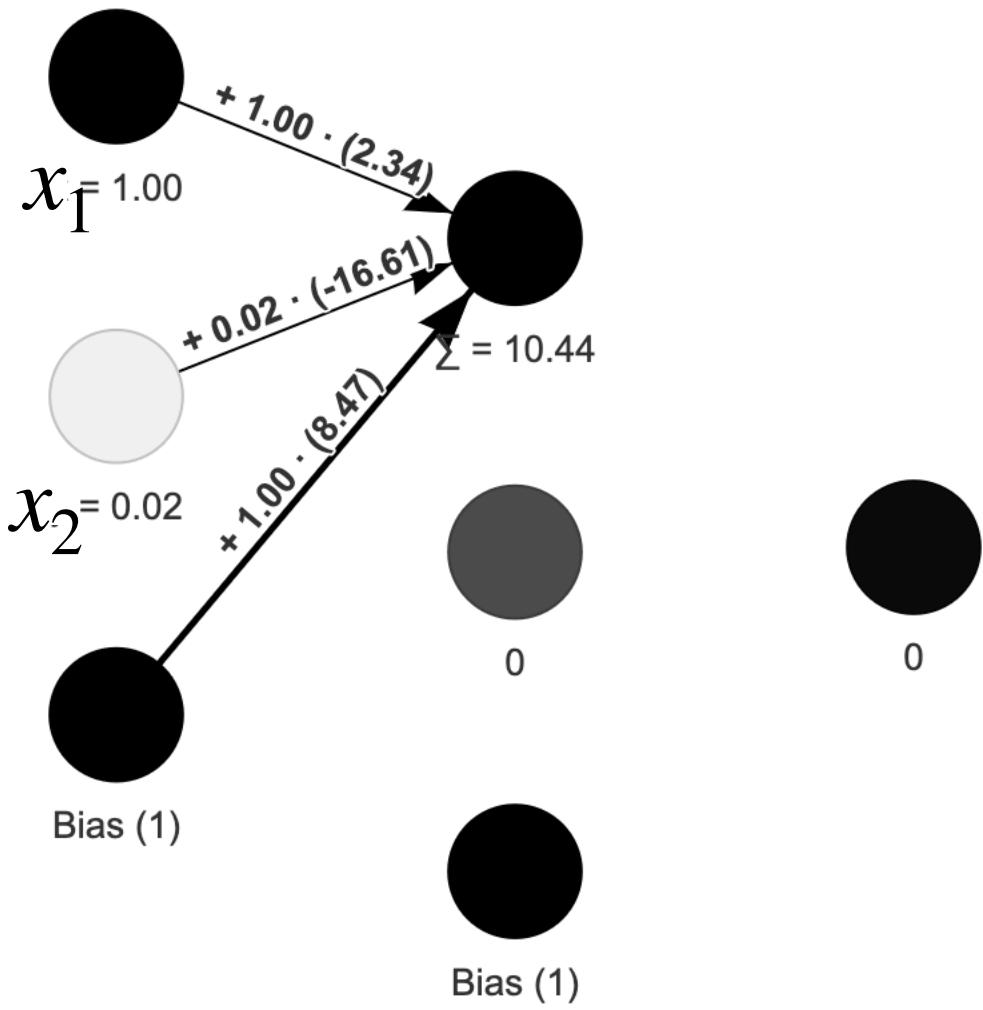


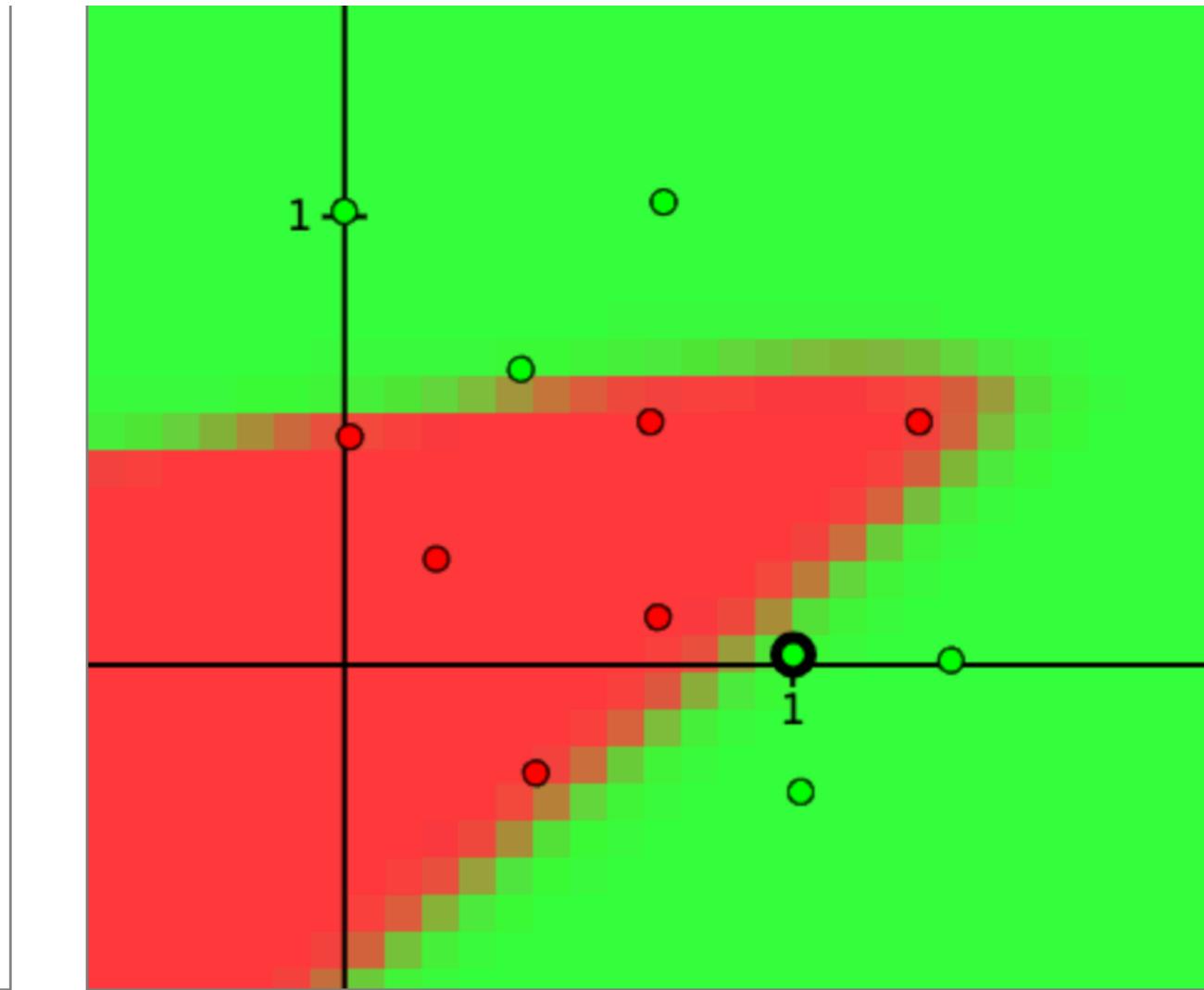
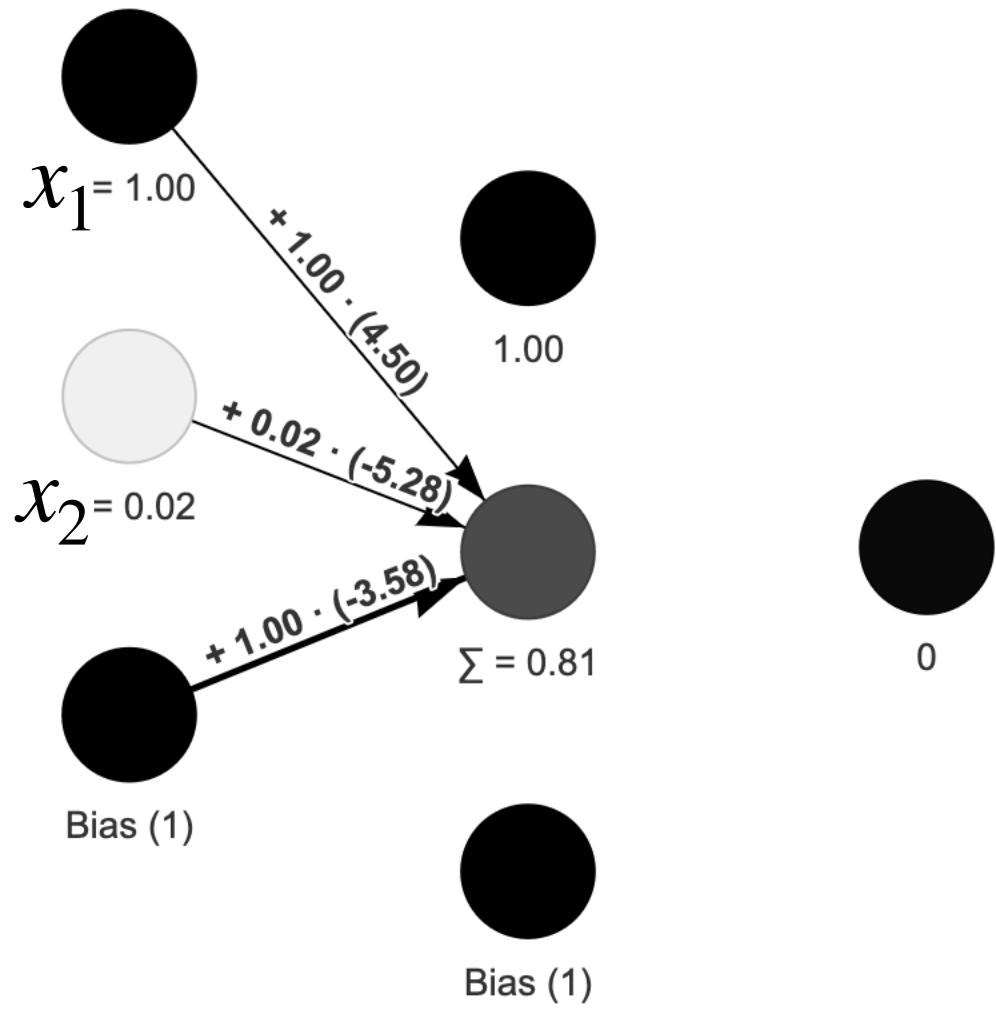
0

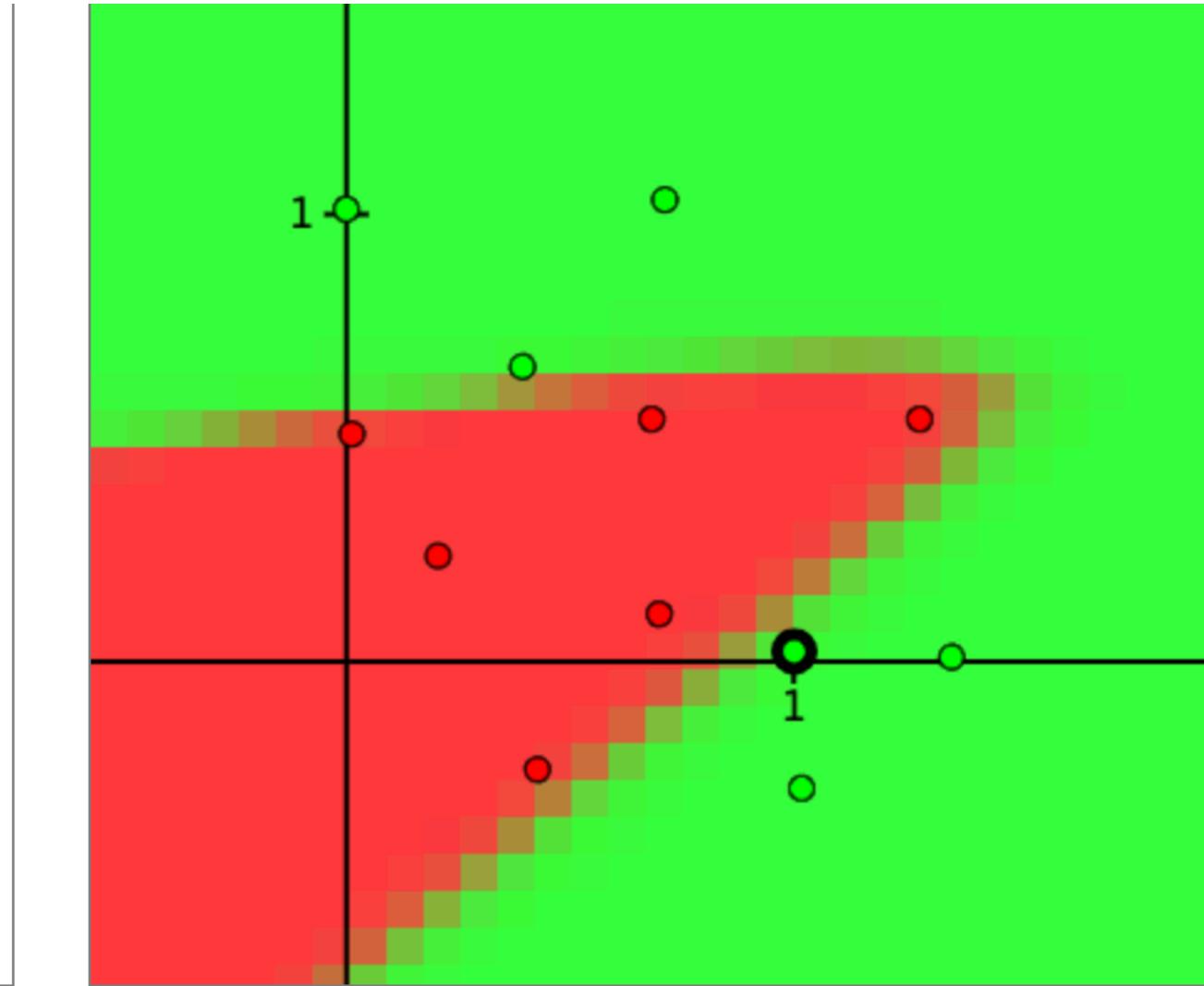
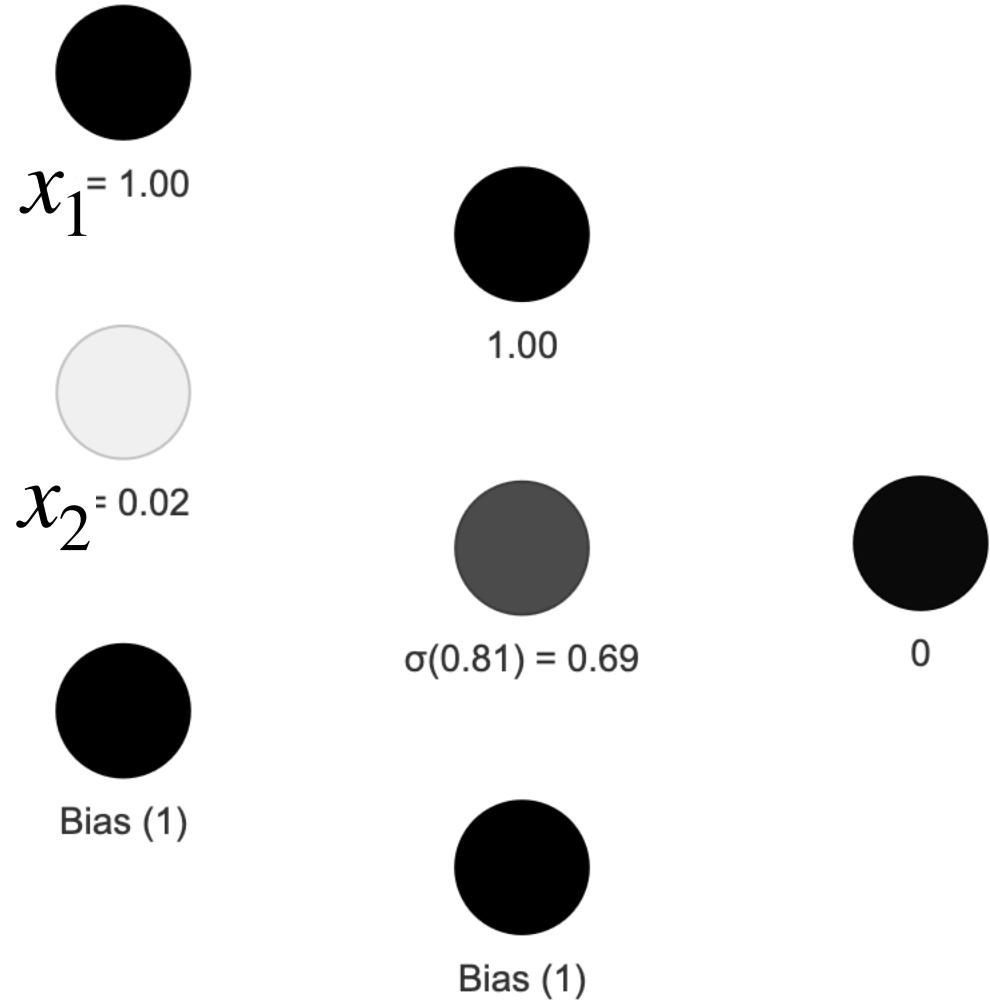


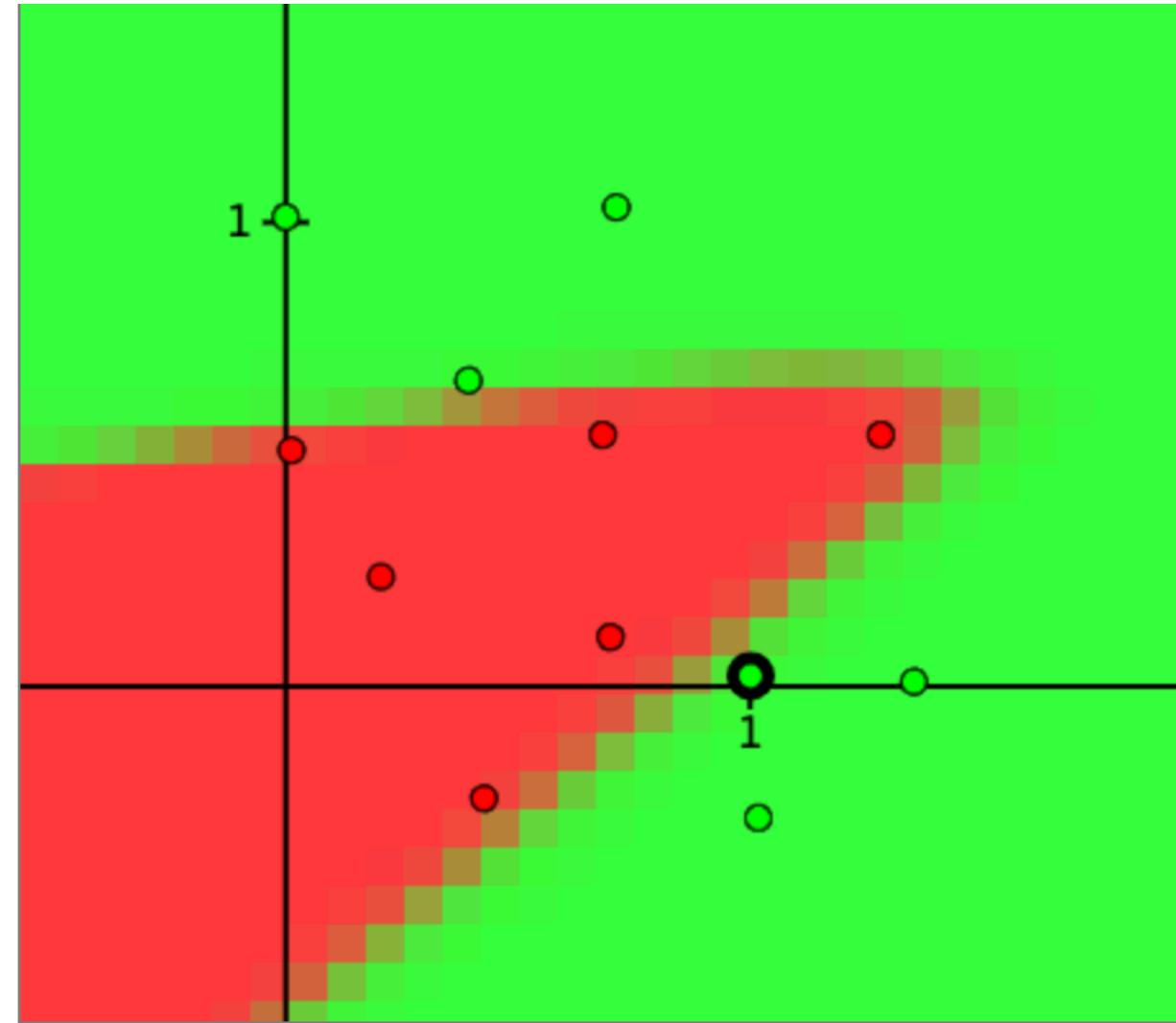
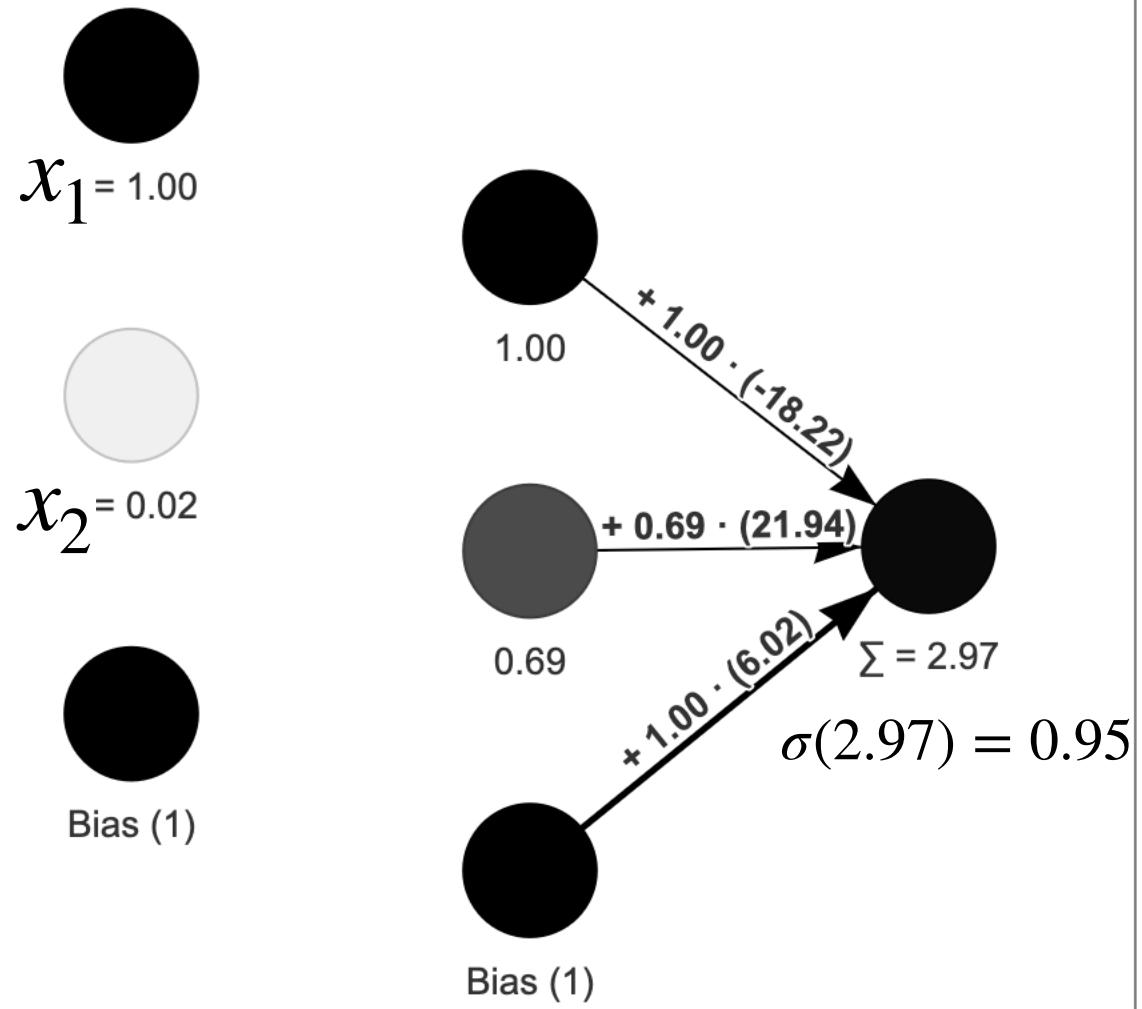
Bias (1)

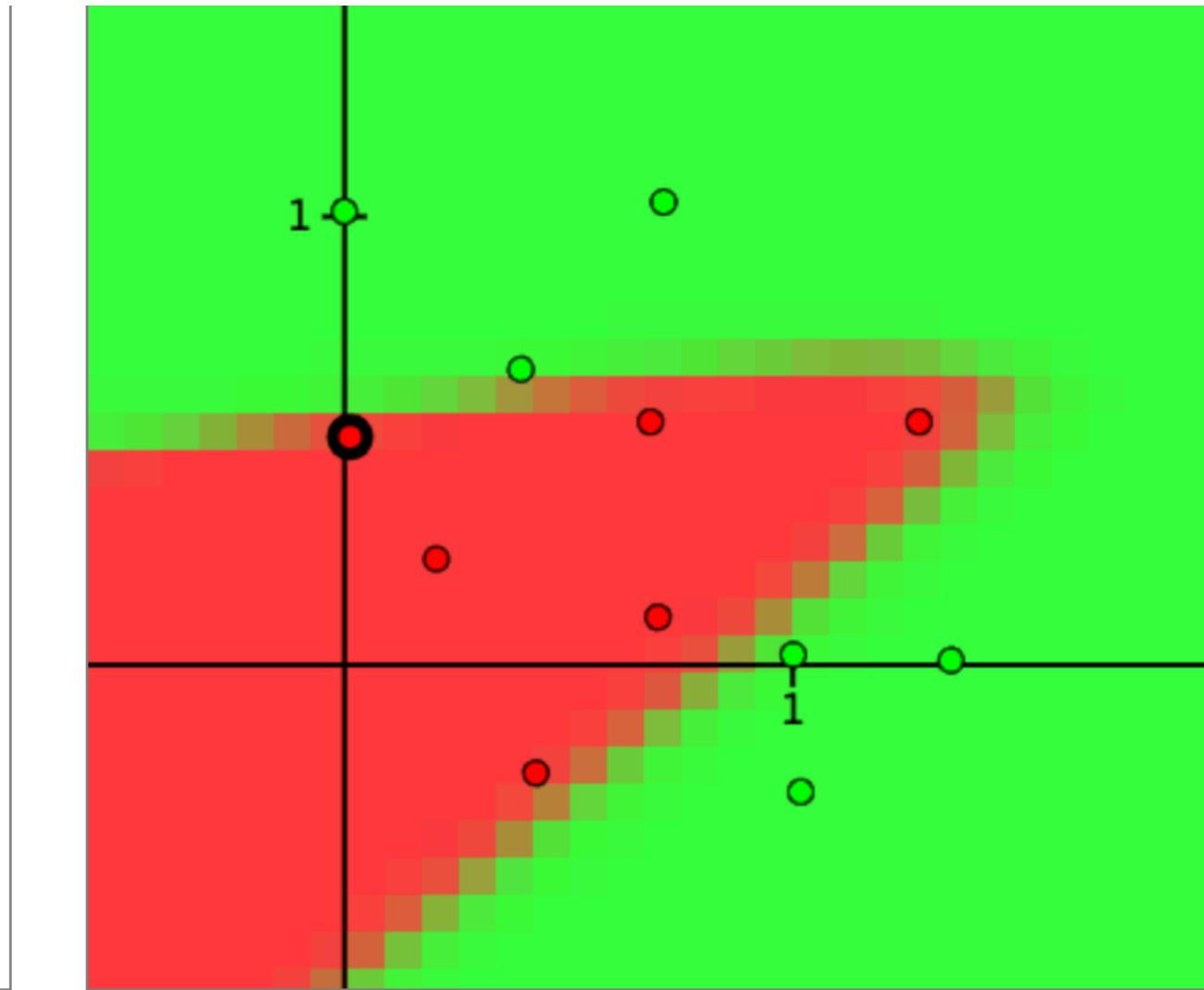
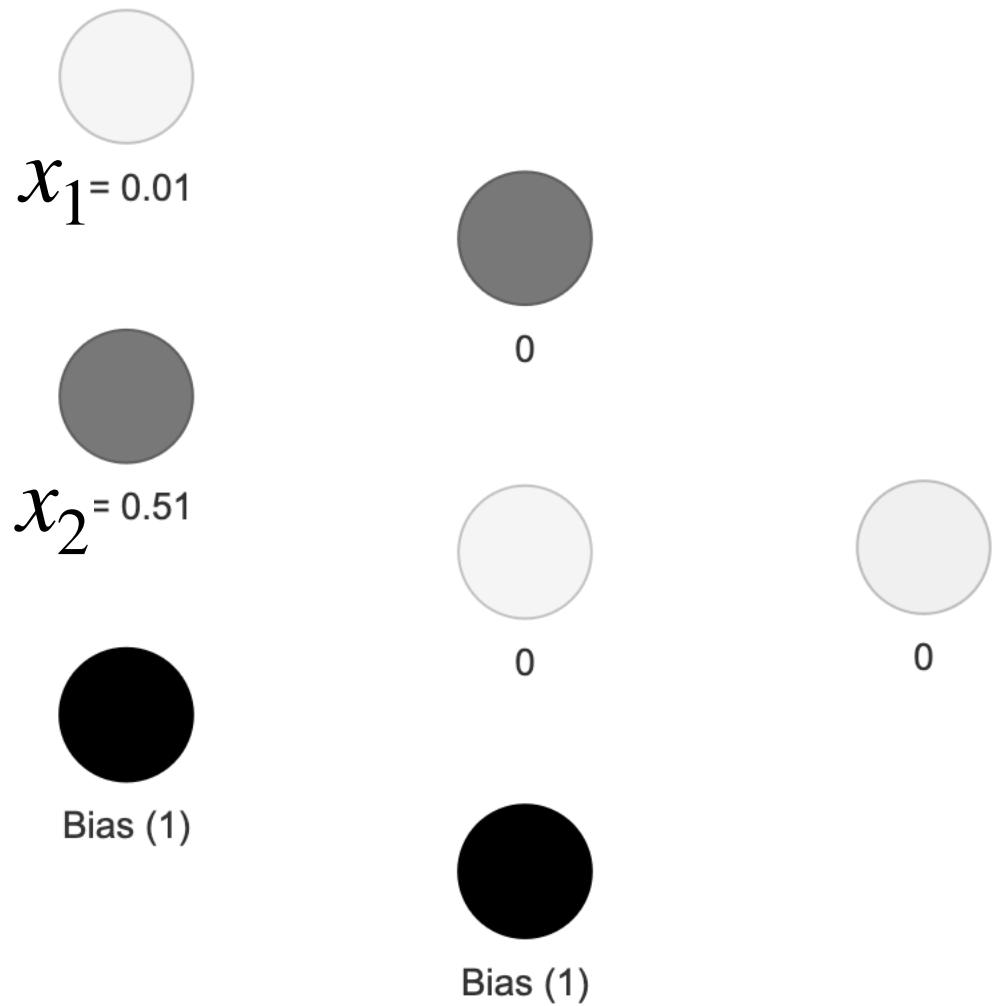


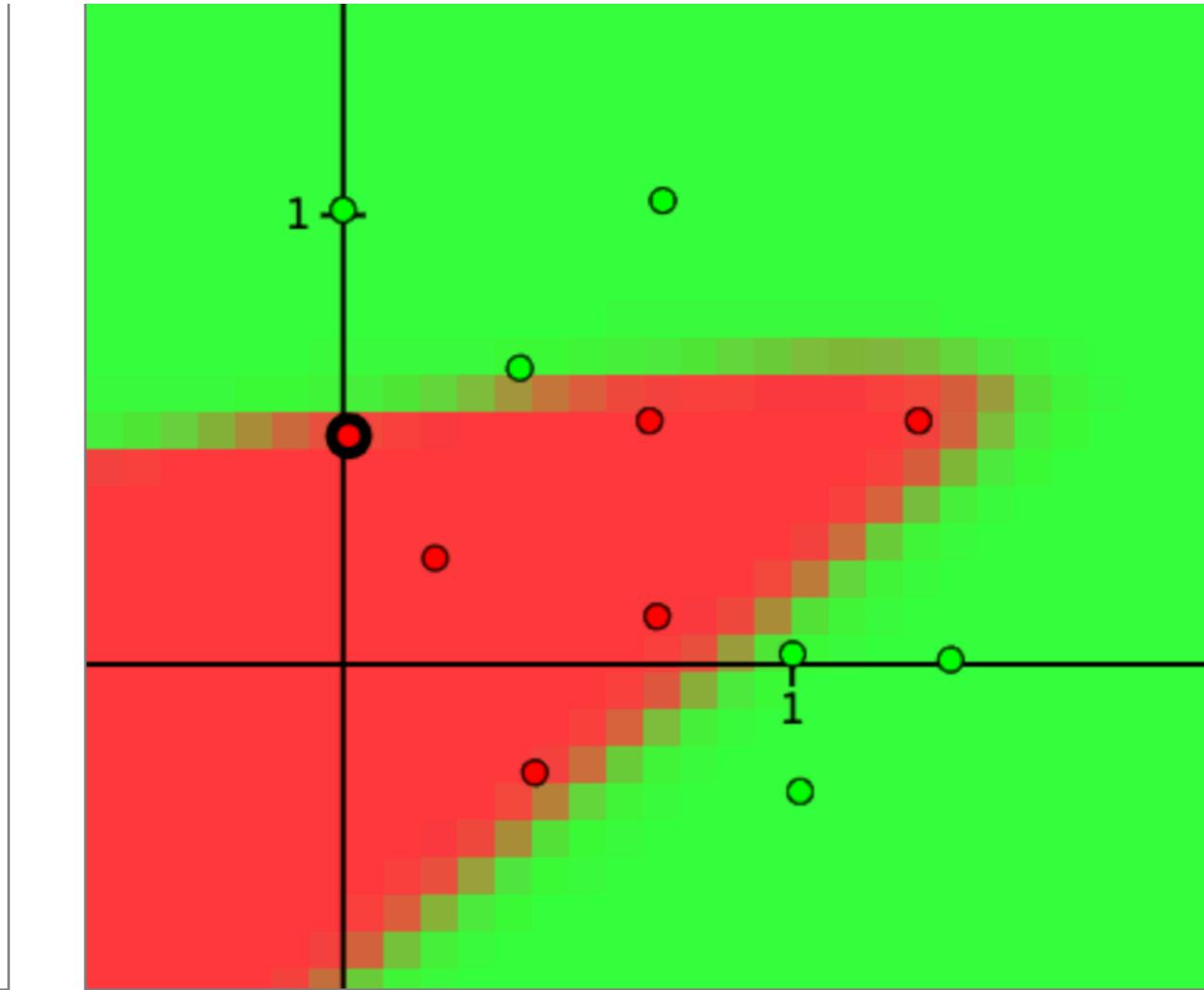
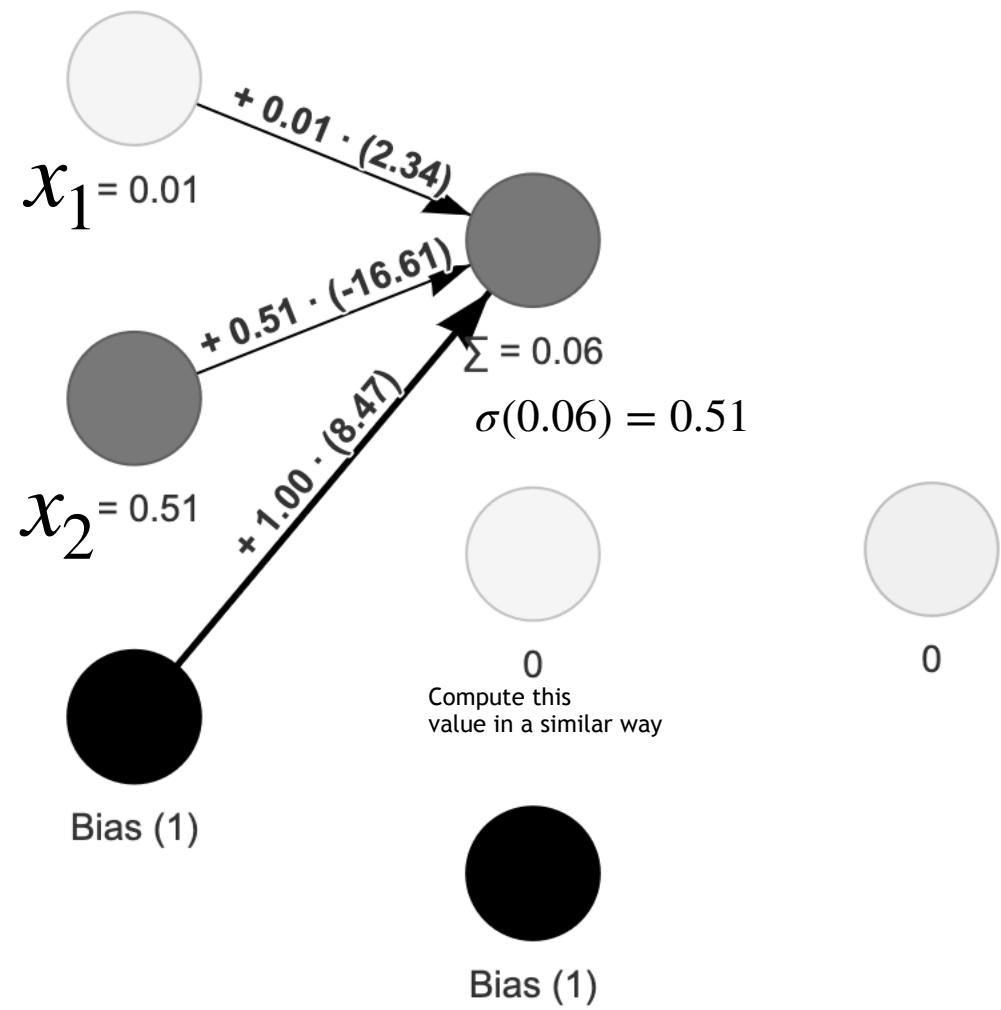


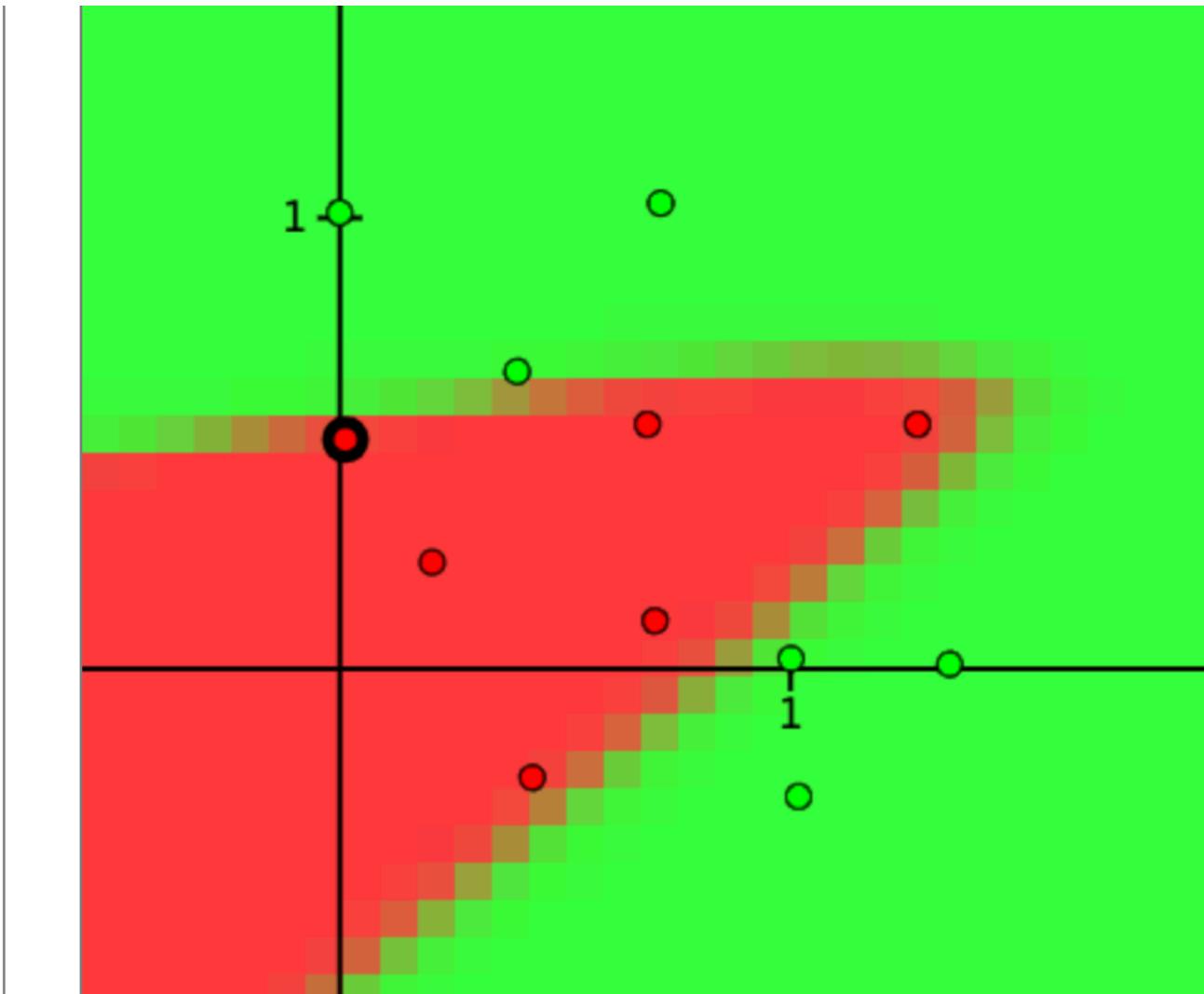
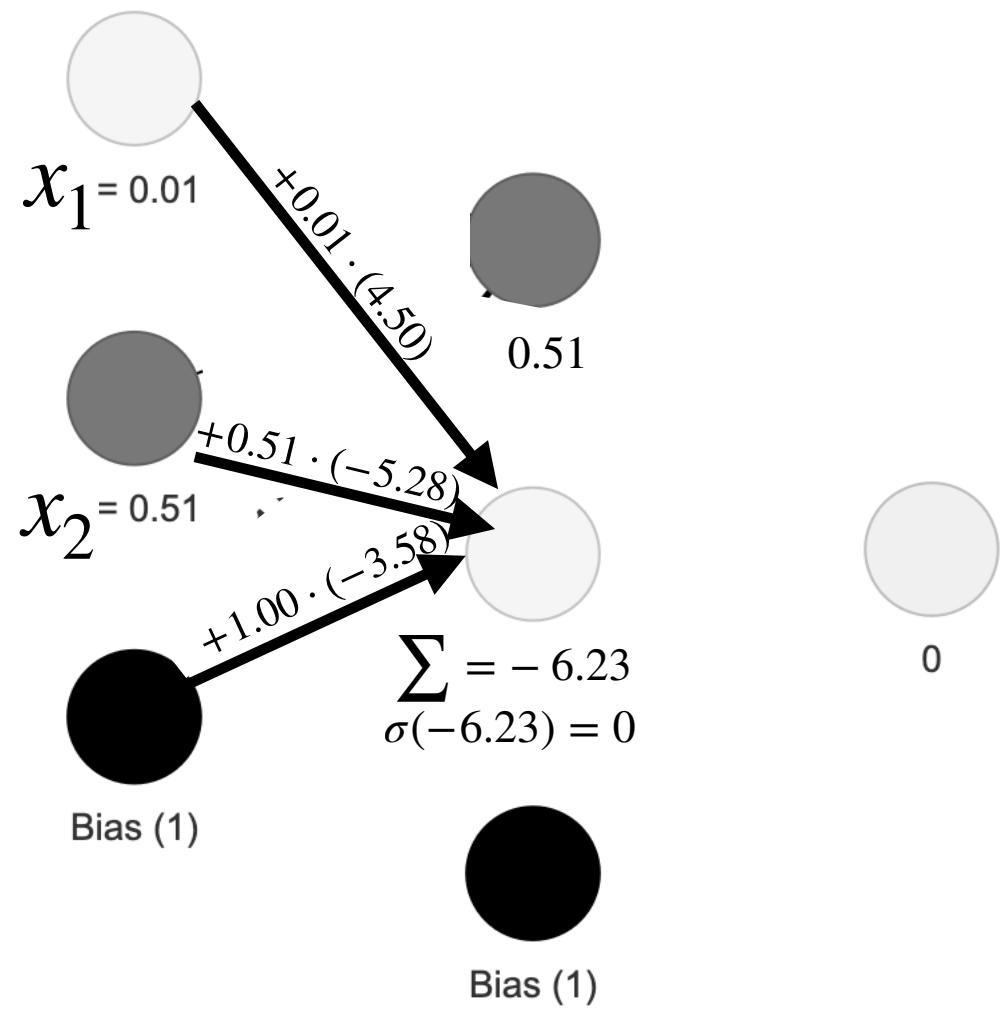


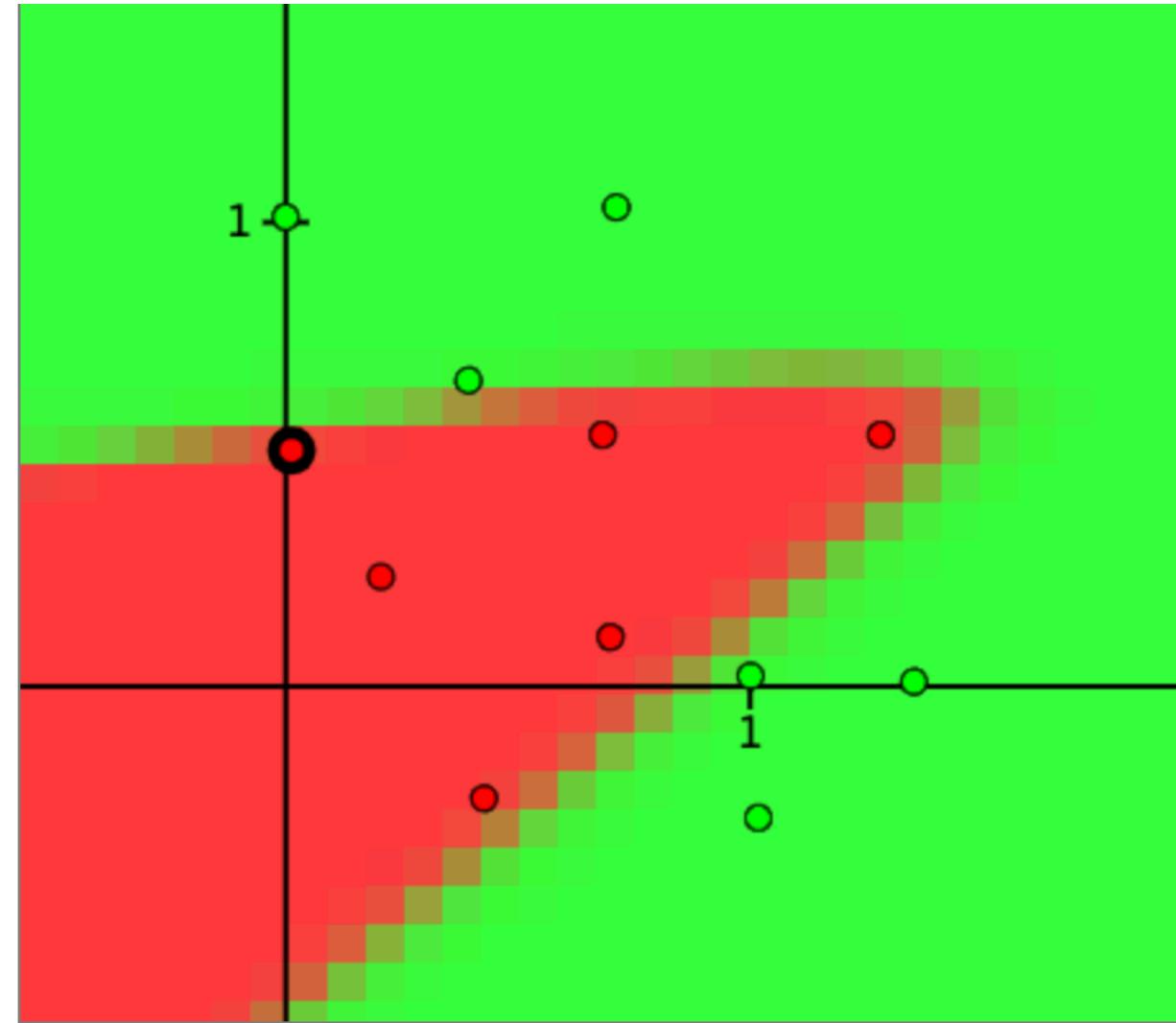
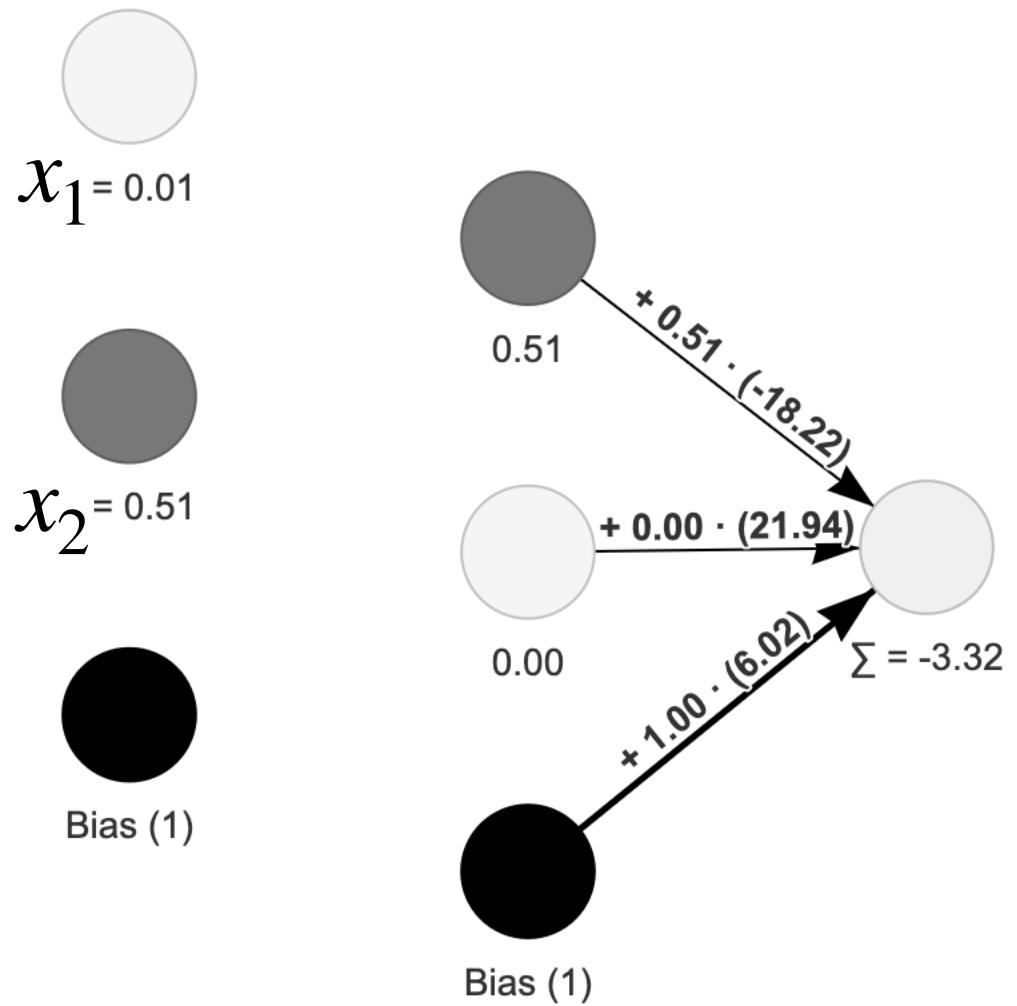














$x_1 = 0.01$



$x_2 = 0.51$



Bias (1)



0.51



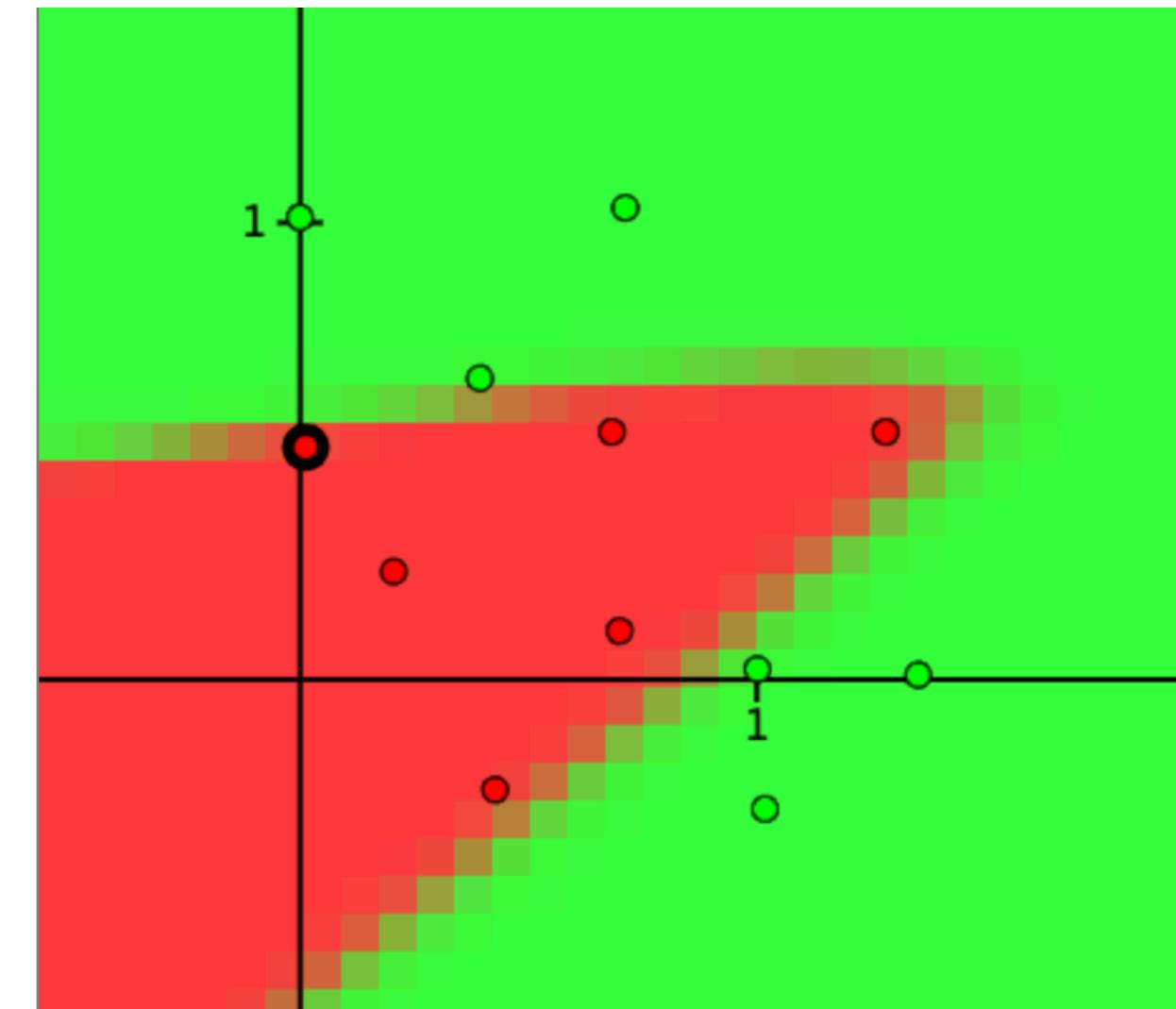
0.00



Bias (1)



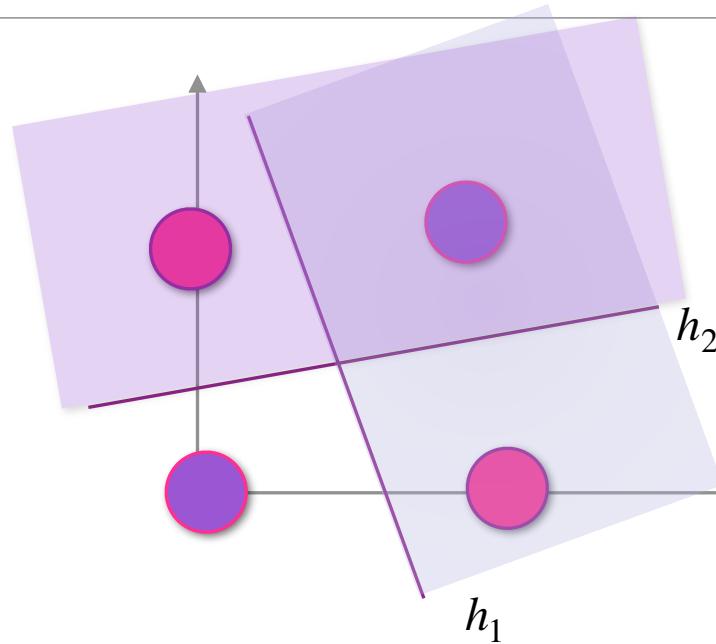
$\sigma(-3.32) = 0.03$



Pair share: can a neural network compute xor?

Boolean functions

x_1	x_2	Xor
0	0	0
0	1	1
1	0	1
1	1	0

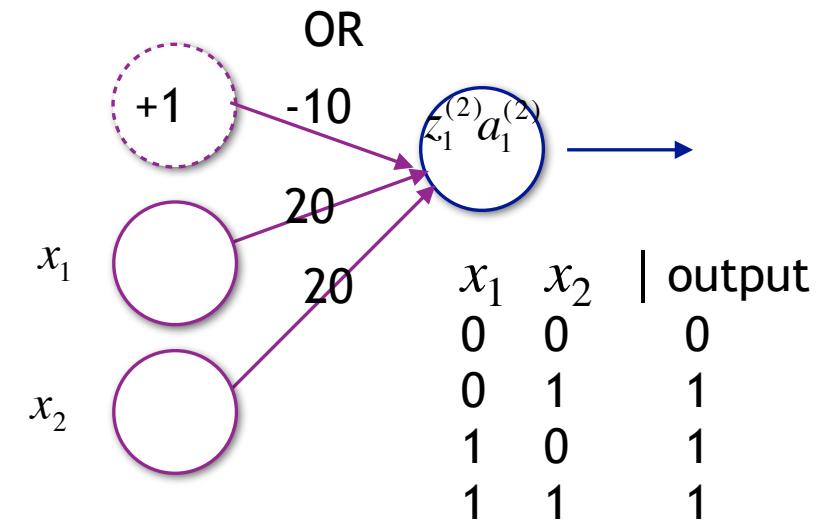
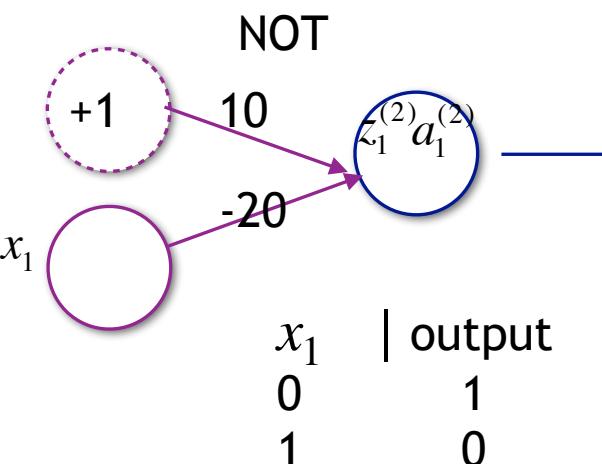
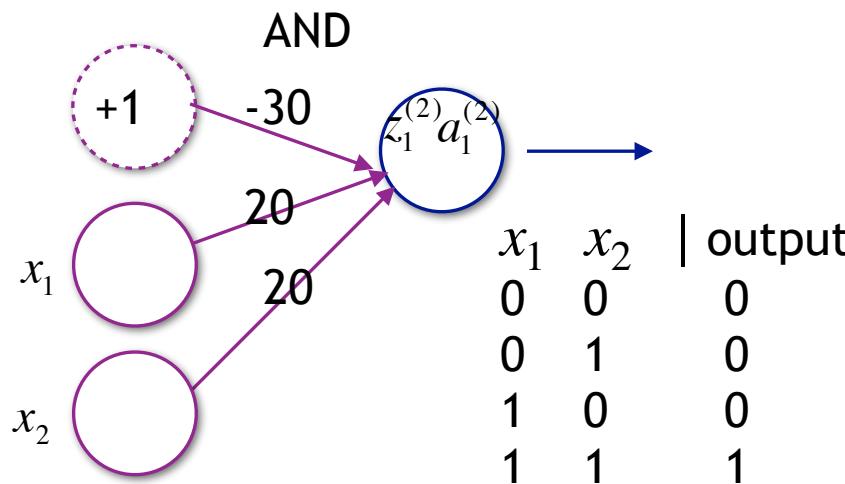


$$h(\mathbf{x}) = \begin{cases} 1 & (h_1(\mathbf{x}) = 1 \text{ and } h_2(x) = 0) \text{ or } (h_1(\mathbf{x}) = 0 \text{ and } h_2(x) = 1) \\ 0 & \text{otherwise} \end{cases}$$

Computing Boolean Functions

It is possible to build a neural network that computes any logical proposition:

$$f(z) = \frac{1}{1 + \exp(-z)}$$



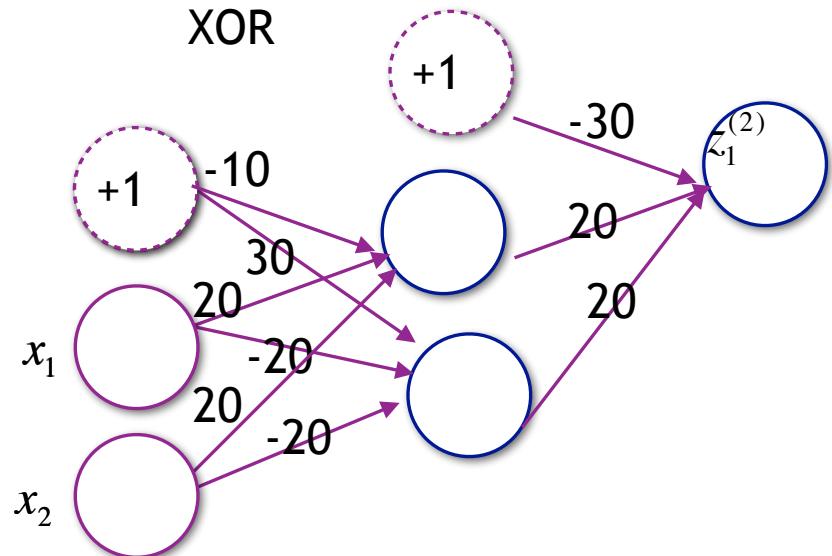
We can create a neural network to approximate XOR using the functions above and the definition:

$$h(\mathbf{x}) = \begin{cases} 1 & (h_1(\mathbf{x}) = 1 \text{ and not } h_2(\mathbf{x}) = 1) \text{ or } (\text{not } h_1(\mathbf{x}) = 1 \text{ and } h_2(\mathbf{x}) = 1) \\ 0 & \text{otherwise} \end{cases}$$

Computing Boolean Functions

$$f(z) = \frac{1}{1 + \exp(-z)}$$

We can approximate XOR using the following:



x_1	x_2	Xor
0	0	0
0	1	1
1	0	1
1	1	0

Every neural network represents an equation

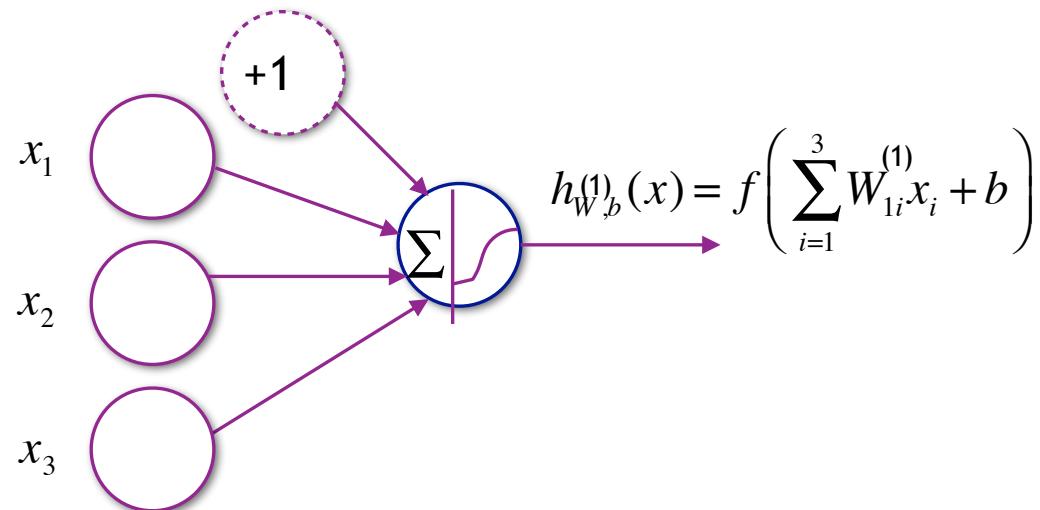
Outline

- ❑ Introduction to neurons
- ❑ Nonlinear classifiers from linear features
-  ❑ Neural networks notation
- ❑ Pseudocode for prediction
- ❑ Training a neural network
- ❑ Implementing gradient descent for neural networks
 - Vectorization
 - Pseudocode
- ❑ Preprocessing
- ❑ Initialization
- ❑ Activations

New Notation

Unfortunately, there will be quite a bit of notation.....

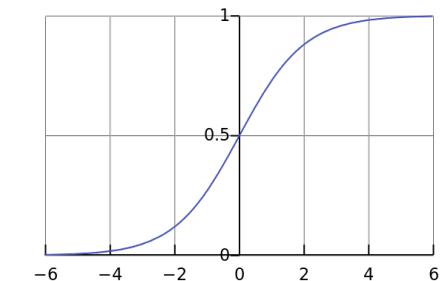
We will follow the notation from http://deeplearning.stanford.edu/wiki/index.php/Neural_Networks



$$W^T = \begin{bmatrix} W_{11}^{(1)} \\ W_{12}^{(1)} \\ W_{13}^{(1)} \end{bmatrix} \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$z = [W_{11}^{(1)}, W_{12}^{(1)}, W_{13}^{(1)}] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + b$$

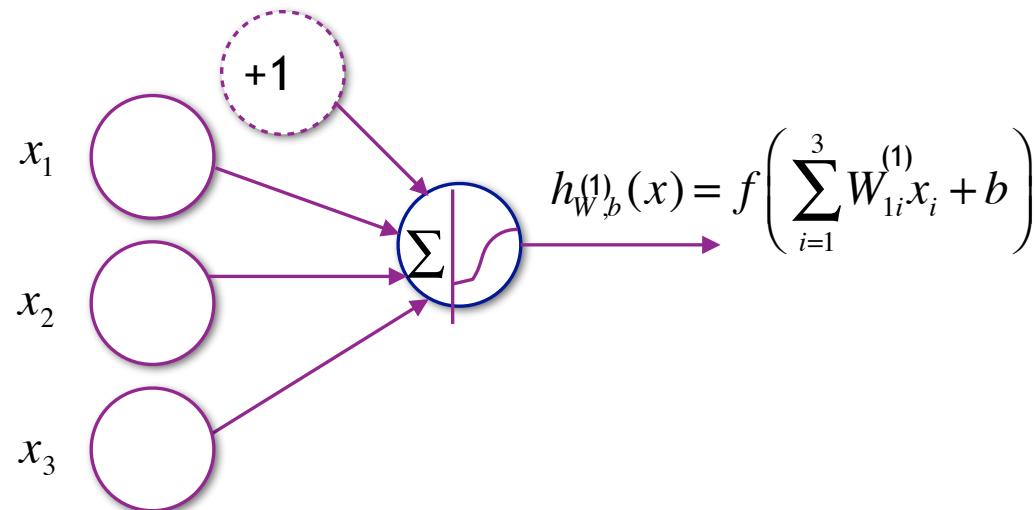
$$f(z) = \frac{1}{1 + \exp(-z)}$$



New Notation

Unfortunately, there will be quite a bit of notation.....

We will follow the notation from http://deeplearning.stanford.edu/wiki/index.php/Neural_Networks

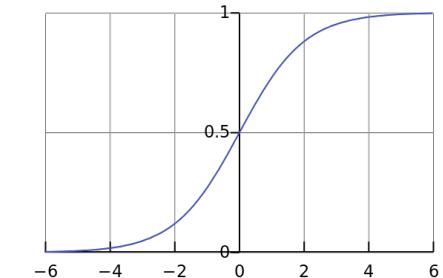


$z = [W_{11}^{(1)}, W_{12}^{(1)}, W_{13}^{(1)}]$ approximating single “neuron” using the sigmoid function

$$W^T = \begin{bmatrix} W_{11}^{(1)} \\ W_{12}^{(1)} \\ W_{13}^{(1)} \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

I added another subscript and superscript - we will need these later

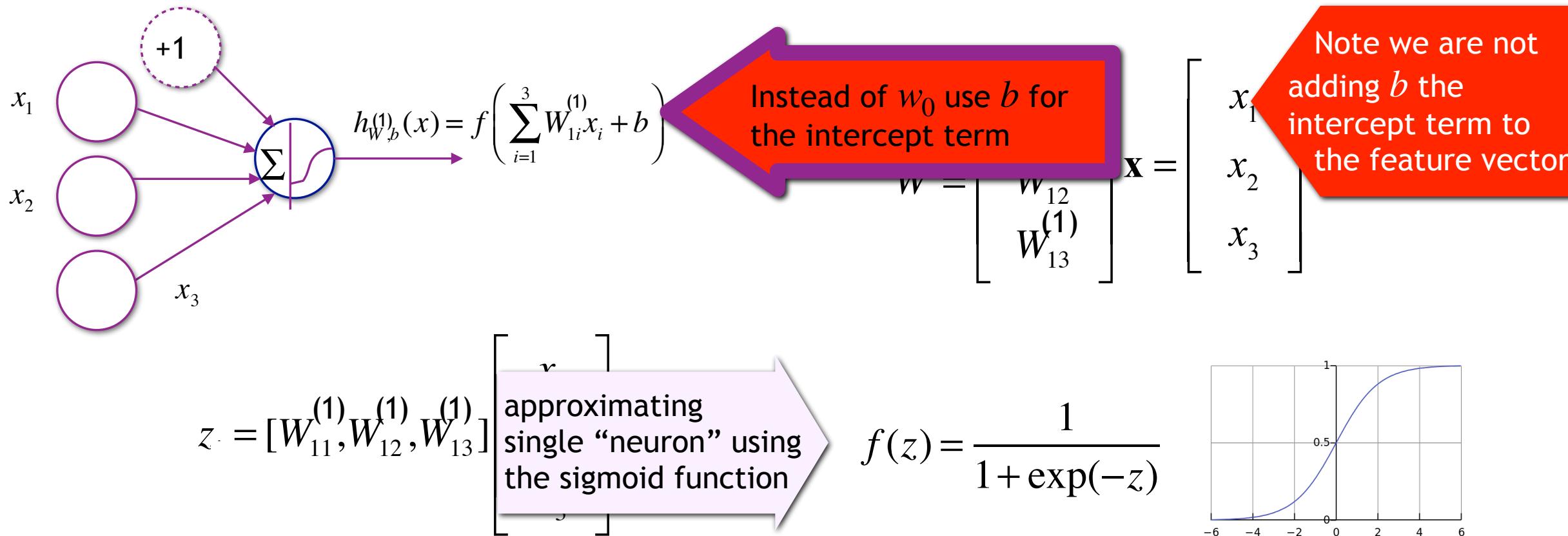
$$f(z) = \frac{1}{1 + \exp(-z)}$$



New Notation

Unfortunately, there will be quite a bit of notation.....

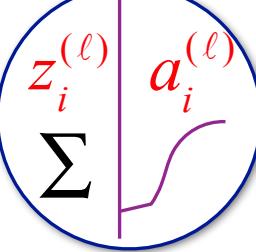
We will follow the notation from http://deeplearning.stanford.edu/wiki/index.php/Neural_Networks

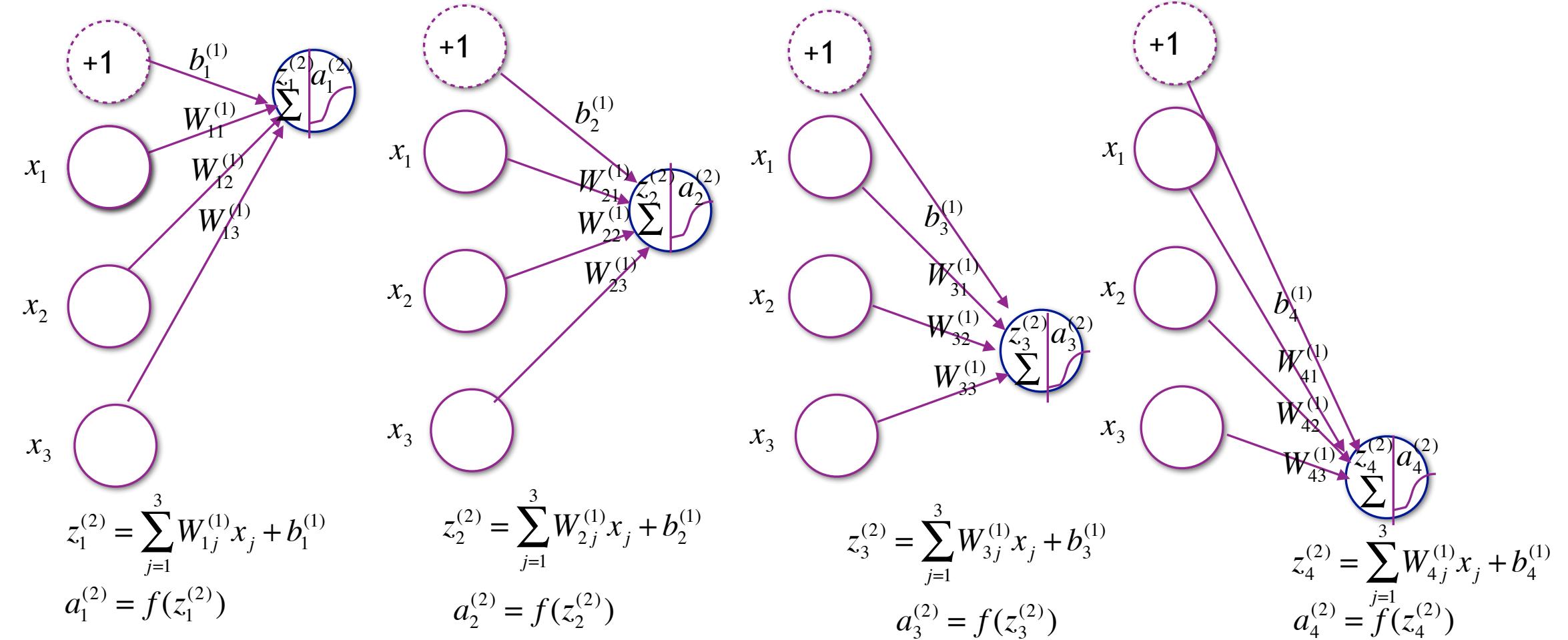


Neural Network Model & Notation

$$f(z) = \frac{1}{1 + \exp(-z)}$$

$$z_i^{(2)} = \sum_{j=1}^3 W_{ij}^{(1)} x_j + b_i^{(1)}$$

$$a_i^{(2)} = f(z_i^{(2)})$$


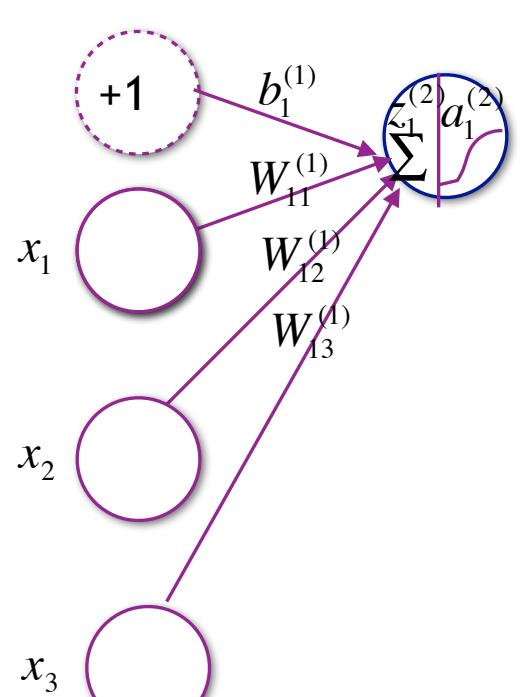


Neural Network Model & Notation

$$z_i^{(2)} = \sum_{j=1}^3 W_j^{(1)} x_j + b_i^{(1)}$$

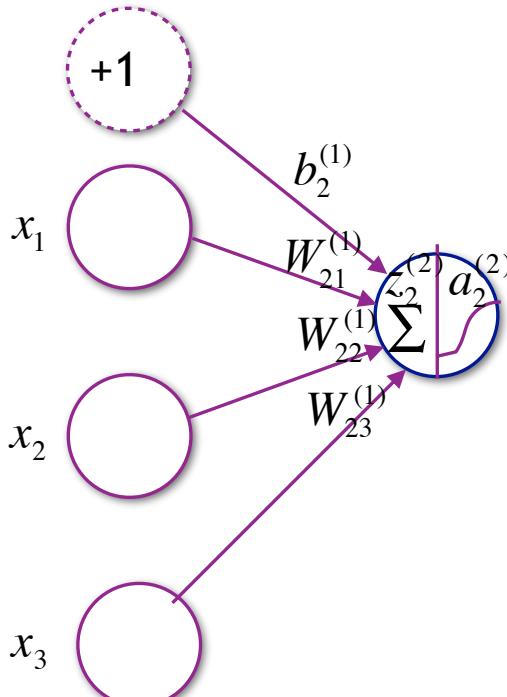
$$a_i^{(2)} = f(z_i^{(2)})$$

$$f(z) = \frac{1}{1 + \exp(-z)}$$



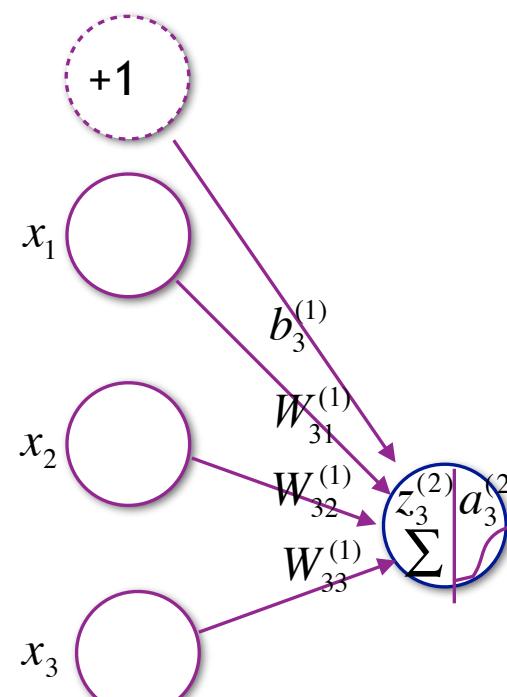
$$z_1^{(2)} = \sum_{j=1}^3 W_{1j} x_j + b_1^{(1)}$$

$$a_1^{(2)} = f(z_1^{(2)})$$



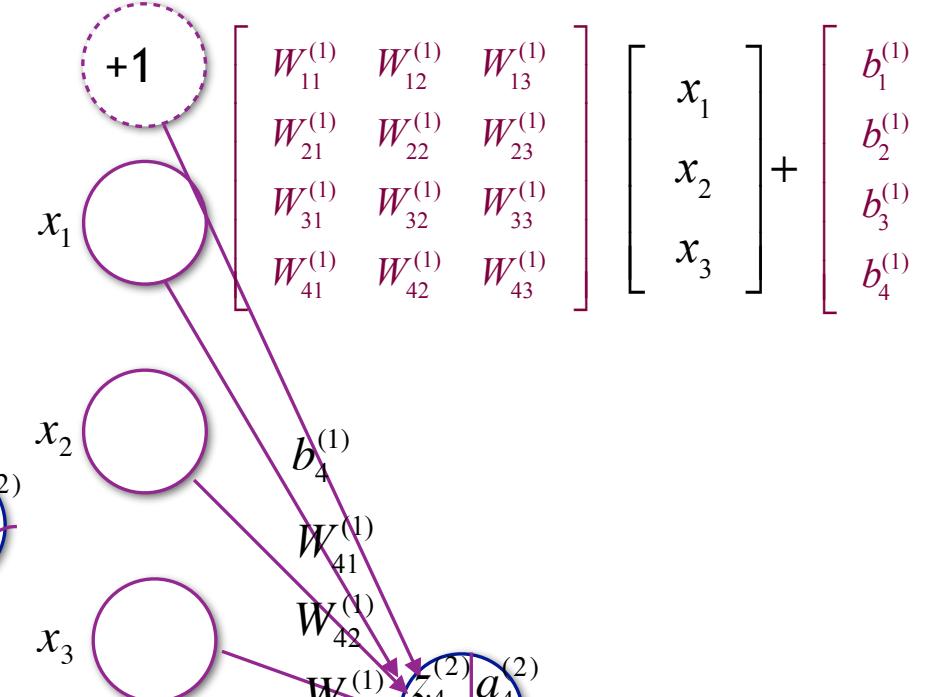
$$z_2^{(2)} = \sum_{j=1}^3 W_{2j} x_j + b_2^{(1)}$$

$$a_2^{(2)} = f(z_2^{(2)})$$



$$z_3^{(2)} = \sum_{j=1}^3 W_{3j} x_j + b_3^{(1)}$$

$$a_3^{(2)} = f(z_3^{(2)})$$



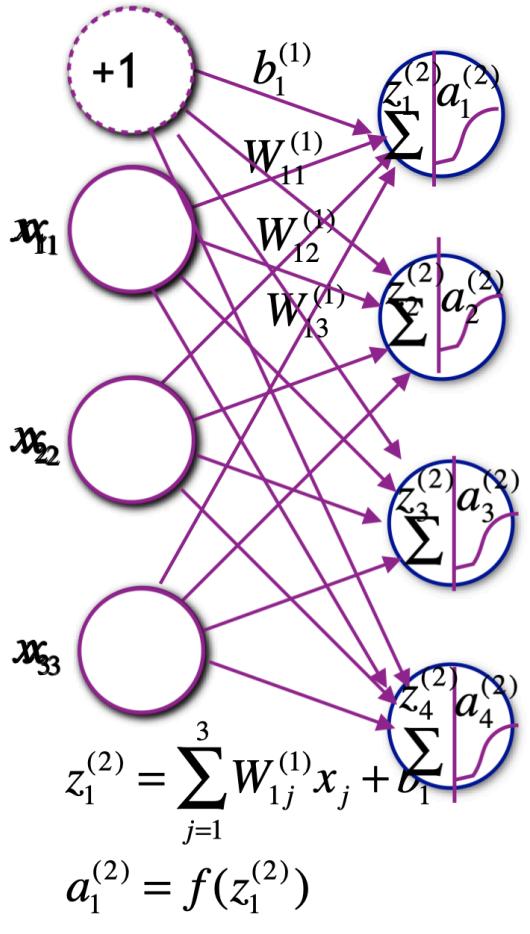
$$z_4^{(2)} = \sum_{j=1}^3 W_{4j} x_j + b_4^{(1)}$$

$$a_4^{(2)} = f(z_4^{(2)})$$

Neural Network Model & Notation

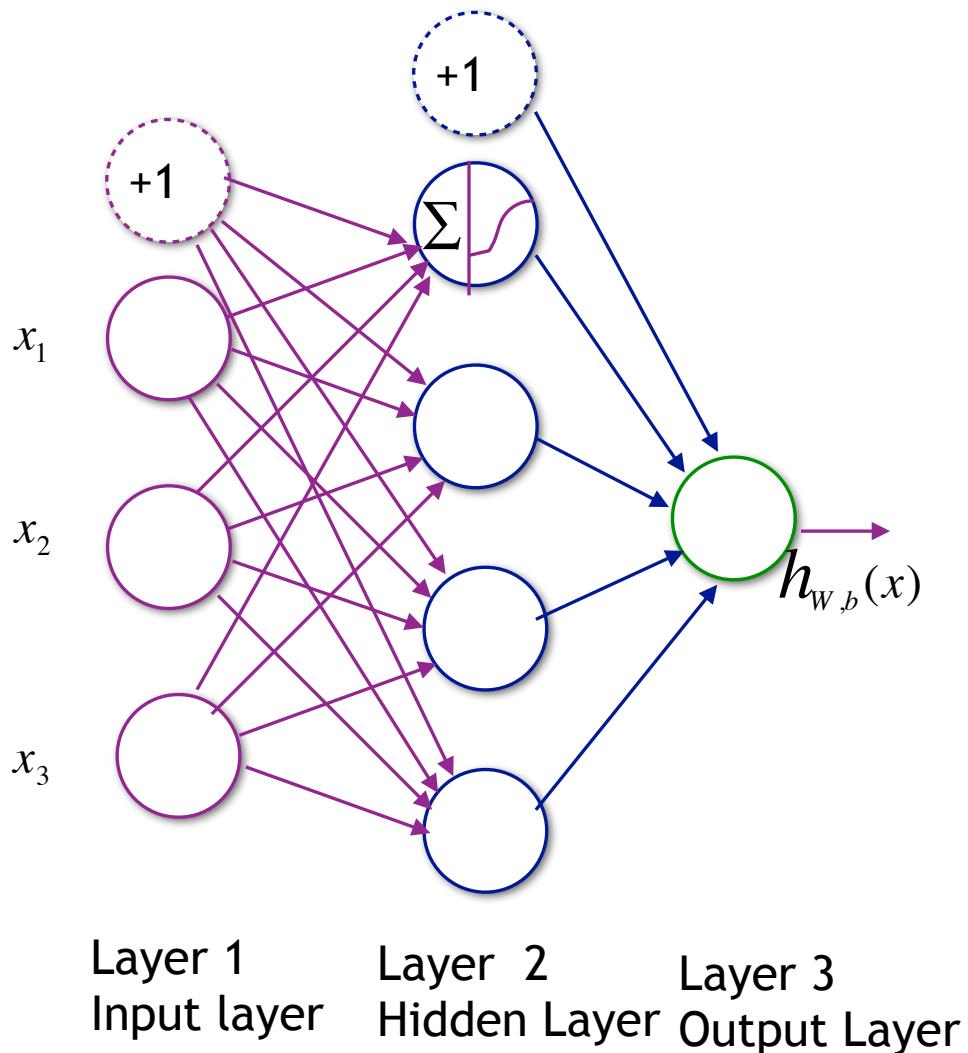
$$z_i^{(2)} = \sum_{j=1}^3 W_j^{(1)} x_j + b_i^{(1)} \quad a_i^{(2)} = f(z_i^{(2)})$$

$$f(z) = \frac{1}{1 + \exp(-z)}$$



$$\begin{bmatrix} W_{11}^{(1)} & W_{12}^{(1)} & W_{13}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} & W_{23}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} & W_{33}^{(1)} \\ W_{41}^{(1)} & W_{42}^{(1)} & W_{43}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}$$

Neural Network Model & Notation

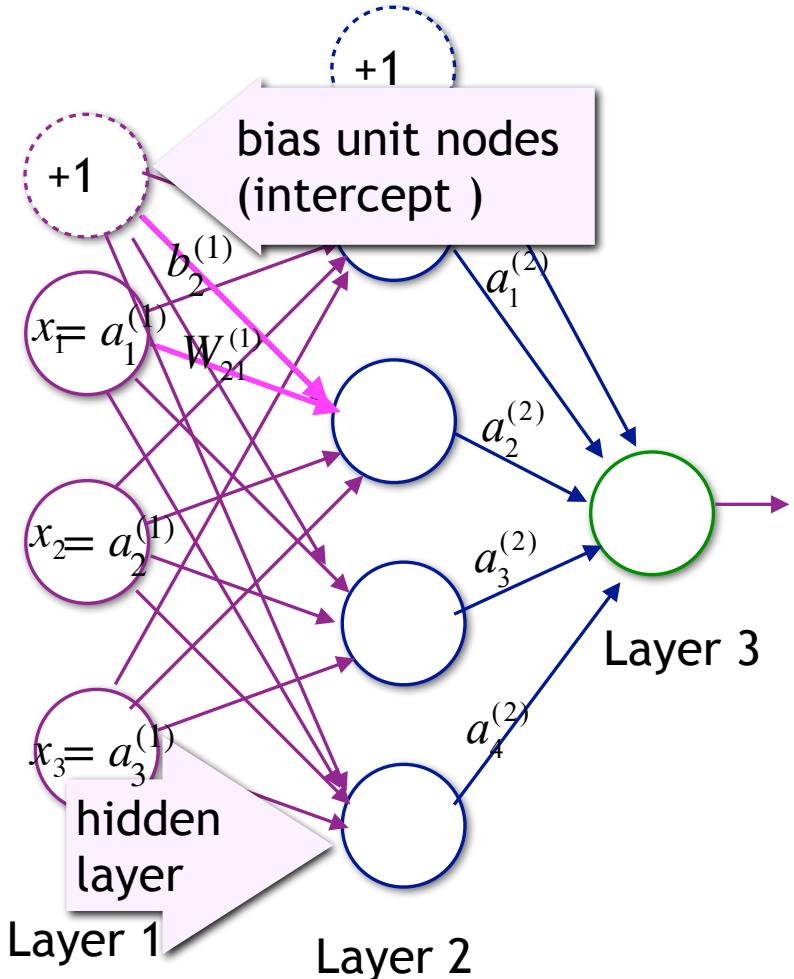


- We can combine many of the simple “neurons” so the output of one is the input of another neuron
- Input layer units are set by x
- The input for all other nodes are typically the *weighted* sum of the **activation** on the links connected to this node
- The **activation function** is applied to the input to provide the value (we will use the sigmoid function)

$$f(z) = \frac{1}{1 + \exp(-z)}$$

Neural Network Model & Notation

input layer
from input units



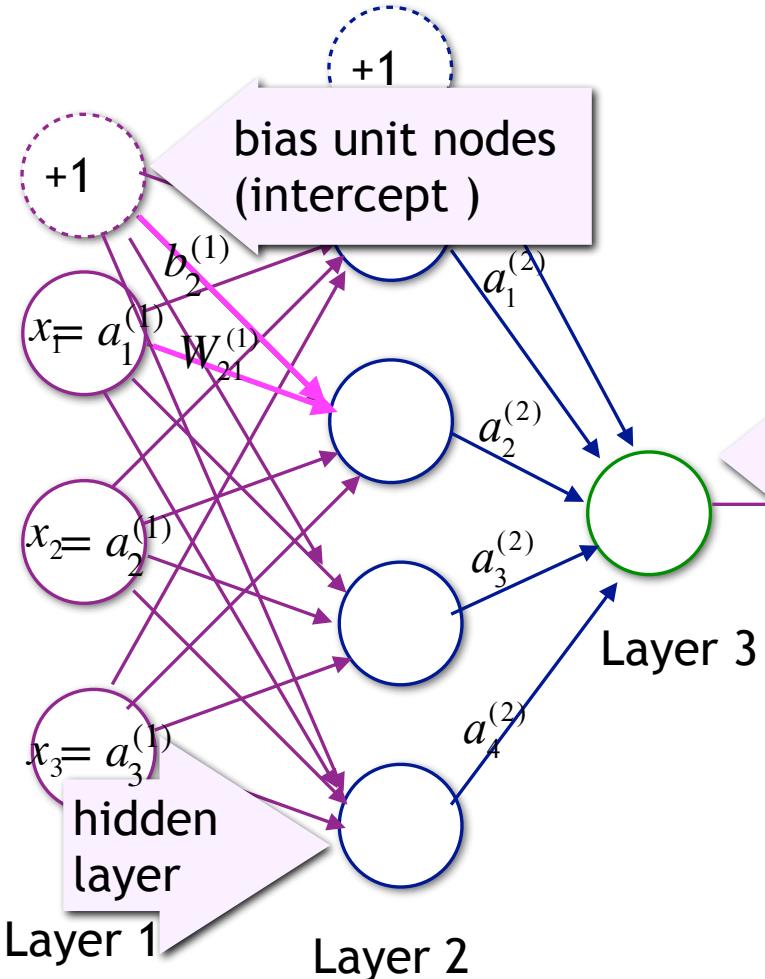
$$W_{ij}^{(\ell)} \quad \left\{ \begin{array}{lll} 1 \leq \ell < n_l & \text{layers} \\ 1 \leq j \leq s_\ell & \text{inputs} \\ 1 \leq i \leq s_{\ell+1} & \text{outputs} \end{array} \right.$$

- ❑ The network has 3 input units (not 4), 4 hidden units (not 5) and 1 output unit
- ❑ n_ℓ is the number of layers in the network
- ❑ Each “neuron” has its own weights. We will store all the weights for one level together in a matrix $W^{(\ell)}$ such that $W_{ij}^{(\ell)}$ is the weight for leaving node j at level 1 and entering node i at level 2
- ❑ $b_i^{(1)}$ is the bias of unit i in level 2

$$W^{(1)} \in \mathbb{R}^{4 \times 3} \quad W^{(2)} \in \mathbb{R}^{1 \times 4}$$

Neural Network Model & Notation

input layer from input units



$$W_{ij}^{(\ell)} \quad \left\{ \begin{array}{lll} 1 \leq \ell < n_l & \text{layers} \\ 1 \leq j \leq s_\ell & \text{inputs} \\ 1 \leq i \leq s_{\ell+1} & \text{outputs} \end{array} \right.$$

- ❑ The network has 3 input units (not 4), 4 hidden units (not 5) and 1 output unit

The number of layers is $n_\ell=3$. We are counting the input layer - many sources do not count this layer

store all the weights for one level together in a matrix $W^{(\ell)}$ such that $W_{ij}^{(1)}$ is the weight for leaving node j at level 1 and entering node i at level 2

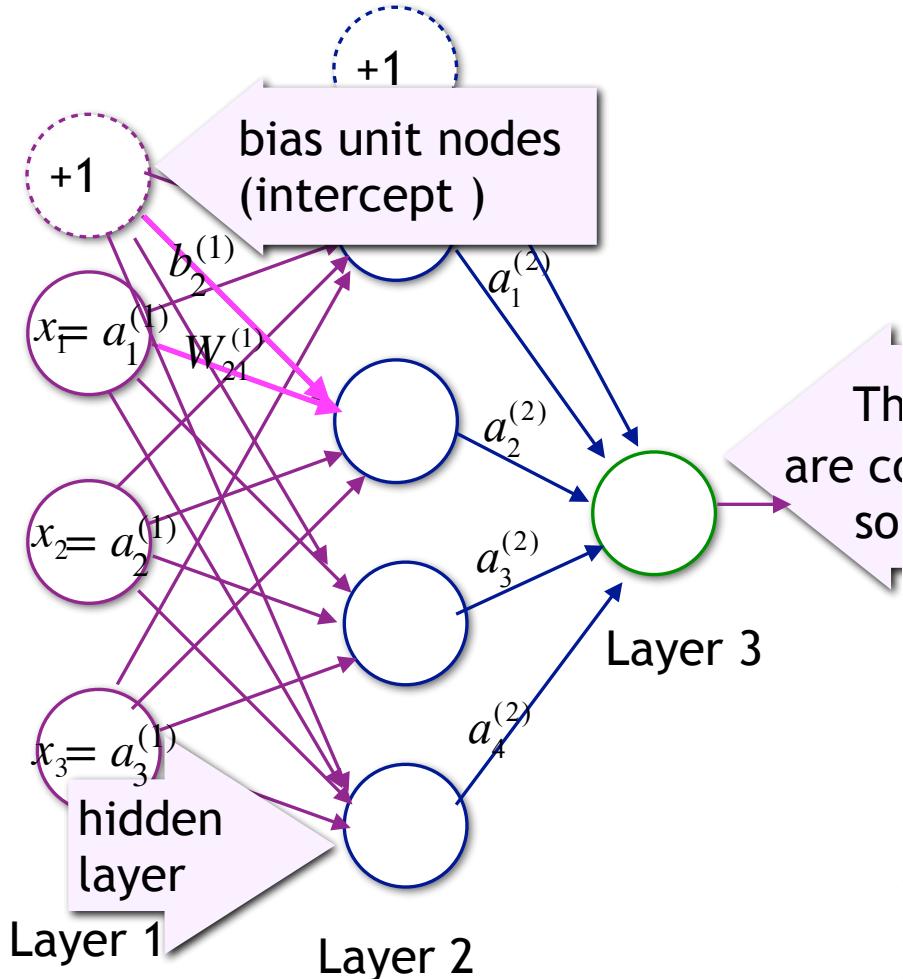
- ❑ $b_i^{(1)}$ is the bias of unit i in level 2

$$W^{(1)} \in \mathbb{R}^{4 \times 3}$$

$$W^{(2)} \in \mathbb{R}^{1 \times 4}$$

Neural Network Model & Notation

input layer
from input units



1. How many parameters to train in this neural network?

(a) 0-3

(b) 4-6

(c) 7-10

(d) 11-14

(e) 5-18

(f) 19-22

(g) 23-26

(h) more than 26

layers
inputs
outputs

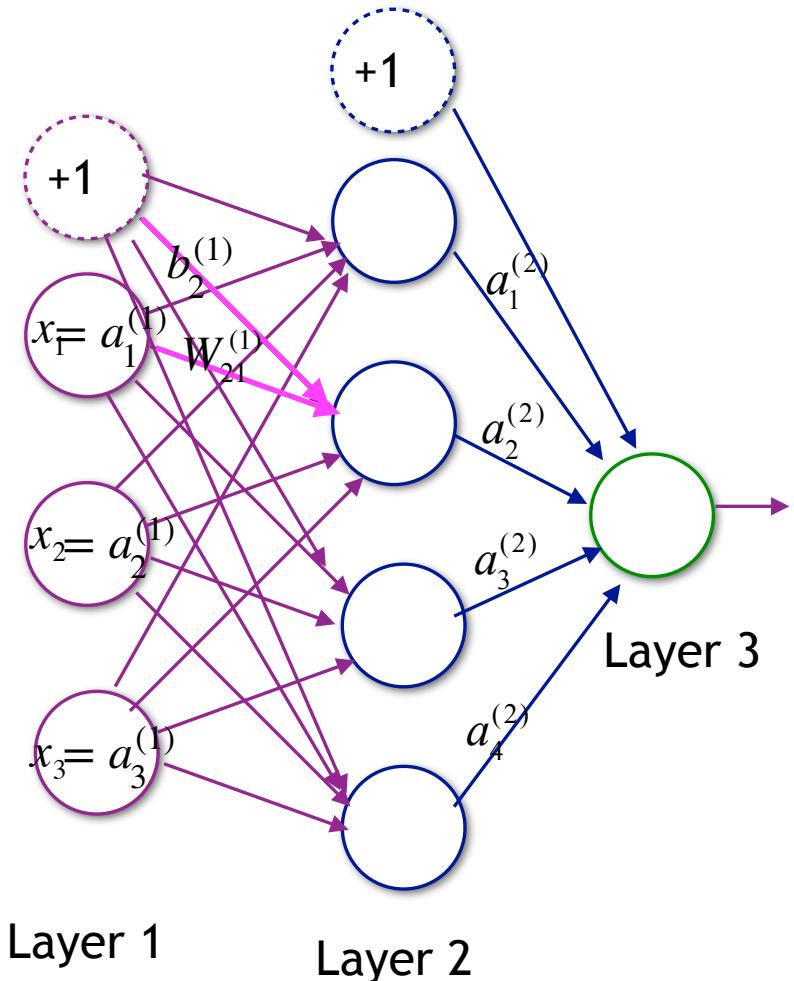
4 hidden

ork

? will
her in a
t for
ode i at

Neural Network Model & Notation

input layer
from input units



Let s_j be the number of nodes at layer j where we don't include the bias unit

$$W_{ij}^{(\ell)} \quad \left\{ \begin{array}{ll} 1 \leq \ell < n_l & \text{layers} \\ 1 \leq j \leq s_\ell & \text{inputs} \\ 1 \leq i \leq s_{\ell+1} & \text{outputs} \end{array} \right.$$

- ❑ The network has 3 input units (not 4), 4 hidden units (not 5) and 1 output unit
- ❑ n_ℓ is the number of layers in the network
- ❑ Each “neuron” has its own weights. We will store all the weights for one level together in a matrix $W^{(1)}$ such that $W_{ij}^{(1)}$ is the weight for leaving node j at level 1 and entering node i at level 2
- ❑ $b_i^{(1)}$ is the bias of unit i in level 2

$$W^{(1)} \in \mathbb{R}^{4 \times 3}$$

$$W^{(2)} \in \mathbb{R}^{1 \times 4}$$

Neural Network Model & Notation

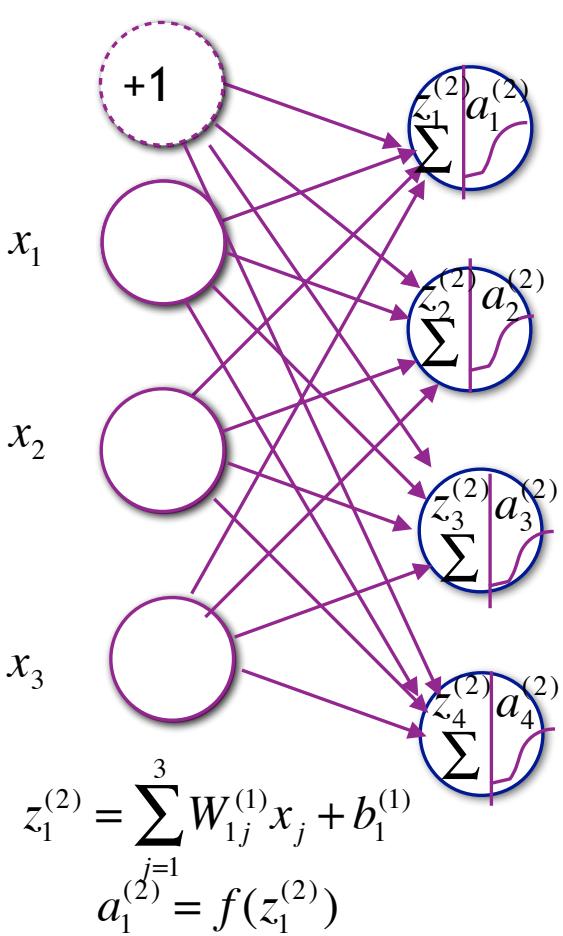
$$z_i^{(2)} = \sum_{j=1}^3 W_{ij}^{(1)} x_j + b_i^{(1)}$$

$$a_i^{(2)} = f(z_i^{(2)})$$

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

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

$$f(z) = \frac{1}{1 + \exp(-z)}$$



$$\begin{bmatrix} W_{11}^{(1)} & W_{12}^{(1)} & W_{13}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} & W_{23}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} & W_{33}^{(1)} \\ W_{41}^{(1)} & W_{42}^{(1)} & W_{43}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix} = \begin{bmatrix} W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 \\ W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 \\ W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 \\ W_{41}^{(1)}x_1 + W_{42}^{(1)}x_2 + W_{43}^{(1)}x_3 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}$$

$$= \begin{bmatrix} W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)} \\ W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)} \\ W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)} \\ W_{41}^{(1)}x_1 + W_{42}^{(1)}x_2 + W_{43}^{(1)}x_3 + b_4^{(1)} \end{bmatrix} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \\ z_4^{(2)} \end{bmatrix}$$

$$a^{(2)} = \begin{bmatrix} f(z_1^{(2)}) \\ f(z_2^{(2)}) \\ f(z_3^{(2)}) \\ f(z_4^{(2)}) \end{bmatrix} = \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \\ a_4^{(2)} \end{bmatrix}$$

$$z_1^{(2)} = \sum_{j=1}^3 W_{1j}^{(1)} x_j + b_1^{(1)}$$

$$a_1^{(2)} = f(z_1^{(2)})$$

$$z_2^{(2)} = \sum_{j=1}^3 W_{2j}^{(1)} x_j + b_2^{(1)}$$

$$a_2^{(2)} = f(z_2^{(2)})$$

$$z_3^{(2)} = \sum_{j=1}^3 W_{3j}^{(1)} x_j + b_3^{(1)}$$

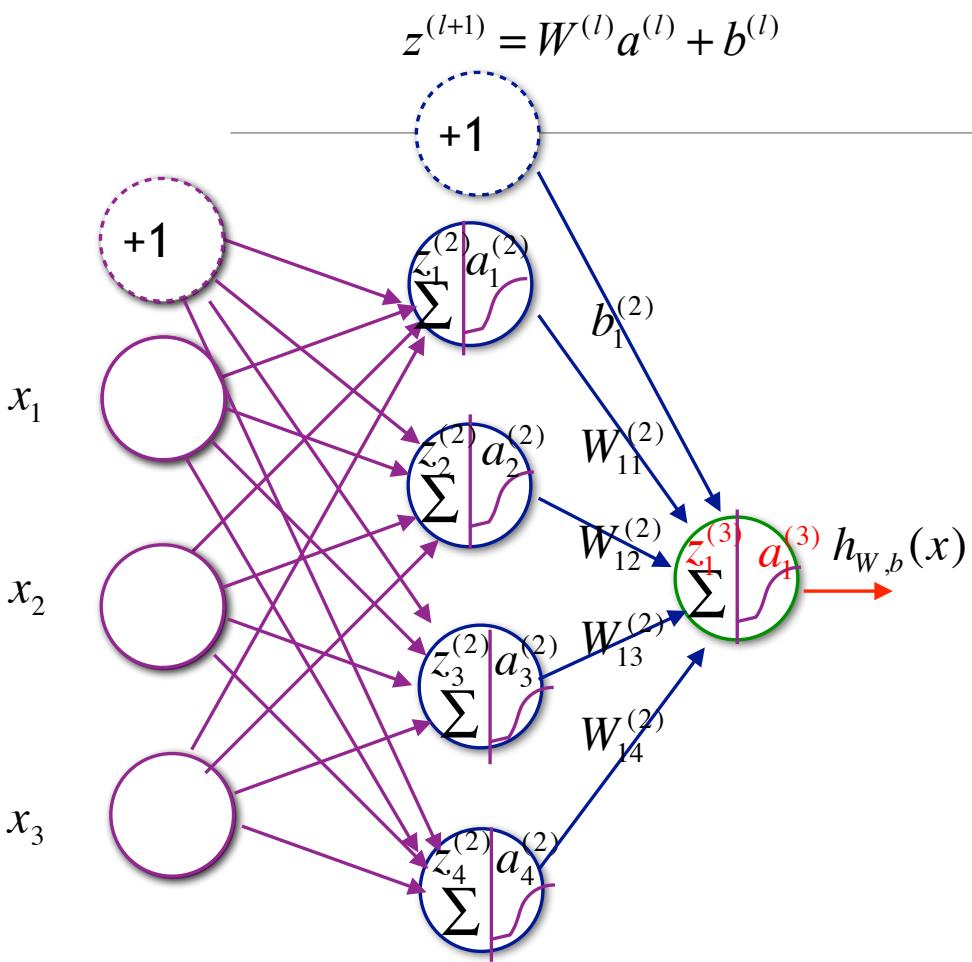
$$a_3^{(2)} = f(z_3^{(2)})$$

$$z_4^{(2)} = \sum_{j=1}^3 W_{4j}^{(1)} x_j + b_4^{(1)}$$

$$a_4^{(2)} = f(z_4^{(2)})$$

Neural Network Model & Notation

$$\Sigma \quad f(z) = \frac{1}{1 + \exp(-z)}$$



$$\begin{bmatrix} W_{11}^{(2)} & W_{12}^{(2)} & W_{13}^{(2)} & W_{14}^{(2)} \end{bmatrix} \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \\ a_4^{(2)} \end{bmatrix} + \begin{bmatrix} b_1^{(2)} \end{bmatrix} = \begin{bmatrix} z_1^{(3)} \end{bmatrix}$$

$$\mathbf{a}^{(3)} = \begin{bmatrix} f(z_1^{(2)}) \end{bmatrix} = \begin{bmatrix} a_1^{(3)} \end{bmatrix}$$

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

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

$$z_1^{(3)} = \sum_{j=1}^4 W_{1j}^{(2)} a_j^{(2)} + b_1^{(2)}$$

$$a_1^{(3)} = f(z_1^{(3)})$$

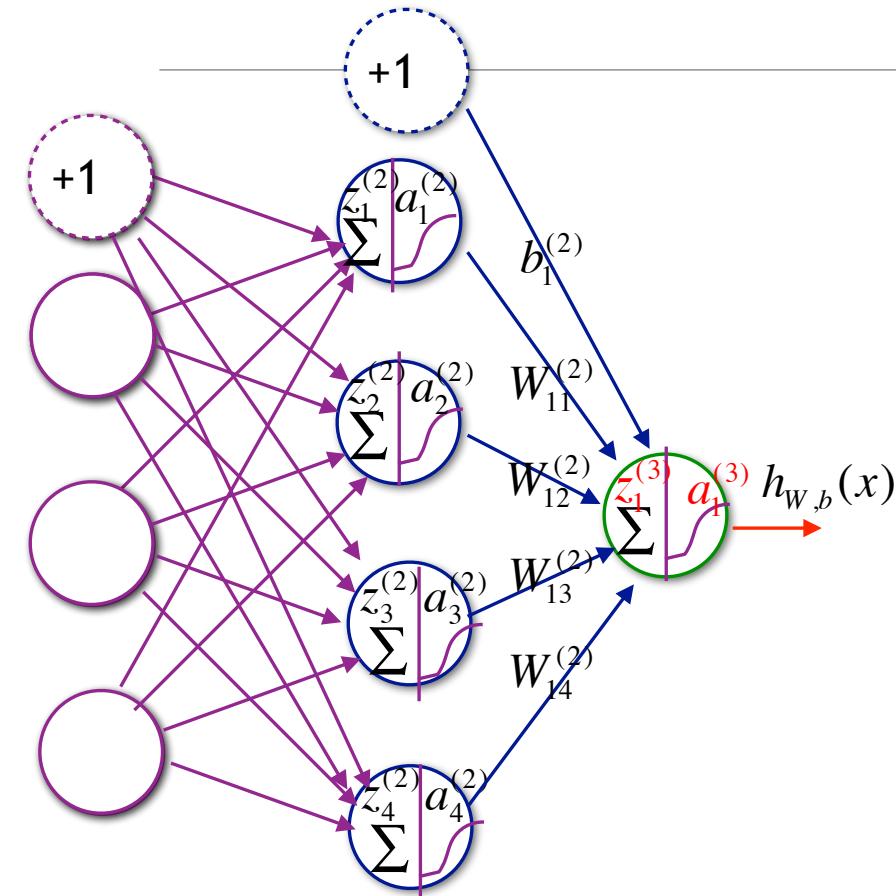
$$h_{W,b} = \mathbf{a}^{(3)} = f(\mathbf{z}^{(3)})$$

Neural Network Model & Notation

$$\Sigma \quad f(z) = \frac{1}{1 + \exp(-z)}$$

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

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



What is $a_3^{(2)}$?

- a) One of the inputs to neurons on layer 2
- b) One of the inputs to neurons on layer 3
- c) activation value of 3rd neuron on layer 2
- d) activation value of 2nd neuron on layer 3
- e) activation value of 3rd neuron on layer 3
- f) activation value of 2nd neuron on layer 3

$$+[b_1^{(2)}] = [z_1^{(3)}]$$

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

Artificial Neural Networks

- ❑ Each node (we call the nodes *neurons*) represents a linear function of its input, similar to logistic regression

$$z_1^{(2)} = W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)}$$

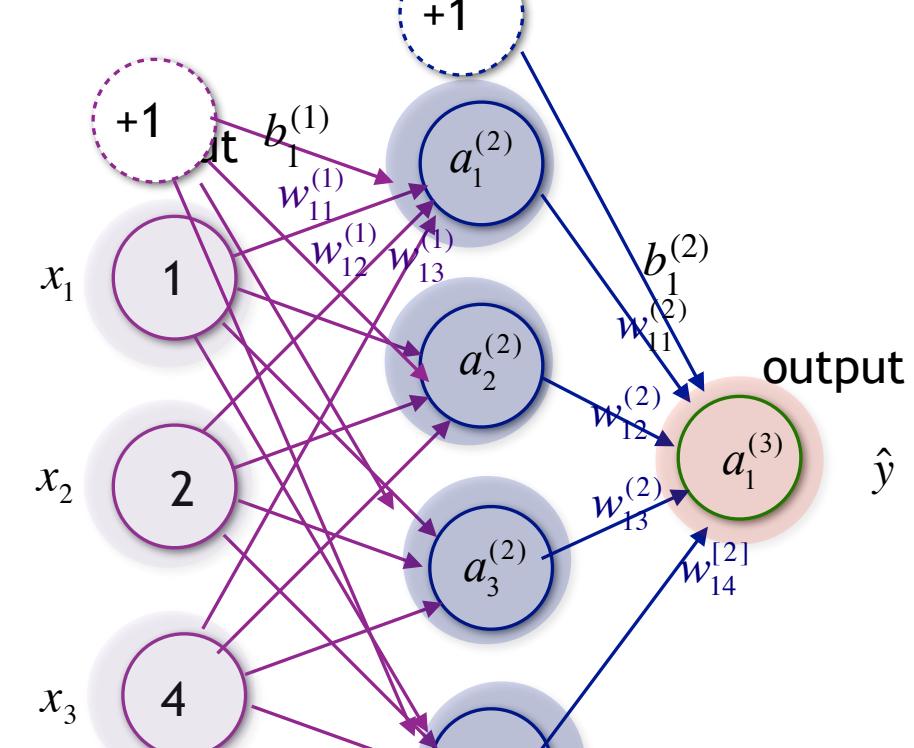
$$a_1^{(2)} = f(z_1^{(2)})$$

- $W_{11}^{(1)} W_{12}^{(1)} W_{13}^{(1)} b_1^{(1)}$ are the **weights** (aka parameters)
- $W^{(1)}$ is the matrix of weights that map values from input layer 1 to the first hidden layer (layer 2) ($W^{(\ell)}$ is the matrix that maps layer ℓ to layer $\ell+1$)
- $b^{(1)}$ is the **vector of bias values** for the neurons on layer 1

- ❑ $f^{(1)}$ is the “activation function” for layer 1. ($f^{(\ell)}$ is the activation function for layer ℓ)

- ❑ $a_j^{(2)}$ is the activation value for the j^{th} neuron of layer 2
 $a_j^{(\ell)}$ is the activation value for the j^{th} neuron of layer the ℓ

Let s_j is the number of nodes at layer j where we don't include the bias unit



$$W^{(1)} = \begin{bmatrix} W_{11}^{(1)} & W_{12}^{(1)} & W_{13}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} & W_{23}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} & W_{33}^{(1)} \\ W_{41}^{(1)} & W_{42}^{(1)} & W_{43}^{(1)} \end{bmatrix} \quad b^{(1)} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}$$

$$W^{(2)} = \begin{bmatrix} W_{11}^{(2)} & W_{12}^{(2)} & W_{13}^{(2)} & W_{14}^{(2)} \end{bmatrix} \quad b^{(2)} = \begin{bmatrix} b_1^{(2)} \end{bmatrix}$$

Artificial Neural Networks

- ❑ Each node (we call the nodes *neurons*) represents a linear function of its input, similar to logistic regression

$$z_1^{(2)} = W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)}$$

$$a_1^{(2)} = f(z_1^{(2)})$$

- $W_{11}^{(1)} W_{12}^{(1)} W_{13}^{(1)} b_1^{(1)}$ are the **weights** (aka parameters)
- $W^{(1)}$ is the matrix of weights that map values from input layer 1 to the first hidden layer (layer 2) ($W^{(\ell)}$ is the matrix that maps layer ℓ to layer $\ell+1$)
- $b^{(1)}$ is the **vector of bias values** for the neurons on layer 1

- ❑ $f^{(1)}$ is the “activation function” for layer 1. ($f^{(\ell)}$ is the activation function for layer ℓ)

- ❑ $a_j^{(2)}$ is the activation value for the j^{th} neuron of layer 2
 $a_j^{(\ell)}$ is the activation value for the j^{th} neuron of layer the ℓ

Let s_j is the number of nodes at layer j where we don't include the bias unit

$$W^{(2)} = \begin{bmatrix} W_{11}^{(2)} & W_{12}^{(2)} & W_{13}^{(2)} & W_{14}^{(2)} \end{bmatrix} \quad b^{(2)} = \begin{bmatrix} b_1^{(2)} \end{bmatrix}$$
51

