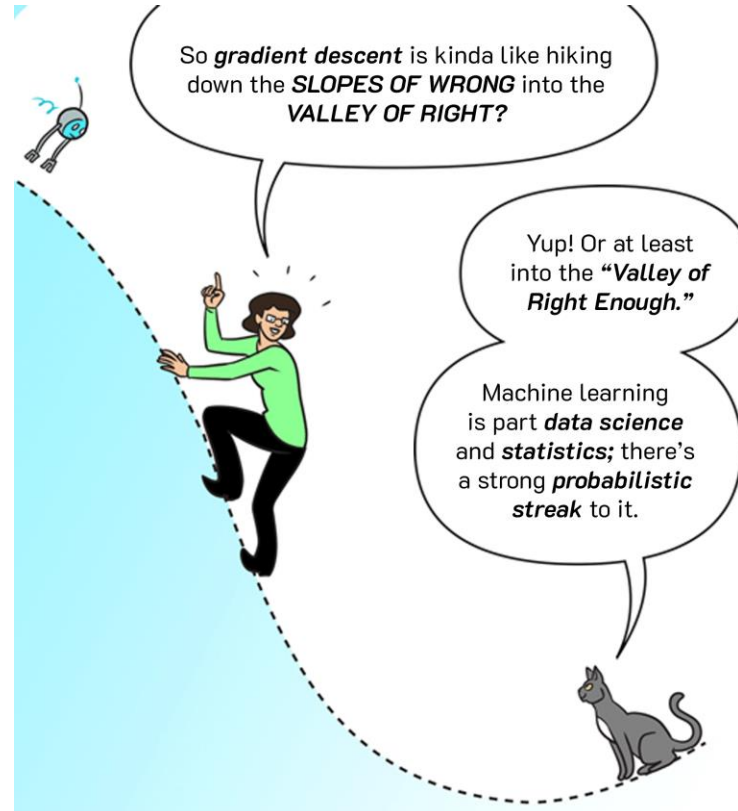




Neural Networks



Perceptron learning algorithm

repeat until convergence (or for some # of iterations):

for each training example (f_1, f_2, \dots, f_n , label):

$$\text{prediction} = b + \sum_{i=1}^n w_i f_i$$

if $\text{prediction} * \text{label} \leq 0$: // they don't agree

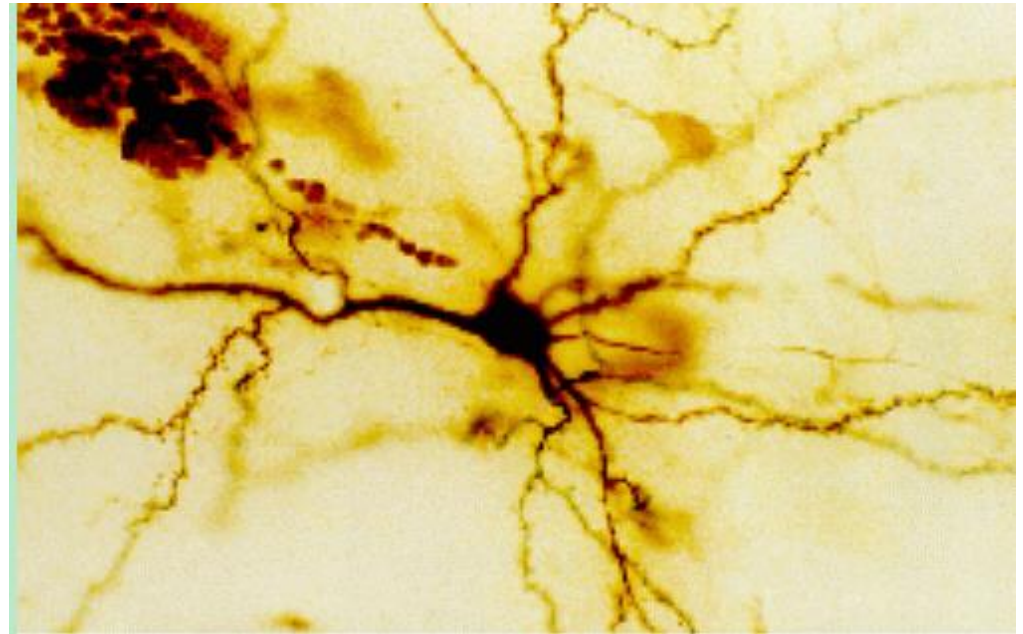
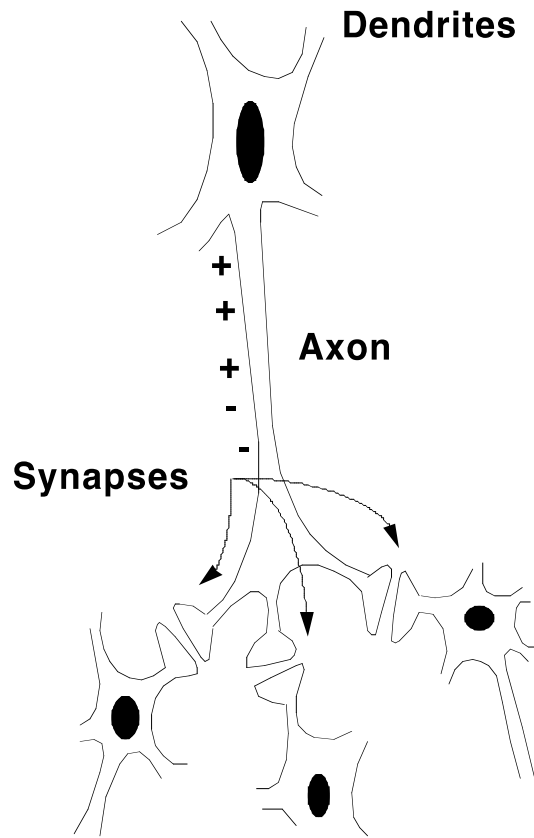
for each w_i :

$$w_i = w_i + f_i * \text{label}$$

$$b = b + \text{label}$$

Why is it called the “perceptron” learning algorithm if what it learns is a line? Why not “line learning” algorithm?

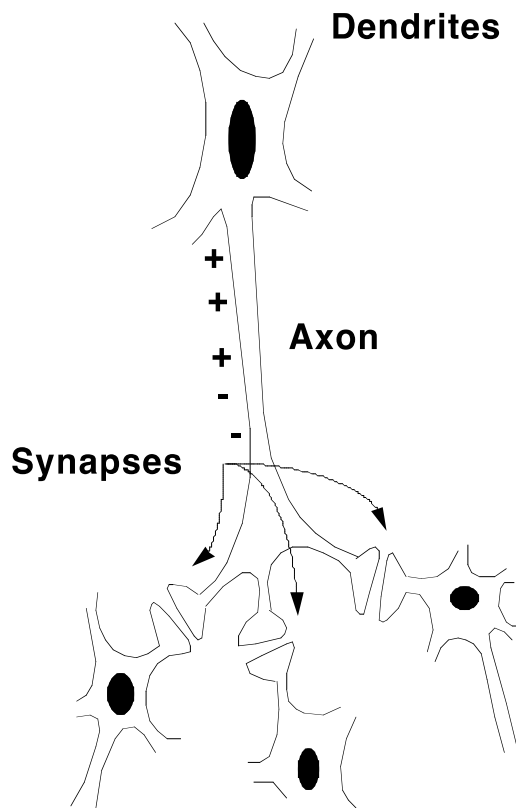
Our Nervous System



Neuron

What do you know?

Our nervous system: the computer science view

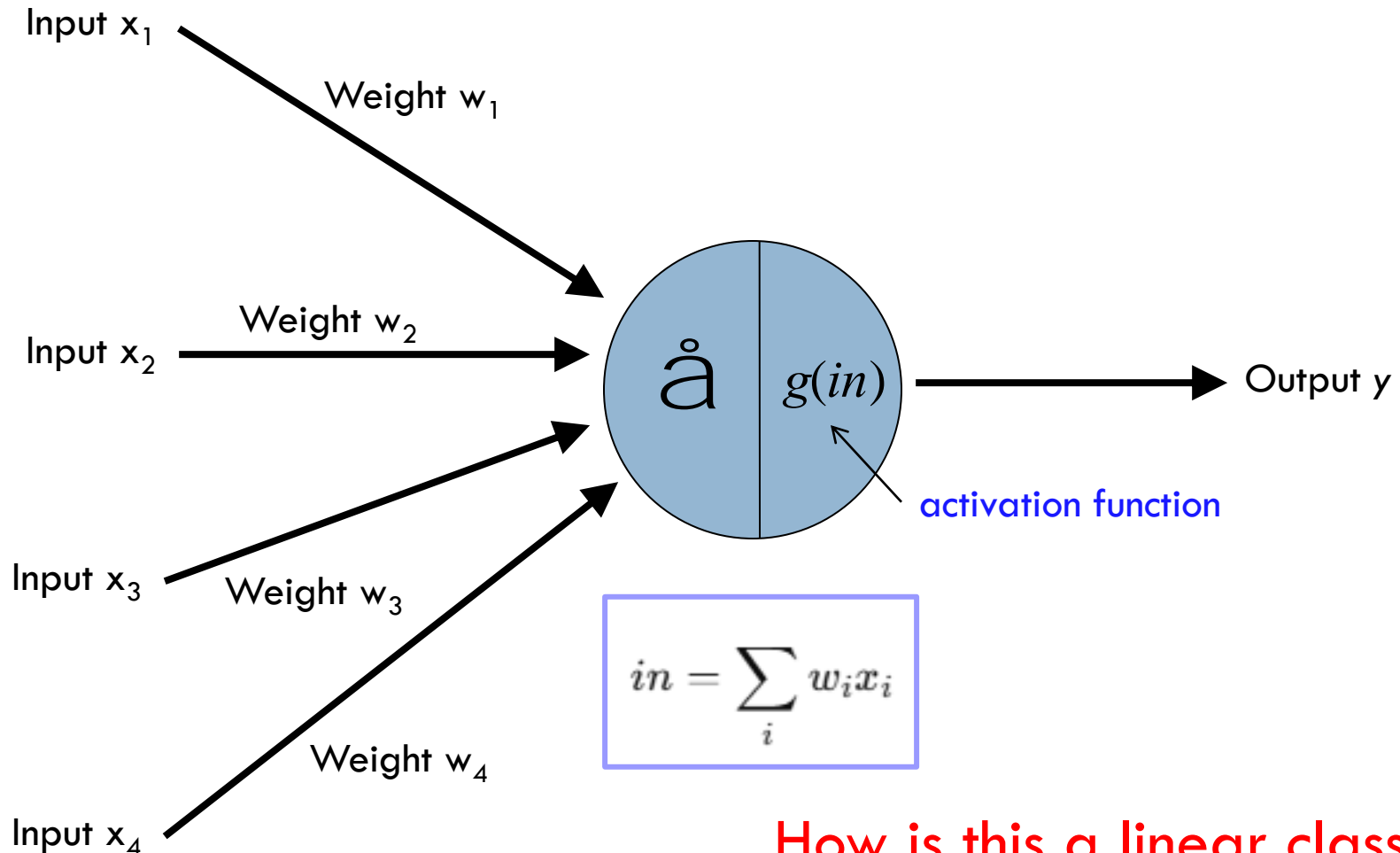


the human brain is a large collection of interconnected neurons

a **NEURON** is a brain cell

- ▣ they collect, process, and disseminate electrical signals
- ▣ they are connected via synapses
- ▣ they **FIRE** depending on the conditions of the neighboring neurons

A neuron/perceptron

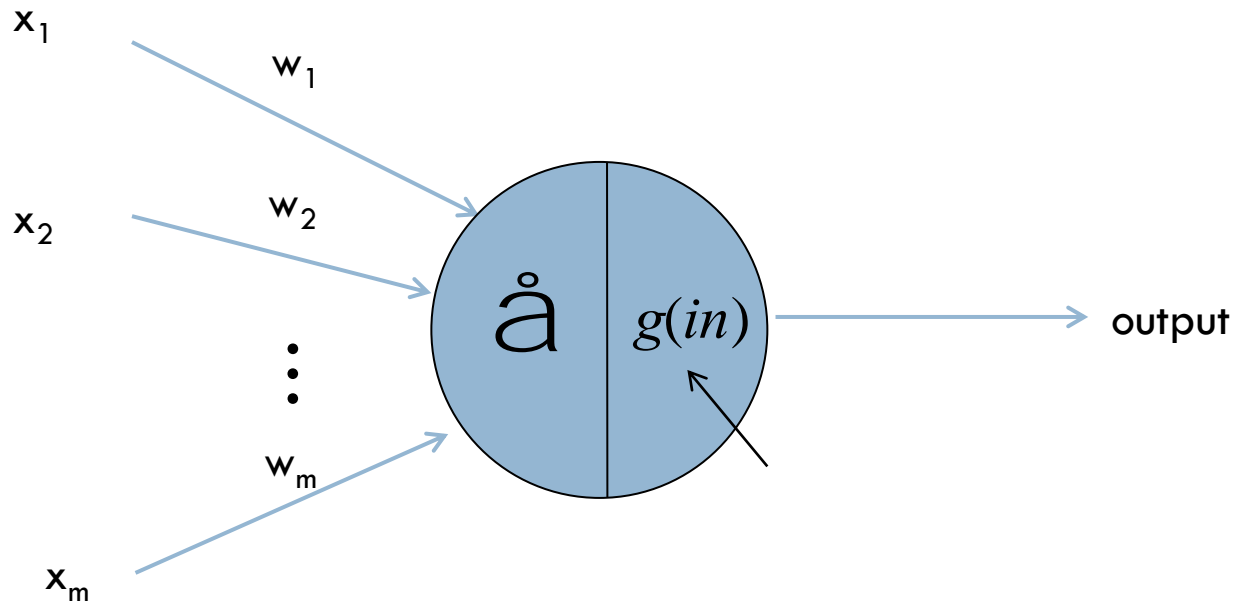


How is this a linear classifier
(i.e. perceptron)?

Hard threshold = linear classifier

hard threshold:

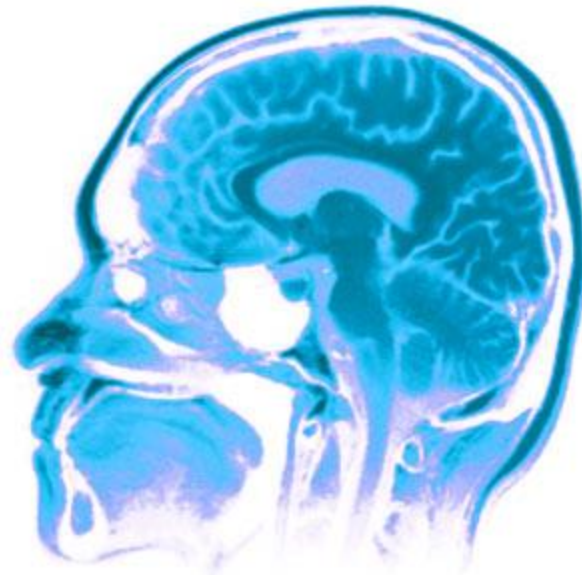
$$g(in) = \begin{cases} 1 & \text{if } in > -b \\ 0 & \text{otherwise} \end{cases} \quad output = \begin{cases} 1 & \text{if } \sum_i w_i x_i + b > 0 \\ 0 & \text{otherwise} \end{cases}$$



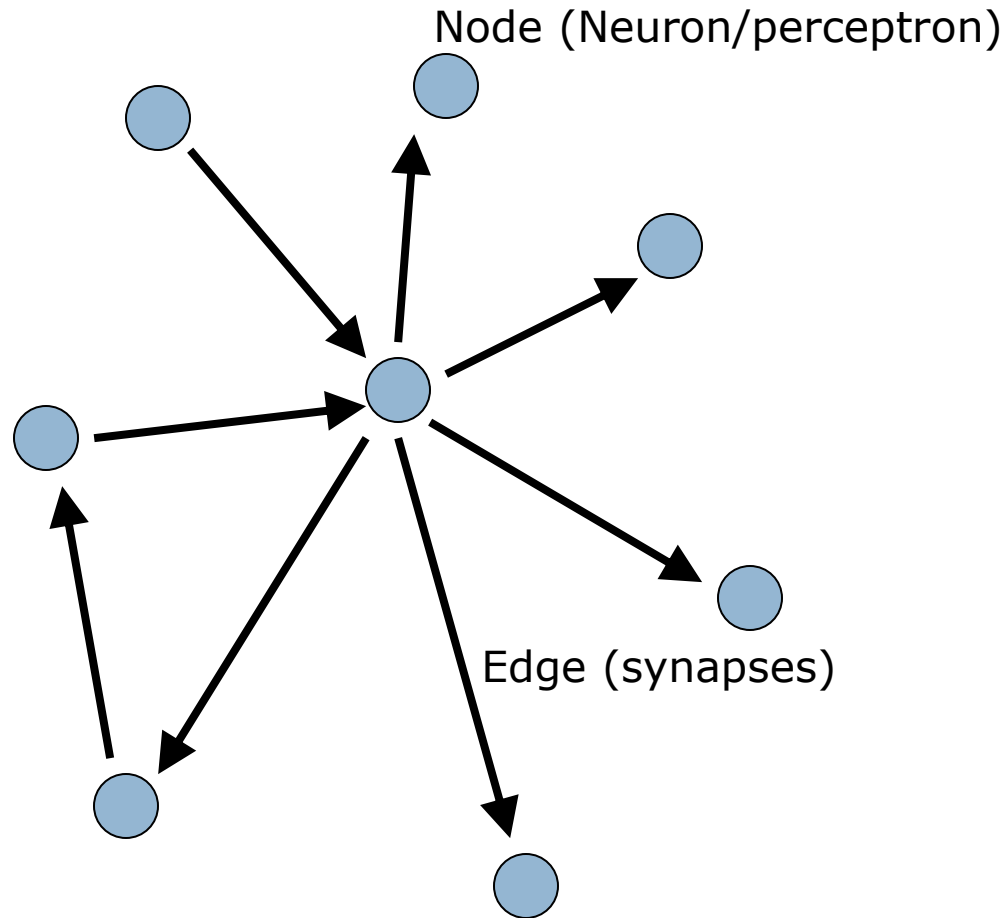
Neural Networks

Neural Networks try to mimic the structure and function of our nervous system

People like biologically motivated approaches



Artificial Neural Networks

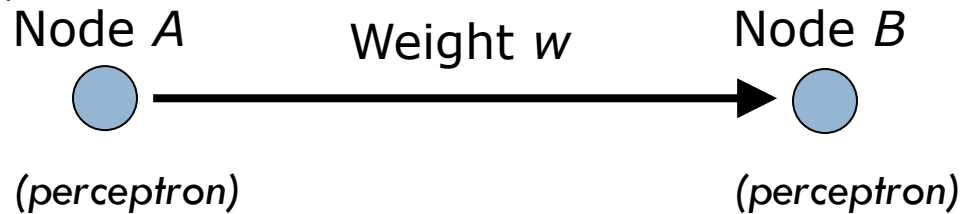


our approximation

When computing the input to neuron B:

$$in_B = \sum_i w_i x_i$$

- If $x_i = 1$ (A fires), and:
 - $w_i > 0$: this **adds** to $in_B \rightarrow$ B is **more likely** to fire
 - $w_i < 0$: this **subtracts** from $in_B \rightarrow$ B is **less likely** to fire



- Neuron A send a signal to Neuron B.
- w = strength (and type) of the connection from A to B.

Then:

- If **A fires** (i.e., outputs a 1 or sends a spike), and:
 - $w > 0 \rightarrow$ A **stimulates** B (excitation)
 - $w < 0 \rightarrow$ A **inhibits** B (suppression)
 - $w = 0 \rightarrow$ A has **no effect** on B

$$output = \begin{cases} 1 & \text{if } \sum_i w_i x_i + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

W is the strength of signal sent between A and B.

If A fires and w is **positive**, then A **stimulates** B.

If A *fires* and w is **negative**, then A **inhibits** B.

Other activation functions

hard threshold:

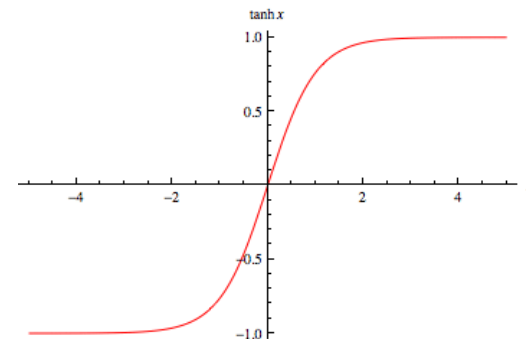
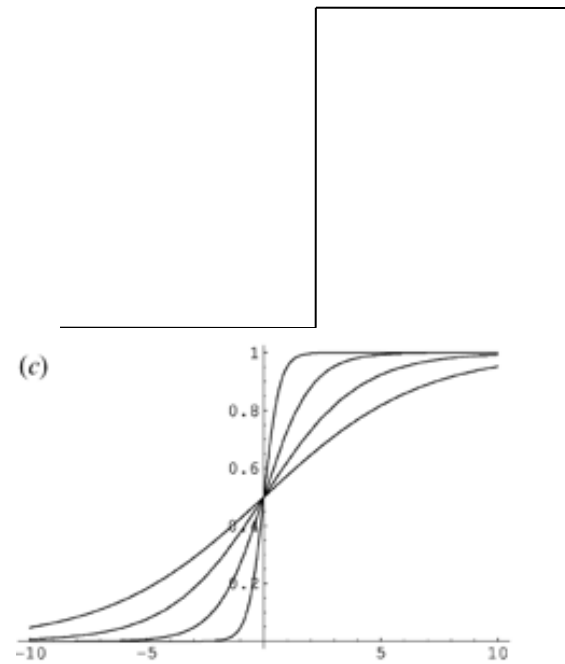
$$g(in) = \begin{cases} 1 & \text{if } in > -b \\ 0 & \text{otherwise} \end{cases}$$

sigmoid

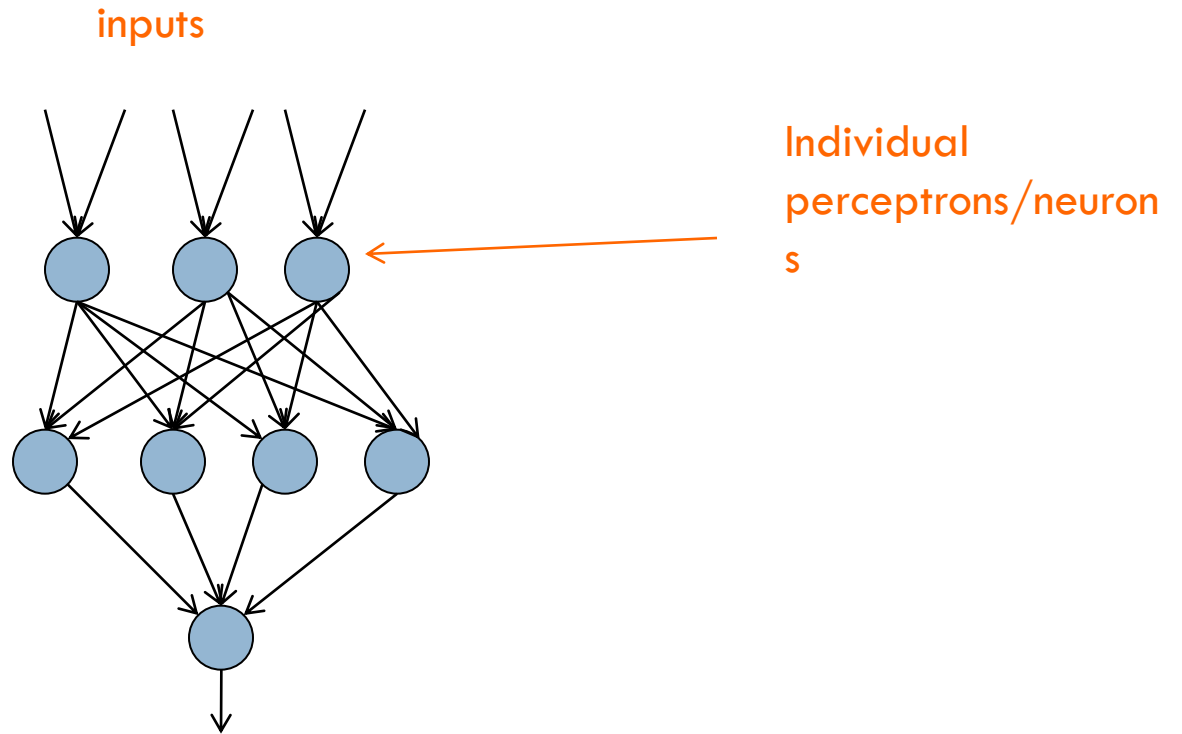
$$g(x) = \frac{1}{1 + e^{-ax}}$$

$\tanh x$

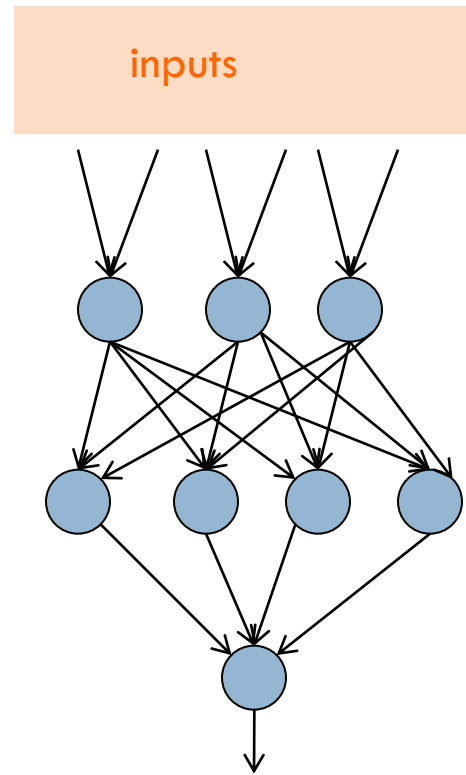
why other threshold functions?



Neural network



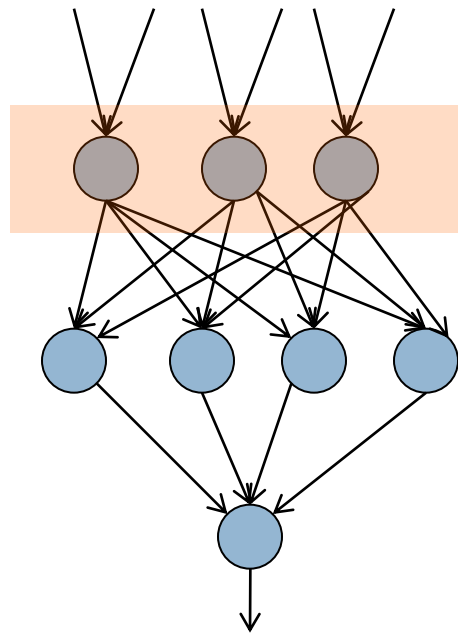
Neural network



some inputs are
provided/entered

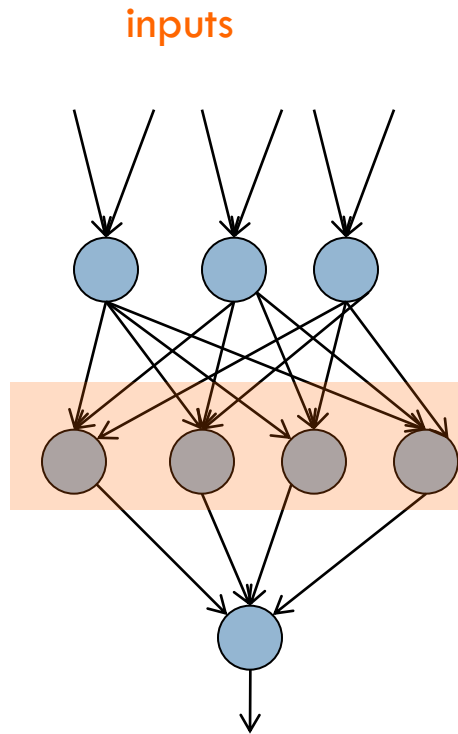
Neural network

inputs



each perceptron computes and calculates an answer

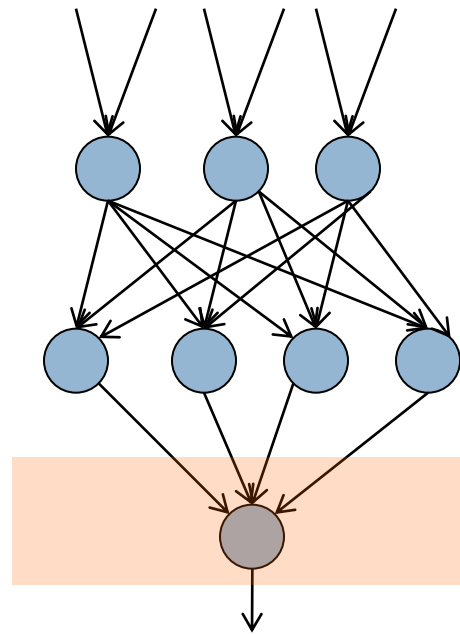
Neural network



those answers become inputs
for the next level

Neural network

inputs



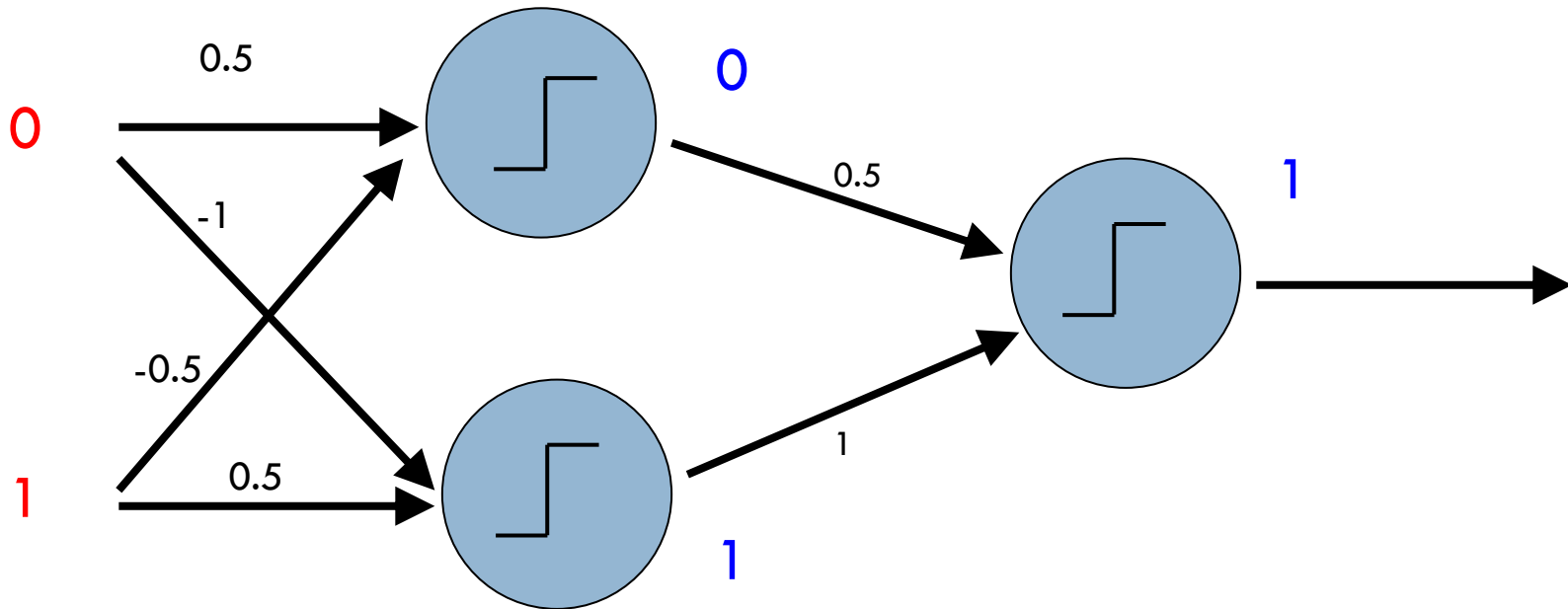
finally get the answer after all levels compute

Activation spread



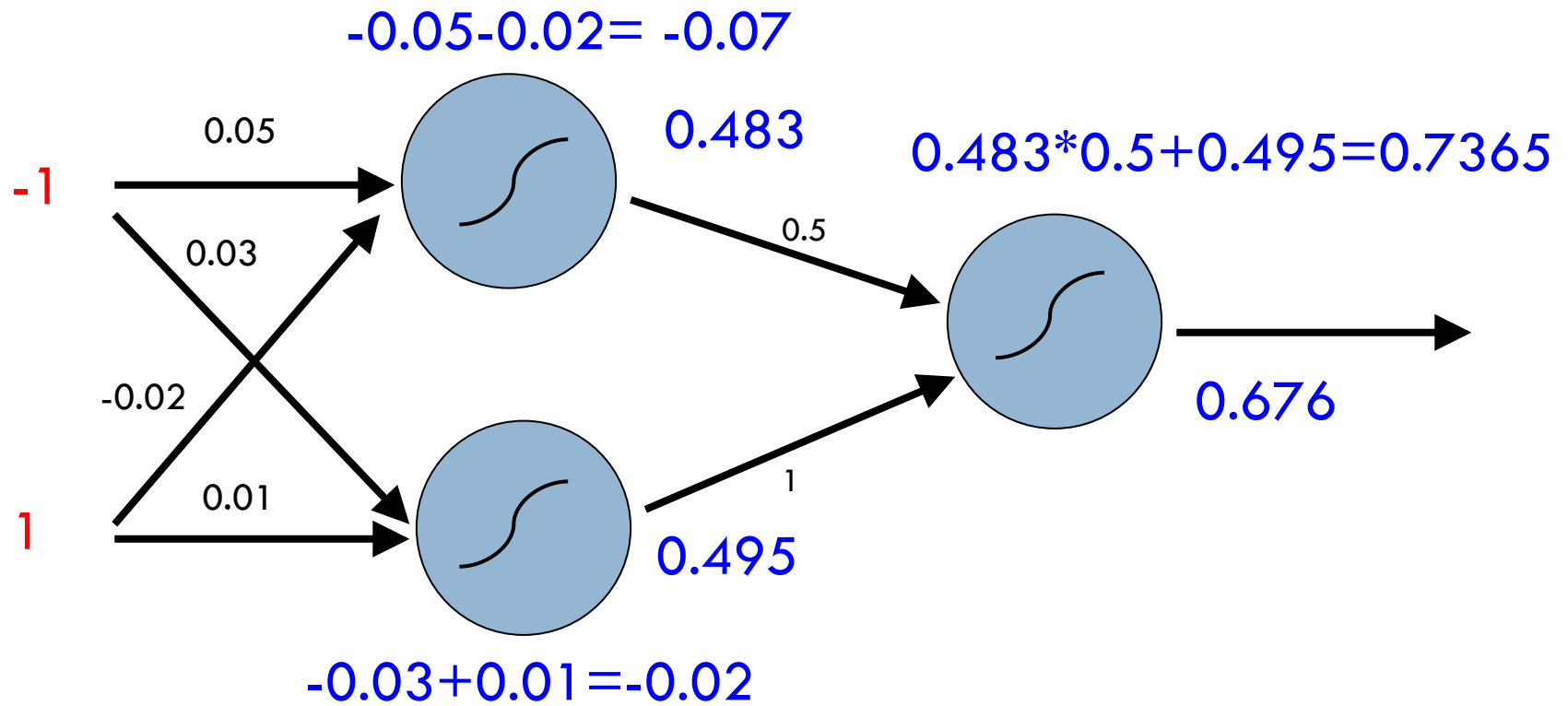
<http://www.youtube.com/watch?v=Yq7d4ROvZ6I>

Computation (assume 0 bias)



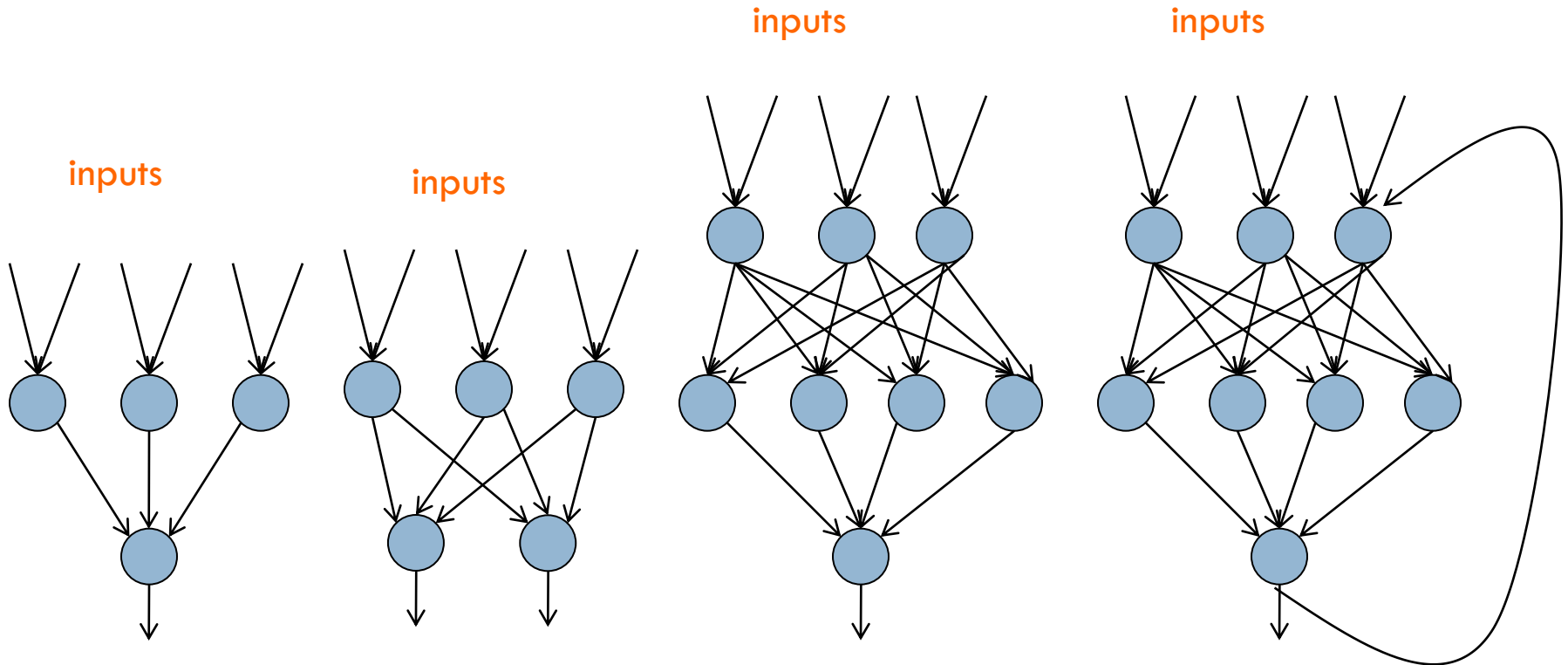
$$g(in) = \begin{cases} 1 & \text{if } in > -b \\ 0 & \text{otherwise} \end{cases}$$

Computation



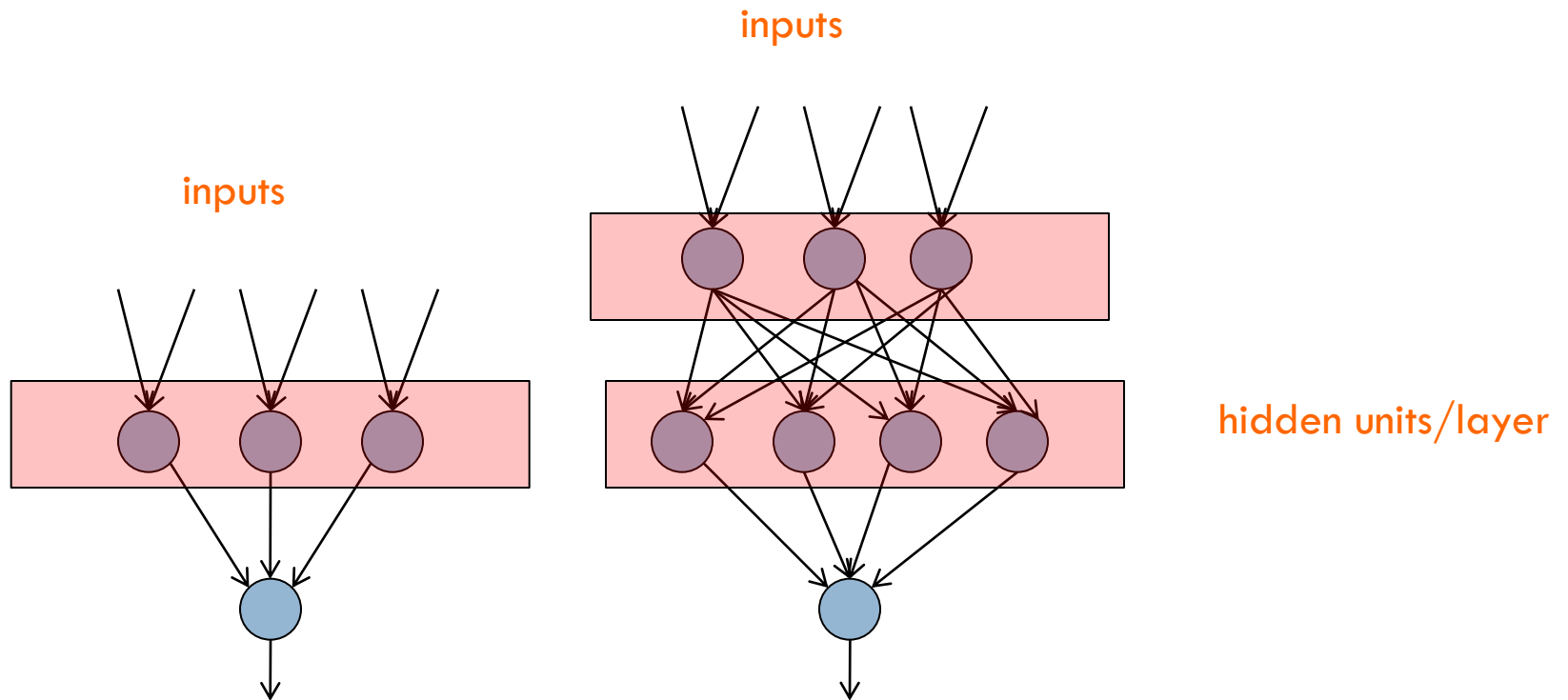
Neural networks

Different kinds/characteristics of networks



How are these different?

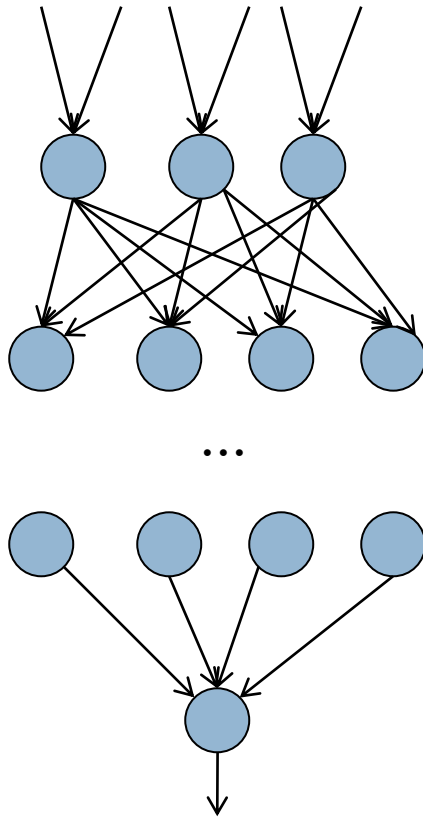
Hidden units/layers



Feed forward networks

Hidden units/layers

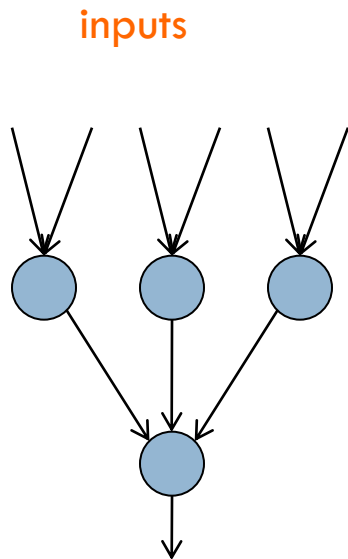
inputs



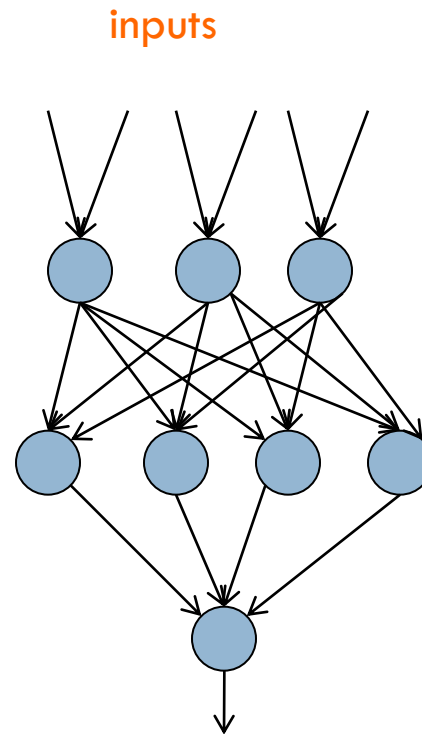
Can have many layers of hidden units of differing sizes

To count the number of layers, you count all but the inputs

Hidden units/layers



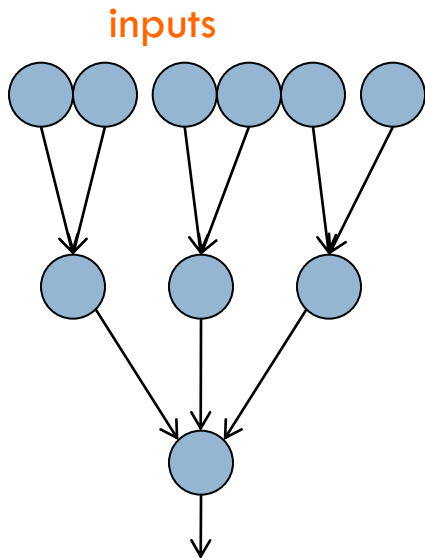
2-layer network



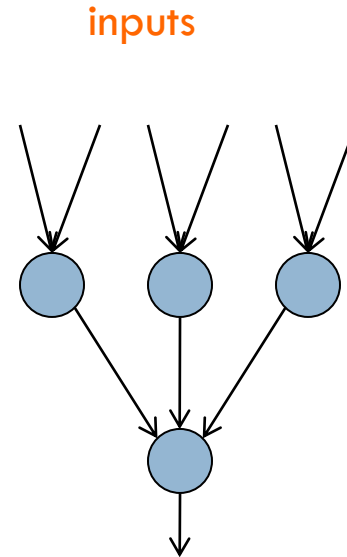
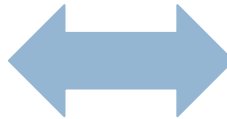
3-layer network

Alternate ways of visualizing

Sometimes the input layer will be drawn with nodes as well

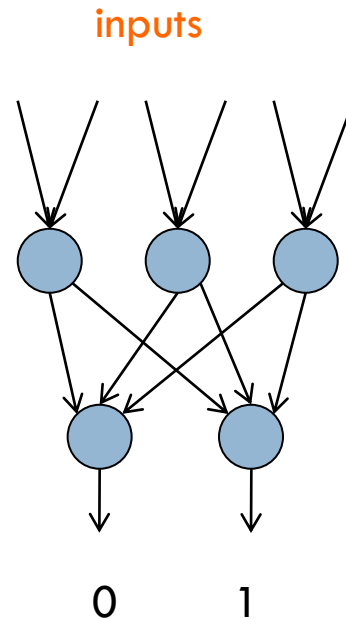


2-layer network



2-layer network

Multiple outputs



Can be used to model multiclass datasets or more interesting predictors, e.g. images

Multiple outputs



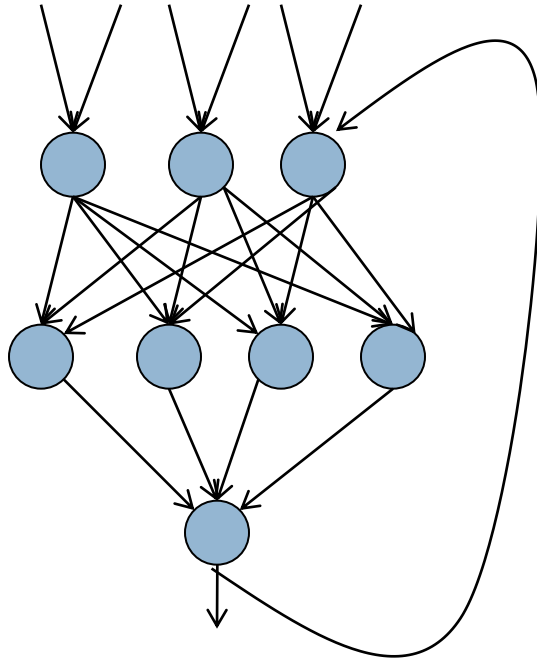
input



output
(edge detection)

Neural networks

inputs



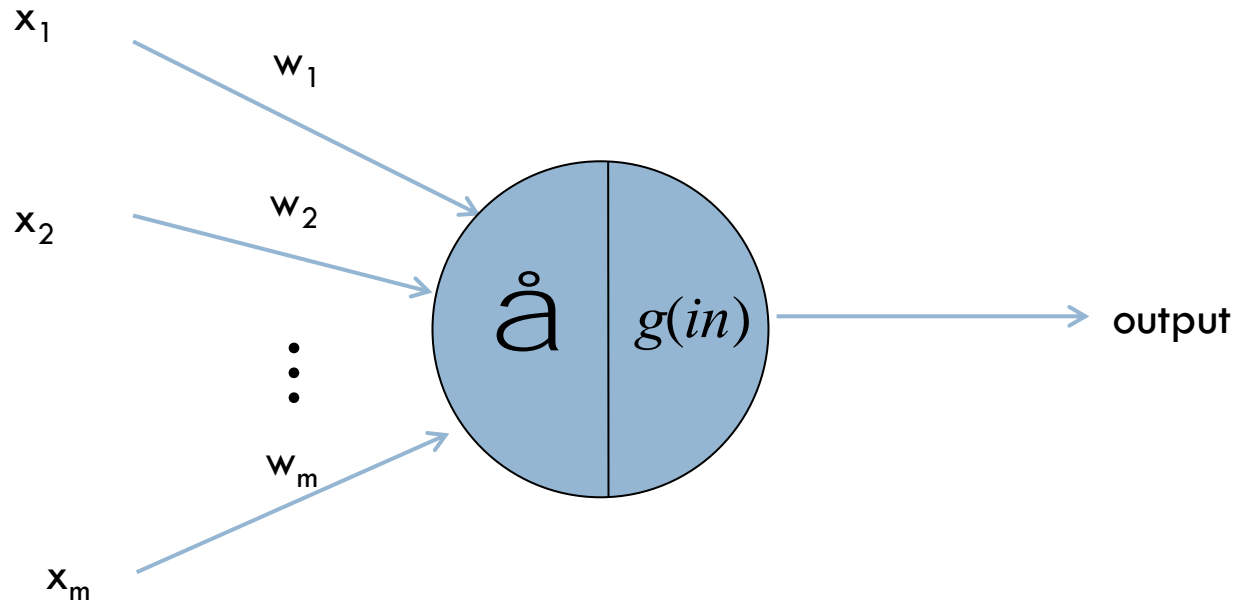
Recurrent network

Output is fed back to input

Can support memory!

Good for temporal/sequential data

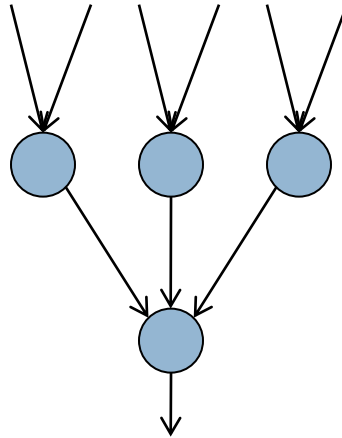
NN decision boundary



What does the decision boundary of a perceptron look like?

Line (linear set of weights)

NN decision boundary

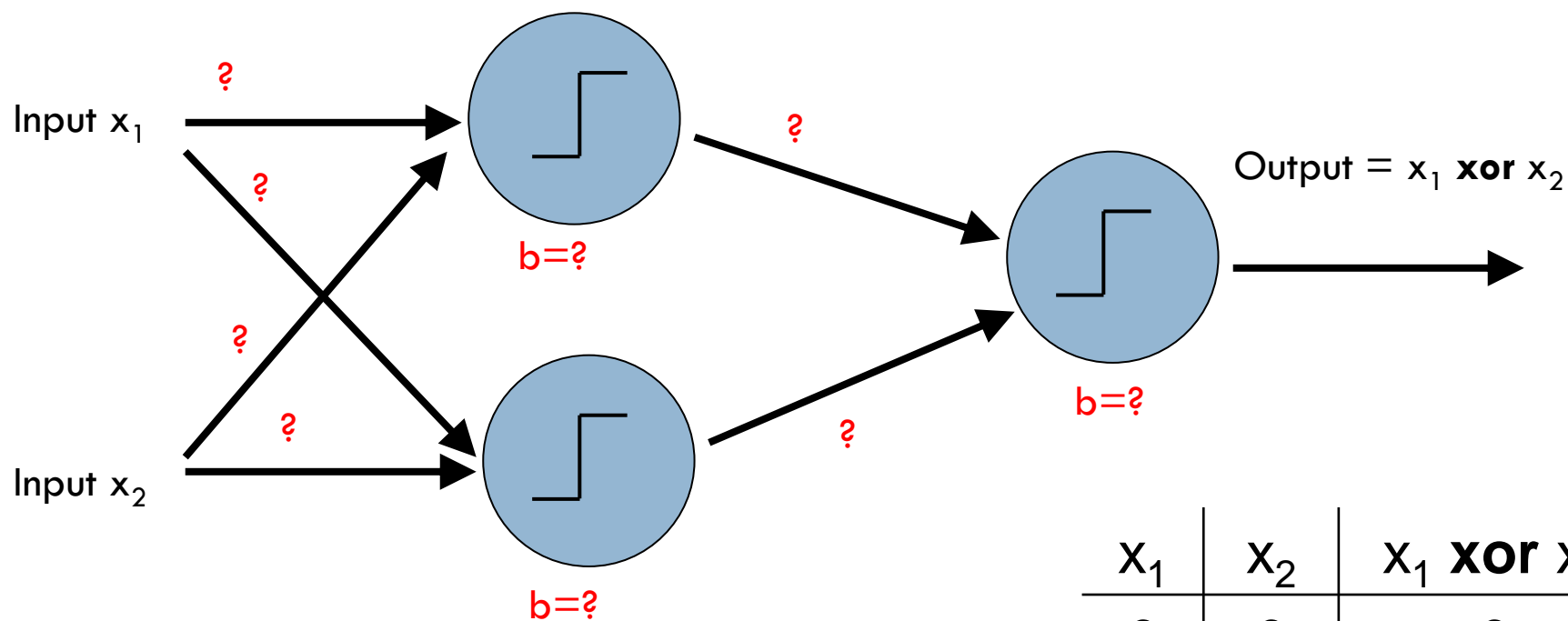


What does the decision boundary of a 2-layer network look like?

Is it linear?

What types of things can and can't it model?

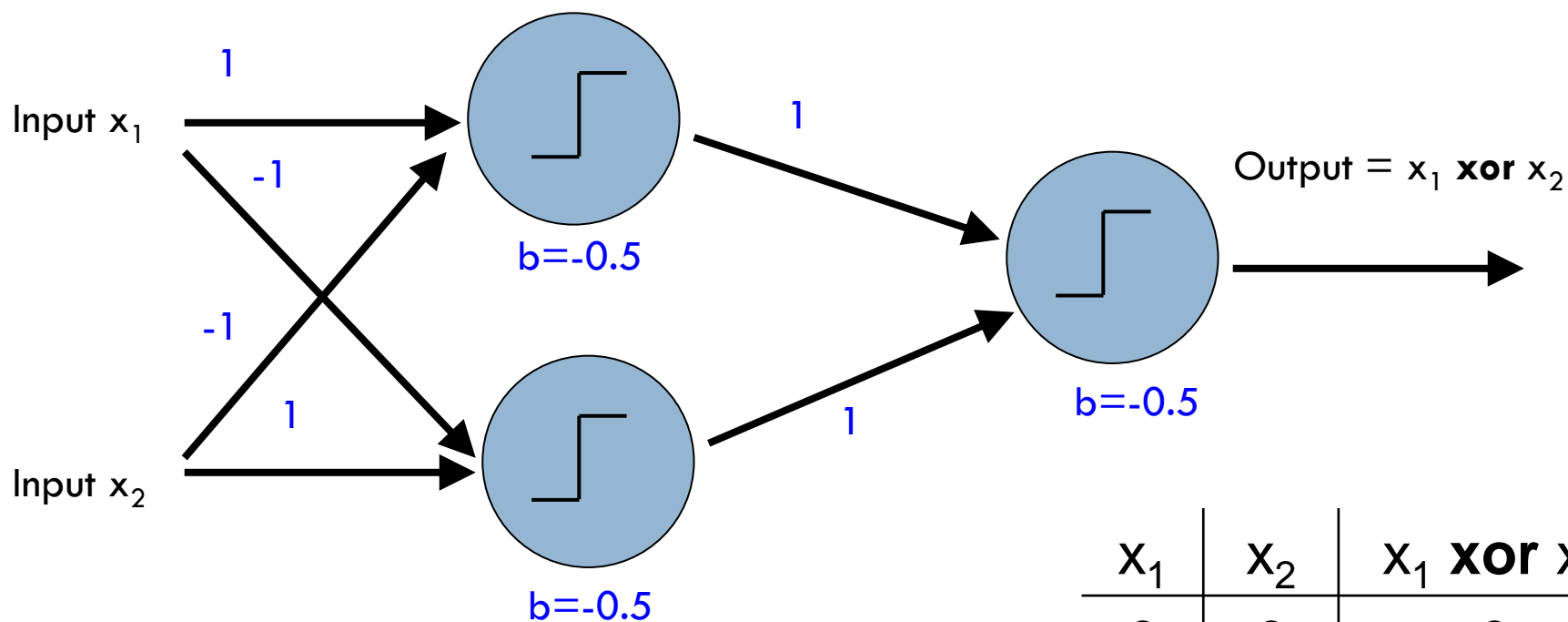
XOR



$$output = \begin{cases} 1 & \text{if } \sum_i w_i x_i + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

x_1	x_2	x_1 xor x_2
0	0	0
0	1	1
1	0	1
1	1	0

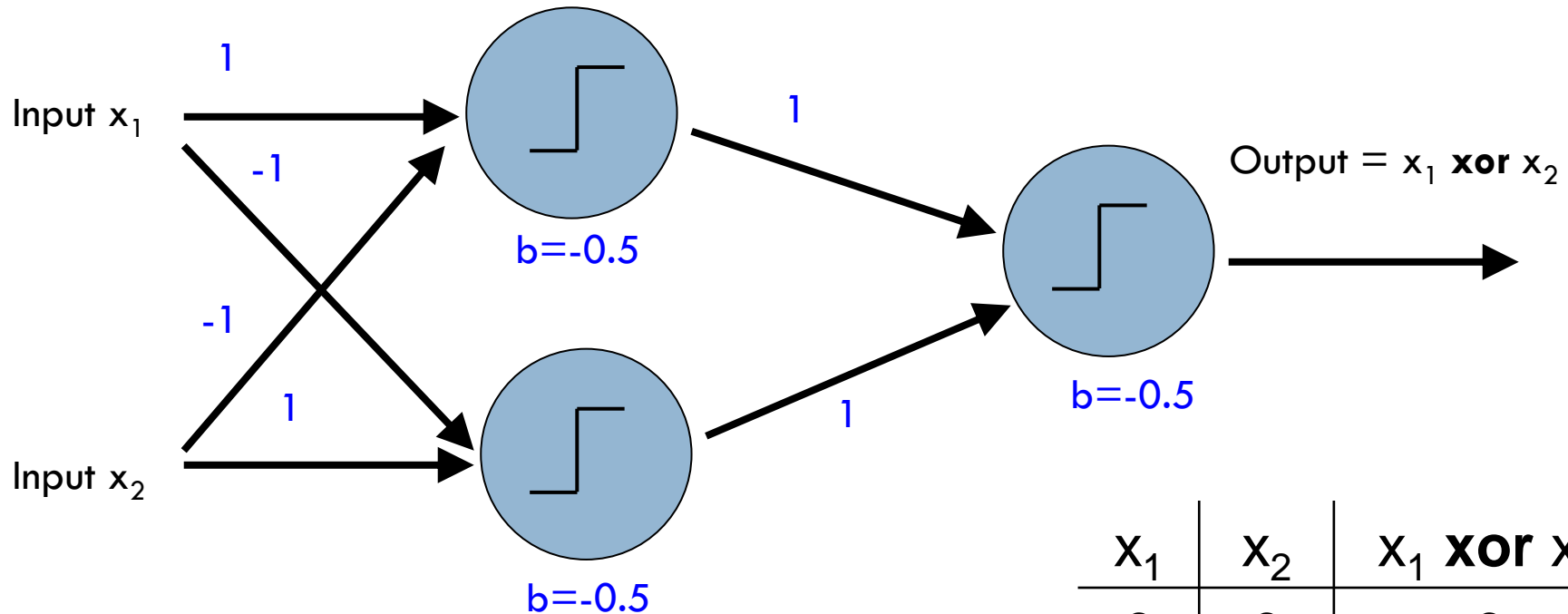
XOR



$$output = \begin{cases} 1 & \text{if } \sum_i w_i x_i + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

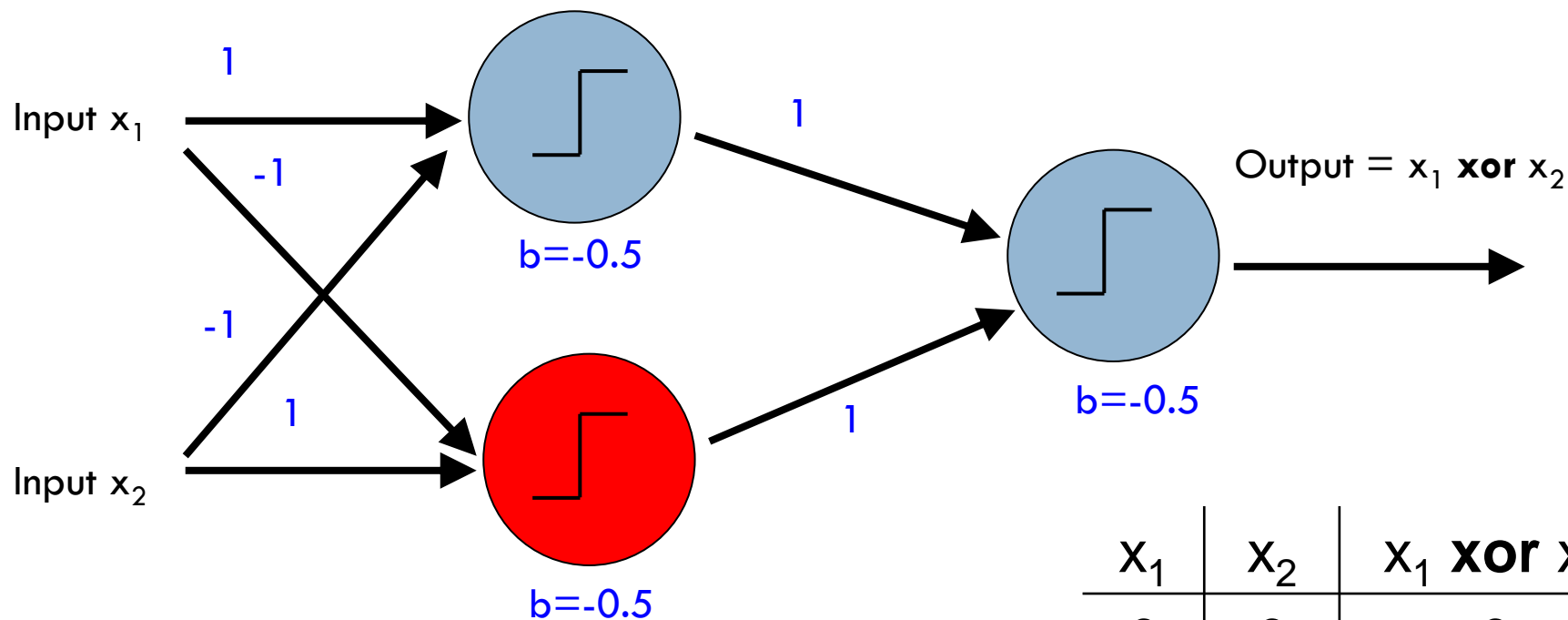
x_1	x_2	x_1 xor x_2
0	0	0
0	1	1
1	0	1
1	1	0

What does the decision boundary look like?



x_1	x_2	x_1 xor x_2
0	0	0
0	1	1
1	0	1
1	1	0

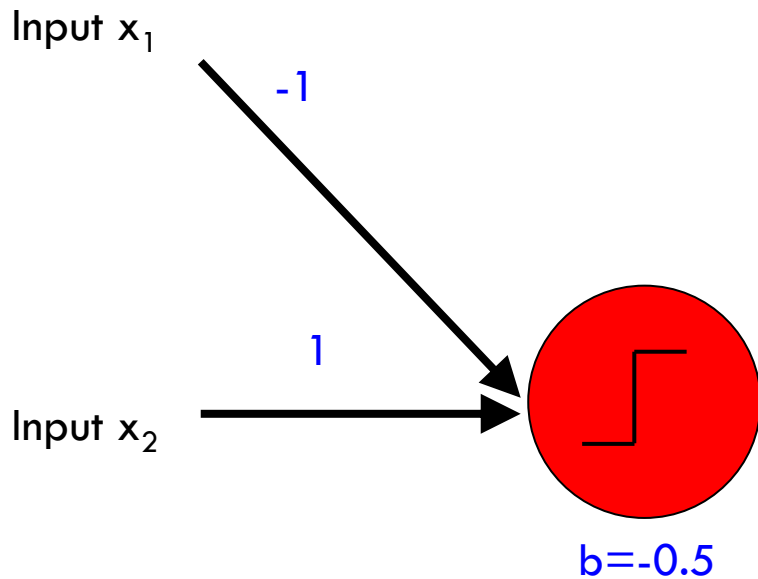
What does the decision boundary look like?



What does this perceptron's decision boundary look like?

x_1	x_2	x_1 xor x_2
0	0	0
0	1	1
1	0	1
1	1	0

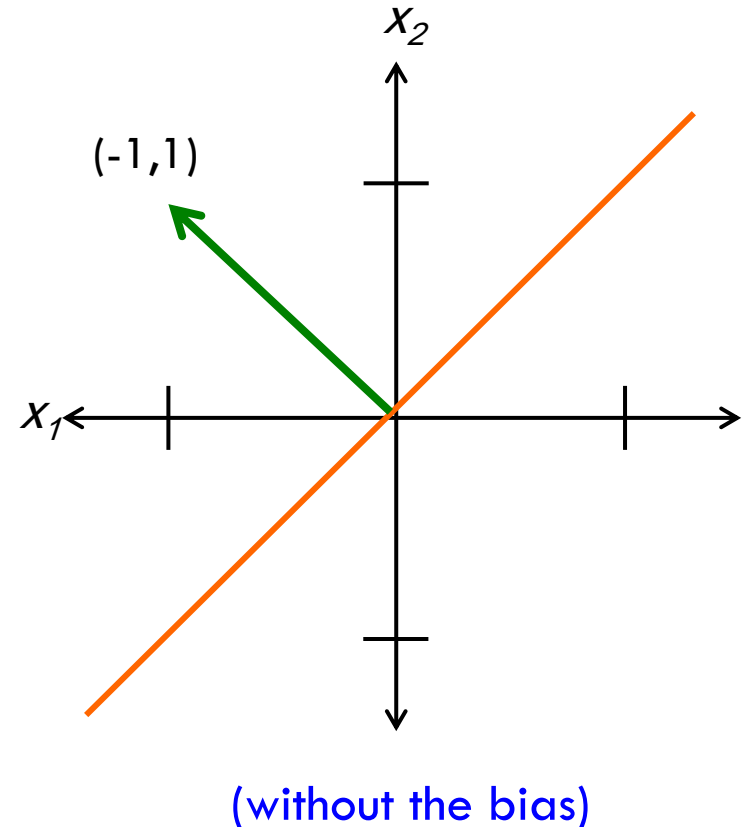
NN decision boundary



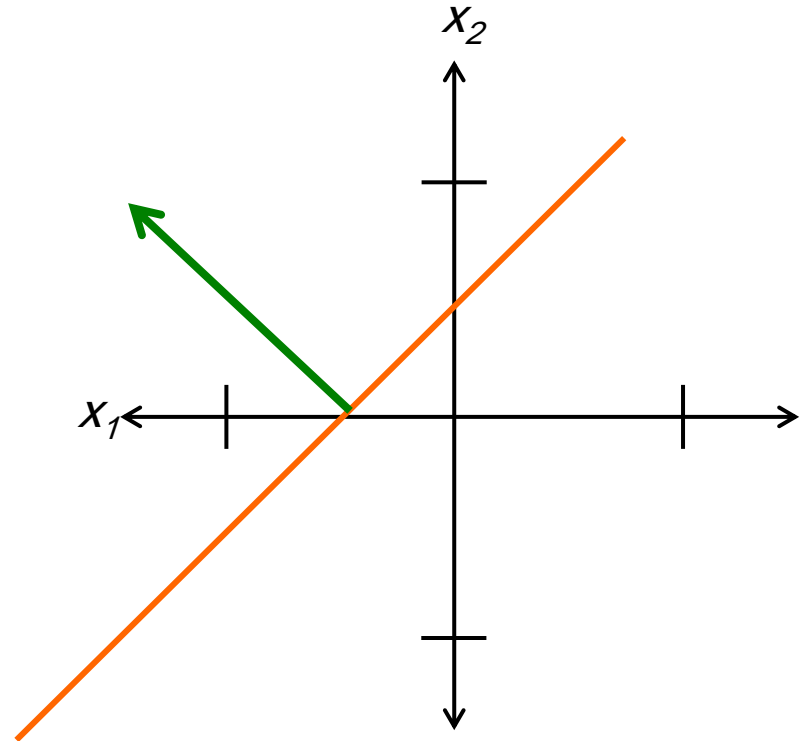
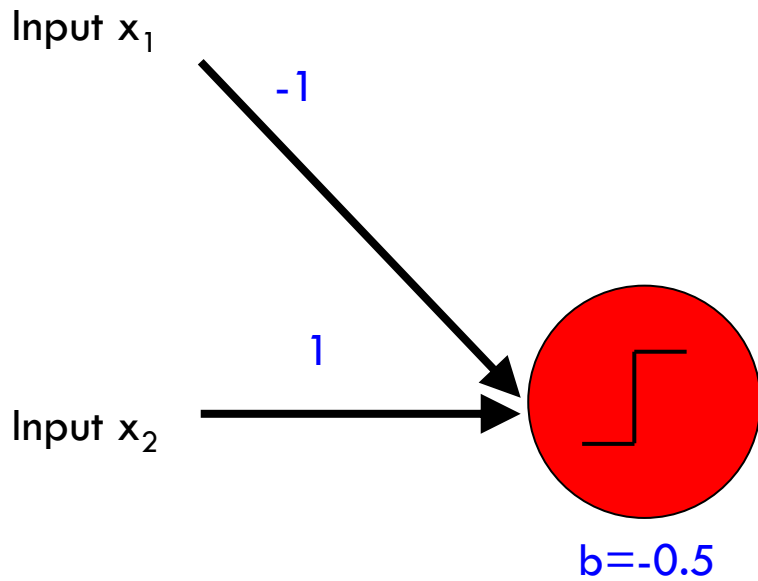
Let $x_2 = 0$, then:

$$-x_1 - 0.5 = 0$$

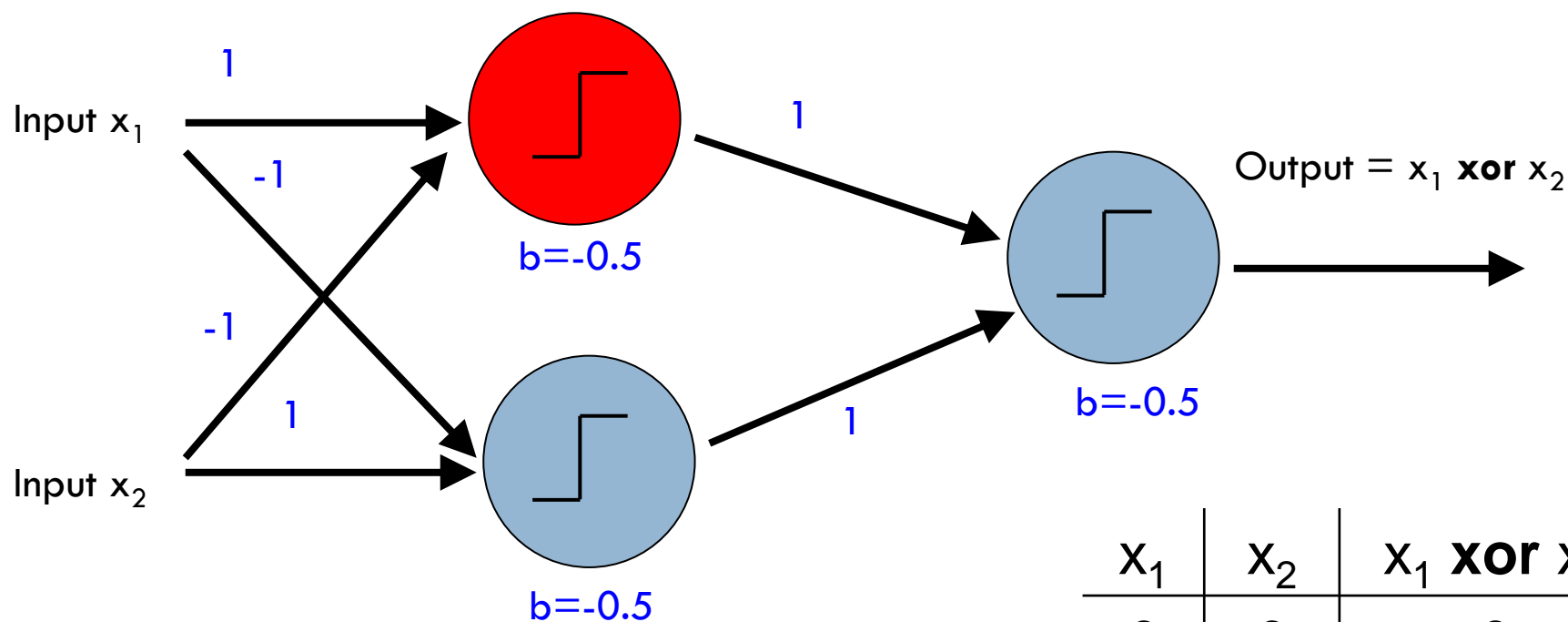
$$x_1 = -0.5$$



NN decision boundary



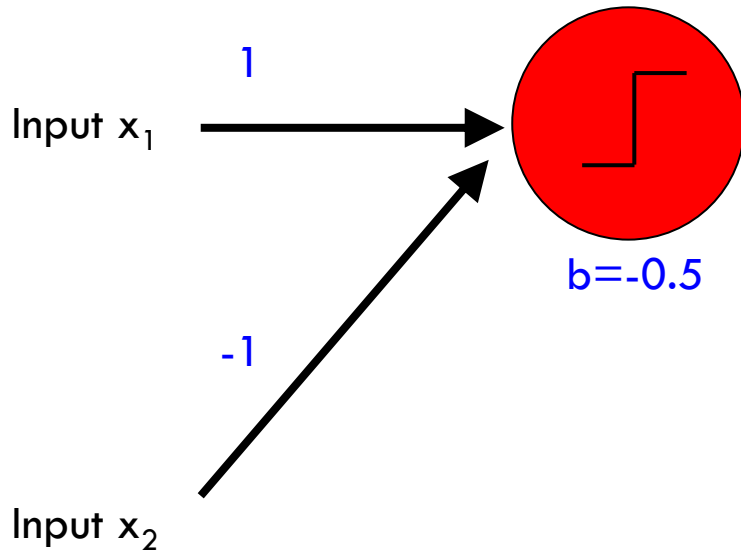
What does the decision boundary look like?



What does this perceptron's decision boundary look like?

x_1	x_2	x_1 xor x_2
0	0	0
0	1	1
1	0	1
1	1	0

NN decision boundary



Let $x_2 = 0$, then:

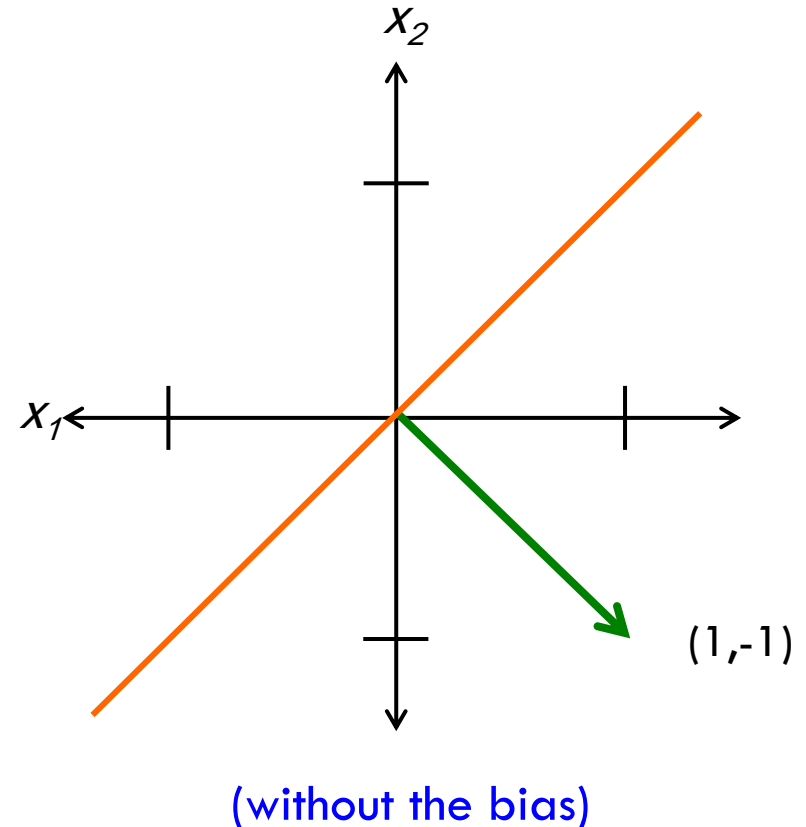
$$-x_1 - 0.5 = 0$$

$$x_1 = 0.5$$

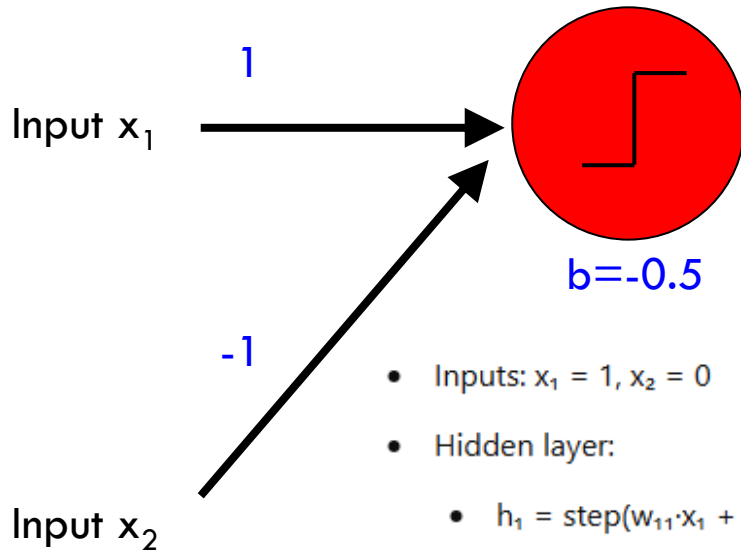
Q1: A neural network has the following configuration:

- Inputs: $x_1 = 1, x_2 = 1$
- Hidden Layer (2 neurons):
 - $h_1 = \text{step}(x_1 - x_2 + 0.5)$
 - $h_2 = \text{step}(-x_1 + x_2 + 0.5)$
- Output: $y = \text{step}(h_1 + h_2 - 0.5)$

Compute the output y .



NN decision boundary

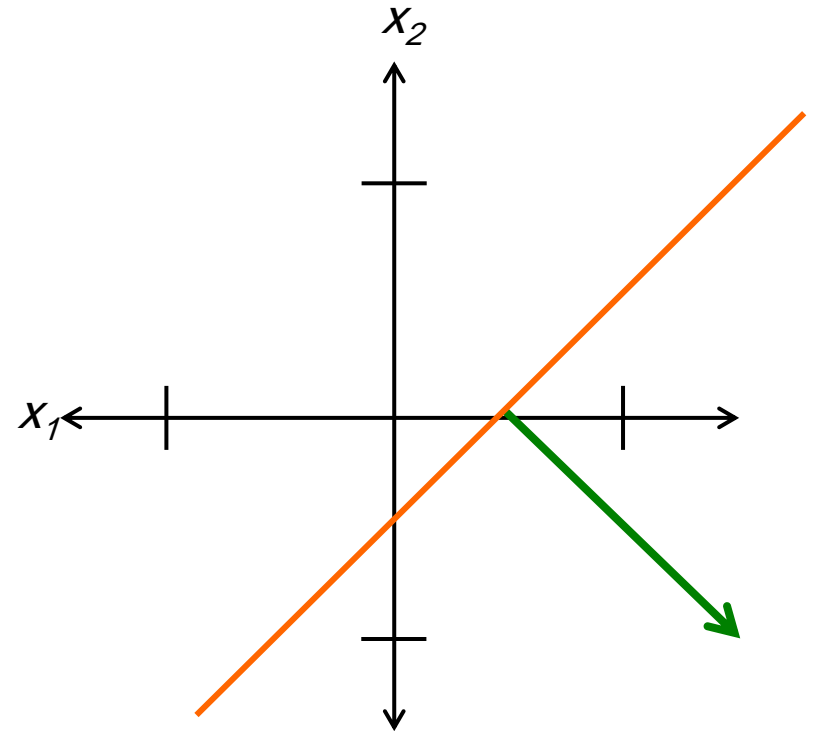


- Inputs: $x_1 = 1, x_2 = 0$
- Hidden layer:
 - $h_1 = \text{step}(w_{11} \cdot x_1 + w_{12} \cdot x_2 - \theta_1)$
 - $h_2 = \text{step}(w_{21} \cdot x_1 + w_{22} \cdot x_2 - \theta_2)$
- Output: $y = \text{step}(v_1 \cdot h_1 + v_2 \cdot h_2 - \theta_3)$

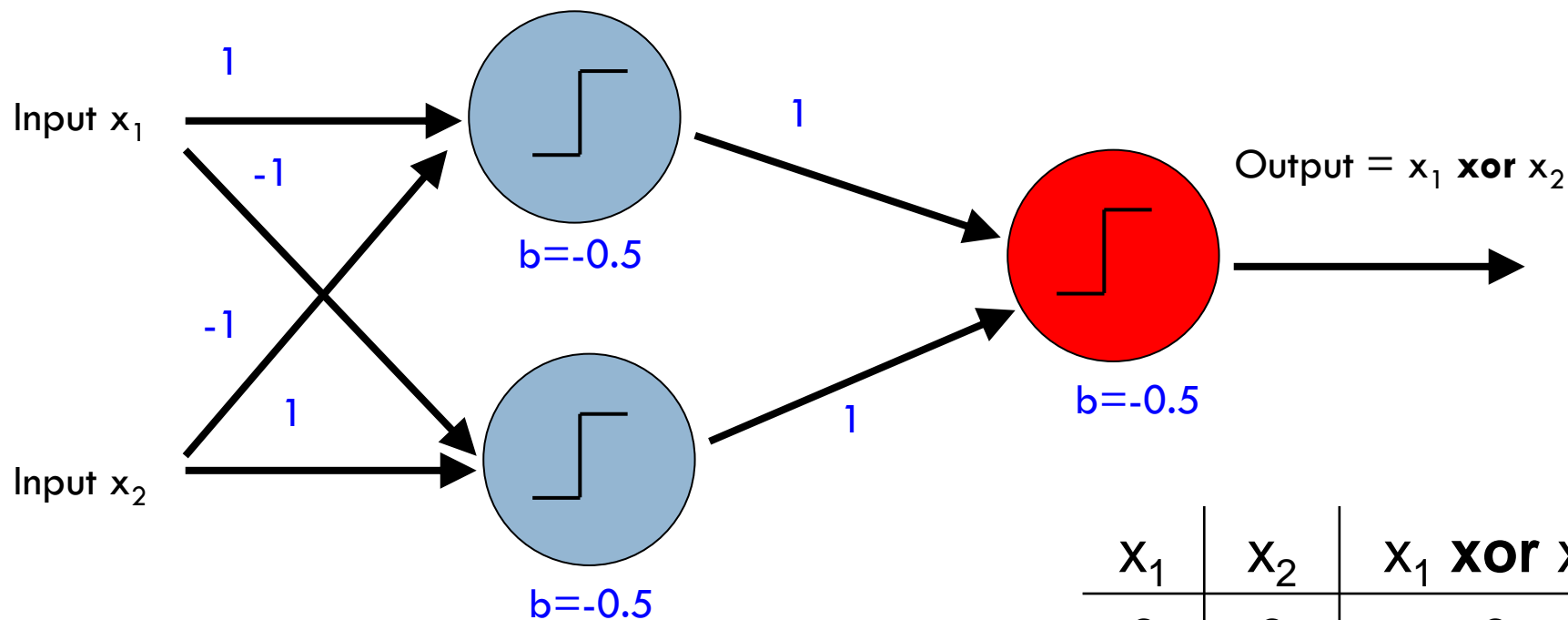
Given:

- $w_{11} = 1, w_{12} = 1, \theta_1 = 1$
- $w_{21} = 1, w_{22} = 1, \theta_2 = 0$
- $v_1 = 1, v_2 = -2, \theta_3 = 0.5$

Compute the output y .



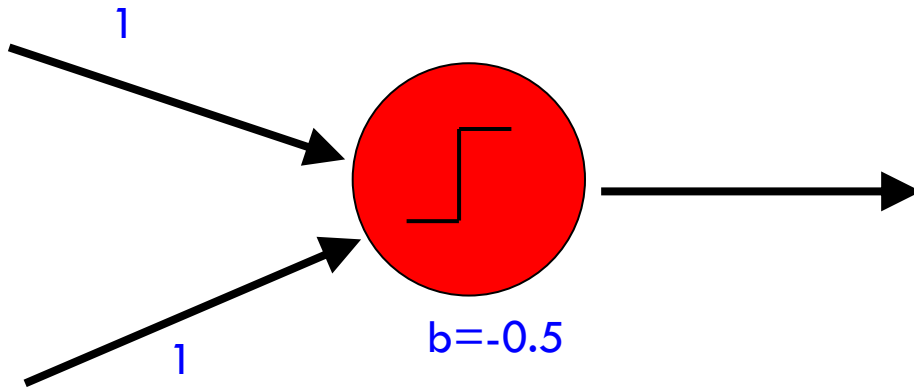
What does the decision boundary look like?



What operation does this perceptron perform on the result?

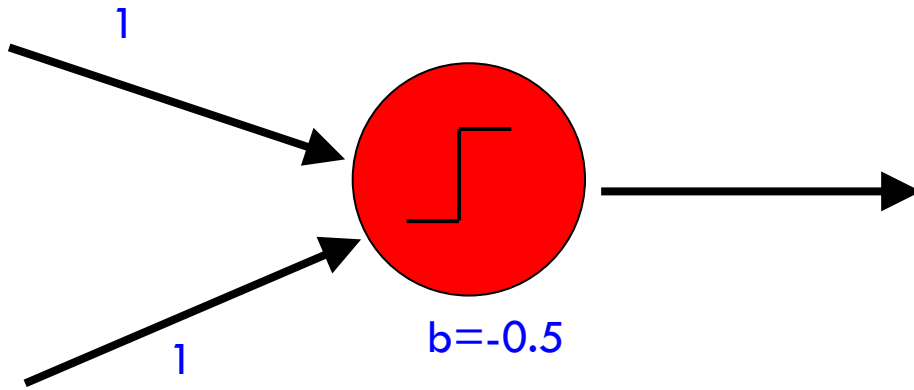
x_1	x_2	x_1 xor x_2
0	0	0
0	1	1
1	0	1
1	1	0

Fill in the truth table



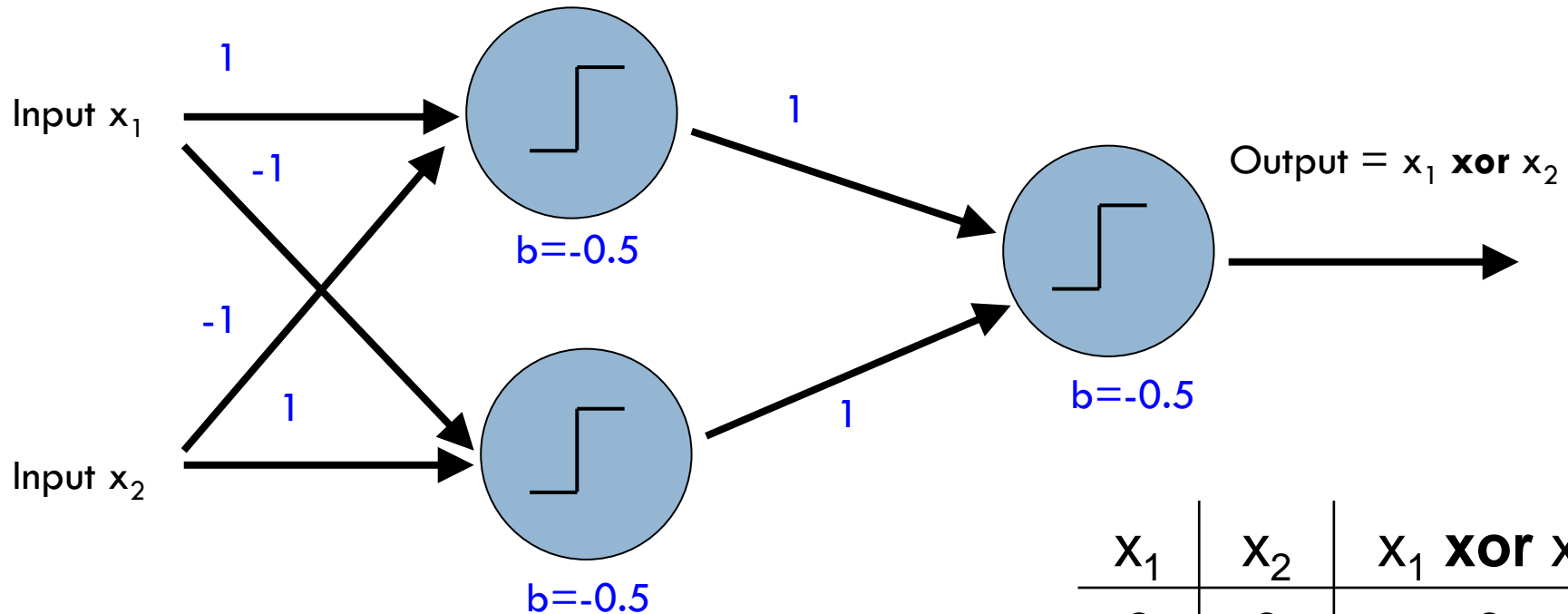
out1	out2	
0	0	?
0	1	?
1	0	?
1	1	?

OR

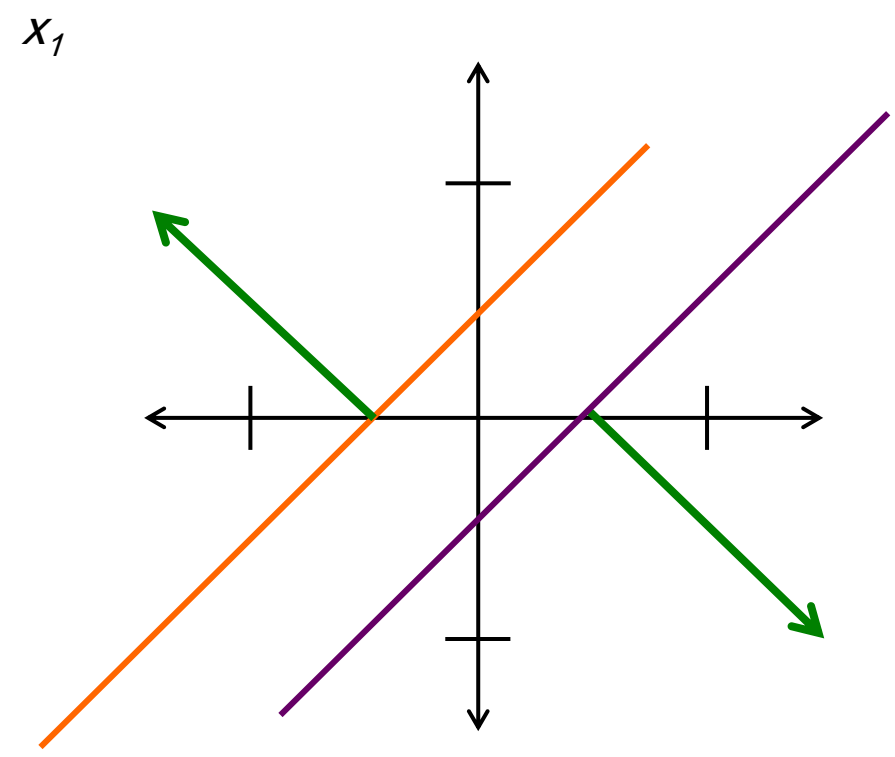
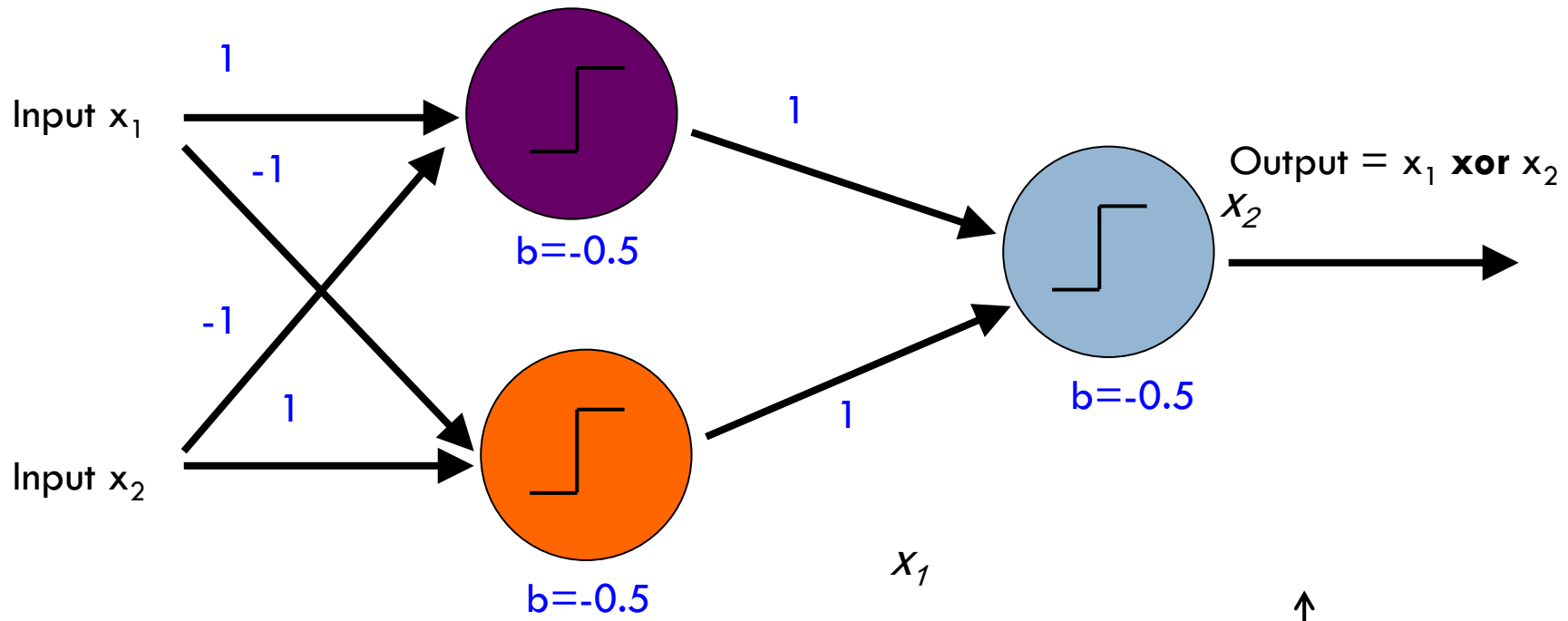


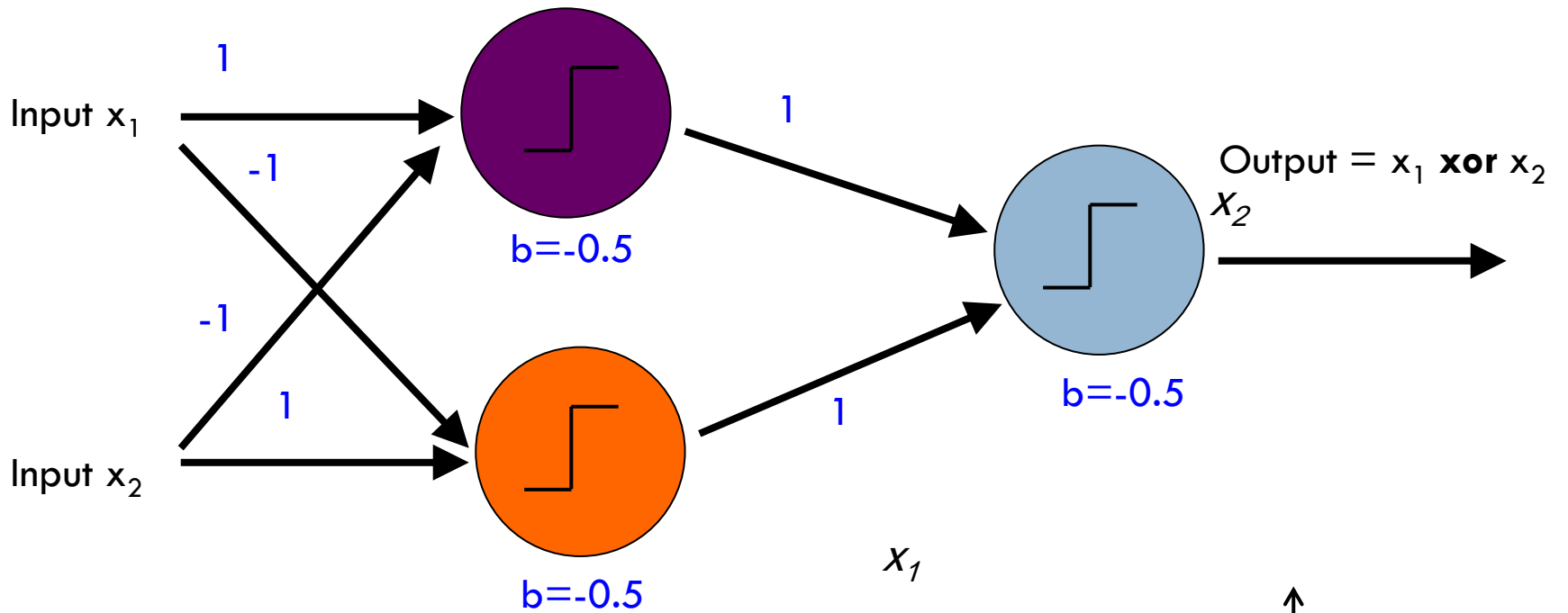
out1	out2	
0	0	0
0	1	1
1	0	1
1	1	1

What does the decision boundary look like?

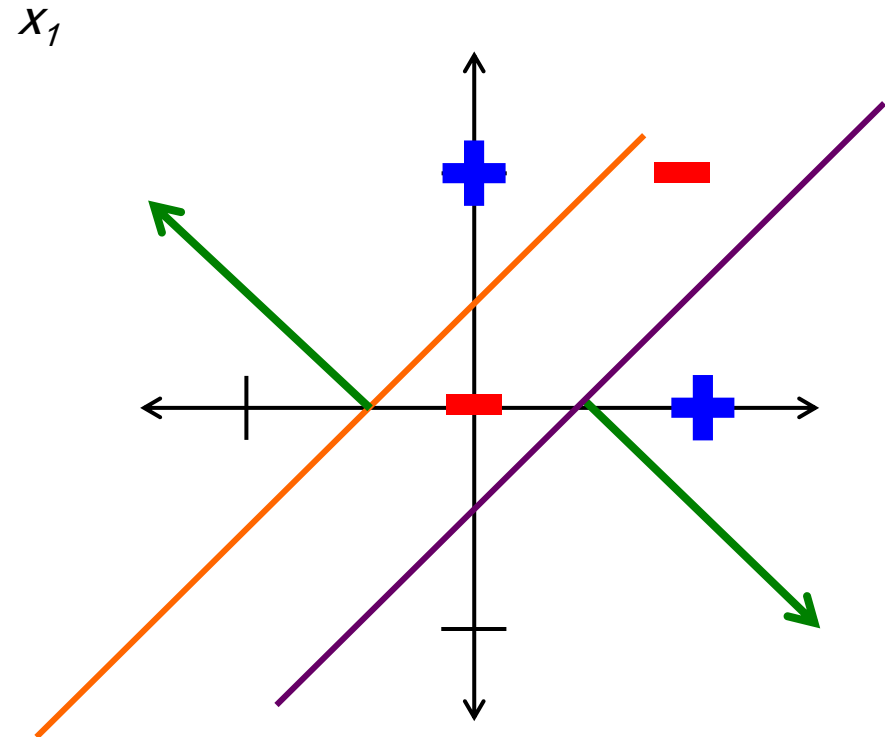


x_1	x_2	x_1 xor x_2
0	0	0
0	1	1
1	0	1
1	1	0

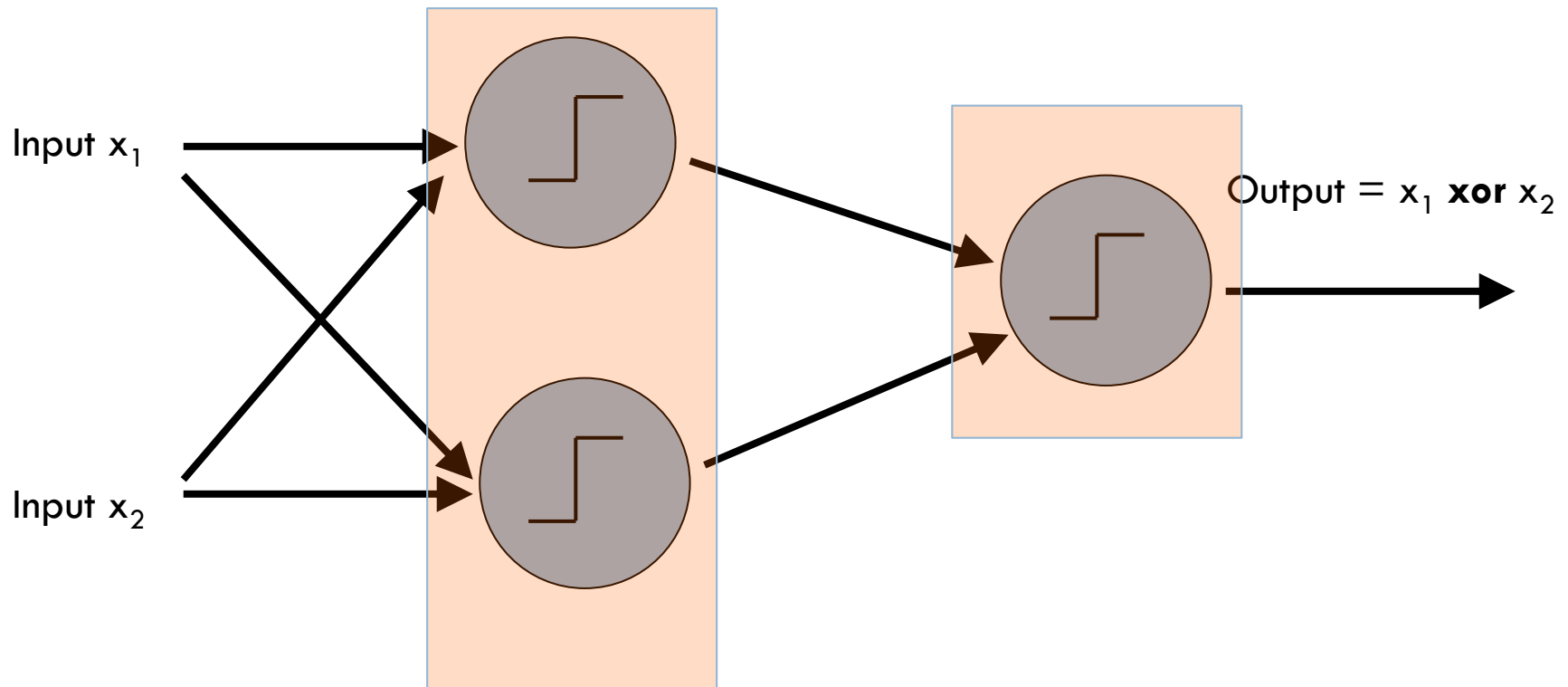




x_1	x_2	$x_1 \text{ XOR } x_2$
0	0	0
0	1	1
1	0	1
1	1	0



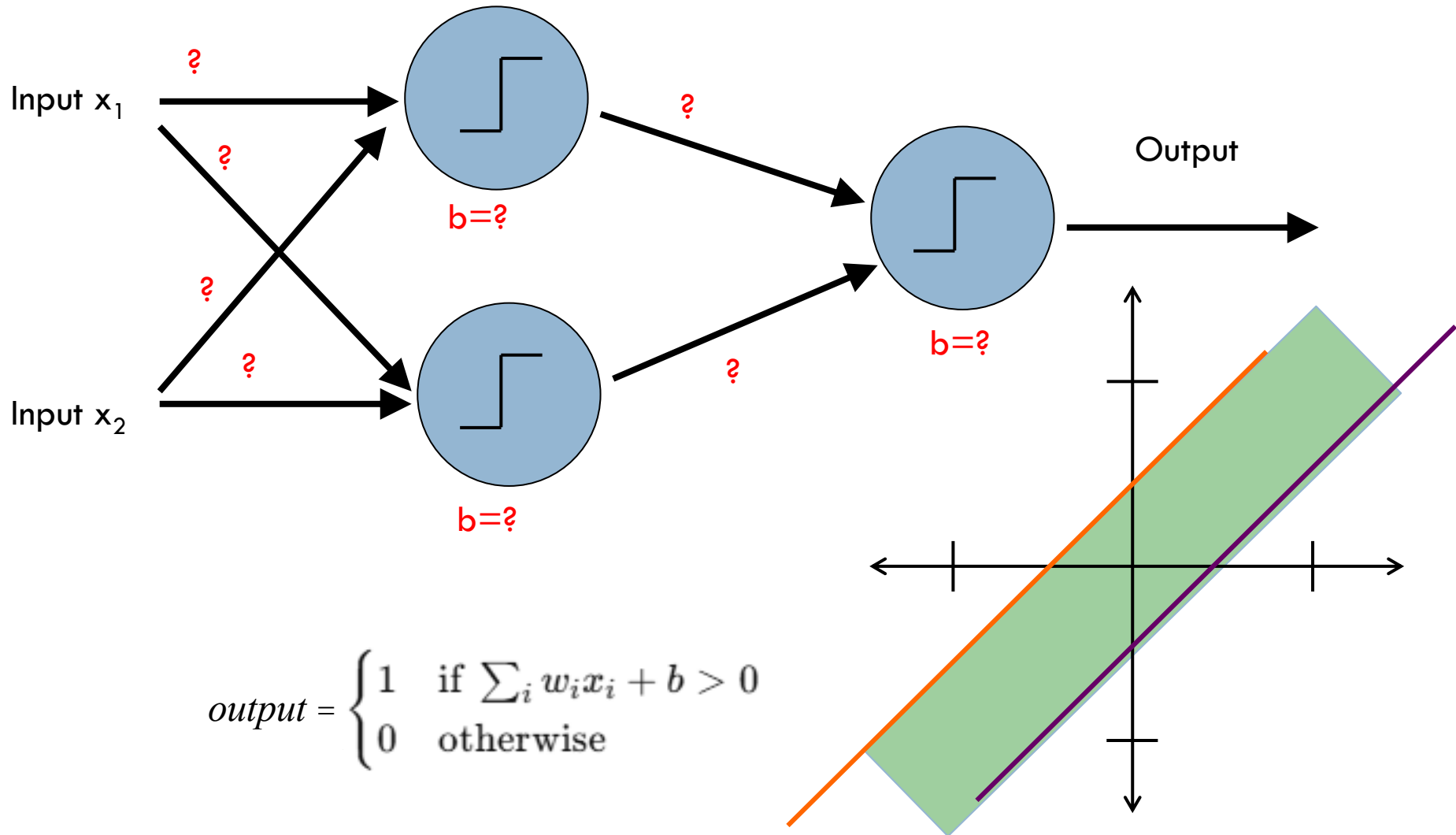
What does the decision boundary look like?



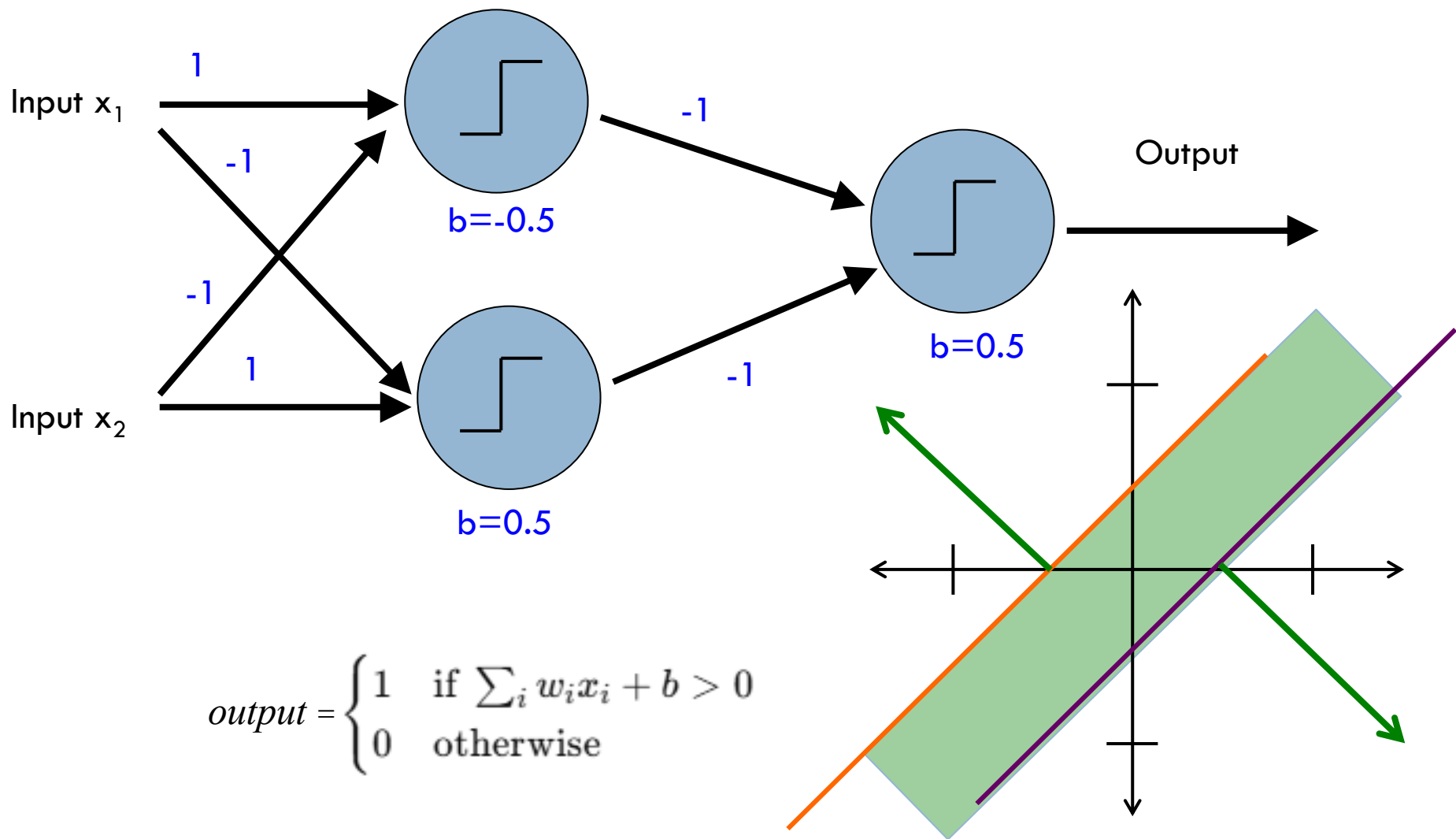
linear splits of the
feature space

combination of
these linear spaces

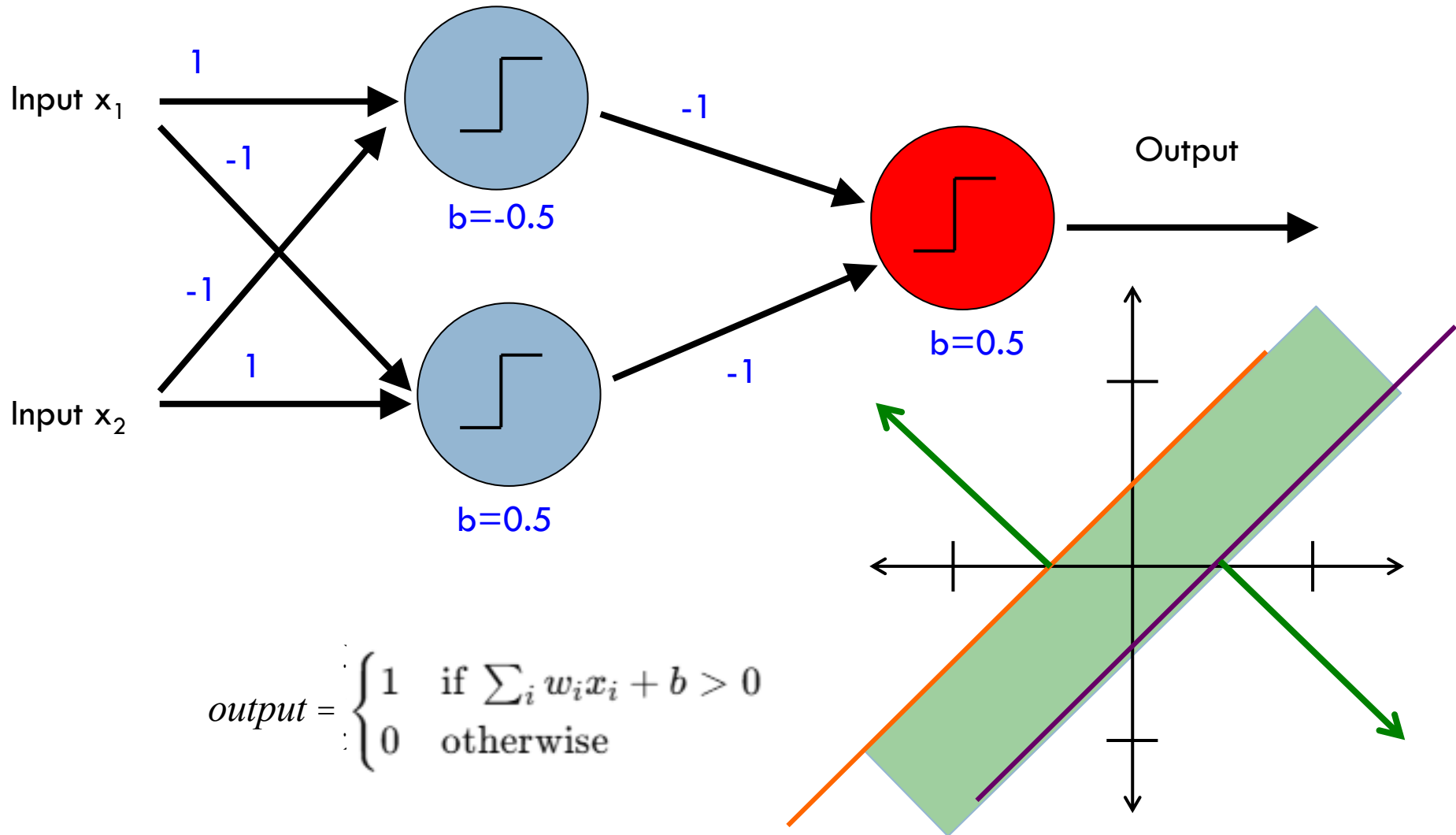
This decision boundary?

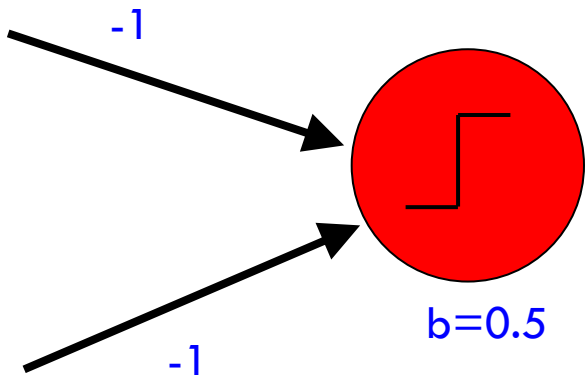


This decision boundary?



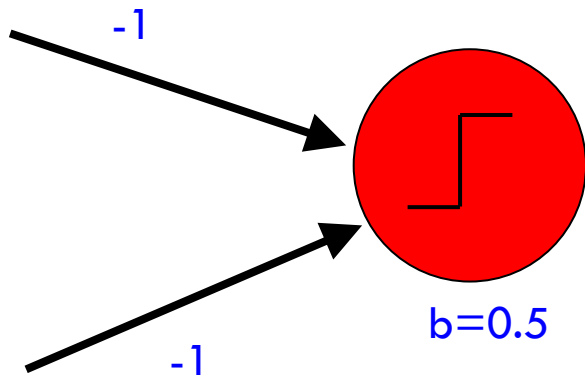
This decision boundary?





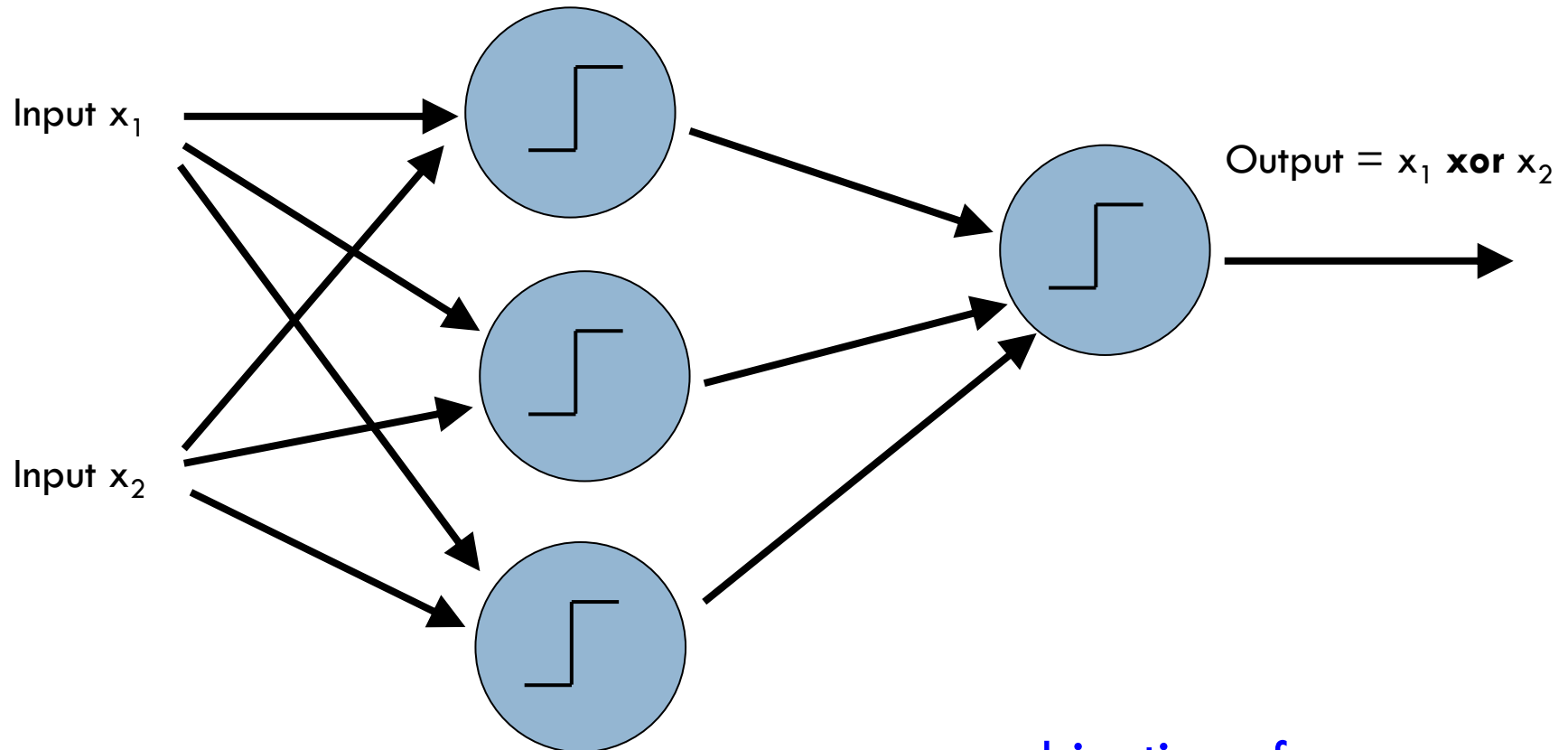
out1	out2	
0	0	?
0	1	?
1	0	?
1	1	?

NOR



out1	out2	
0	0	1
0	1	0
1	0	0
1	1	0

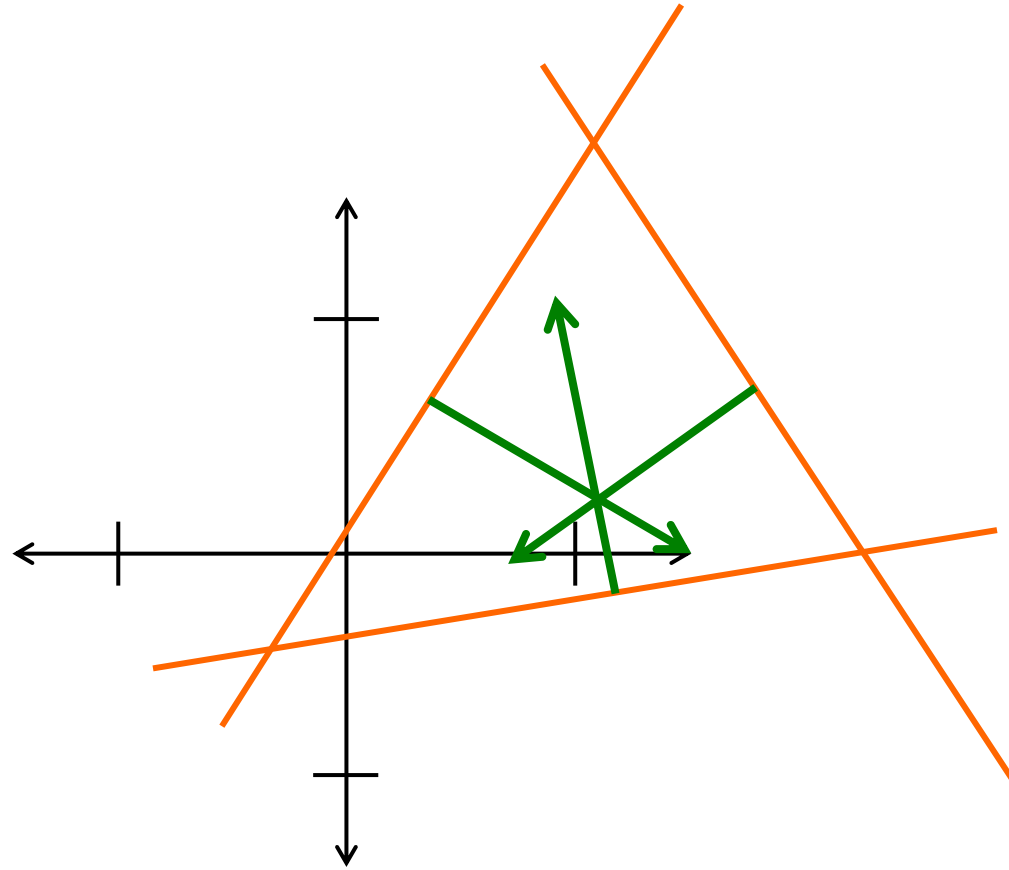
What does the decision boundary look like?



linear splits of the
feature space

combination of
these linear spaces

Three hidden nodes



NN decision boundaries

For DT, as the tree gets larger, the model gets more complex

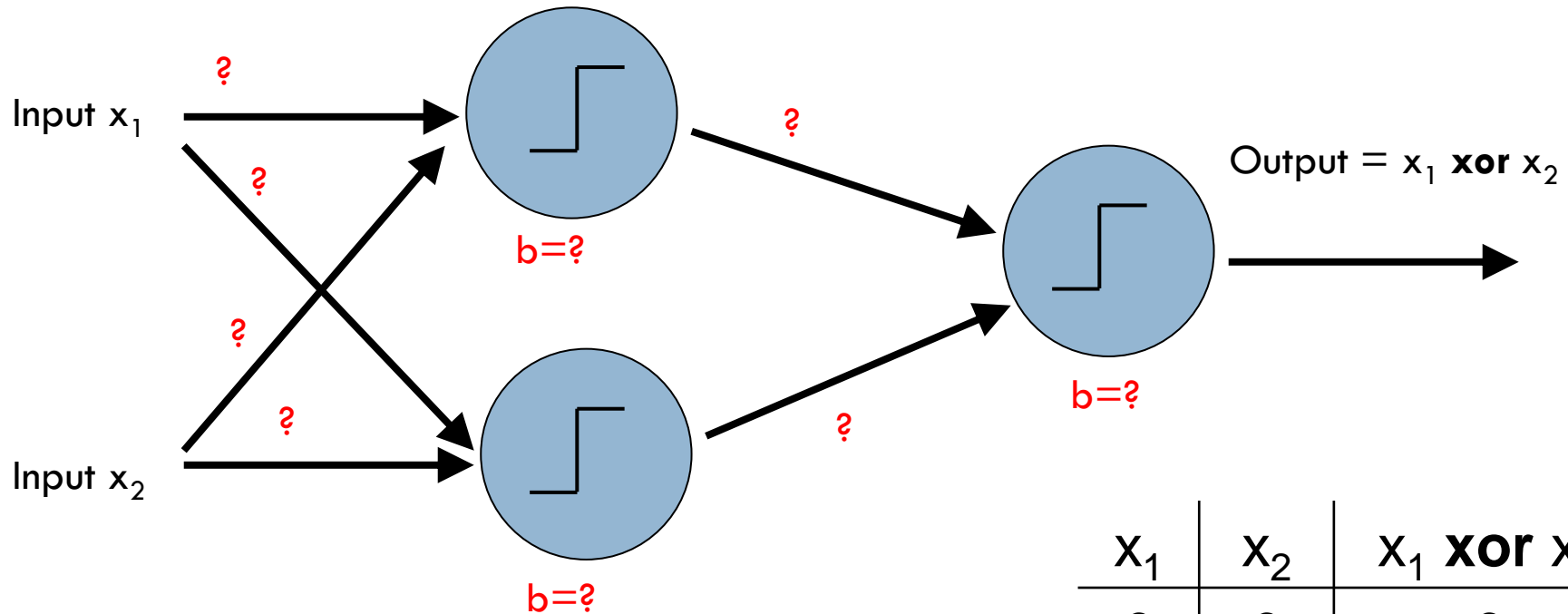
The same is true for neural networks:
more hidden nodes = more complexity

Adding more layers adds even more complexity (and much more quickly)

Good rule of thumb:

$$\text{number of 2-layer hidden nodes} \leq \frac{\text{number of examples}}{\text{number of dimensions}}$$

Training



x_1	x_2	$x_1 \text{ XOR } x_2$
0	0	0
0	1	1
1	0	1
1	1	0

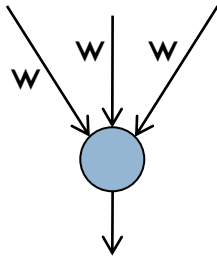
How do we learn the weights?

Training multilayer networks

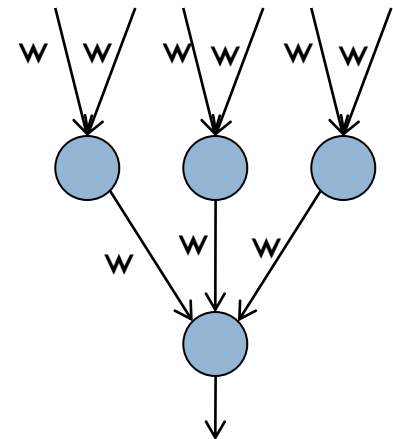
perceptron learning: if the perceptron's output is different than the expected output, update the weights

gradient descent: compare output to label and adjust based on loss function

Any other problem with these for general NNs?



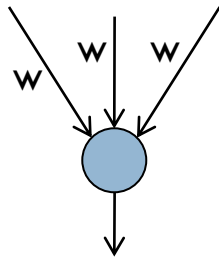
perceptron/
linear model



neural network

Learning in multilayer networks

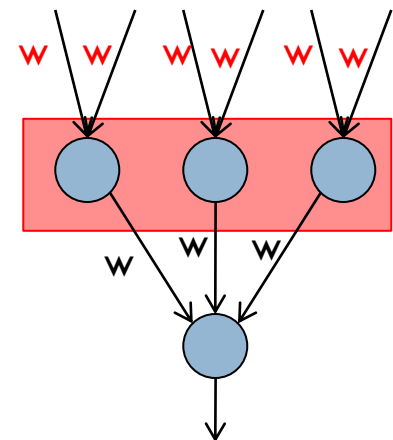
Challenge: for multilayer networks, we don't know what the expected output/error is for the internal nodes!



perceptron/
linear model

how do we learn these weights?

expected output?



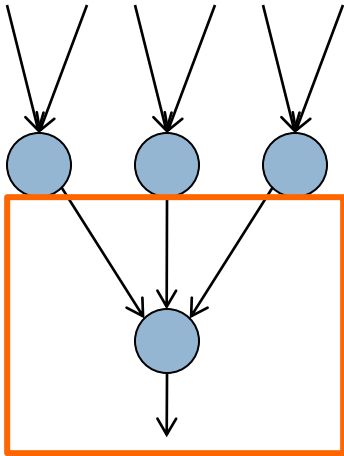
neural network

Backpropagation: intuition

Gradient descent method for learning weights by optimizing a loss function

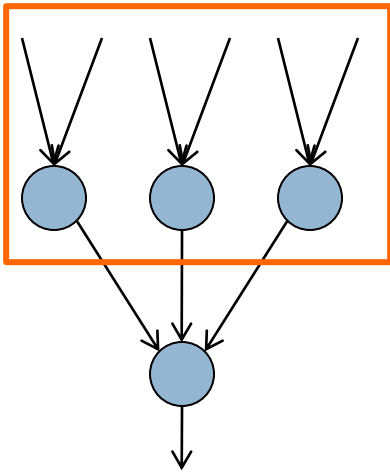
1. calculate output of all nodes
2. calculate the weights for the output layer based on the error
3. “backpropagate” errors through hidden layers

Backpropagation: intuition



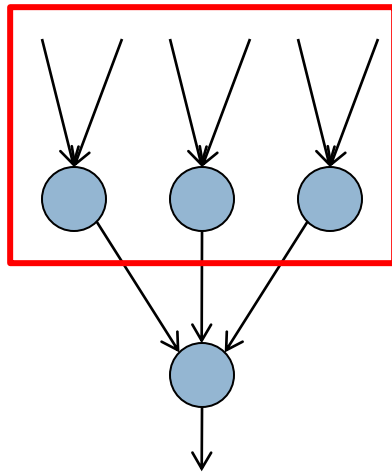
We can calculate the actual error here

Backpropagation: intuition



Key idea: propagate the error back to this layer

Backpropagation: intuition

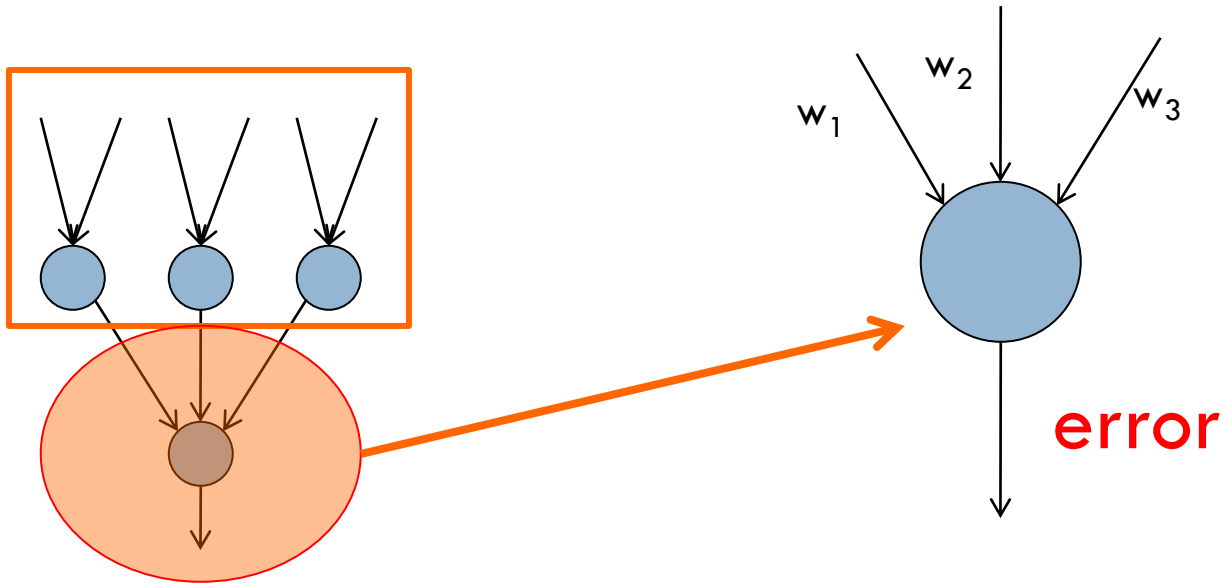


“backpropagate” the error:

Assume all of these nodes were responsible for some of the error

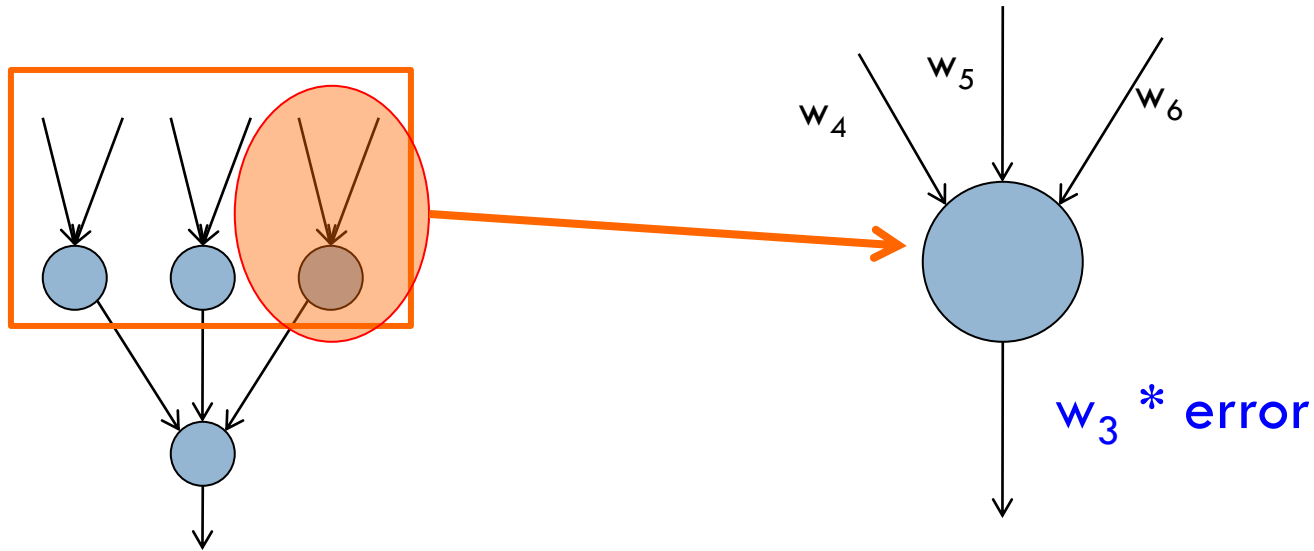
How can we figure out how much they were responsible for?

Backpropagation: intuition



error for node is $\sim w_i * \text{error}$

Backpropagation: intuition



Calculate as normal using this as the error

Backpropagation: the details

Gradient descent method for learning weights by optimizing a **loss function**

1. calculate output of all nodes
2. calculate the updates directly for the output layer
3. “backpropagate” errors through hidden layers

What loss function?

Backpropagation: the details

Gradient descent method for learning weights by optimizing a **loss function**

1. calculate output of all nodes
2. calculate the updates directly for the output layer
3. “backpropagate” errors through hidden layers

$$loss = \sum_x \frac{1}{2} (y - \hat{y})^2 \quad \text{squared error}$$