# Neural Networks & Deep Learning
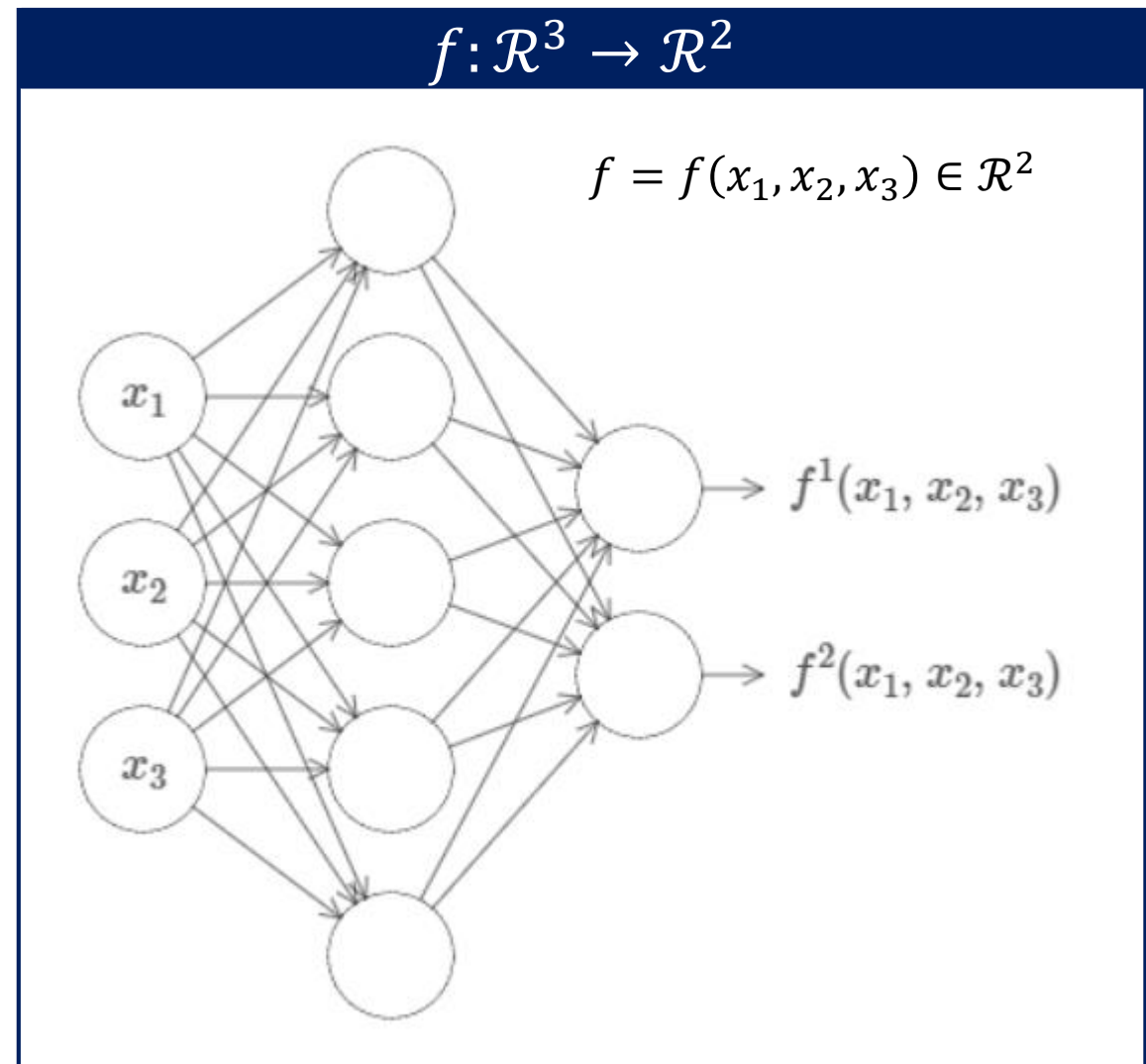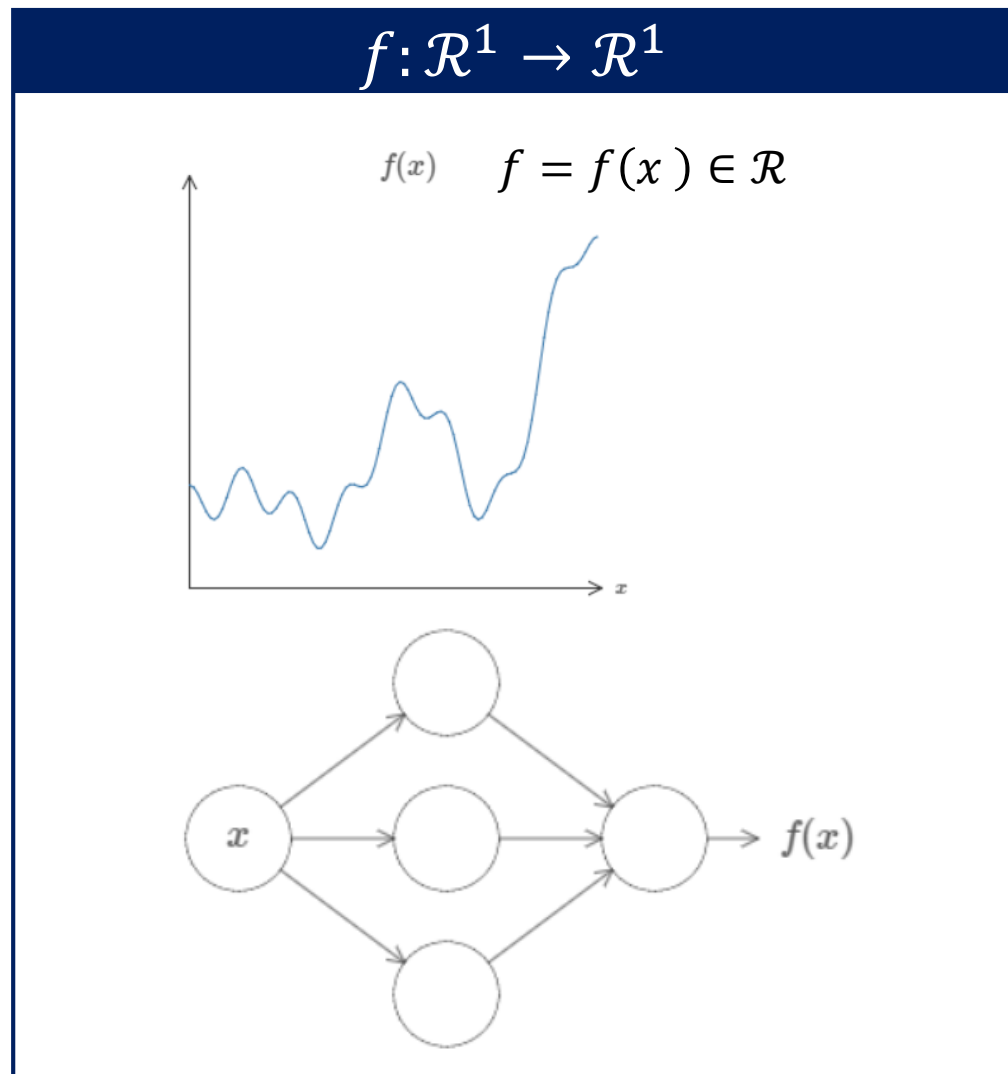
http://neuralnetworksanddeeplearning.com/chap4.html

# The universality theorem

**neural networks can compute any function at all.**



$f: \mathcal{R}^1 \rightarrow \mathcal{R}^1$

$f = f(x) \in \mathcal{R}$

$f: \mathcal{R}^3 \rightarrow \mathcal{R}^2$

$f = f(x_1, x_2, x_3) \in \mathcal{R}^2$

$f^1(x_1, x_2, x_3)$

$f^2(x_1, x_2, x_3)$

# The universality theorem

## Two caveats

# 1.approximation

First, this doesn't mean that a network can be used to *exactly* compute any function.
Rather, we can get an *approximation* that is as good as we want.
By increasing the number of hidden neurons we can improve the approximation.
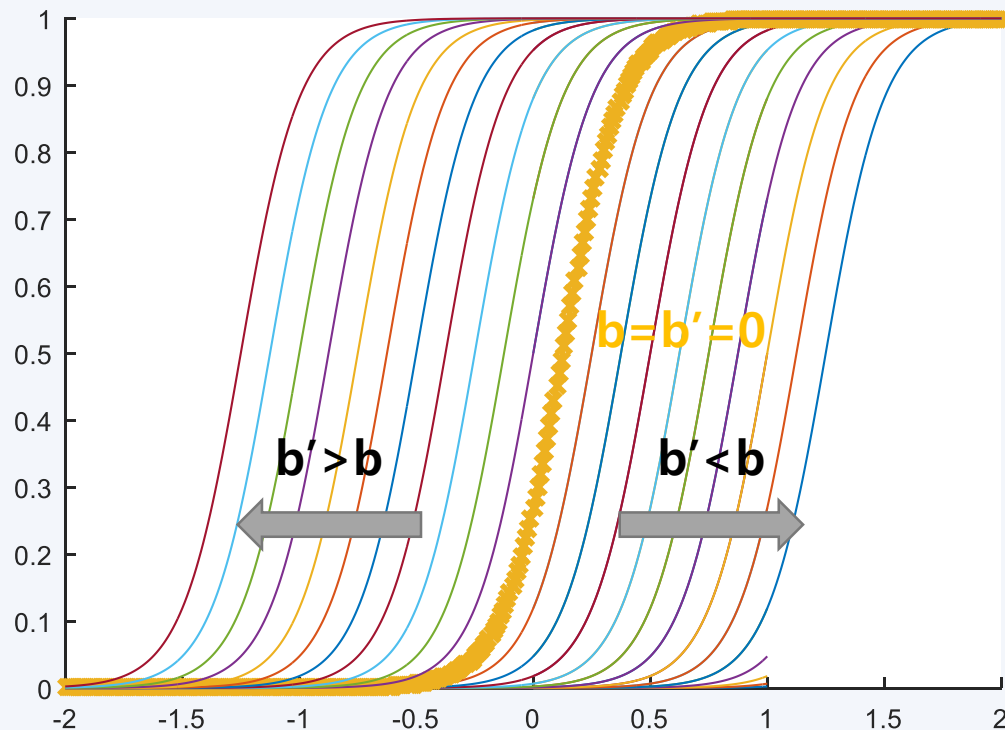
# 2.continuous function

The second caveat is that the class of functions which can be approximated in the way described are the *continuous* functions. If a function is discontinuous, i.e., makes sudden, sharp jumps, then it won't in general be possible to approximate using a neural net.

# Universality with one input and one output
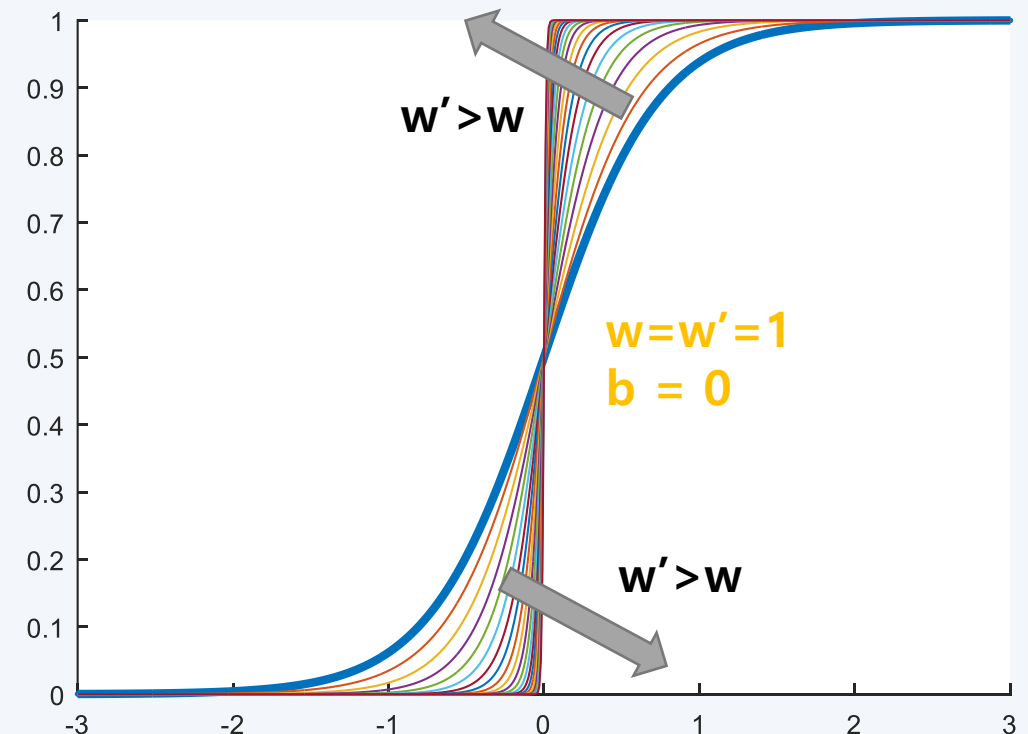
## Shape change by weight and bias

- As bias term increases, graph for output from hidden neuron goes left without shape change
- As bias term decreases, graph for output from hidden neuron goes right without shape change

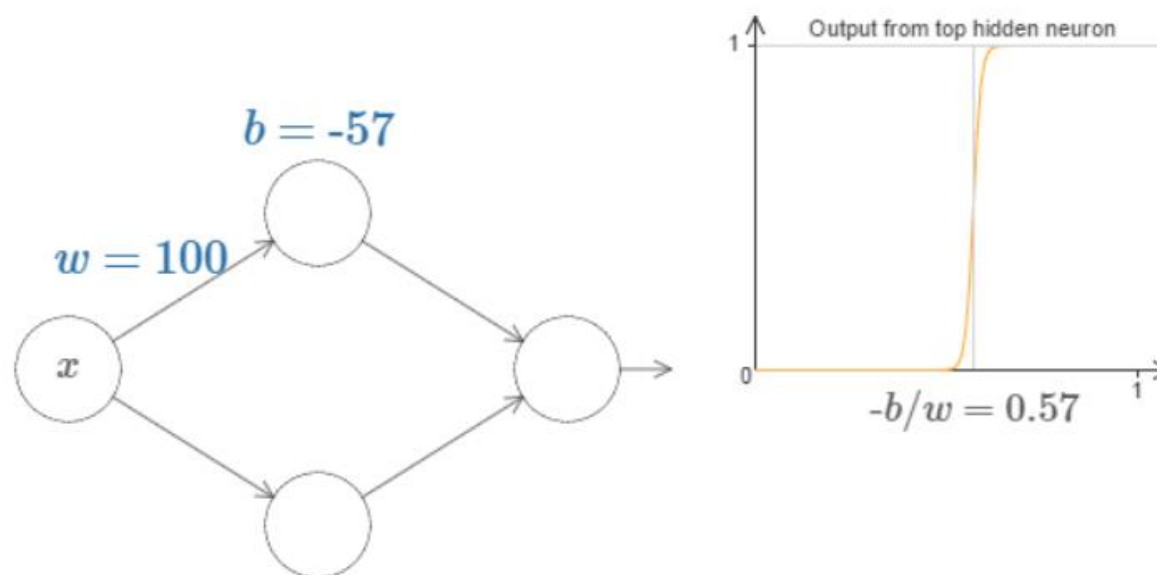$$\sigma(wx + b) = \sigma(z) \rightarrow \sigma(wx + b')$$

- As weight term increases, graph for output from hidden neuron changes its shape.
- The curve gets steeper, until eventually it begins to look like a step function

$$\sigma(wx + b) = \sigma(z) \rightarrow \sigma(w'x + b)$$

# Universality with one input and one output

## One hidden layer with two neurons



Output from top hidden neuron

$b = \text{-}57$

$w = 100$

$x$

$\text{-}b/w = 0.57$

- It's actually quite a bit easier to work with step functions than general sigmoid functions.
- The reason is that in the output layer we add up contributions from all the hidden neurons.
- For this, we use 'w' with a very large – big enough that the step function is a very good approximation.
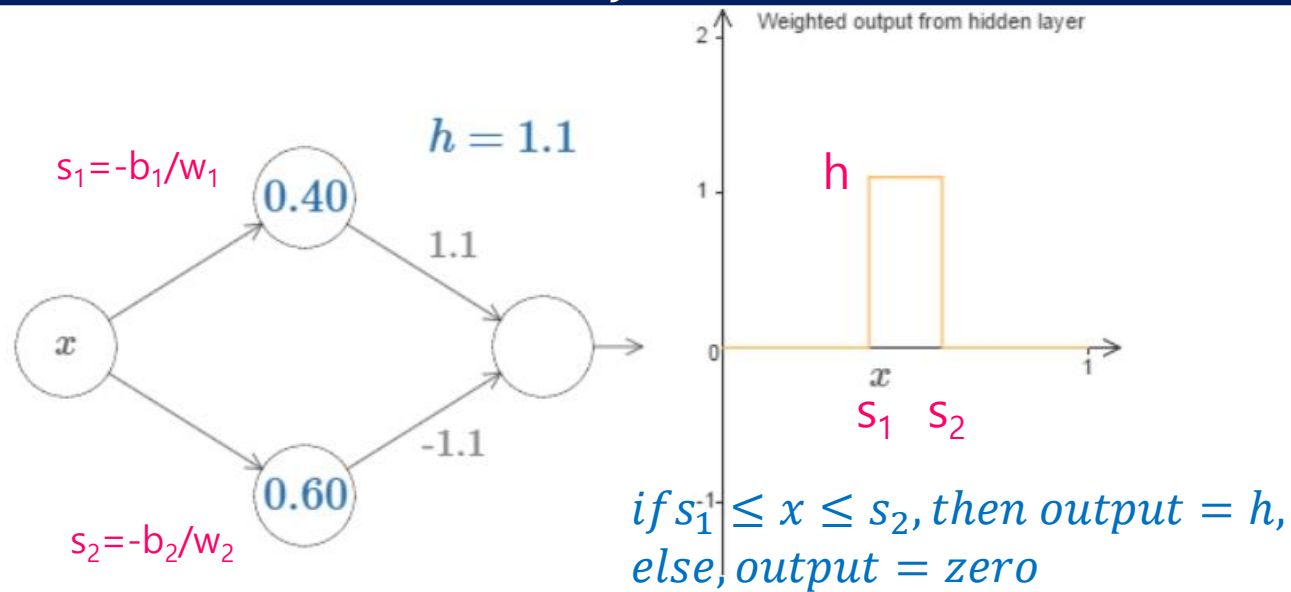- If so, at what value of x does the step occur?

$$\text{x} = \text{-b/w} = \text{s}$$

# Universality with one input and one output

**One hidden layer with two neurons**

A pair of step-function like neurons with weights summing up to zero gives "bump" function

| Bumb function | Delta function |
|---|---|



$s_1 = -b_1/w_1$

$h = 1.1$

0.40

1.1

$x$

-1.1

0.60

$s_2 = -b_2/w_2$

Weighted output from hidden layer

h

$s_1$   $s_2$

$if\ s_1^{-1} \le x \le s_2,\ then\ output = h,$
$else,\ output = zero$

"bump" function is also called if-then-else statement

$h = -0.7$

0.54

-0.7

$x$       $s_2 - s_1 \rightarrow 0$

0.7

0.55

Weighted output from hidden layer

We can approximate any functions with may pairs of delta-function

# Universality with one input and one output

**Example**

$$f(x) = 0.2 + 0.4x^2 + 0.3x\sin(15x) + 0.05\cos(50x)$$



- 5 pairs of neurons
- A combination of 5 step functions
- Only 'h' values are learned

$$f(x) = \sigma(z)$$

$$z = \sigma^{-1} \circ f(x)$$

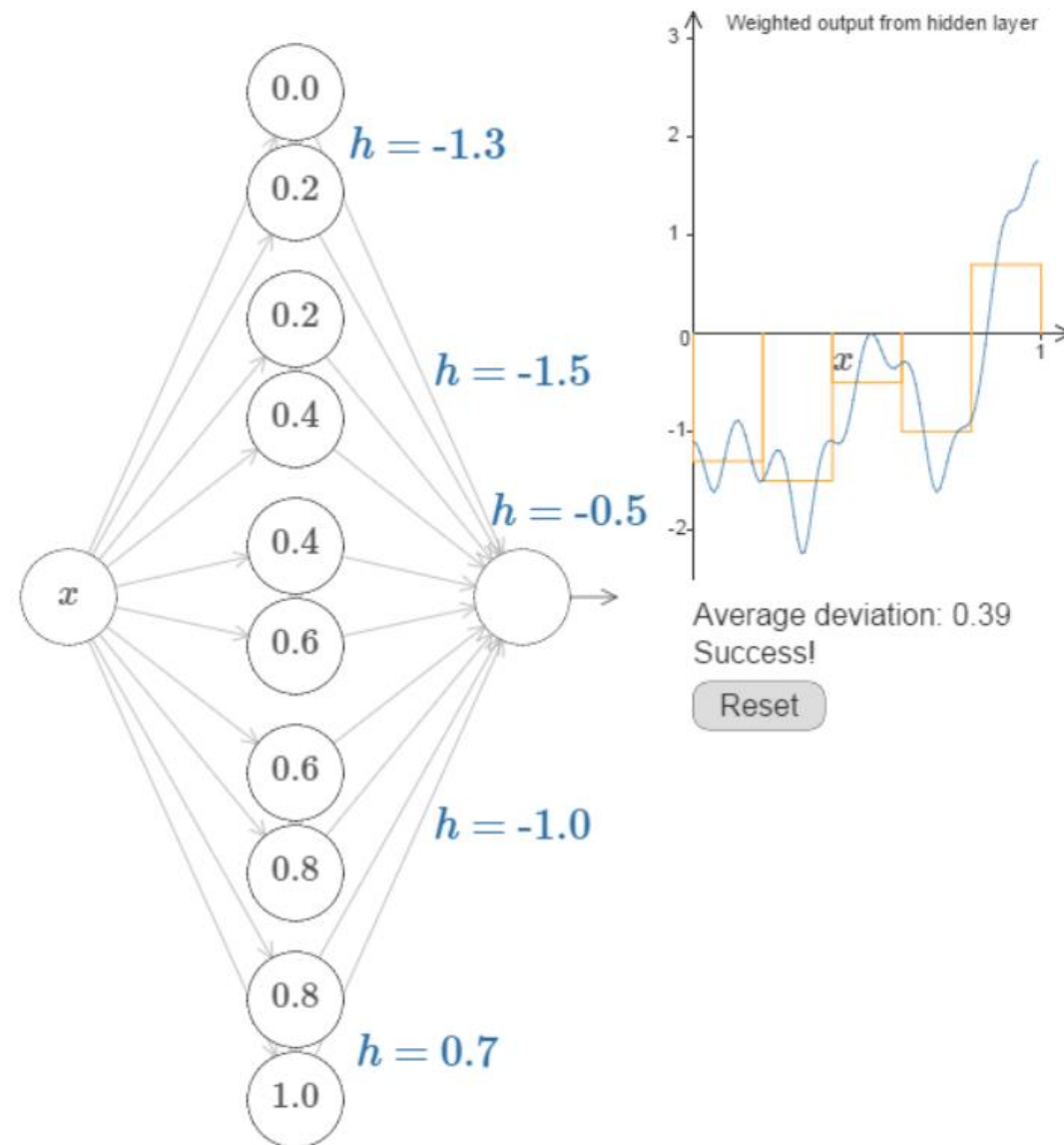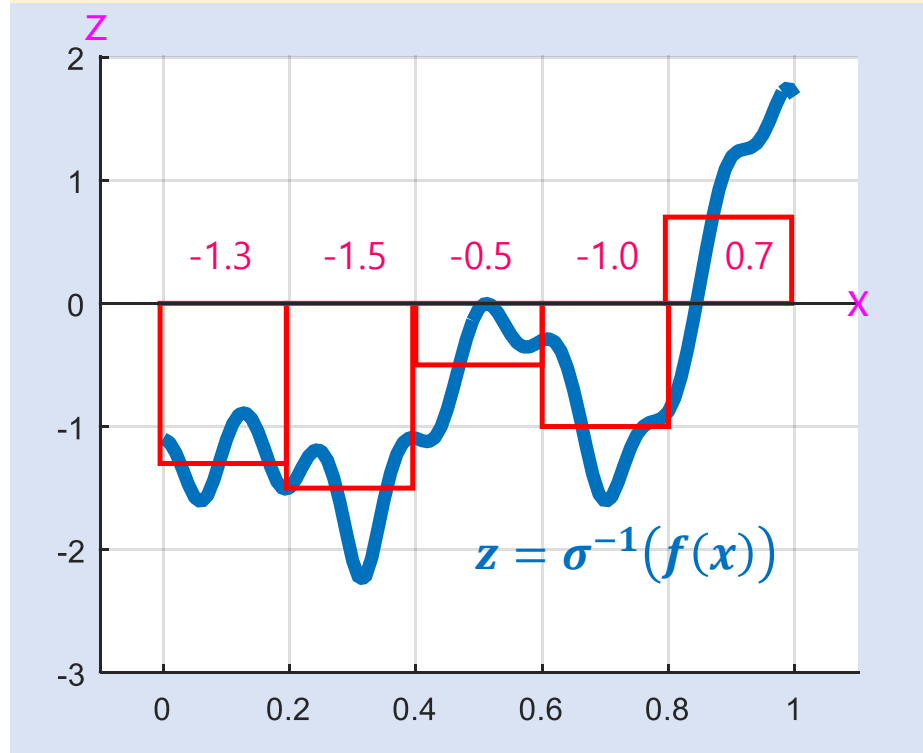For the simplicity, we are focusing on 'z'

# Universality with one input and one output

**Example**

$$f(x) = 0.2 + 0.4x^2 + 0.3x\sin(15x) + 0.05\cos(50x)$$
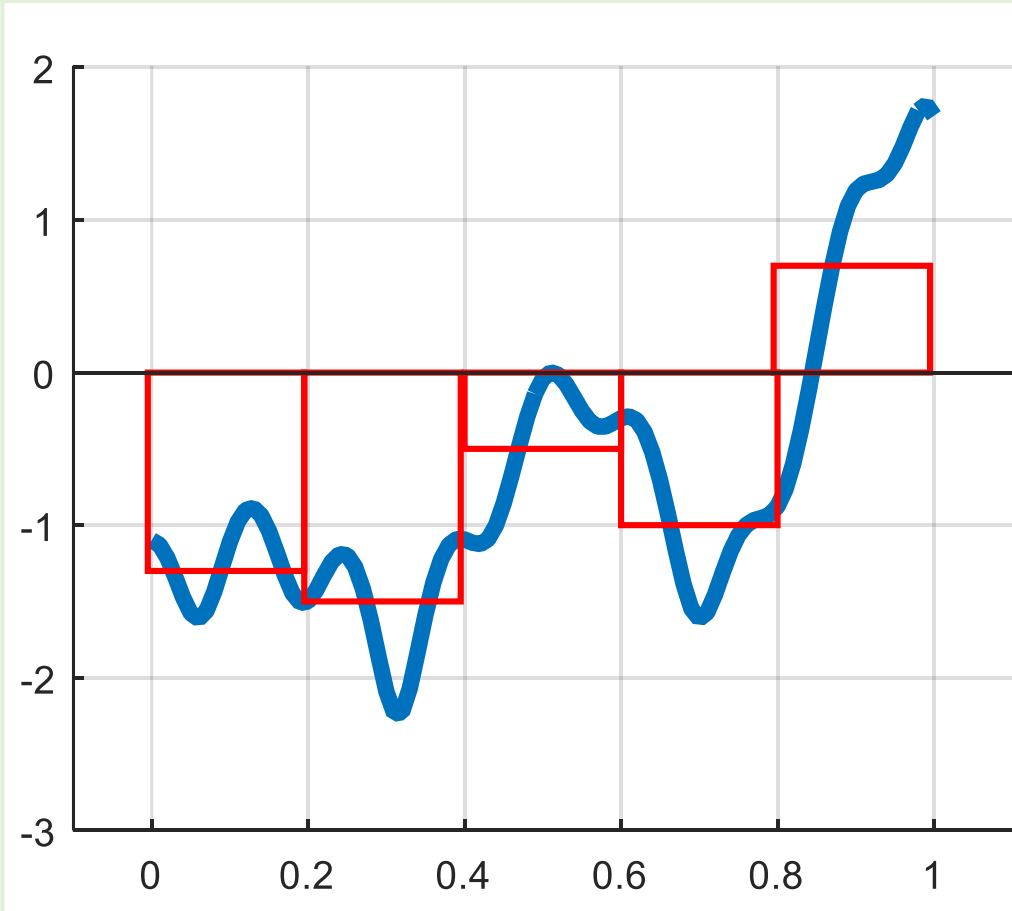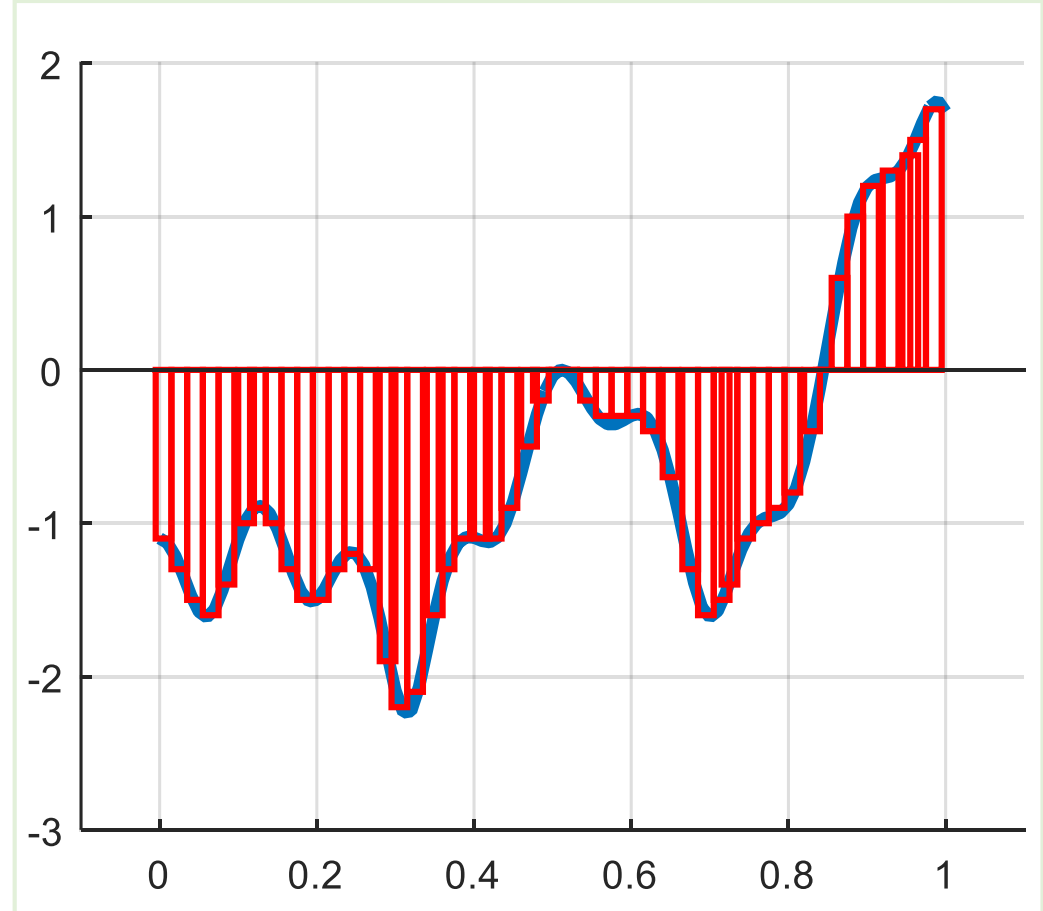
$$z = \sigma^{-1}(f(x)) = \sigma^{-1}(y) = \ln\left(\frac{y}{1-y}\right)$$

# Universality with one input and one output

**Example**

# Universality with one input and one output

**Example**

How about adding a convolution layer??



**Magnified View**

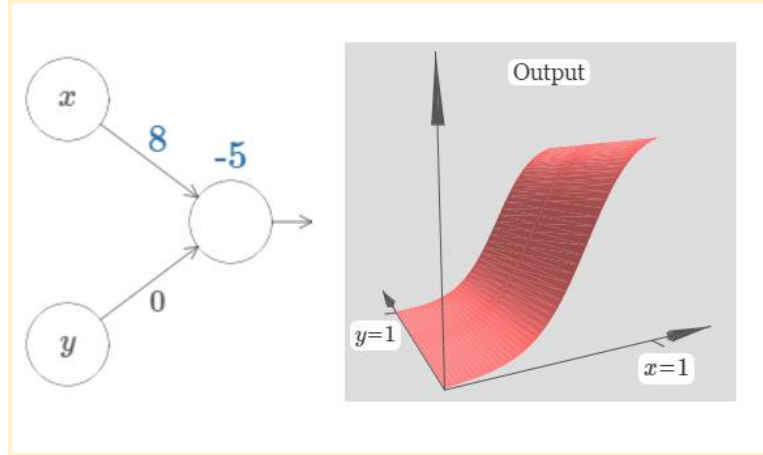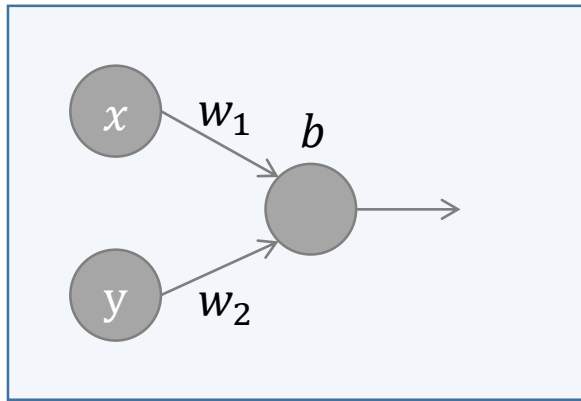# Universality with one input and one output

**Example**

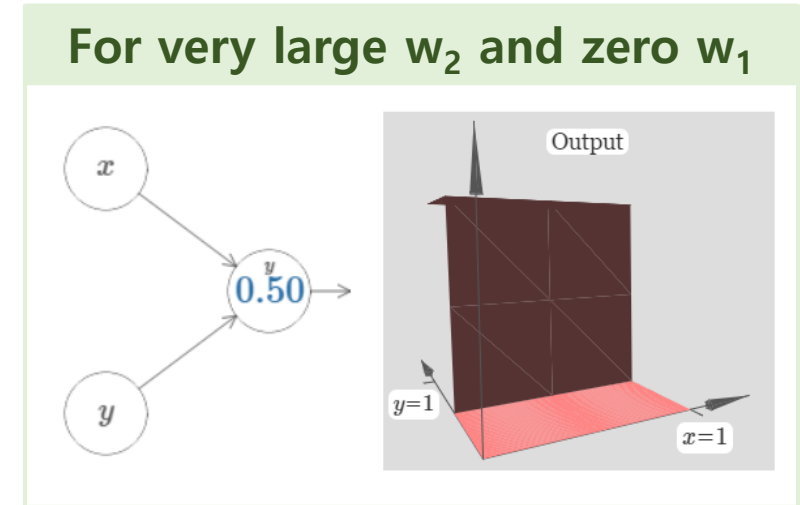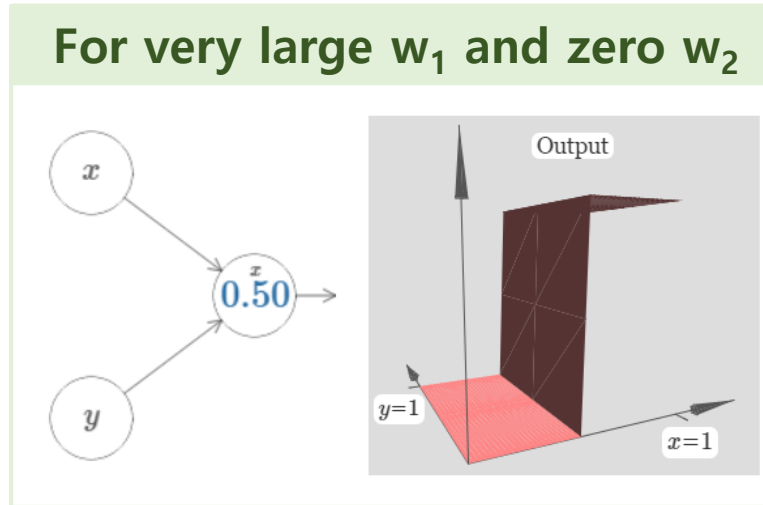How about adding a convolution layer??

# Many input variables

## Two input variables



As bias term increases, graph for output from hidden neuron goes left on the x-axis without shape change
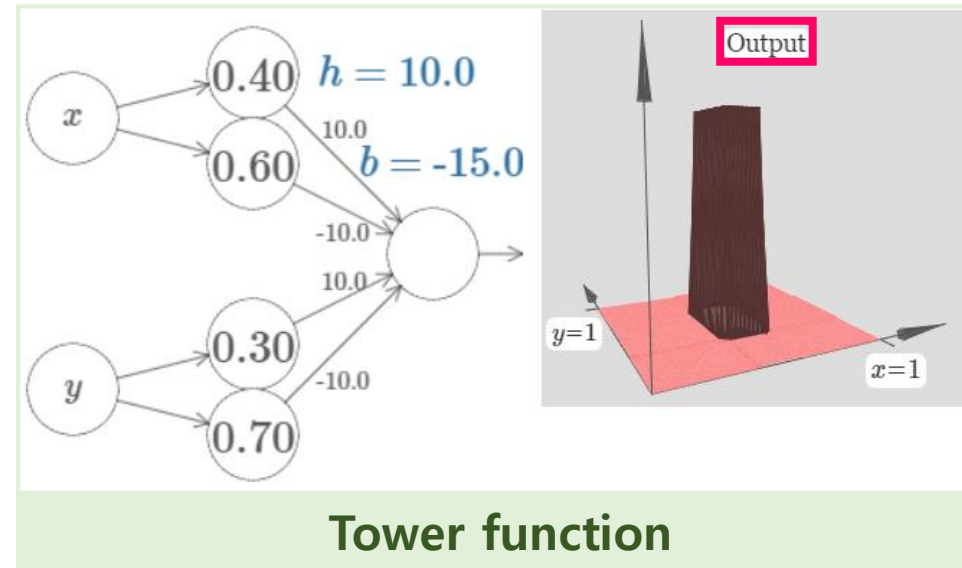As bias term decreases, graph for output from hidden neuron goes right on the x-axis without shape change
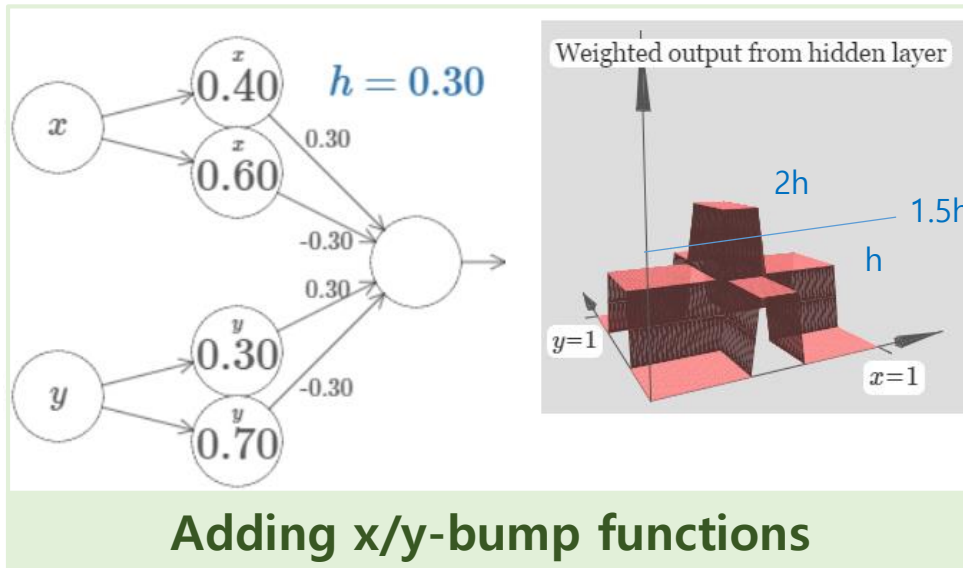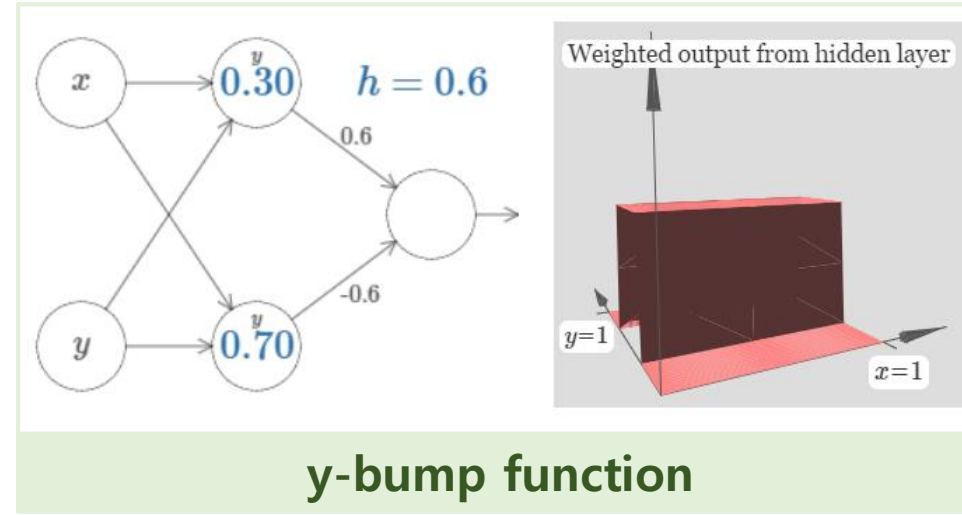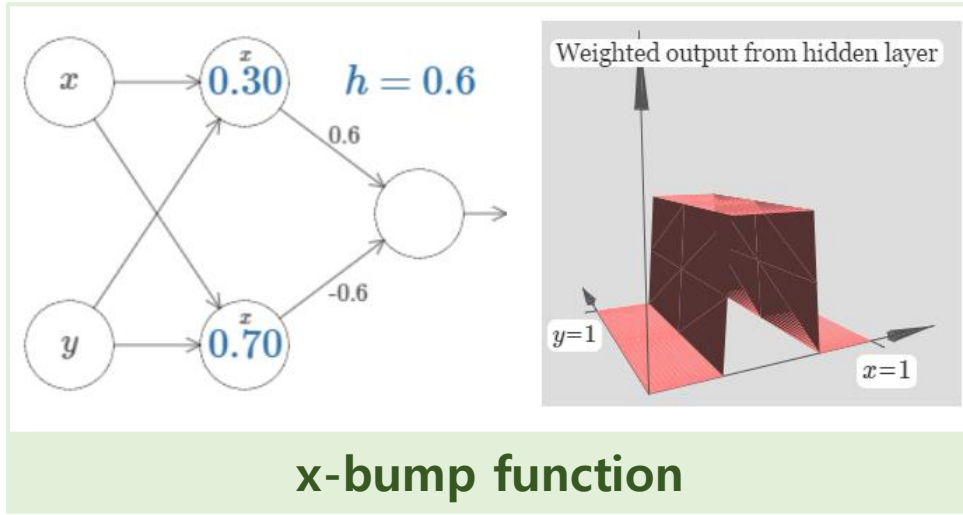
As weight term increases, graph for output from hidden neuron changes its shape.
( The curve gets steeper, until eventually it begins to look like a step function)

Same properties as in one input variable are observed!

### For very large $w_1$ and zero $w_2$



### For very large $w_2$ and zero $w_1$

# Many input variables

## Two input variables



**x-bump function**



**y-bump function**



**Adding x/y-bump functions**



**Tower function**
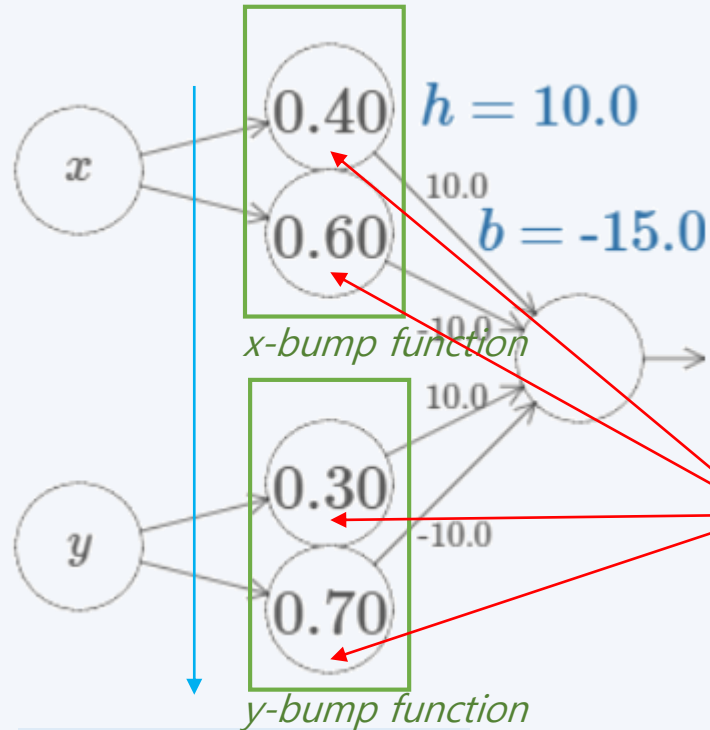
Output = σ(weighted output from hidden layer + b)
Tower function can be obtained from using threshold value of 1.5h. It means b should be -1.5h

(1) To get the output neuron to show the right kind of if-then-else behaviour, we need the input weights (all $h$ or $-h$) to be large

(2) the value of $b$ determines the scale of the if-then-else threshold.
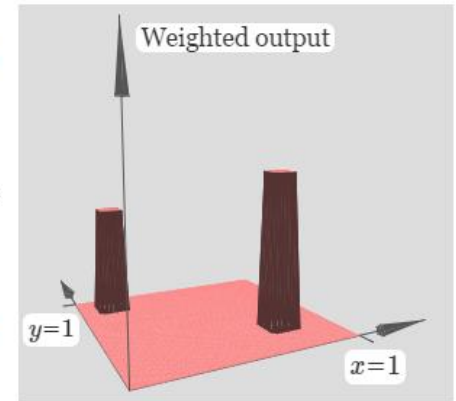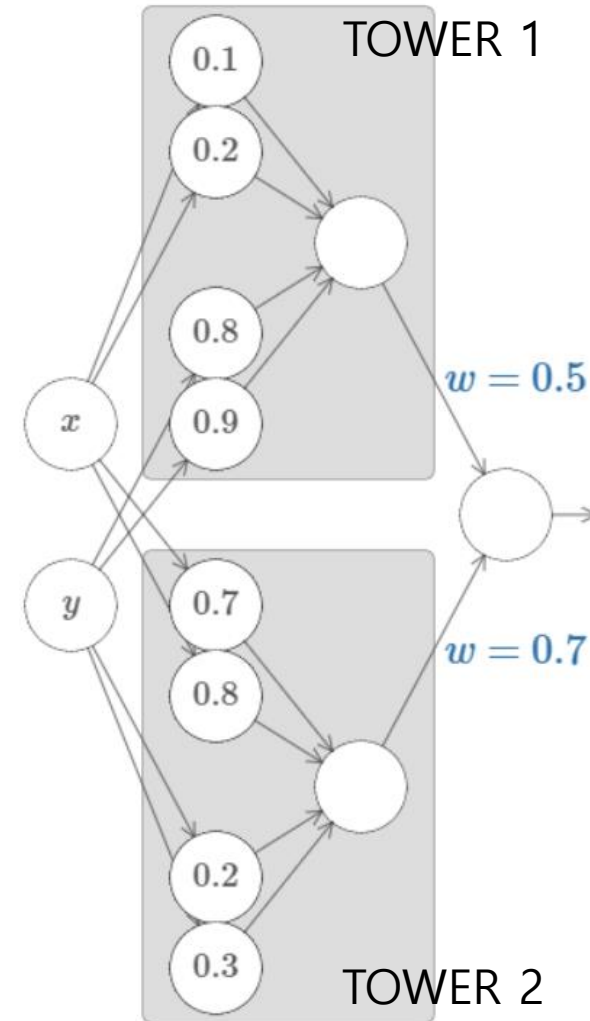
# Many input variables

**Two input variables**



-b/w determines location and width of bump functions

$h = 10.0$

$b = -15.0$

*x-bump function*

*y-bump function*

Large weights for step functions

*h* must be large and *b* is set to properly (-1.5h) for tower function

These define the position of tower function
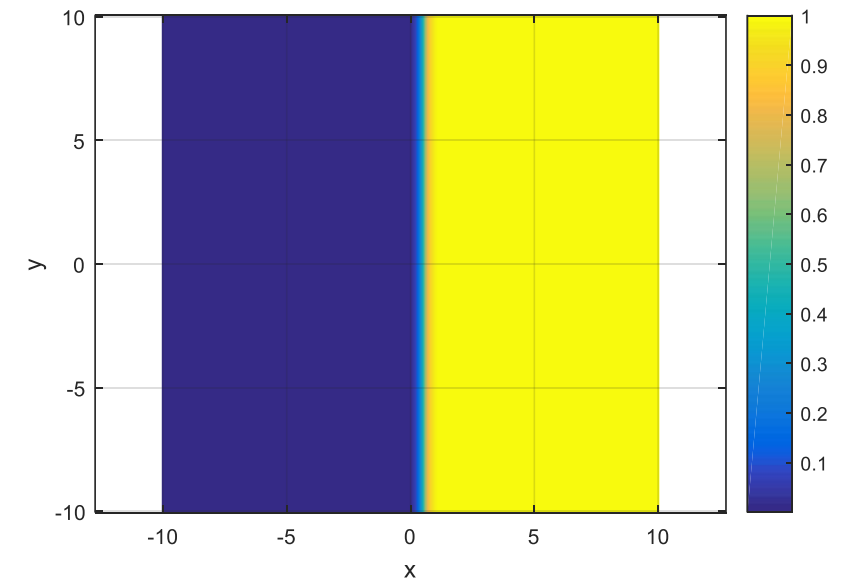
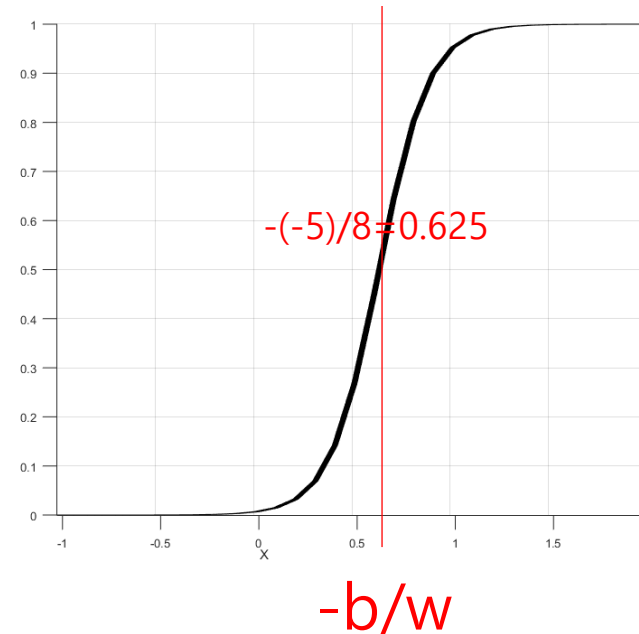TOWER 1

$w = 0.5$

$w = 0.7$

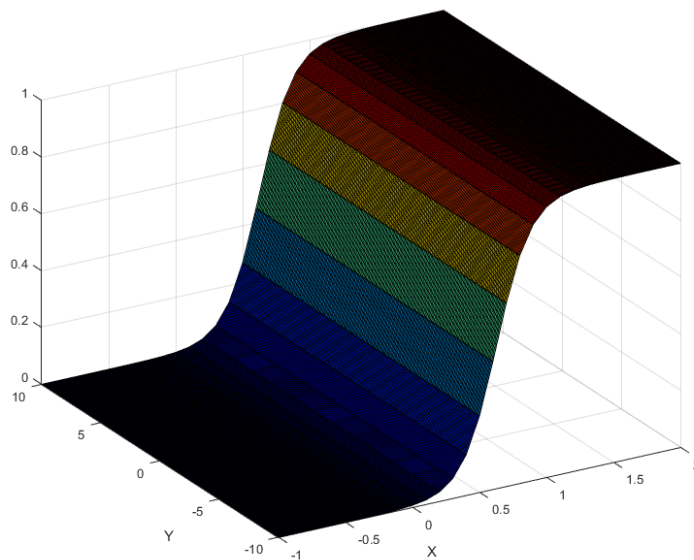Weighted output

TOWER 2

# Many input variables

## Two input variables



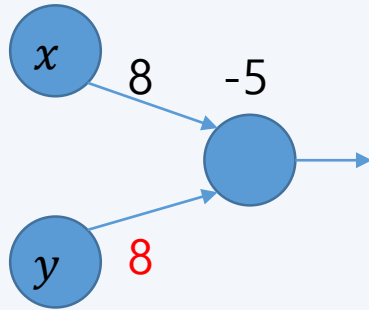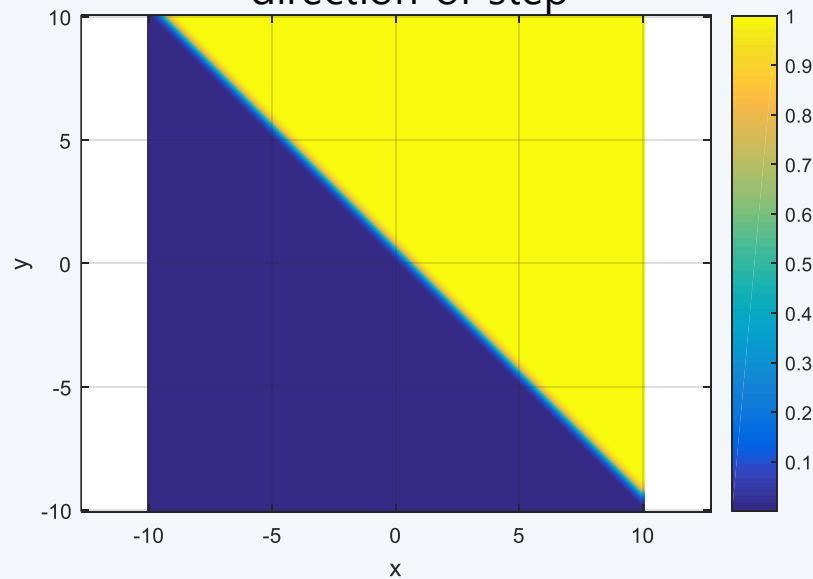We've dealt with one weight zero cases.



-(-5)/8=0.625

-b/w
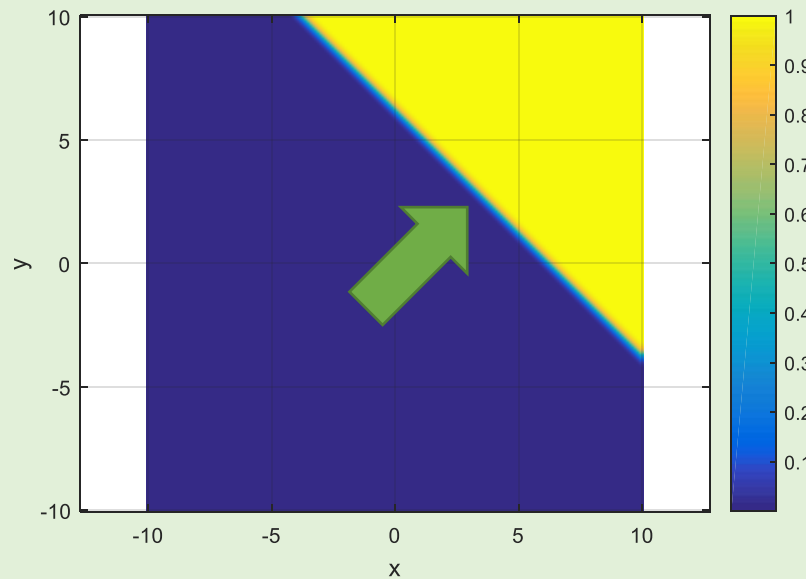
# Many input variables

## Two input variables

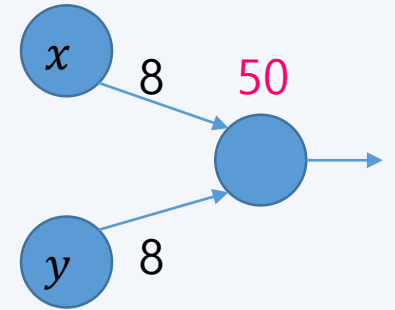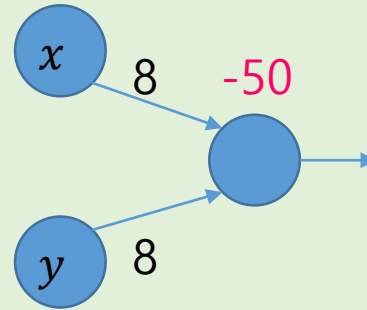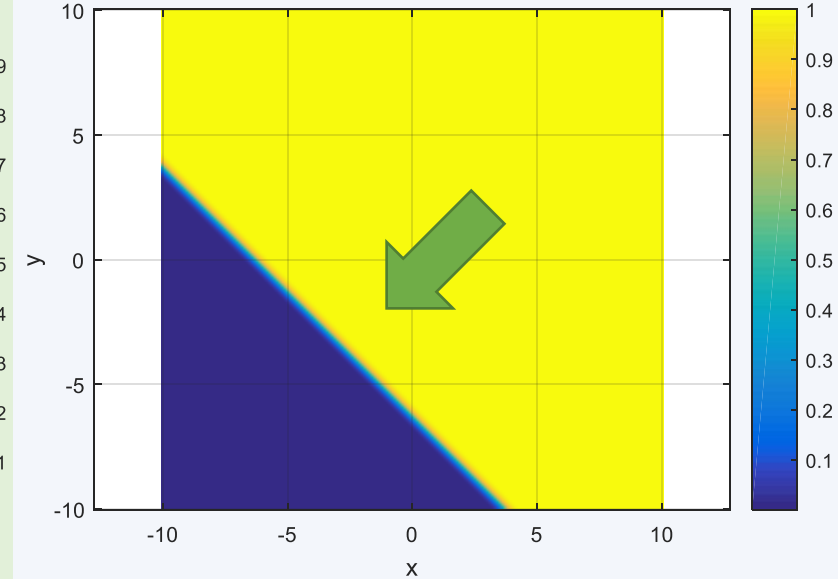If neither weights are zero, it gives rotation effect to step function



Rotation!! (ratio of weights defines the direction of step
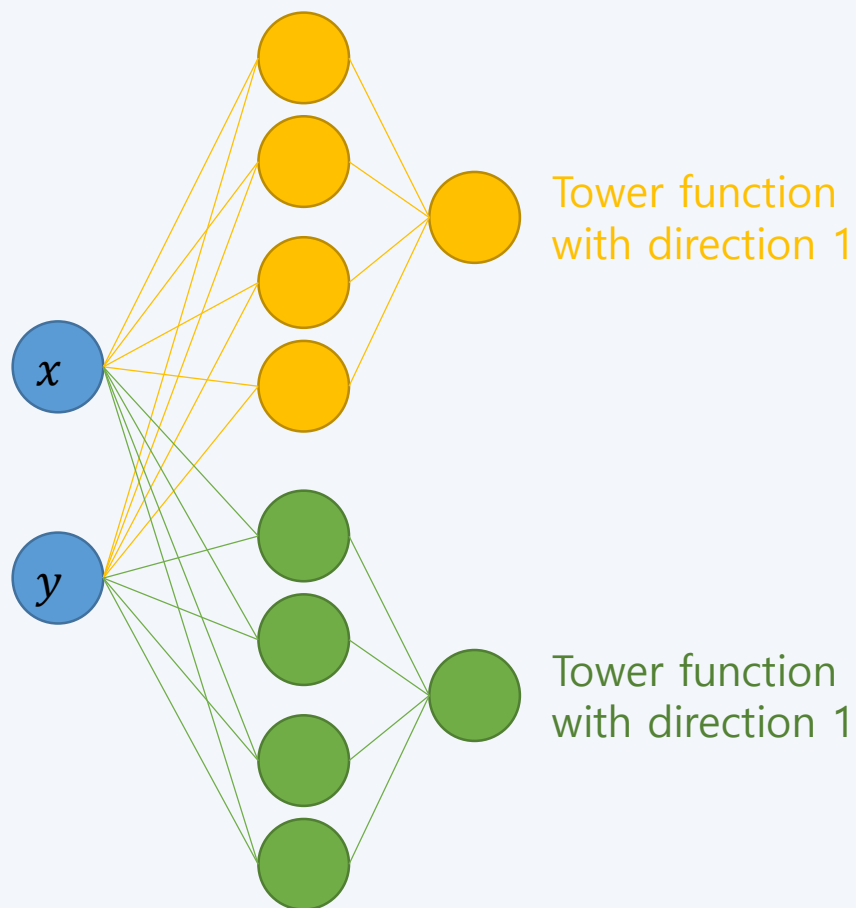


Shift up along the direction



Shift down along the direction

# Many input variables
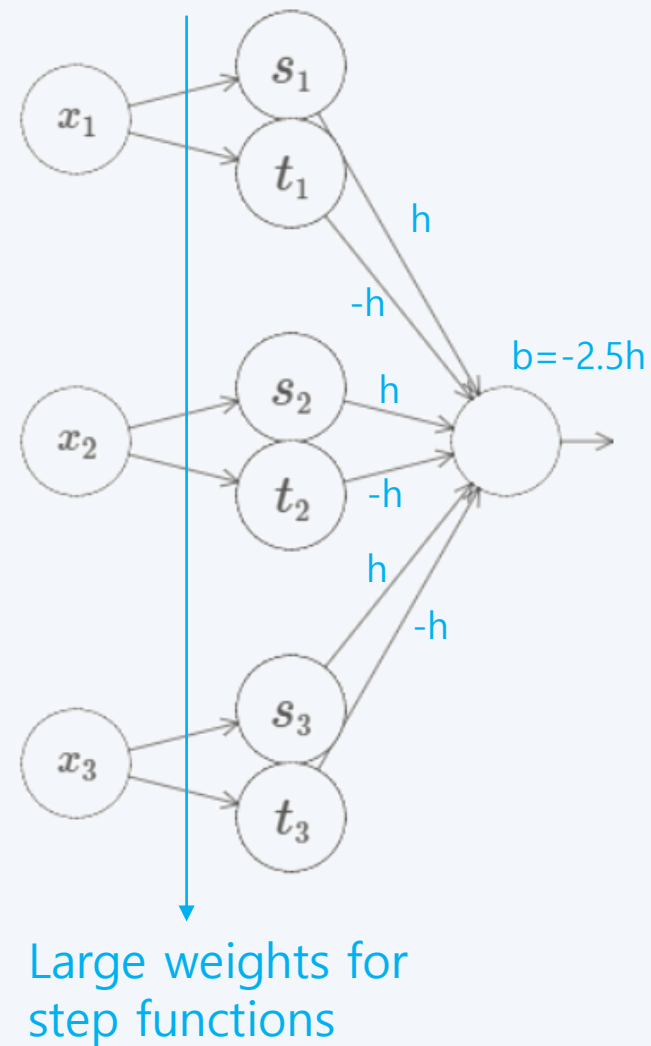


**Two** input variables

$x$

$y$

Tower function with direction 1

Tower function with direction 1

**Three** input variables

$s_1$

$x_1$

$t_1$

h

-h

$x_2$

$s_2$ h

$t_2$ -h

b=-2.5h

h

-h

$x_3$

$s_3$

$t_3$

Large weights for step functions

# Extension beyond sigmoid neurons

**Are rectified linear units universal for computation? YES!!**