# Deep Learning

Lecture 6

# Activation functions

Neural Network

Linear classifiers

Activation Functions
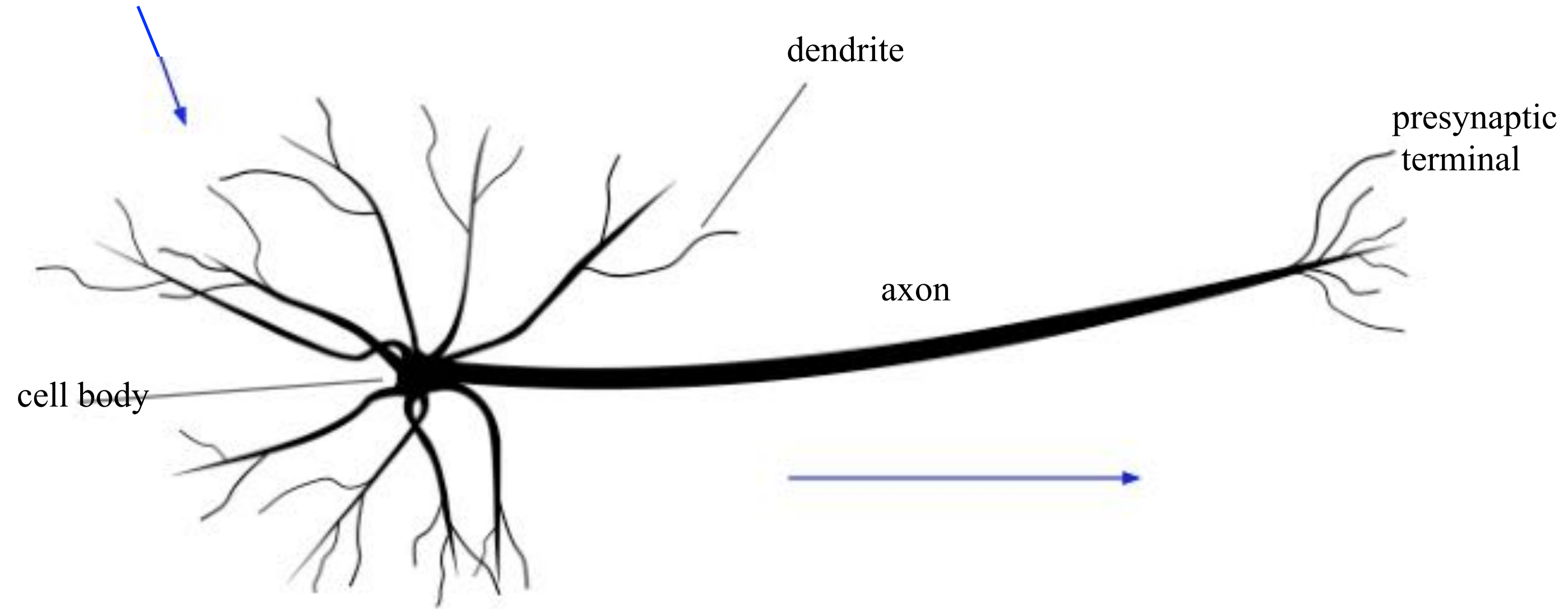
This image by Fotis Bobolas is
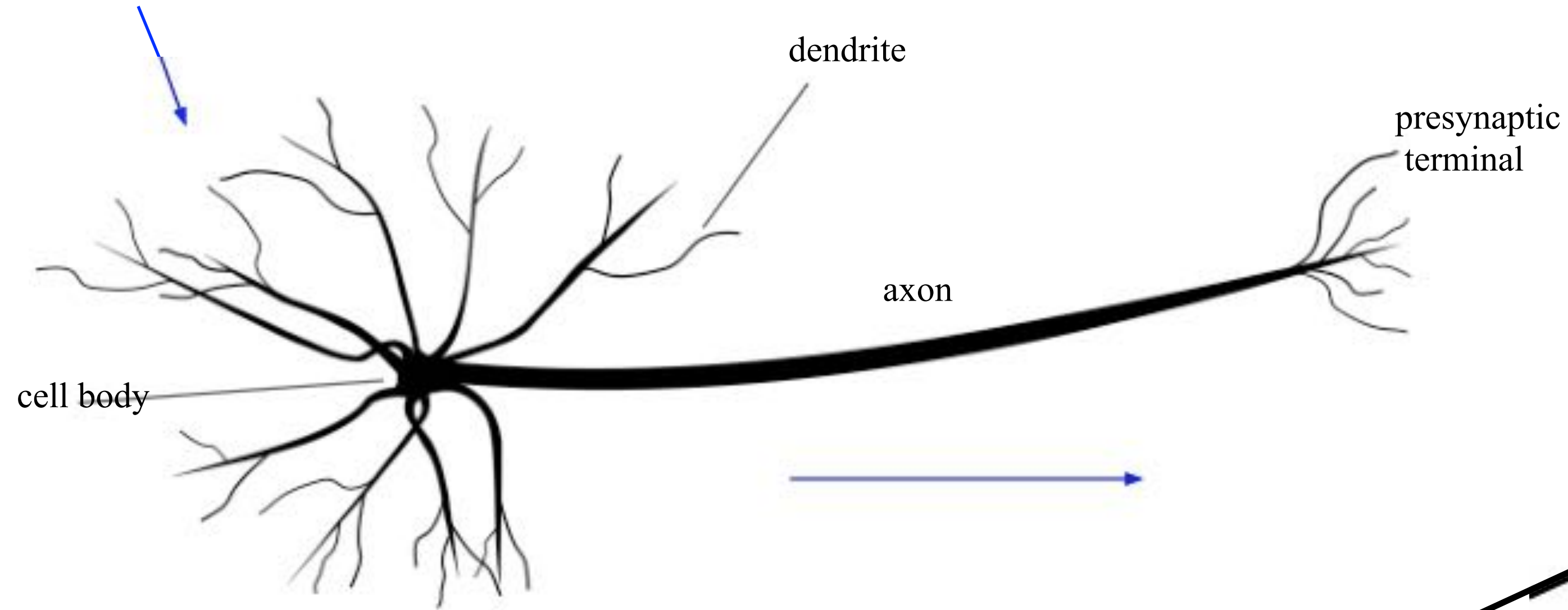licensed under CC-BY 2.0

Impulses carried toward cell body

dendrite

presynaptic
terminal

axon

cell body

Impulses carried toward cell body

dendrite

presynaptic
terminal

axon

cell body

$x_0$

$w_0$

axon from a neuron

synapse

$w_0 x_0$

dendrite

cell body

$w_1 x_1$

$\sum_i w_i x_i + b$  $f$

$f\left(\sum_i w_i x_i + b\right)$

output axon

activation
function

$w_2 x_2$

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

$$x_0$$

$$w_0$$

synapse

axon from a neuron

$$w_0 x_0$$

dendrite

$$w_1 x_1$$

cell body

$$\sum_i w_i x_i + b$$ $$f$$

$$f\left(\sum_i w_i x_i + b\right)$$

output axon

$$w_2 x_2$$

activation function

sigmoid activation function

$$\frac{1}{1+e^{-x}}$$

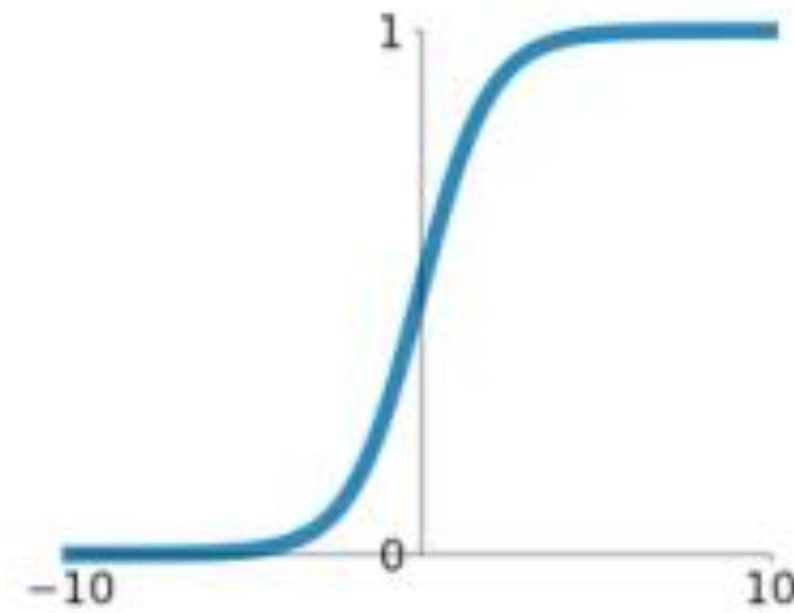# Be very careful with your brain analogies!

**Biological Neurons:**
- Many different types
- Dendrites can perform complex non-linear computations
- Synapses are not a single weight but a complex non-linear dynamical system
- Rate code may not be adequate

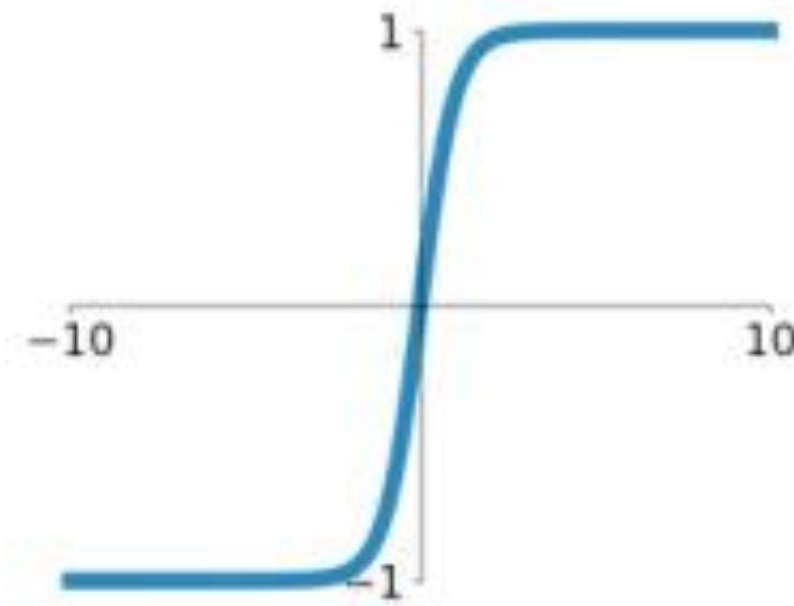[Dendritic Computation. London and Hausser]

# Activation functions

**Sigmoid**

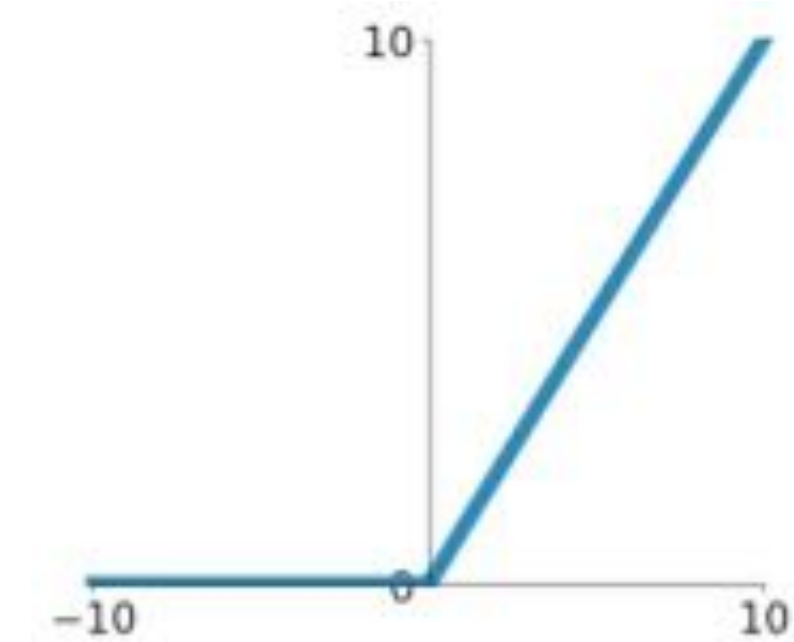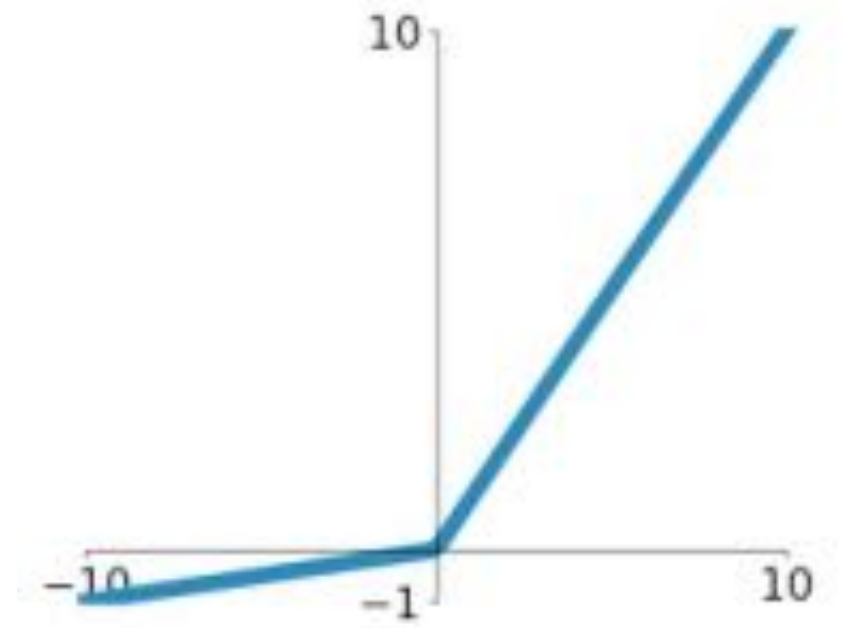$$\sigma(x) = \frac{1}{1+e^{-x}}$$



**tanh**

$$\tanh(x)$$



**ReLU**

$$\max(0, x)$$



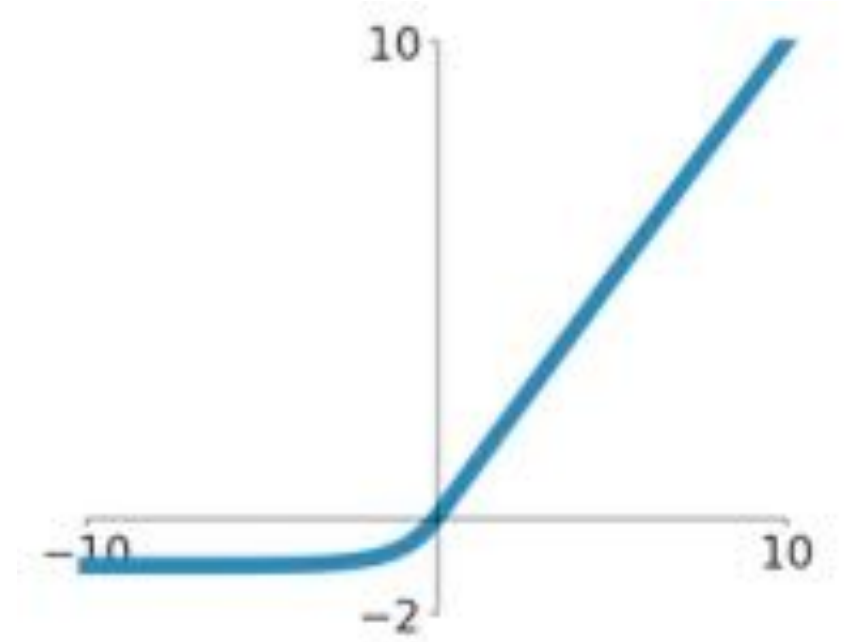**Leaky ReLU**

$$\max(0.1x, x)$$



**Maxout**

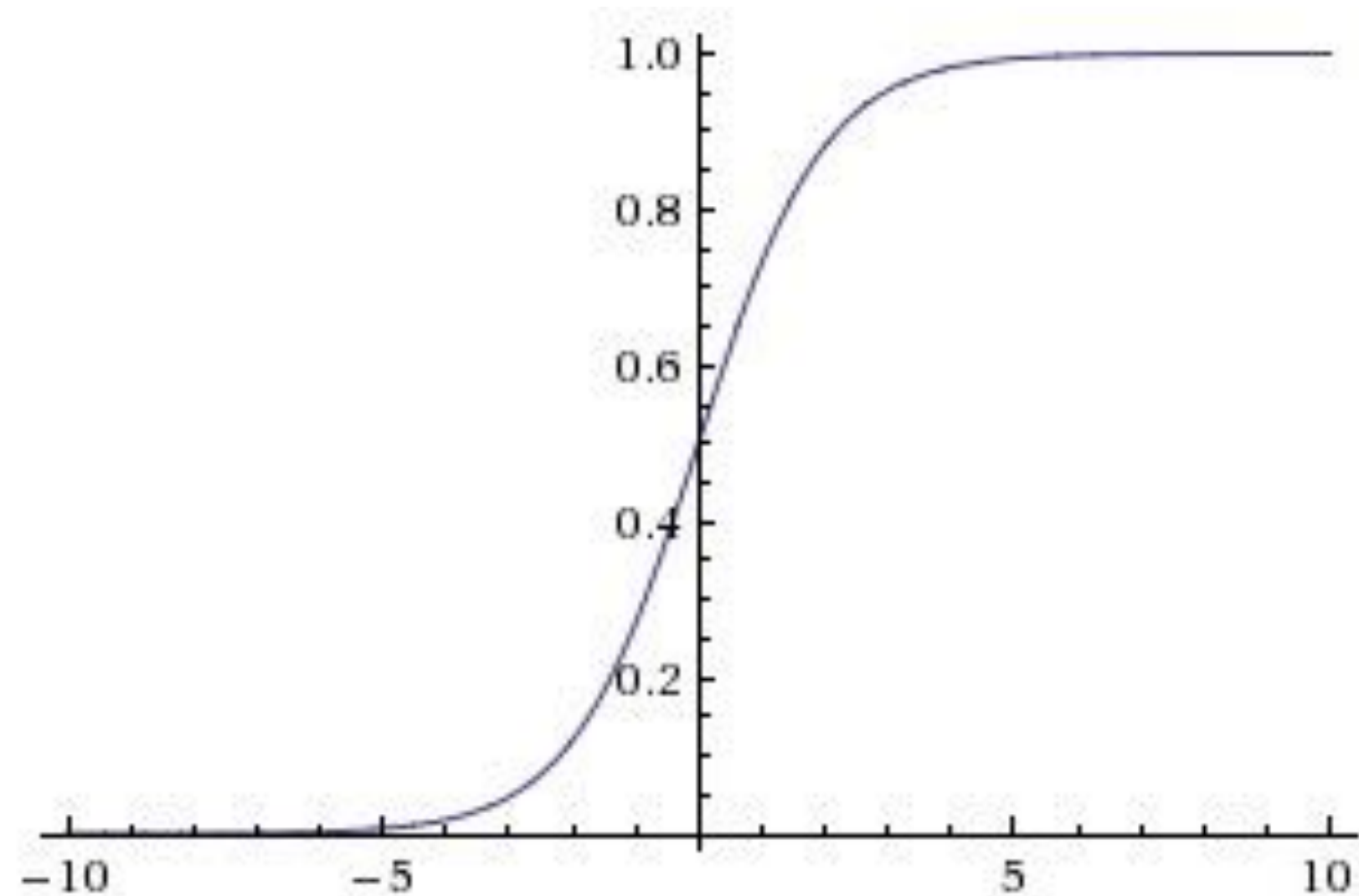$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

# Activation functions

$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron



**Sigmoid**

# Activation functions

$$\sigma(x) = 1/(1 + e^{-x})$$



**Sigmoid**

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

3 problems:

1. Saturated neurons "kill" the gradients

x

$$\frac{\partial \sigma}{\partial x}$$

sigmoid gate

$$\sigma(x) = 1/(1 + e^{-x})$$

$$\frac{\partial L}{\partial x} = \frac{\partial \sigma}{\partial x} \frac{\partial L}{\partial \sigma}$$

$$\frac{\partial L}{\partial \sigma}$$

What happens when x = -10?
What happens when x = 0?
What happens when x = 10?

# Activation functions

$$\sigma(x) = 1/(1 + e^{-x})$$



**Sigmoid**
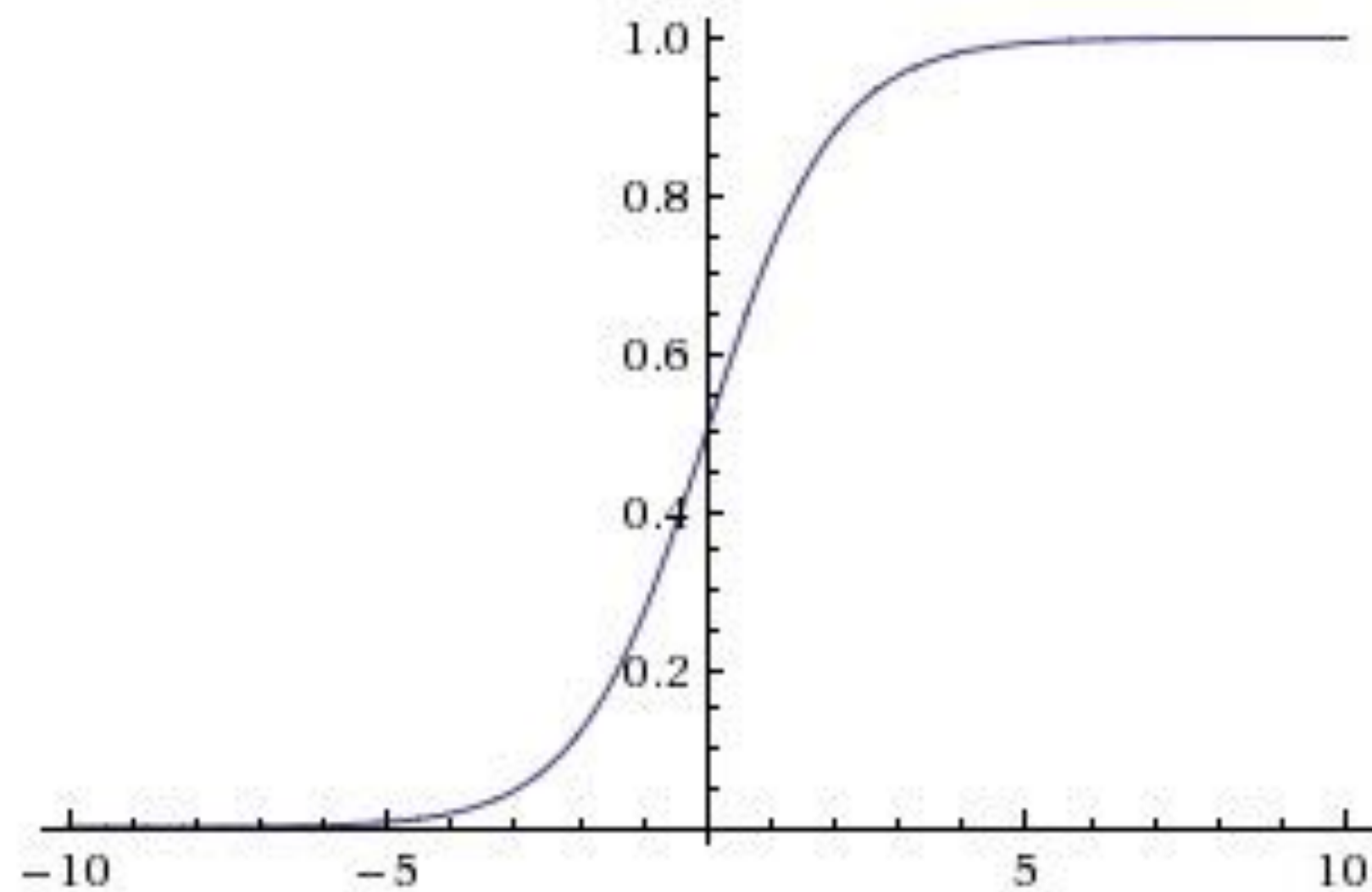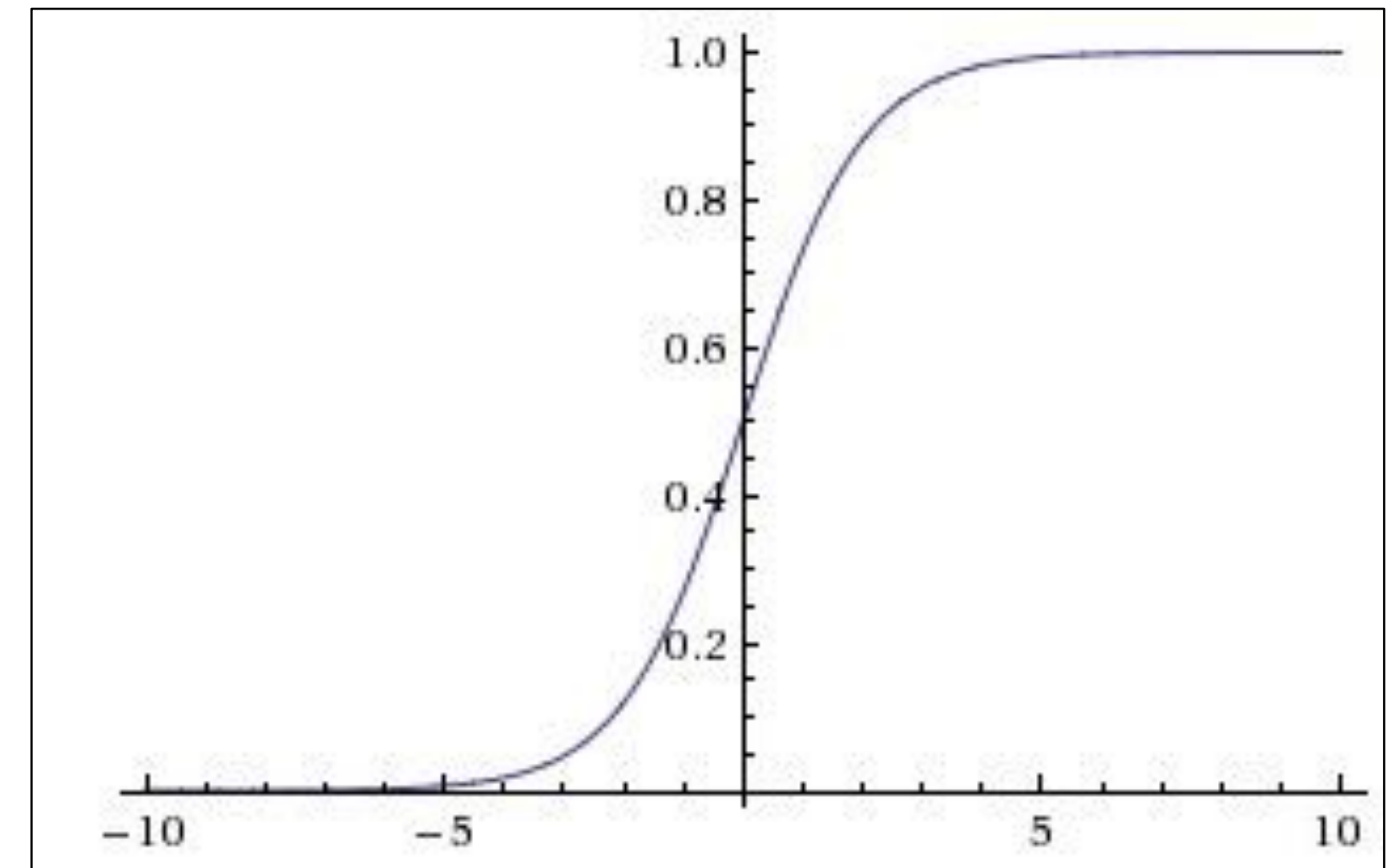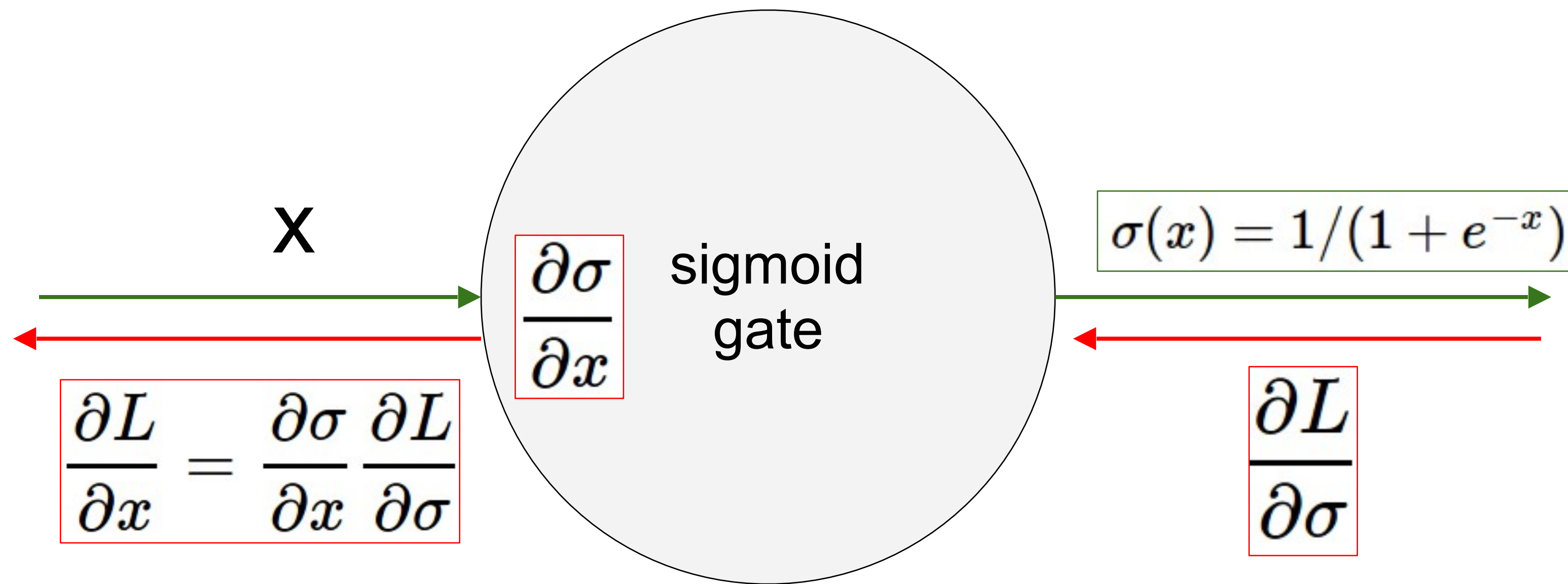
- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron
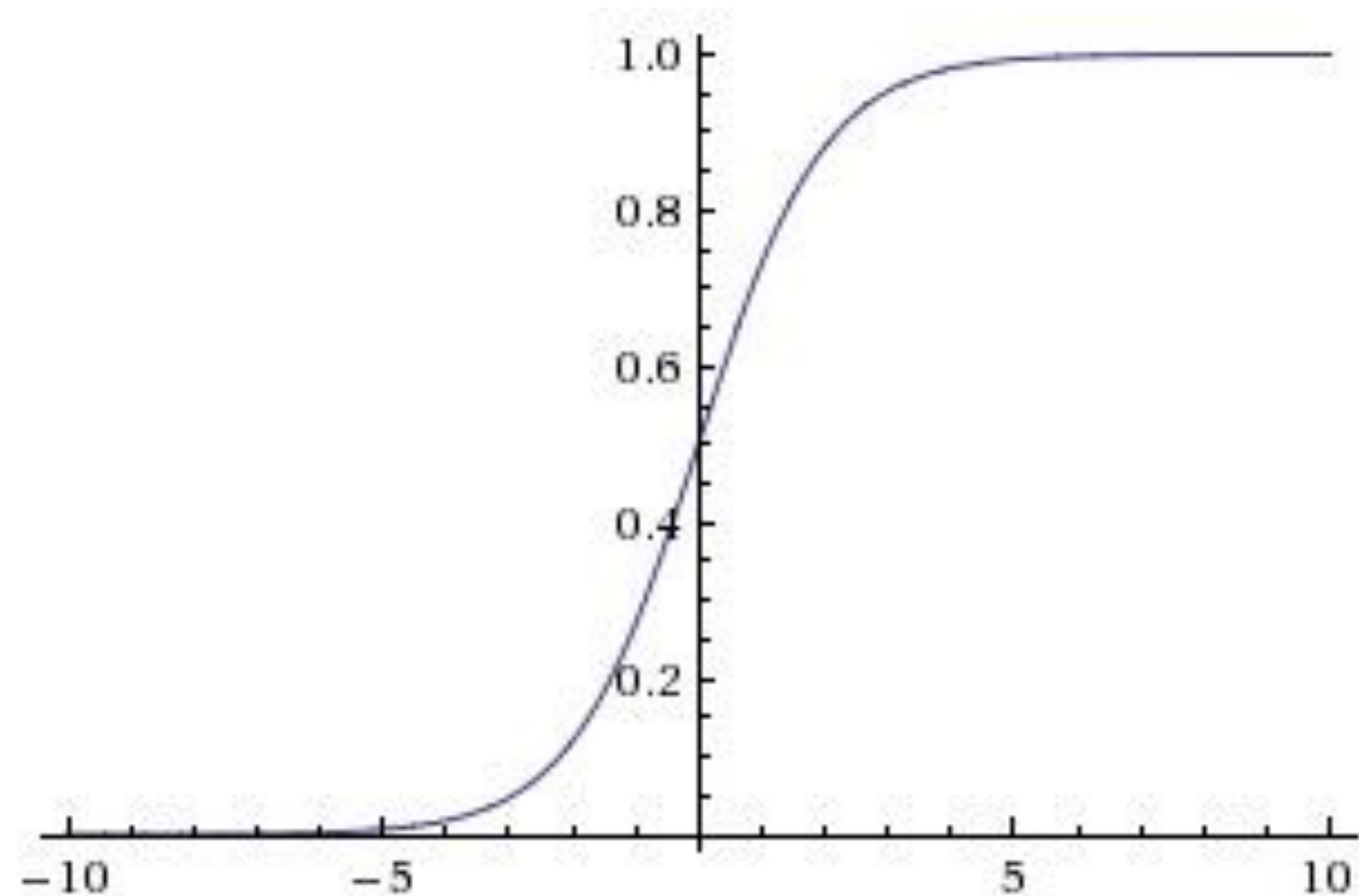
3 problems:

1. Saturated neurons "kill" the gradients
2. Sigmoid outputs are not zero-centered

Consider what happens when the input to a neuron (x)
is always positive:



$$f\left(\sum_i w_i x_i + b\right)$$

What can we say about the gradients on **w**?

Consider what happens when the input to a neuron is always positive...

$$f\left(\sum_i w_i x_i + b\right)$$

allowed gradient update directions

allowed gradient update directions

zig zag path

hypothetical optimal w vector

What can we say about the gradients on **w**?
Always all positive or all negative :(
(this is also why you want zero-mean data!)

# Activation functions

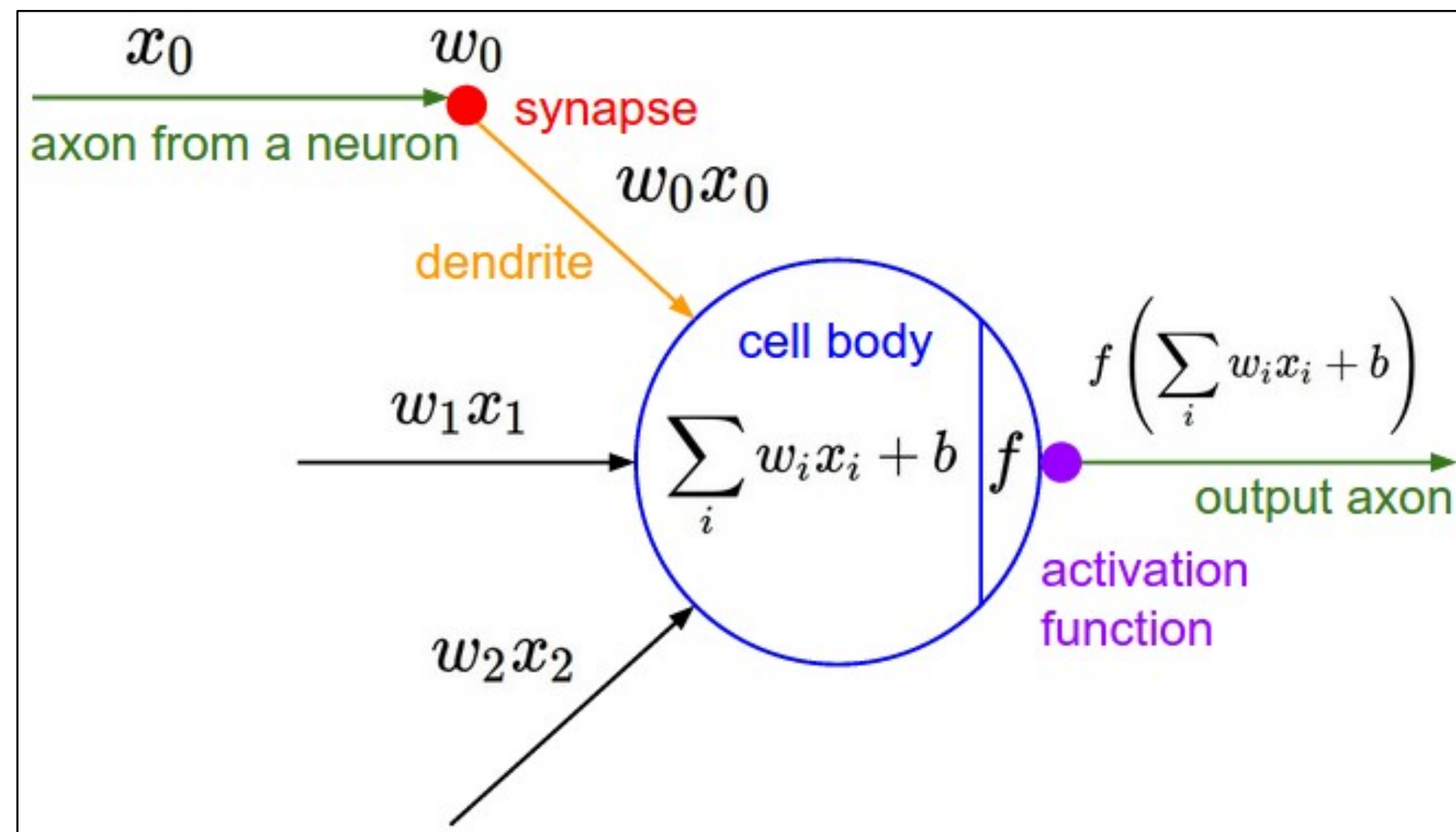$$\sigma(x) = 1/(1 + e^{-x})$$



**Sigmoid**

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

3 problems:

1. Saturated neurons "kill" the gradients
2. Sigmoid outputs are not zero-centered
3. exp() is a bit compute expensive

# Activation functions



**tanh(x)**

- Squashes numbers to range [-1,1]
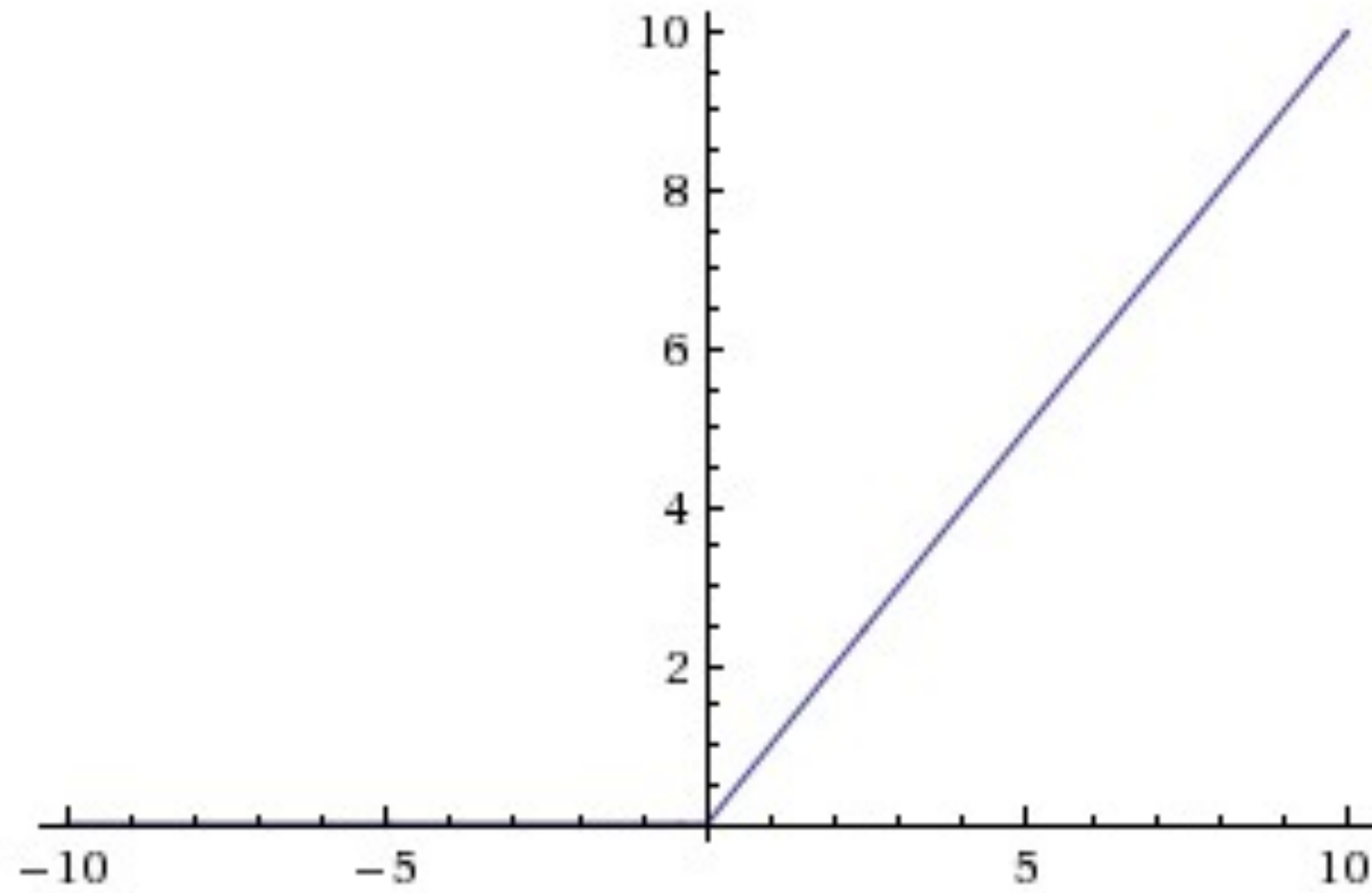- zero centered (nice)
- still kills gradients when saturated :(

[LeCun et al., 1991]

# Activation functions

- Computes **f(x) = max(0,x)**

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)



**ReLU**
(Rectified Linear Unit)

[Krizhevsky et al., 2012]

# Activation functions



**ReLU**
(Rectified Linear Unit)

- Computes **f(x) = max(0,x)**

- Does not saturate (in +region)
- Very computationally efficient
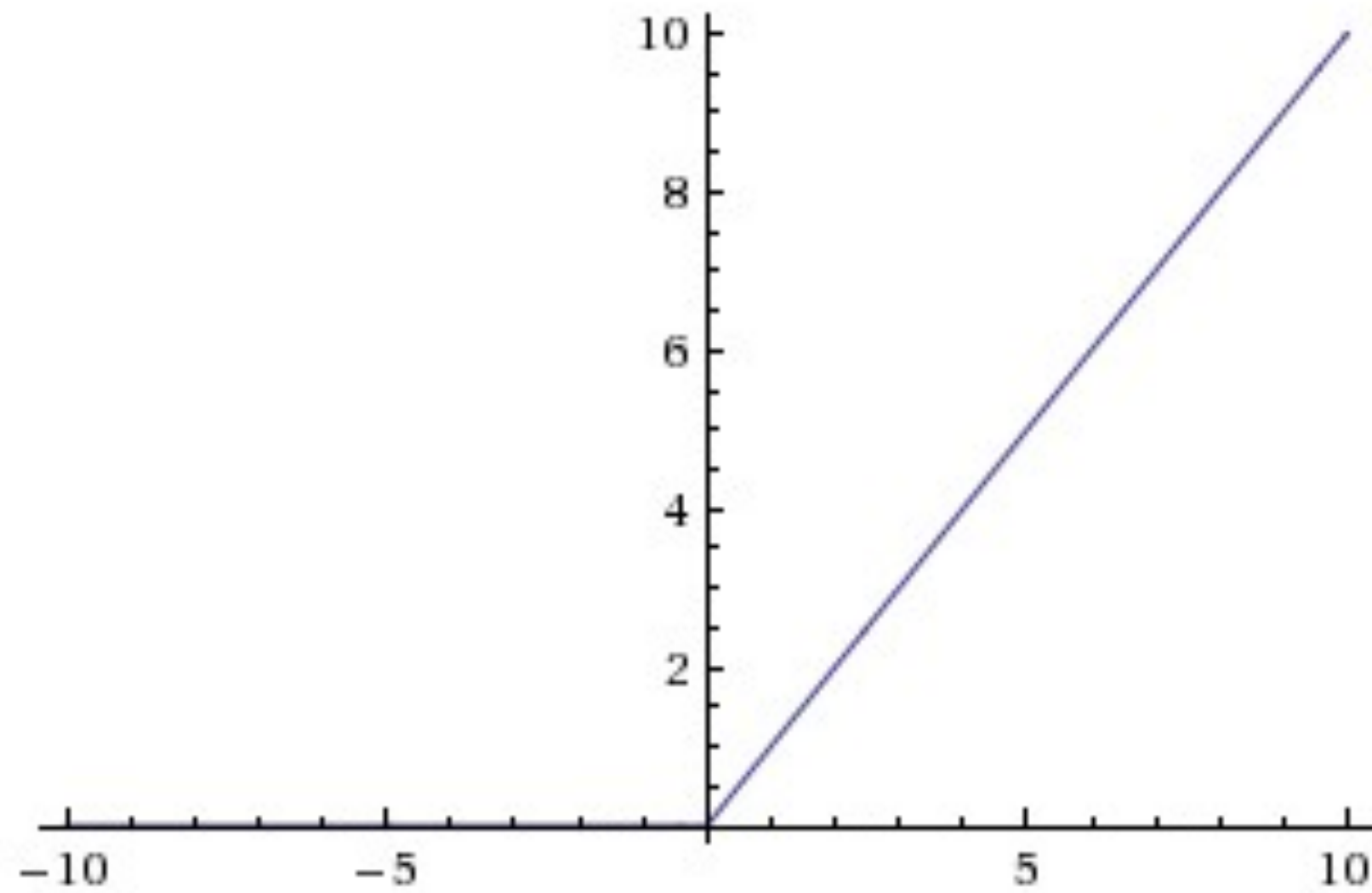- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

- Not zero-centered output

- An annoyance:

hint: what is the gradient when x < 0?

x

$$\frac{\partial \sigma}{\partial x}$$

ReLU gate

$$\sigma(x) = \max(0, x)$$

$$\frac{\partial L}{\partial x} = \frac{\partial \sigma}{\partial x}\frac{\partial L}{\partial \sigma}$$

$$\frac{\partial L}{\partial \sigma}$$



What happens when x = -10?
What happens when x = 0?
What happens when x = 10?

**DATA CLOUD**

active ReLU

dead ReLU
will never activate
=> never update

**DATA CLOUD**

active ReLU

dead ReLU
will never activate
=> never update

=> people like to initialize ReLU neurons with slightly positive biases (e.g. 0.01)

# Activation functions

[Mass et al., 2013]
[He et al., 2015]



- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not "die".**

**Leaky ReLU**

$$f(x) = \max(0.01x, x)$$

# Activation functions

[Mass et al., 2013]
[He et al., 2015]



- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not "die".**

**Leaky ReLU**

$$f(x) = \max(0.01x, x)$$
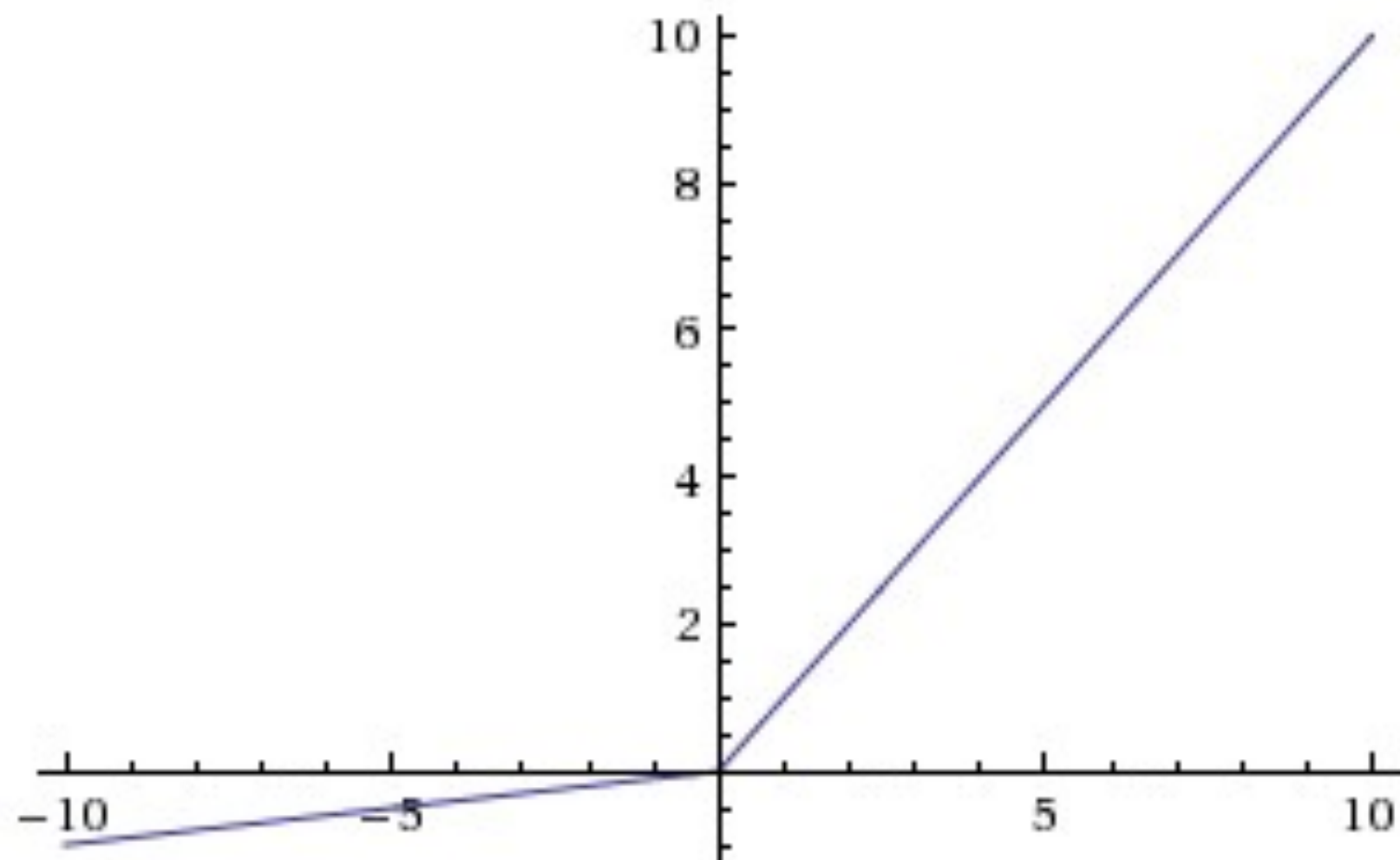
**Parametric Rectifier (PReLU)**

$$f(x) = \max(\alpha x, x)$$

backprop into \alpha (parameter)

# Activation functions

## Exponential Linear Units (ELU)



- All benefits of ReLU
- Does not die
- Closer to zero mean outputs

- Computation requires exp()

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha \left( \exp(x) - 1 \right) & \text{if } x \leq 0 \end{cases}$$

Used with permission from David Wingate

# Maxout "neuron"

- Does not have the basic form of dot product -> nonlinearity
- Generalizes ReLU and Leaky ReLU
- Linear Regime! Does not saturate! Does not die!

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

Problem: doubles the number of parameters/neuron :(

# TLDR: In practice, you should

- Use ReLU. Be careful with your learning rates
- Try out Leaky ReLU / Maxout / ELU
- Try out tanh but don't expect much
- Don't use sigmoid

# Fully-connected Neural Networks

# Neural networks: without the brain stuff

( **Before** ) Linear score function:   $f = Wx$

# Neural networks: without the brain stuff

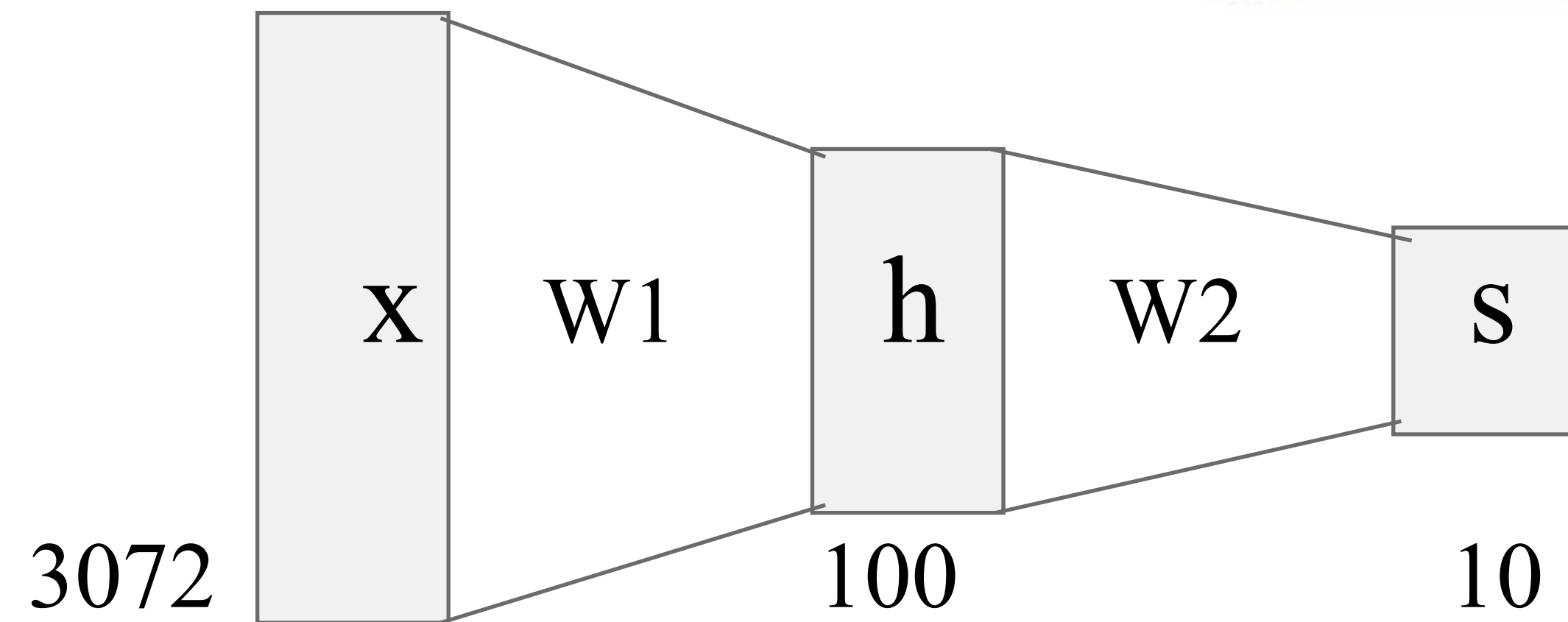( **Before** ) Linear score function:     $f = Wx$

( **Now** ) 2-layer Neural Network     $f = W_2 \max(0, W_1 x)$

# Neural networks: without the brain stuff

( **Before** ) Linear score function:

$$f = Wx$$

( **Now** ) 2-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$



3072    x  W1    h  W2    s    100    10

# Neural networks: without the brain stuff

( **Before** ) Linear score function:

$$f = Wx$$

( **Now** ) 2-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$



3072      x   W1     h   W2     s     100     10

| plane | car | bird | cat | deer | dog | frog | horse | ship | truck |

# Neural networks: without the brain stuff

( **Before** ) Linear score function:     $f = Wx$

( **Now** ) 2-layer Neural Network     $f = W_2 \max(0, W_1 x)$
      or 3-layer Neural Network

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

# Full implementation of training a 2-layer Neural Network needs ~20 lines:

```python
import numpy as np
from numpy.random import randn


N, D_in, H, D_out = 64, 1000, 100, 10
x, y = randn(N, D_in), randn(N, D_out)
w1, w2 = randn(D_in, H), randn(H, D_out)


for t in range(2000):
  h = 1 / (1 + np.exp(-x.dot(w1)))
  y_pred = h.dot(w2)
  loss = np.square(y_pred - y).sum()
  print(t, loss)


  grad_y_pred = 2.0 * (y_pred - y)
  grad_w2 = h.T.dot(grad_y_pred)
  grad_h = grad_y_pred.dot(w2.T)
  grad_w1 = x.T.dot(grad_h * h * (1 - h))


  w1 -= 1e-4 * grad_w1
  w2 -= 1e-4 * grad_w2
```

# In HW: Writing a 2-layer net

```python
# receive W1,W2,b1,b2 (weights/biases), X (data)
# forward pass:
h1 = #... function of X,W1,b1
scores = #... function of h1,W2,b2
loss = #... (several lines of code to evaluate Softmax loss)
# backward pass:
dscores = #...
dh1,dW2,db2 = #...
dW1,db1 = #...
```
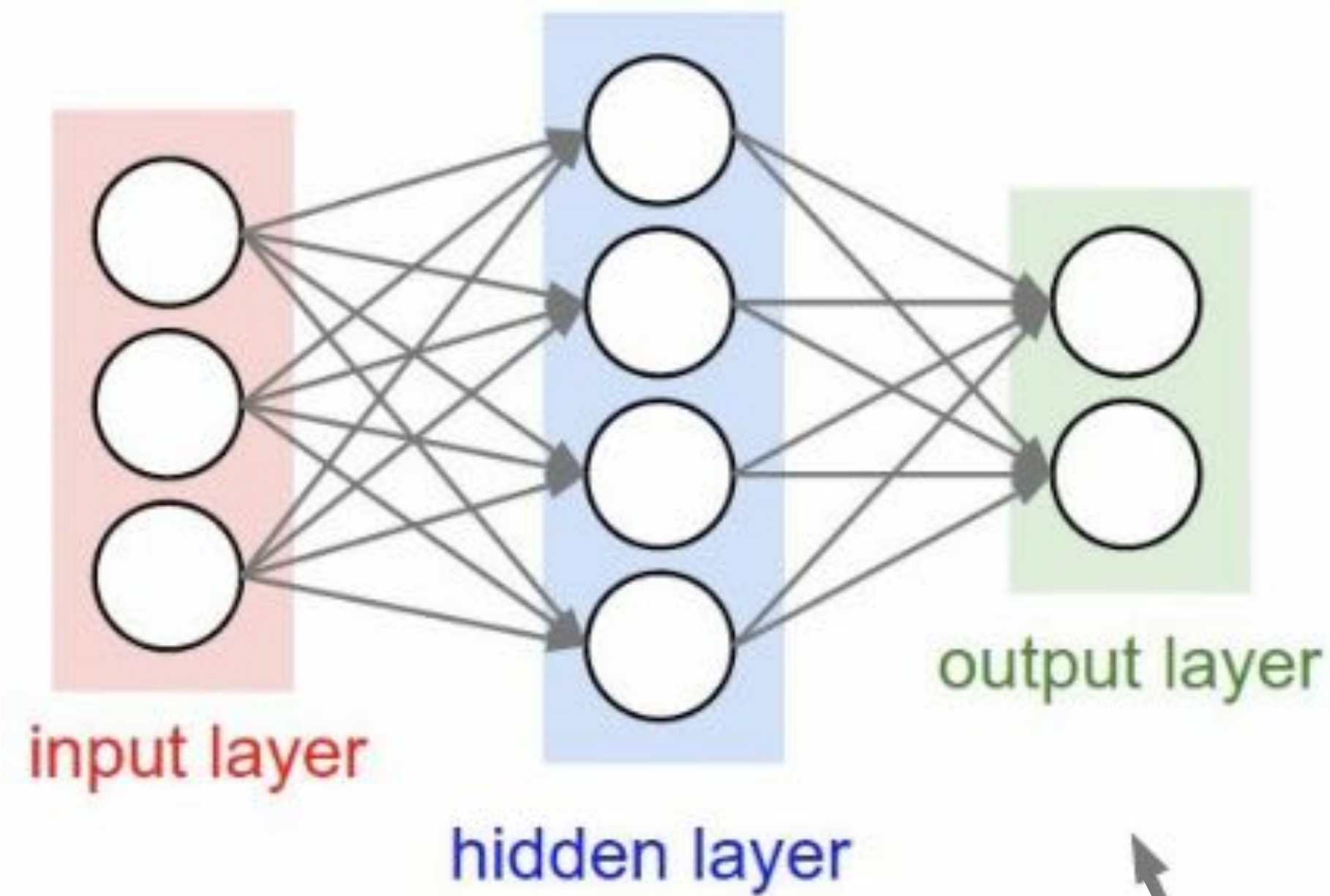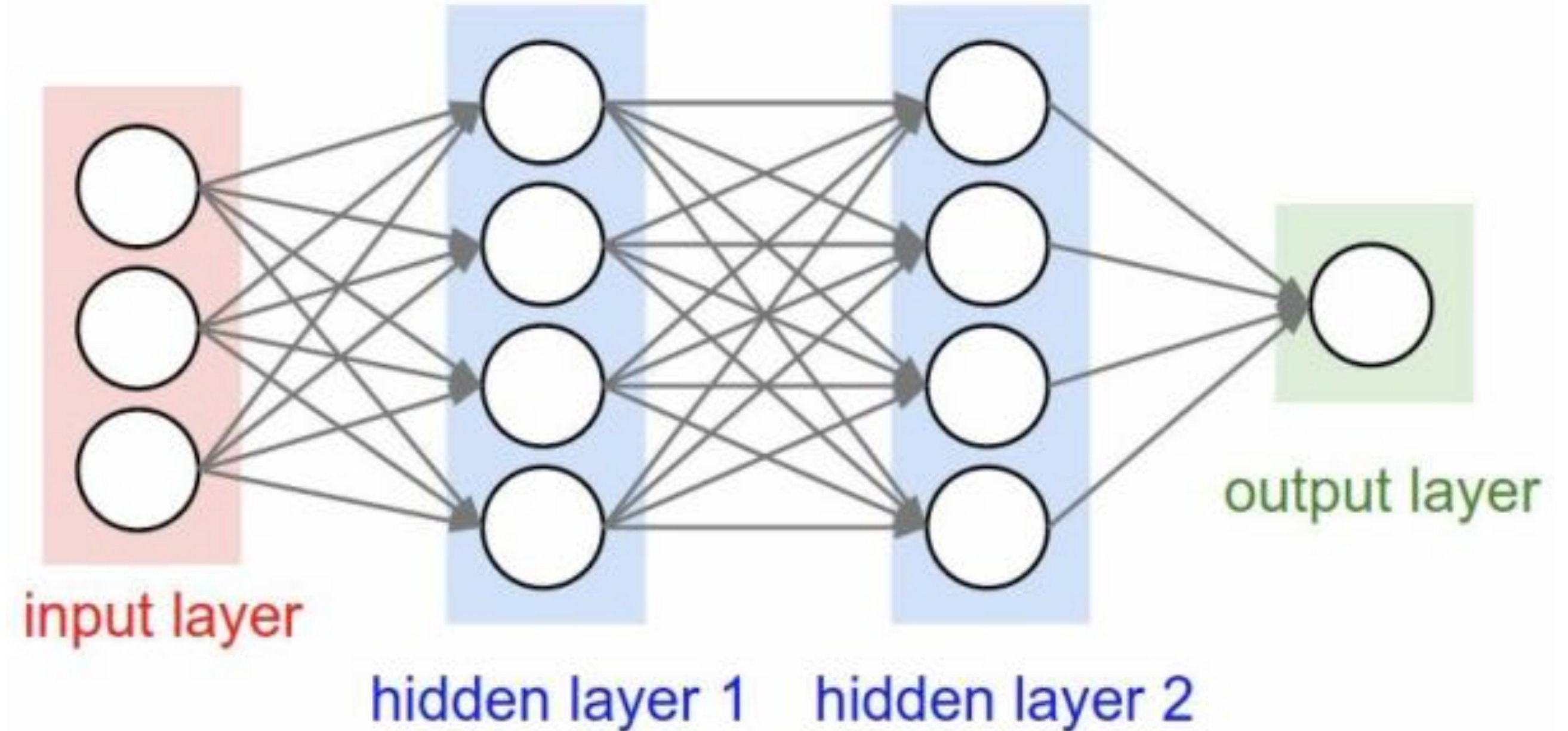
# Neural networks: Architectures



"2-layer Neural Net", or
"1-hidden-layer Neural Net"

**"Fully-connected" layers**

"3-layer Neural Net", or
"2-hidden-layer Neural Net"
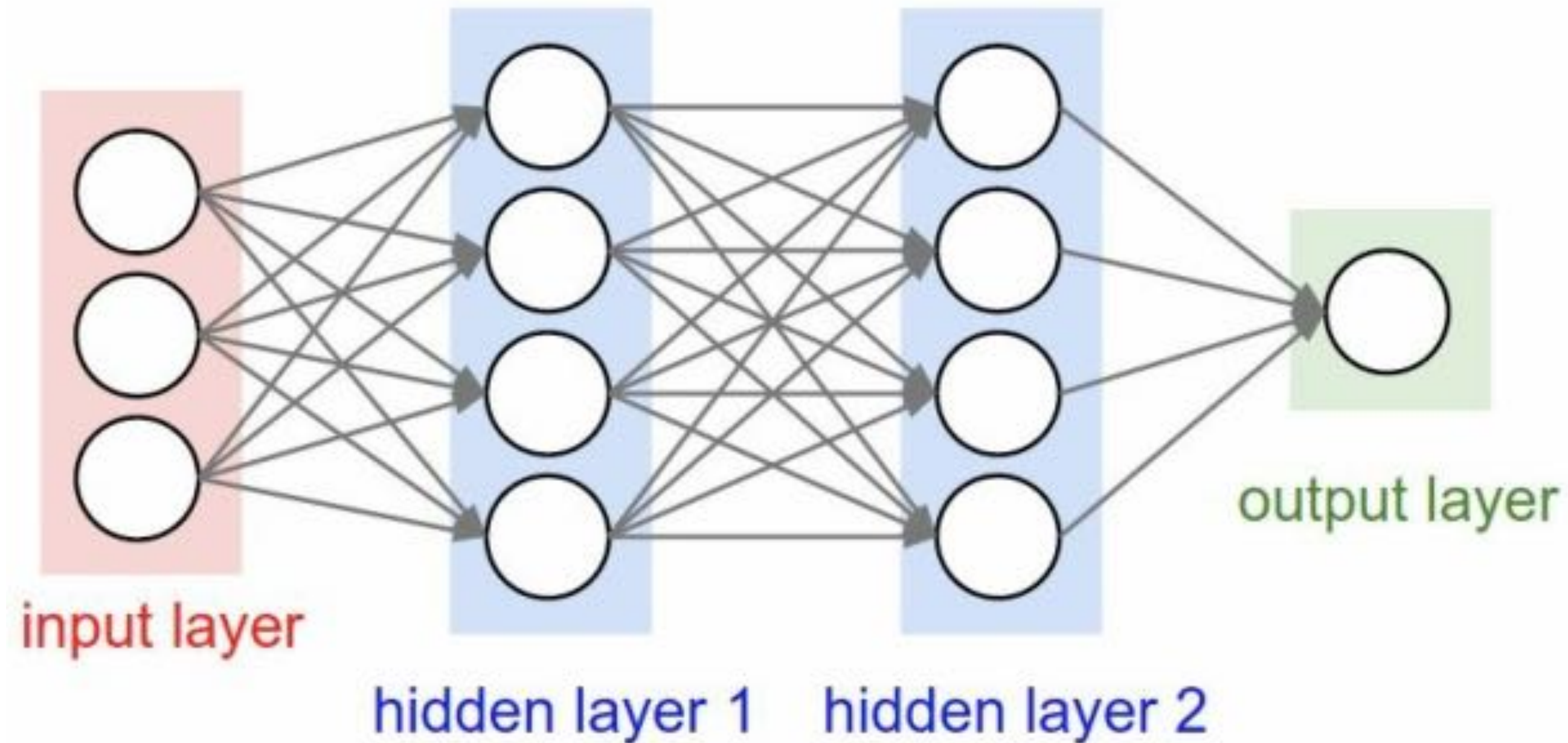
# Example feed-forward computation of a neural network

```python
class Neuron:
    # ...
    def neuron_tick(inputs):
        """ assume inputs and weights are 1-D numpy arrays and bias is a number """
        cell_body_sum = np.sum(inputs * self.weights) + self.bias
        firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation function
        return firing_rate
```

We can efficiently evaluate an entire layer of neurons.

# Example feed-forward computation of a neural network



input layer
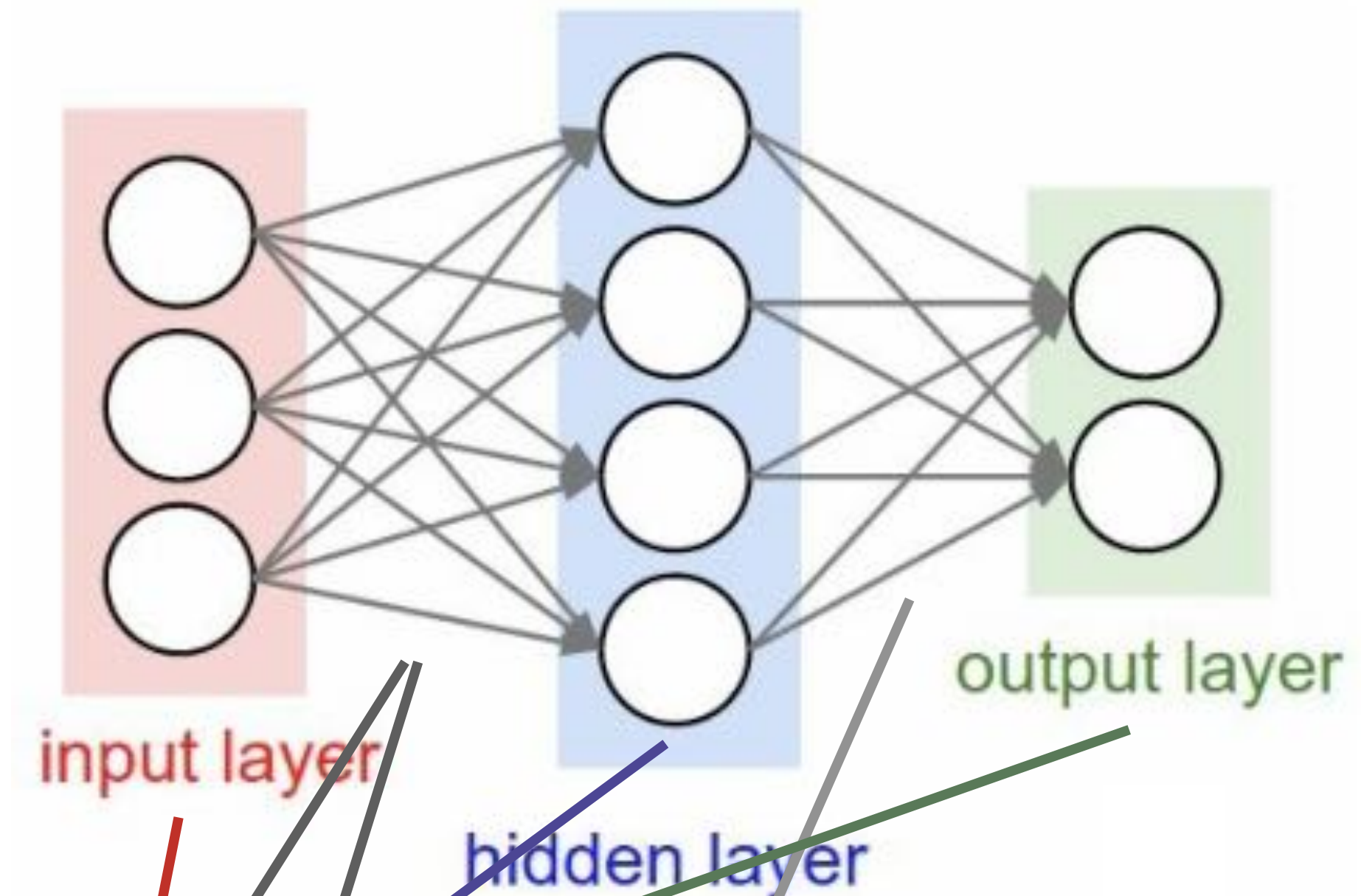hidden layer 1    hidden layer 2
output layer

```
# forward-pass of a 3-layer neural network:
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

# Universal approximation theorem (Cybenko 1989)



**Definition.** We say that $\sigma$ is *sigmoidal* if

$$\sigma(t) \to \begin{cases} 1 & \text{as} \quad t \to +\infty, \\ 0 & \text{as} \quad t \to -\infty. \end{cases}$$

**Theorem 2.** *Let $\sigma$ be any continuous sigmoidal function. Then finite sums of the form*
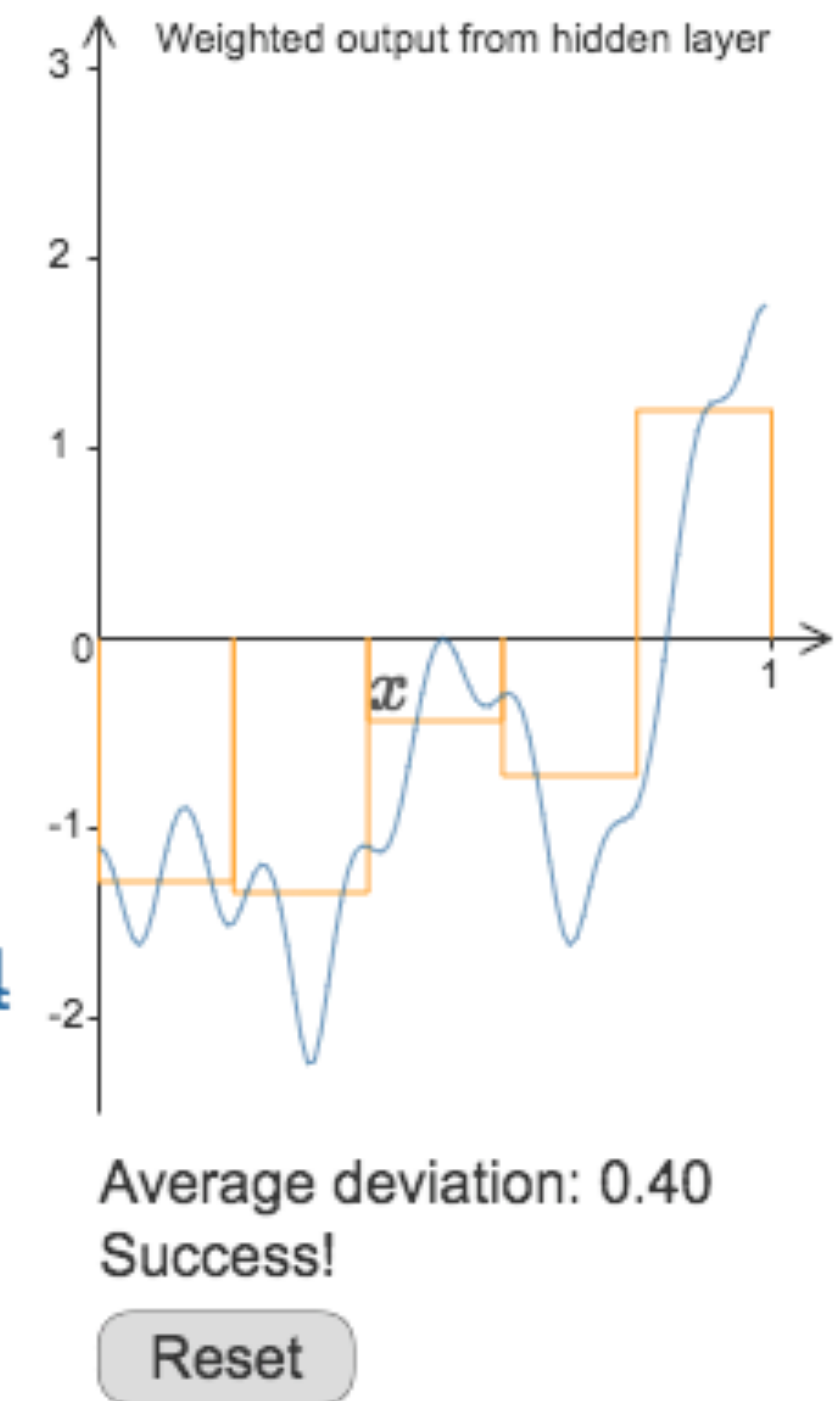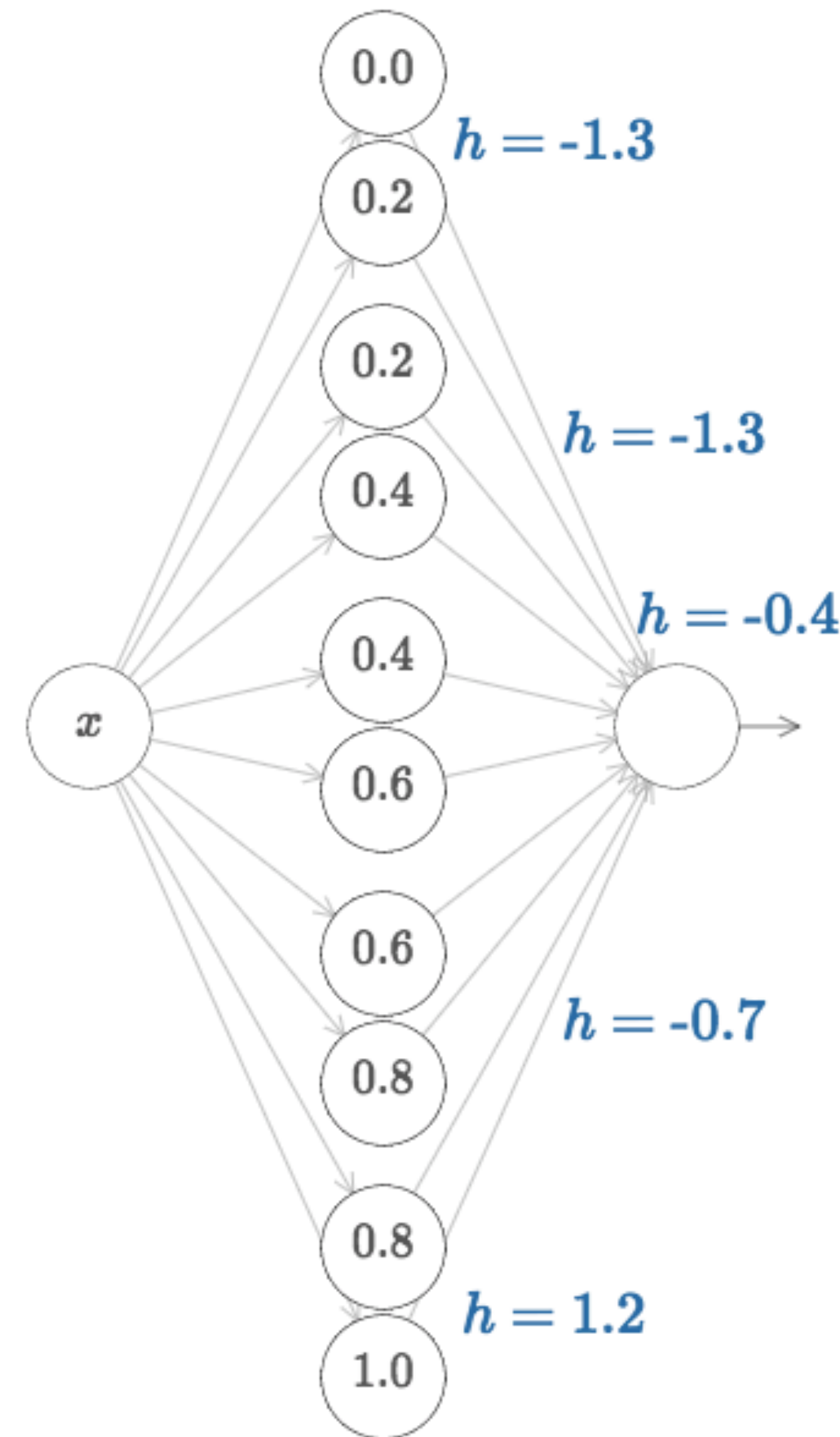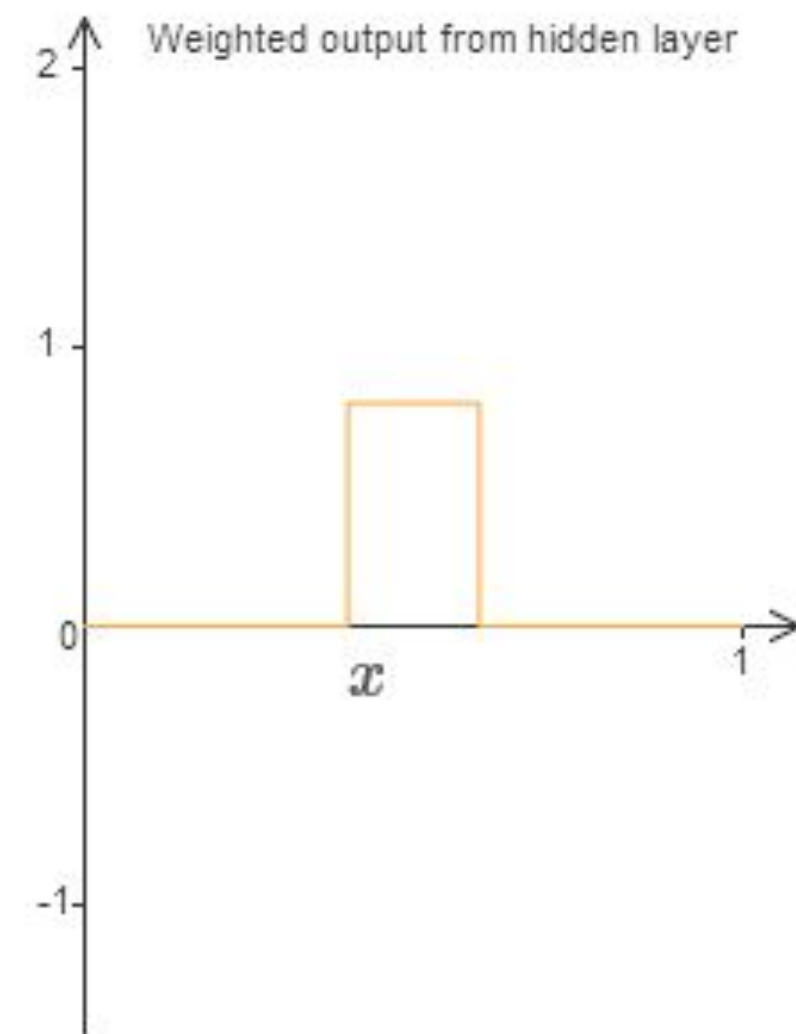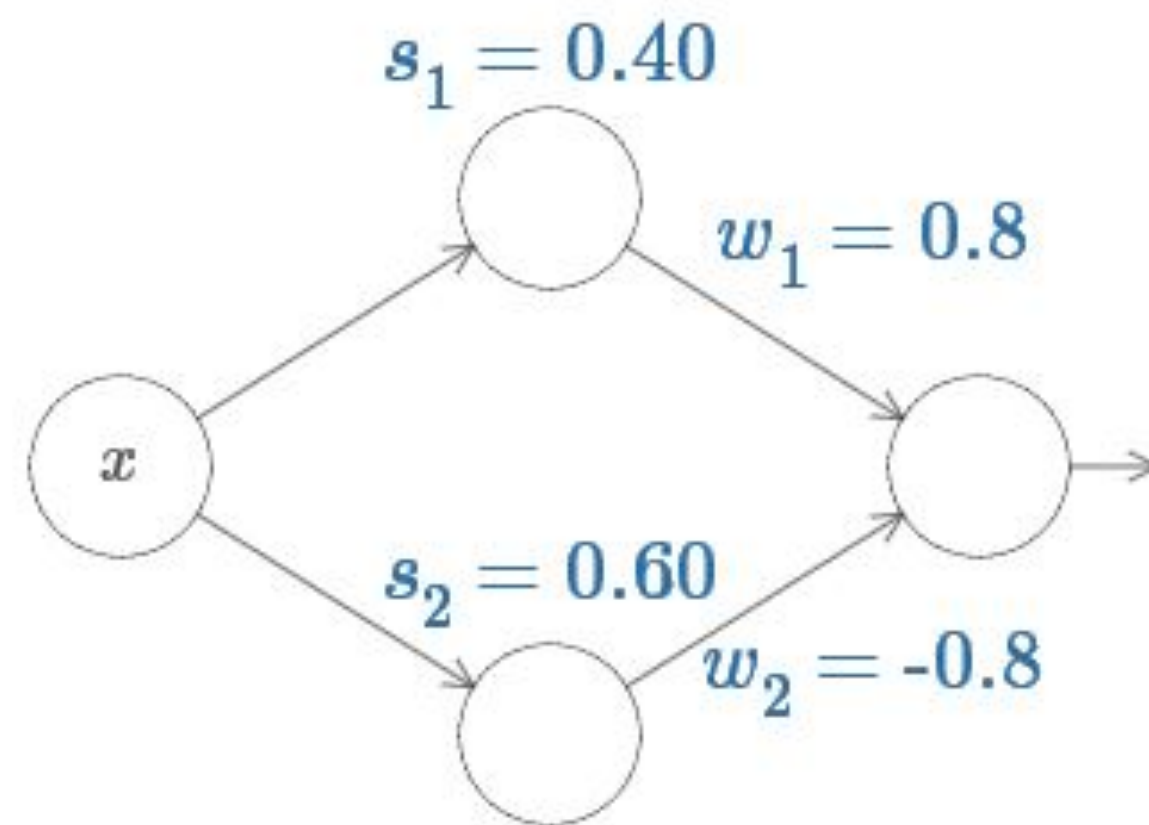
$$G(x) = \sum_{j=1}^{N} \alpha_j \sigma(y_j^T x + \theta_j)$$

*are dense in $C(I_n)$. In other words, given any $f \in C(I_n)$ and $\varepsilon > 0$, there is a sum, $G(x)$, of the above form, for which*

$$|G(x) - f(x)| < \varepsilon \qquad \text{for all} \quad x \in I_n.$$

Not every function, but a lot!
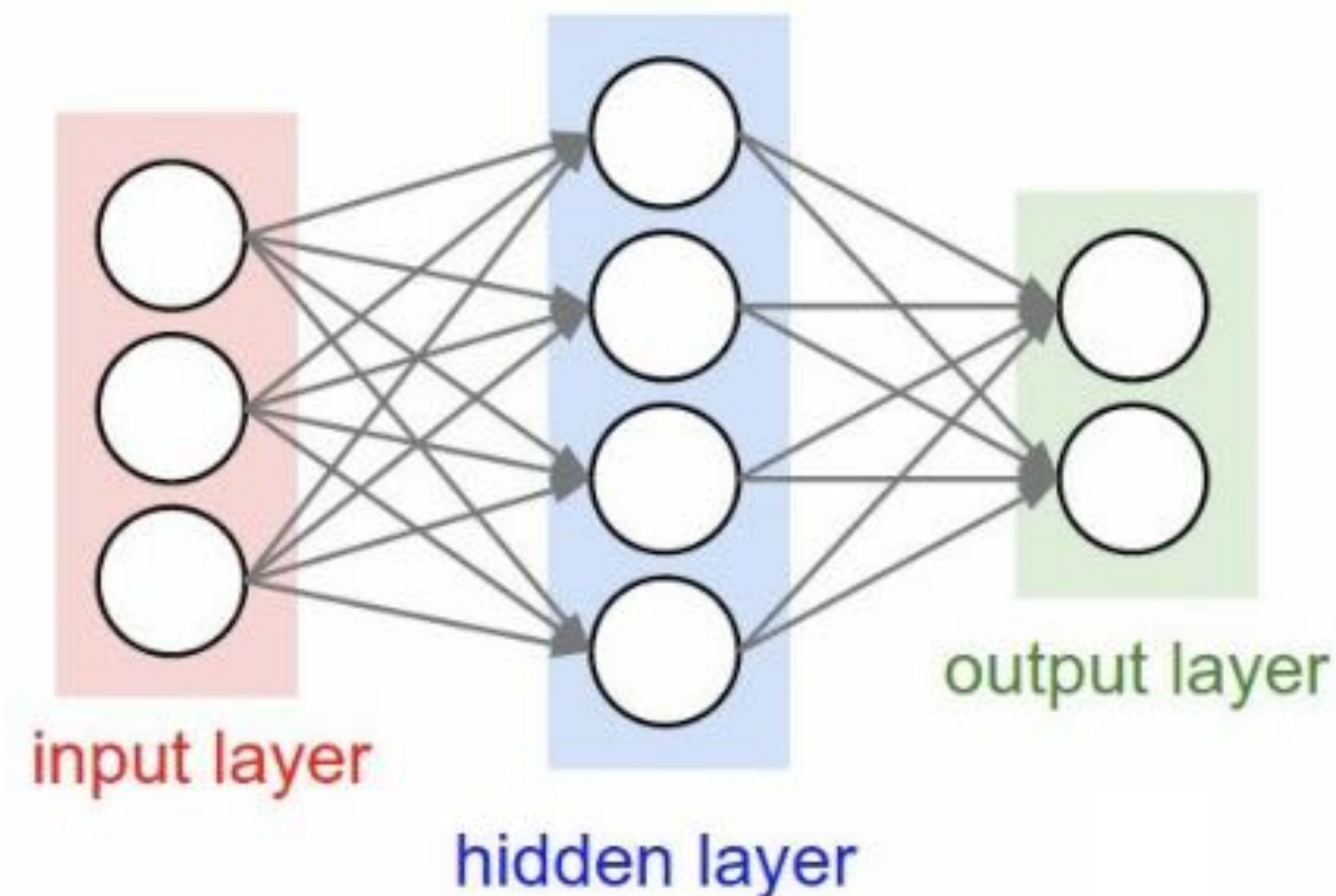
# Universal approximation theorem (Cybenko 1989)

# Universal approximation theorem (Cybenko 1989)

Exciting result right?

Not really…requires one neuron for every small volume of the space (so number of neurons grows exponentially with dimension!)

Overcome this with *depth*!



**Theorem 2.** *Let σ be any continuous sigmoidal function. Then finite sums of the form*

$$G(x) = \sum_{j=1}^{N} \alpha_j \sigma(y_j^T x + \theta_j)$$

*are dense in $C(I_n)$. In other words, given any $f \in C(I_n)$ and $\varepsilon > 0$, there is a sum, $G(x)$, of the above form, for which*

$$|G(x) - f(x)| < \varepsilon \qquad for\ all \quad x \in I_n.$$

# Tensorflow Playground

# Summary

- We arrange neurons into fully-connected layers
- The abstraction of a **layer** has the nice property that it allows us to use efficient vectorized code (e.g. matrix multiplies)
- Neural networks are not really *neural*
- Fully-connected neural networks are *universal* (but inefficient)
- Overcome inefficiency with *depth*!