# Neural Network
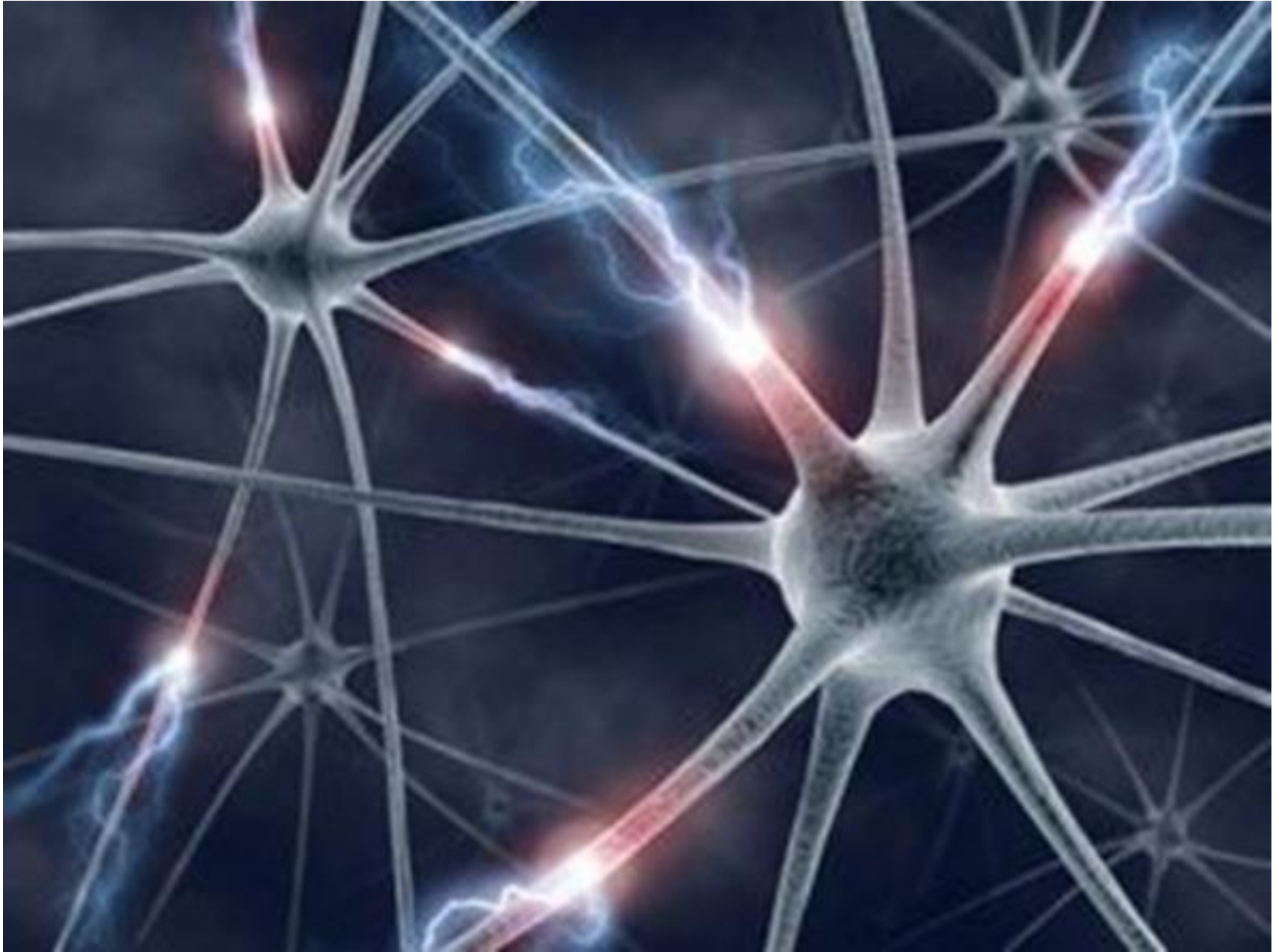
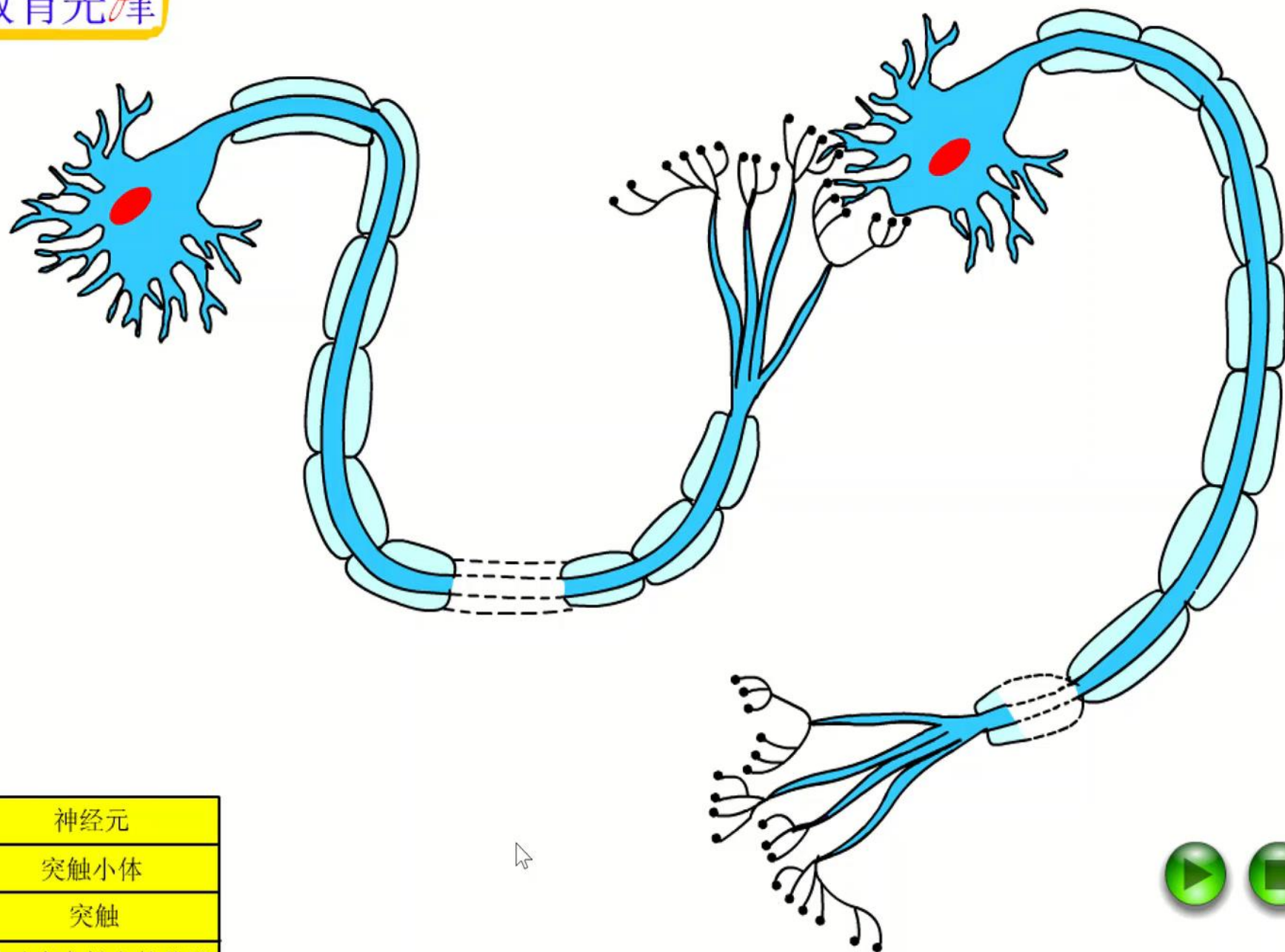# Outline

❖ Modeling one neuron

❖ Activation functions

❖ Fully connected feed-forward network

❖ How to train a multi-layer network

❖ Representational power of NN

# Biology

- Neurons respond slowly
  - $10^{-3}$ s compared to $10^{-9}$ s for electrical circuits
- The brain uses massively parallel computation
  - $\approx 10^{11}$ neurons in the brain
  - $\approx 10^{4}$ connections per neuron

教育先锋
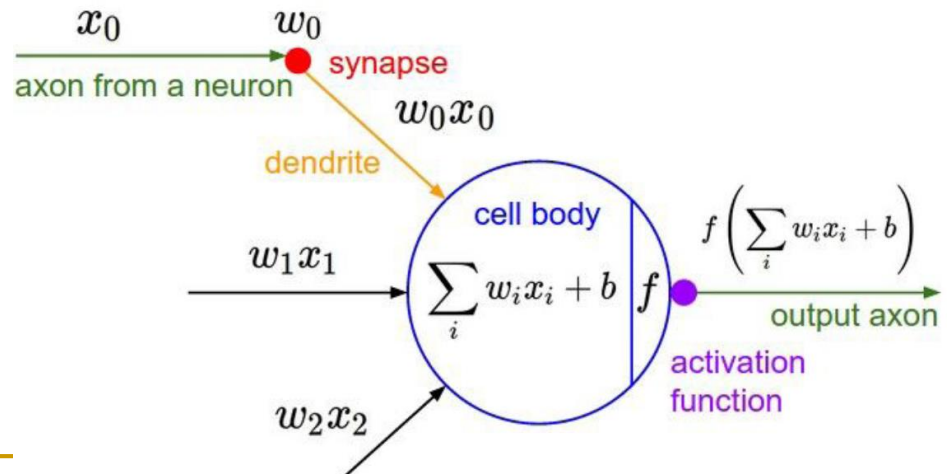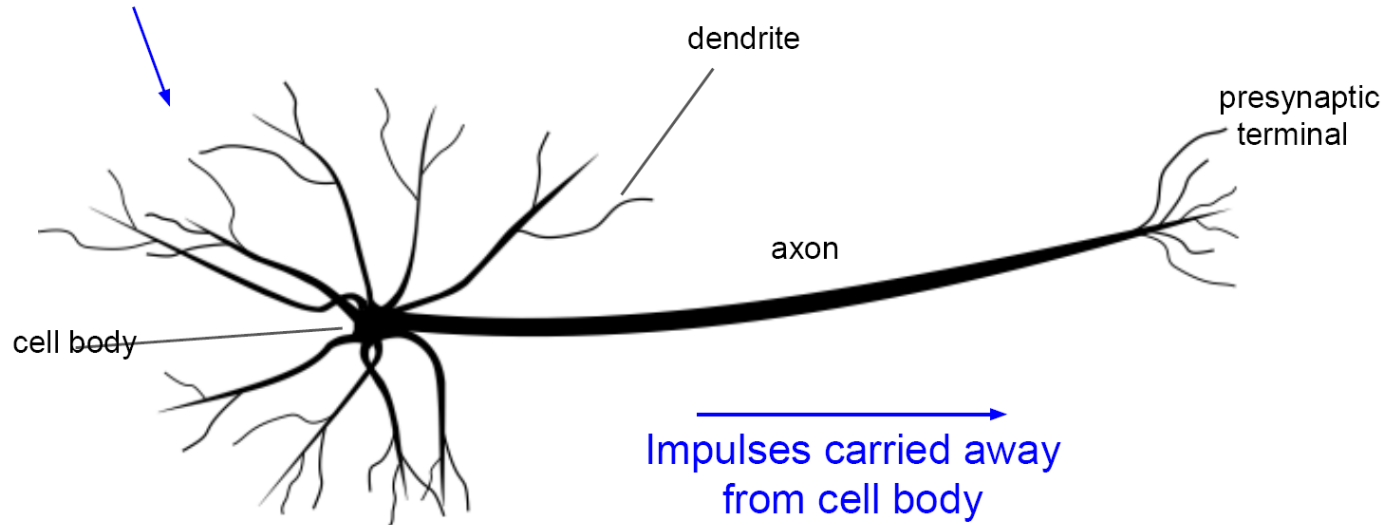
| 神经元 |
| 突触小体 |
| 突触 |
| 冲动在突触上的传递 |

# Modeling one neuron

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from 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

$w_2 x_2$

# Neuron

# Single neuron as a linear classifier

❖ With an appropriate loss function on the neuron's output, a single neuron can be turned into a linear classifier.

❖ A single neuron can be used to implement a binary classifier (e.g. binary Softmax or binary SVM classifiers).

# Single neuron as a linear classifier

❖ **Binary Softmax classifier**

  ◆ Interpreting:

    ▫ $\sigma(\sum_i w_i x_i + b)$  to be the probability of one of the classes

      $P(y_i = 1 \mid x_i; w)$

    ▫ The probability of the other class to be

      $P(y_i = 0 \mid x_i; w) = 1 - P(y_i = 1 \mid x_i; w)$

  ◆ With this interpretation, we can formulate the <span style="color:red">cross-entropy loss</span> ,  and optimizing it would lead to a **binary Softmax classifier** (also known as logistic regression)

❖ Alternatively, we could attach a <span style="color:blue">max-margin hinge loss</span> to the output of the neuron and train it to become a **binary Support Vector Machine**.
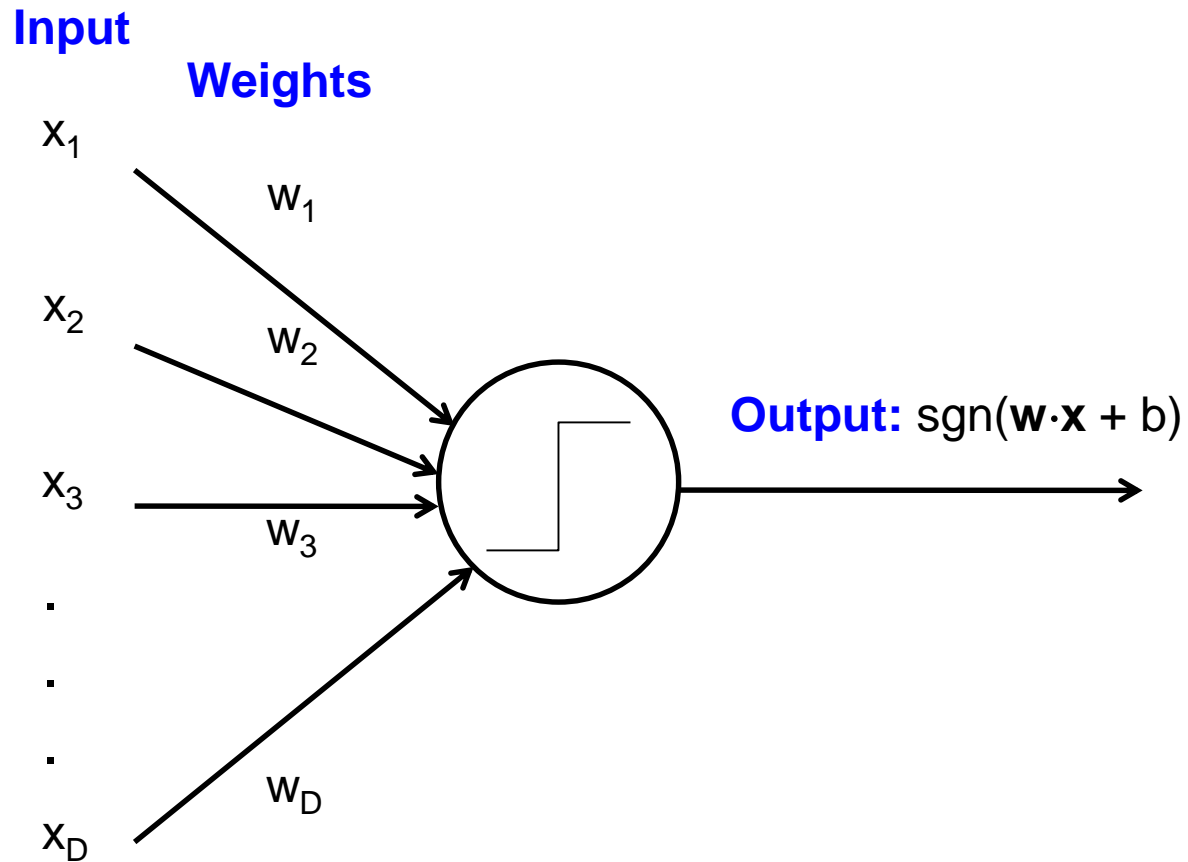
# Outline

❖ Modeling one neuron

❖ Activation functions

❖ Fully connected feed-forward network

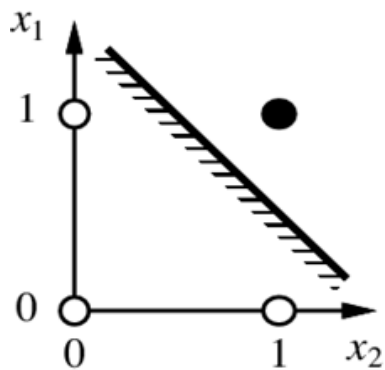❖ How to train a multi-layer network

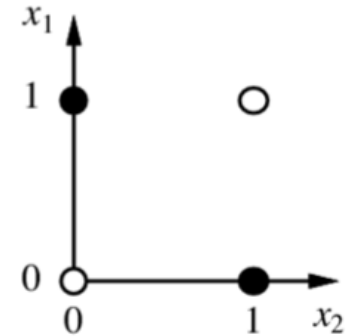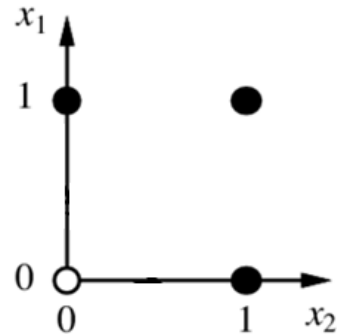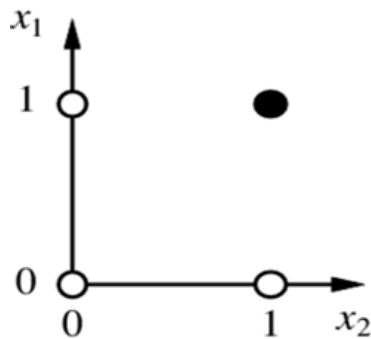❖ Representational power of NN

# Activation Functions

❖ Neural networks are used to implement complex functions, and <span style="color:red">non-linear</span> Activation Functions enable them to approximate arbitrarily complex functions.

❖ Without the non-linearity introduced by the Activation Function：

  ◆ a neuron is equivalent to a linear classifier

  ◆ a multi-layer neural network is equivalent to a single layer neural network

# Perceptron

**Input**
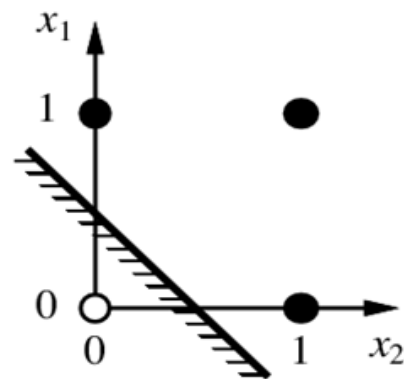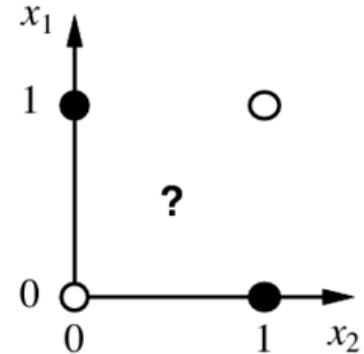
**Weights**

$x_1$

$w_1$

$x_2$

$w_2$

$x_3$

$w_3$

$.$

$.$

$.$

$w_D$

$x_D$

**Output:** $\text{sgn}(\mathbf{w}\cdot\mathbf{x} + b)$

# Linear separability



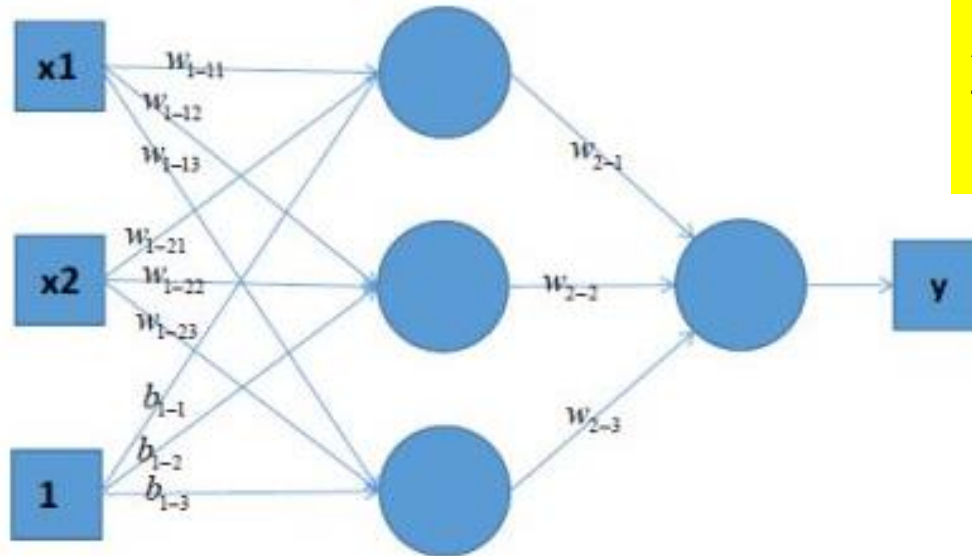$x_1$ **and** $x_2$          $x_1$ **or** $x_2$          $x_1$ **xor** $x_2$

# Multilayer Perceptron（MLP）

❖ A MLP with a single hidden layer



Add nonlinearity though
activation functions!

$$y = w_{2-1}(w_{1-11}x_1 + w_{1-21}x_2 + b_{1-1}) + w_{2-2}(w_{1-12}x_1 + w_{1-22}x_2 + b_{1-2}) + w_{2-3}(w_{1-13}x_1 + w_{1-23}x_2 + b_{1-3})$$

$$y = x_1(w_{2-1}w_{1-11} + w_{2-2}w_{1-12} + w_{2-3}w_{1-13}) + x_2(w_{2-1}w_{1-21} + w_{2-2}w_{1-22} + w_{2-3}w_{1-23}) + w_{2-1}b_{1-1} + w_{2-2}b_{1-2} + w_{2-3}b_{1-3}$$

Still a linear classifier

# Some 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
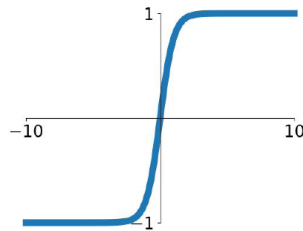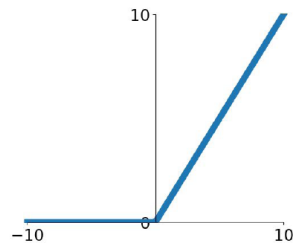
❖ **Sigmoid**  $\sigma(x) = 1/(1 + e^{-x})$

♦ Squashes numbers to range [0,1]

♦ Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

♦ Problems:

□ Saturated neurons "kill" the gradients

□ Sigmoid outputs are not zero-centered

□ exp() is a bit compute expensive



Sigmoid Activation Function

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



Derivative of Sigmoid Activation Function

$\sigma(x) = \sigma(x)(1 - \sigma(x))$

# Activation Functions

$$tanh x = \frac{sinh x}{cosh x} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

❖ **Tanh：**

**tanh(x)=2σ(2x)−1**

◆ Squashes numbers to range [-1,1]

◆ Zero centered

◆ Problem:

  ❑ kills gradients when saturated



Tanh Activation Function



Derivative of Tanh Activation Function

# Activation Functions

❖ **Relu： f(x) = max(0,x)**

（**Rectified Linear Unit**）

◆ Does not saturate (in +region)

◆ Very computationally efficient

◆ Converges much faster than sigmoid/tanh in practice (e.g. 6x)

◆ Actually more biologically plausible than sigmoid

◆ Problem:

▫ Not zero-centered output



ReLU Activation Function

max(0,x)

Derivative of ReLU Activation Function

# Activation Functions:

❖ **Leaky ReLU**

♦ 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)

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

# Activation Functions

❖ **Exponential Linear Units (ELU)**

♦ All benefits of ReLU

♦ Closer to zero mean outputs

♦ Negative saturation regime compared with Leaky ReLU

♦ Adds some robustness to noise

Exponential Linear Units (ELU)

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

# Activation Functions:

❖ **Maxout** : $\max(w_1^T x + b_1, w_2^T x + b_2)$

♦ Does not have the basic form of dot product -> nonlinearity

♦ Generalizes ReLU and Leaky ReLU

♦ Linear Regime! Does not saturate! Does not die!

♦ Problem:

▫ doubles the number of parameters/neuron

# Choosing the right Activation Function

❖ Good or bad – there is no rule of thumb

❖ Depending upon the properties of the problem, we might be able to make a better choice for **easy and quicker** convergence of the network

- ♦ Sigmoid functions and their combinations generally work better **in the case of classifiers**

- ♦ Sigmoid and tanh functions: sometimes to be avoided due to the vanishing gradient problem

- ♦ ReLU function : a general activation function, used in most cases these days, but only be used in the hidden layers. Be careful with the learning rates

- ♦ In the case of dead neurons in networks, try out the Leaky ReLU / Maxout / ELU

❖ Usually can begin with using ReLU function and then move over to other activation functions in case ReLU doesn't provide with optimum results

# Outline

❖ Modeling one neuron

❖ Activation functions

❖ Fully connected feed-forward network

❖ How to train a multi-layer network

❖ Representational power of NN

# Fully connected Feedforward Network

# Fully connected Feedforward Network



Sigmoid Function $\sigma(z)$

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

# Fully connected Feedforward Network

# Fully connected Feedforward Network



This is a function.
Input vector, output vector

$$f\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix} \quad f\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.51 \\ 0.85 \end{bmatrix}$$

Given network structure, we define ***a function set***

# Matrix Operation



$$\sigma\left( \begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

$$\begin{bmatrix} 4 \\ -2 \end{bmatrix}$$

# Neural Network

# Neural Network



$$\boxed{y} = f(\boxed{x})$$

Using parallel computing techniques to speed up matrix operation

$$= \sigma(\boxed{W^L} \cdots \sigma(\boxed{W^2} \sigma(\boxed{W^1} \boxed{x} + \boxed{b^1}) + \boxed{b^2}) \cdots + \boxed{b^L})$$

# Output Layer
# as a Multi-Class Classifier

Feature extractor replacing
feature engineering



$x$

**Input
Layer**

**Hidden Layers**

Softmax

$x_1$

$x_2$

$x_K$

**Output
Layer**

= Multi-class
Classifier

y₁

y₂

y_M

# Example Application



Input    Layer 1    Layer 2    Layer L    Output

A function set containing the candidates for
Handwriting Digit Recognition

$x_1$
$x_2$
$x_N$

$y_1$    is 1
$y_2$    is 2
$y_{10}$    is 0

**Input Layer**    **Hidden Layers**    **Output Layer**

Need to decide the network structure to work well on your dataset.

# "Deep" pipeline

Raw input → [ Layer 1 ] → [ Layer 2 ] → [ Layer 3 ] → Output

- Learn a *feature hierarchy*

- Each layer extracts features from the output of previous layer

- All layers are trained jointly

# Outline

❖ Modeling one neuron

❖ Activation functions

❖ Fully connected feed-forward network

❖ How to train a multi-layer network
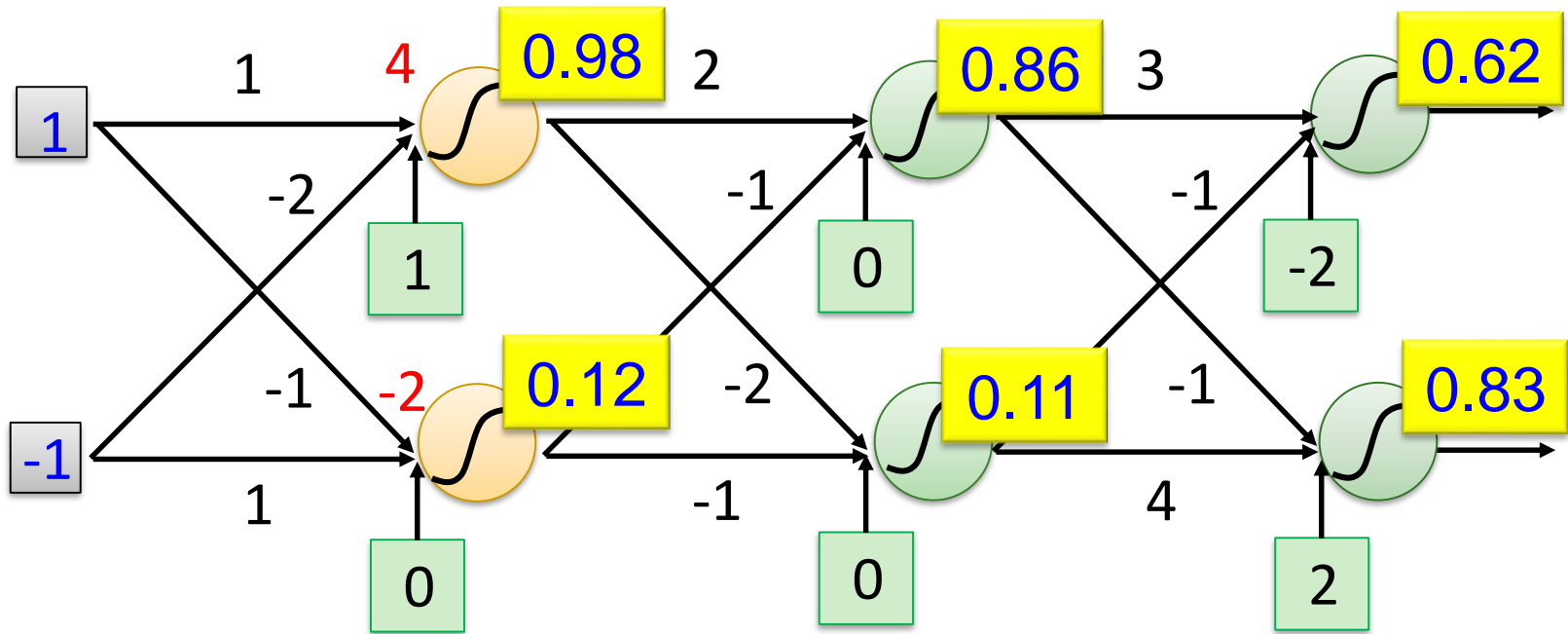
❖ Representational power of NN

# How to train a multi-layer network?



$$x \xrightarrow{\phantom{xx}} f_1(x, w_1) \xrightarrow{h_1} f_2(h_1, w_2) \xrightarrow{h_2} \cdots \xrightarrow{h_{K-1}} f_K(h_{K-1}, w_K) \xrightarrow{h_K} l(h_k, y) \xrightarrow{\phantom{xx}} e$$

input     first layer transformation     second layer transformation     output layer     loss function    error

hidden representation     output

- ❖ We need to find the gradient of the error w. r.t. the parameters of each layer, $\frac{\partial e}{\partial w_k}$, to perform updates     $w_k \leftarrow w_k - \eta \frac{\partial e}{\partial w_k}$

# Computation graph

# Chain Rule

**_Case 1_**     $y = g(x)$     $z = h(y)$

$$\Delta x \to \Delta y \to \Delta z \qquad \frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

**_Case 2_**

$$x = g(s) \qquad y = h(s) \qquad z = k(x, y)$$

$$\frac{dz}{ds} = \frac{\partial z}{\partial x}\frac{dx}{ds} + \frac{\partial z}{\partial y}\frac{dy}{ds}$$

# Chain rule

Let's start with $k = 1$



- $e = l(f_1(x, w_1), y)$

- $\dfrac{\partial}{\partial w_1} l(f_1(x, w_1), y) = \dfrac{\partial e}{\partial h_1} \dfrac{\partial h_1}{\partial w_1}$

- Example:  $e = (y - w_1^T x)^2$

- $h_1 = f_1(x, w_1) = w_1^T x$

- $e = l(h_1, y) = (y - h_1)^2$

$\dfrac{\partial h_1}{\partial w_1} = x$

$\dfrac{\partial e}{\partial h_1} = -2(y - h_1) = -2(y - w_1^T x)$

- $\dfrac{\partial e}{\partial w_1} = \dfrac{\partial e}{\partial h_1} \dfrac{\partial h_1}{\partial w_1} = -2x(y - w_1^T x)$

# Chain rule



$k = 2$

$x$

❖ $e = l(f_2(f_1(x, w_1), w_2))$

$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial w_2} \qquad\qquad \frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_1}$$

❖ Example: $e = -\log\left(\sigma(w_1^T x)\right)$ (assume $y = 1$)

$h_1 = f_1(x, w_1) = w_1^T x \qquad \frac{\partial h_1}{\partial w_1} = x$

$h_2 = f_2(h_1, w_2) = \sigma(h_1) \qquad \frac{\partial h_2}{\partial h_1} = \sigma'(h_1) = \sigma(h_1)(1 - \sigma(h_1))$

$e = l(h_2, 1) = -\log(h_2) \qquad \frac{\partial e}{\partial h_2} = -\frac{1}{h_2}$

$$\frac{\partial e}{\partial w_1} = \frac{\partial e}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_1} = -\frac{1}{\sigma(w_1^T x)} \sigma(w_1^T x)\left(1 - \sigma(w_1^T x)\right) x = \sigma(-w_1^T x) x - x$$

# Chain rule



❖ General case:

❖ $$\frac{\partial e}{\partial w_{\text{i}}} = \boxed{\frac{\partial e}{\partial h_K} \ \frac{\partial h_K}{\partial h_{K-1}} \ ... \ \frac{\partial h_{\text{i}+1}}{\partial h_{\text{i}}}} \ \frac{\partial h_{\text{i}}}{\partial w_{\text{i}}}$$

Upstream gradient
$$\frac{\partial e}{\partial h_{\text{i}}}$$

Local gradient

# Backpropagation summary

$w_k$

Parameter update: $\dfrac{\partial e}{\partial w_k} = \dfrac{\partial e}{\partial h_k}\dfrac{\partial h_k}{\partial w_k}$

Upstream gradient:
$$\frac{\partial e}{\partial h_k}$$

$\dfrac{\partial h_k}{\partial w_k}$ Local gradient

$f_k$

$\dfrac{\partial h_k}{\partial h_{k-1}}$ Local gradient

$h_k$

$h_{k-1}$

Upstream gradient:
$$\frac{\partial e}{\partial h_{k-1}} = \frac{\partial e}{\partial h_k}\frac{\partial h_k}{\partial h_{k-1}}$$

→ Forward pass

→ Backward pass

# What about more general computation graphs?

+ Gradients add at branches

# A detailed example

$$f(x, w) = \frac{1}{1 + \exp[-(w_0 x_0 + w_1 x_1 + w_2)]}$$

# A detailed example

$$f(x, w) = \frac{1}{1 + \exp[-(w_0 x_0 + w_1 x_1 + w_2)]}$$



w0  2.00
x0  -1.00
-2.00
w1  -3.00
x1  -2.00
6.00
+  4.00
w2  -3.00
+  1.00  *-1  -1.00  exp  0.37  +1  1.37  1/x  0.73

Local gradient * Upstream gradient

$$(1/x)' = -1/x^2$$

$$-\frac{1}{1.37^2} * 1 = -0.53$$

# A detailed example

$$f(x, w) = \frac{1}{1 + \exp[-(w_0 x_0 + w_1 x_1 + w_2)]}$$

w0  2.00
    -0.20

x0  -1.00
     0.40

w1  -3.00
    -0.40

x1  -2.00
    -0.60

w2  -3.00
     0.20

\*  -2.00
    0.20

\+  4.00
    0.20

\*  6.00
    0.20

\+  1.00
    0.20

Can simplify computation graph

\*-1   -1.00
       -0.20

exp   0.37
      -0.53

+1   1.37
     -0.53

1/x   0.73
      1.00

Sigmoid gate $\sigma(x) = \frac{1}{1 + \exp(-x)}$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

$$\sigma(1)(1 - \sigma(1)) = 0.73 * (1 - 0.73) = 0.20$$

Source: Stanford 231n

# Patterns in gradient flow

x  3.00

y  -4.00

z  2.00

w  -1.00

# Patterns in gradient flow



Add gate: "gradient distributor"

Multiply gate: "gradient switcher"

Max gate: "gradient router"

Source: Stanford 231n

# Dealing with vectors
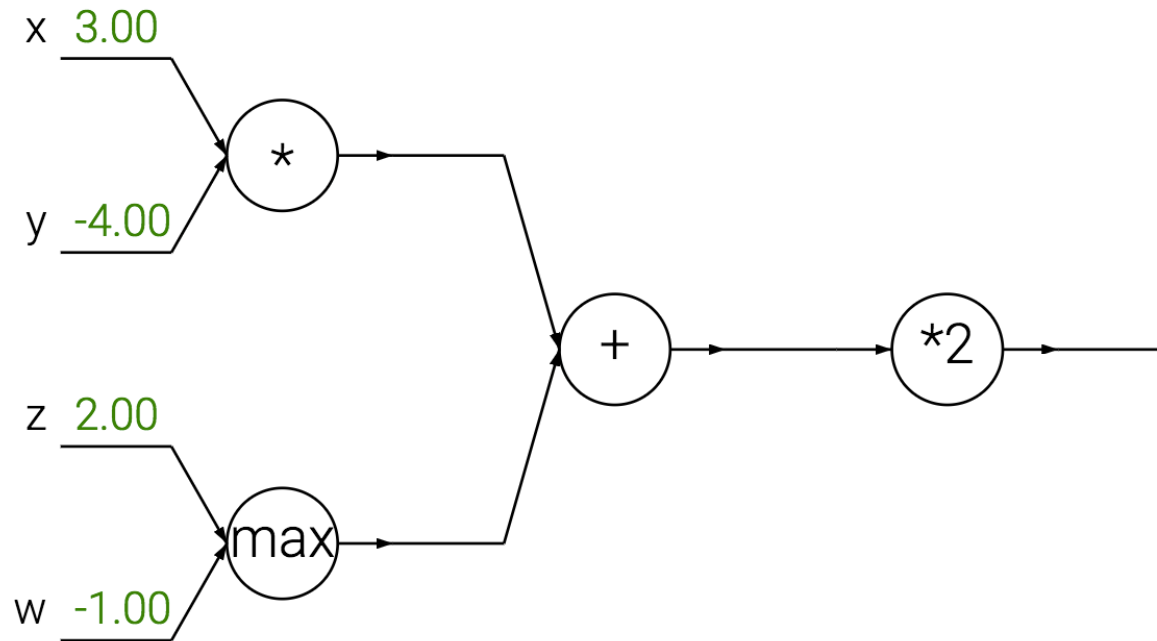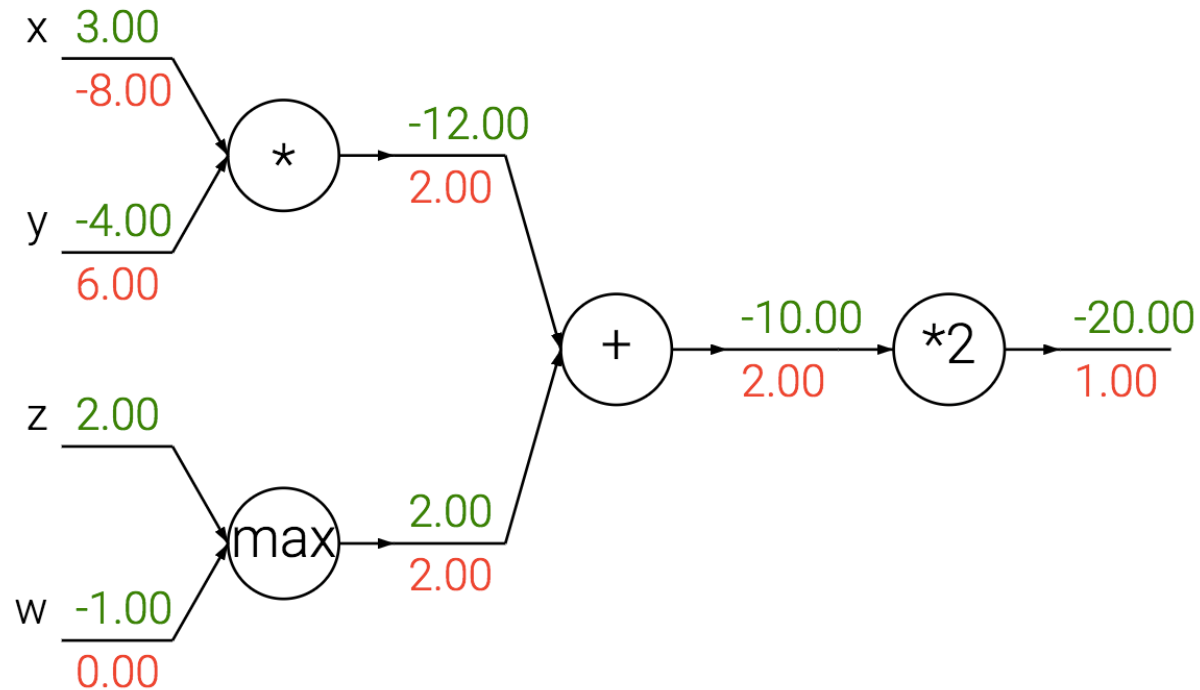
$$\frac{\partial z}{\partial x} = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \cdots & \frac{\partial z_1}{\partial x_M} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_N}{\partial x_1} & \cdots & \frac{\partial z_N}{\partial x_M} \end{pmatrix}$$

$N \times M$
Jacobian

$x$

$f(x)$

$z$

$1 \times M$

$$\frac{\partial e}{\partial x} = \frac{\partial e}{\partial z} \frac{\partial z}{\partial x}$$

$$\frac{\partial e}{\partial z}$$

$1 \times N$

$1 \times M \quad 1 \times N \quad N \times M$

$1 \times N$

# Matrix-vector multiplication

$$\frac{\partial e}{\partial W} = \frac{\partial e}{\partial z} \frac{\partial z}{\partial W}$$

$M \times N \quad 1 \times N \quad N \times (M \times N)$

$\frac{\partial z}{\partial W}$

$W$

$N \times (M \times N)$

$M \times N$

$$f(x, W) = xW$$

$z$

$1 \times N$

$\frac{\partial e}{\partial z}$

$1 \times N$

$\frac{\partial z}{\partial x}$

$x$

$\frac{\partial e}{\partial x} = \frac{\partial e}{\partial z} \frac{\partial z}{\partial x}$

$N \times M$

$1 \times M$

$1 \times M \quad 1 \times N \quad N \times M$

# A vectorized example:

$$f(\mathrm{x}, W) = \sum_{i=1}^{n} (W \cdot \mathrm{x})_i^2$$



$$f(q) = q_1^2 + \cdots + q_n^2$$

$$q = W\mathrm{x} = \begin{pmatrix} w_{11}x_1 + \cdots + w_{1n}x_n \\ \vdots \\ w_{n1}x_1 + \cdots + w_{nn}x_n \end{pmatrix}$$

# A vectorized example: $f(\mathrm{x}, W) = \sum_{i=1}^{n}(W \cdot \mathrm{x})_i^2$

Feed-forward:

$$q = Wx = \begin{pmatrix} w_{11}x_1 + w_{12}x_2 \\ w_{21}x_1 + w_{22}x_2 \end{pmatrix} = \begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$f(q) = \sum_{i=1}^{2} \mathrm{q}_i^2 = 0.116$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} \; \mathrm{W}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} \; \mathrm{x}$$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$*$

L2

0.116

# A vectorized example: $f(x, W) = \sum_{i=1}^{n}(W \cdot x)_i^2$

Backward:

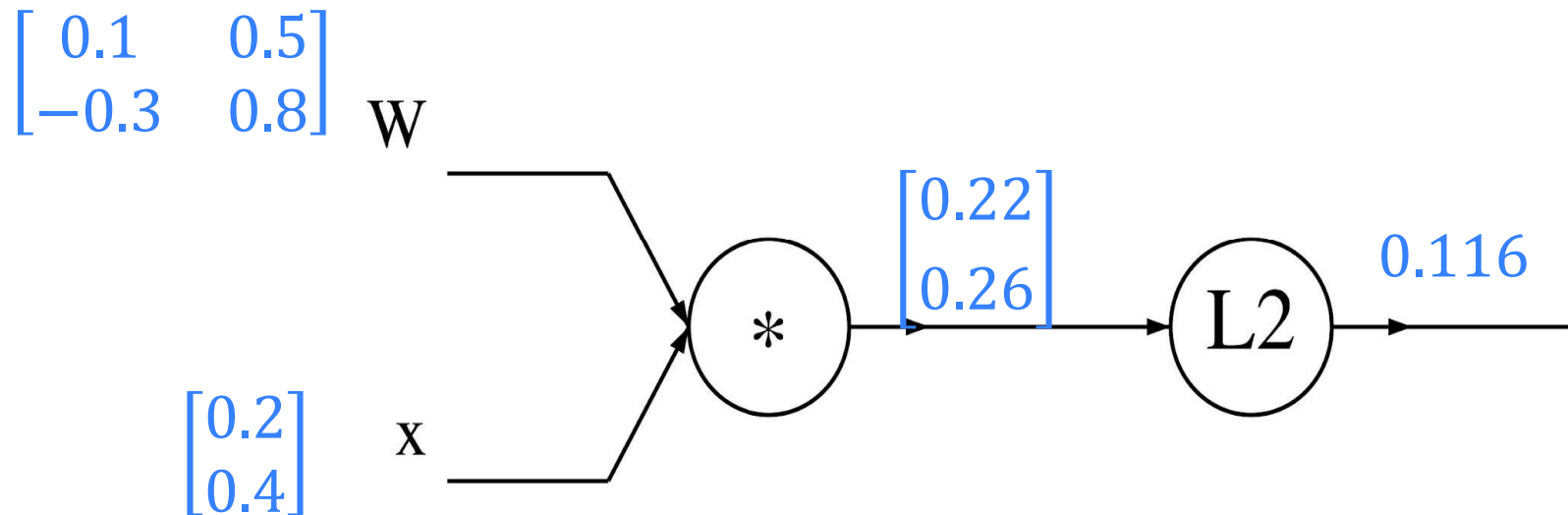$$f(q) = \sum_{i=1}^{2} q_i^2 \implies \frac{\partial f}{\partial q_i} = 2q_i \implies \boxed{\nabla_q f = 2q}$$

# A vectorized example: $f(x, W) = \sum_{i=1}^{n}(W \cdot x)_i^2$

Backward:

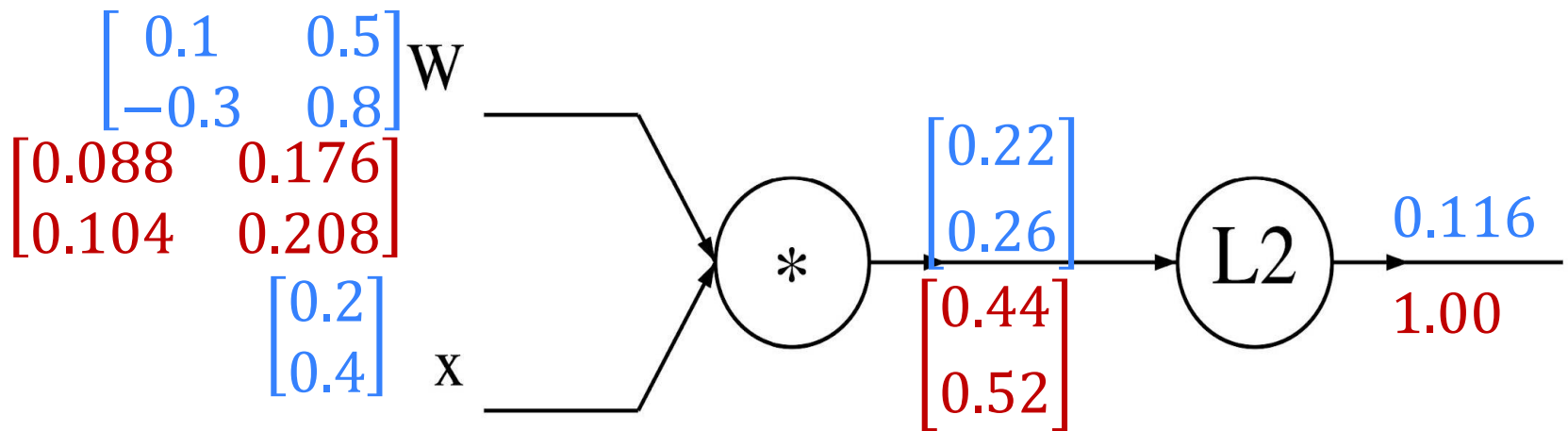$$q = Wx = \begin{pmatrix} w_{11}x_1 + w_{12}x_2 \\ w_{21}x_1 + w_{22}x_2 \end{pmatrix} \implies$$

$$\frac{\partial q}{\partial w_{11}} = x_1 \qquad \frac{\partial q}{\partial w_{12}} = x_2$$

$$\frac{\partial q}{\partial w_{21}} = x_1 \qquad \frac{\partial q}{\partial w_{22}} = x_2$$

$$\frac{\partial f}{\partial w_{ij}} = \sum_k \frac{\partial f}{\partial q_k}\frac{\partial q_k}{\partial w_{ij}} = 2q_i x_j \implies \boxed{\nabla_W f = 2qx^T}$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} x$$

$*$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

# A vectorized example: $f(\mathrm{x}, W) = \sum_{i=1}^{n}(W \cdot \mathrm{x})_i^2$

Backward:

$$q = Wx = \begin{pmatrix} w_{11}x_1 + w_{12}x_2 \\ w_{21}x_1 + w_{22}x_2 \end{pmatrix} \implies$$

$$\frac{\partial q}{\partial x_1} = \begin{pmatrix} w_{11} \\ w_{21} \end{pmatrix}$$

$$\frac{\partial q}{\partial x_2} = \begin{pmatrix} w_{12} \\ w_{22} \end{pmatrix}$$

$$\frac{\partial f}{\partial x_i} = \sum_k \frac{\partial f}{\partial q_k}\frac{\partial q_k}{\partial x_i} = \sum_k 2\, q_k W_{ki} \implies \boxed{\nabla_x\, \mathrm{f} = W^T 2q}$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W$$

$$\begin{bmatrix} 0.088 & 0.176 \\ 0.104 & 0.208 \end{bmatrix}$$

$$\begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} \mathrm{x}$$

$$\begin{bmatrix} -0.112 \\ 0.636 \end{bmatrix}$$

$*$

$$\begin{bmatrix} 0.22 \\ 0.26 \end{bmatrix}$$

$$\begin{bmatrix} 0.44 \\ 0.52 \end{bmatrix}$$

L2

0.116

1.00

# General tips for computation

- Derive error signal (upstream gradient) directly, avoid explicit computation of huge local derivatives

- Write out expression for a single element of the Jacobian, then deduce the overall formula

- Keep consistent indexing conventions, order of operations

- Use dimension analysis

# Outline

❖ Modeling one neuron

❖ Activation functions

❖ Fully connected feed-forward network

❖ How to train a multi-layer network
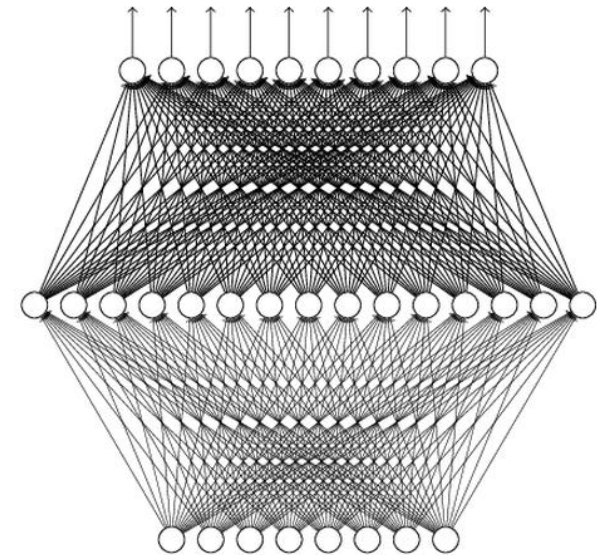
❖ Representational power of NN

# Representational power

❖ Neural Networks with fully-connected layers define a family of functions that are parameterized by the weights of the network.

❖ A Neural Network with at least one hidden layer are *universal approximators, which* means that it can approximate any continuous function.

# Universality Theorem
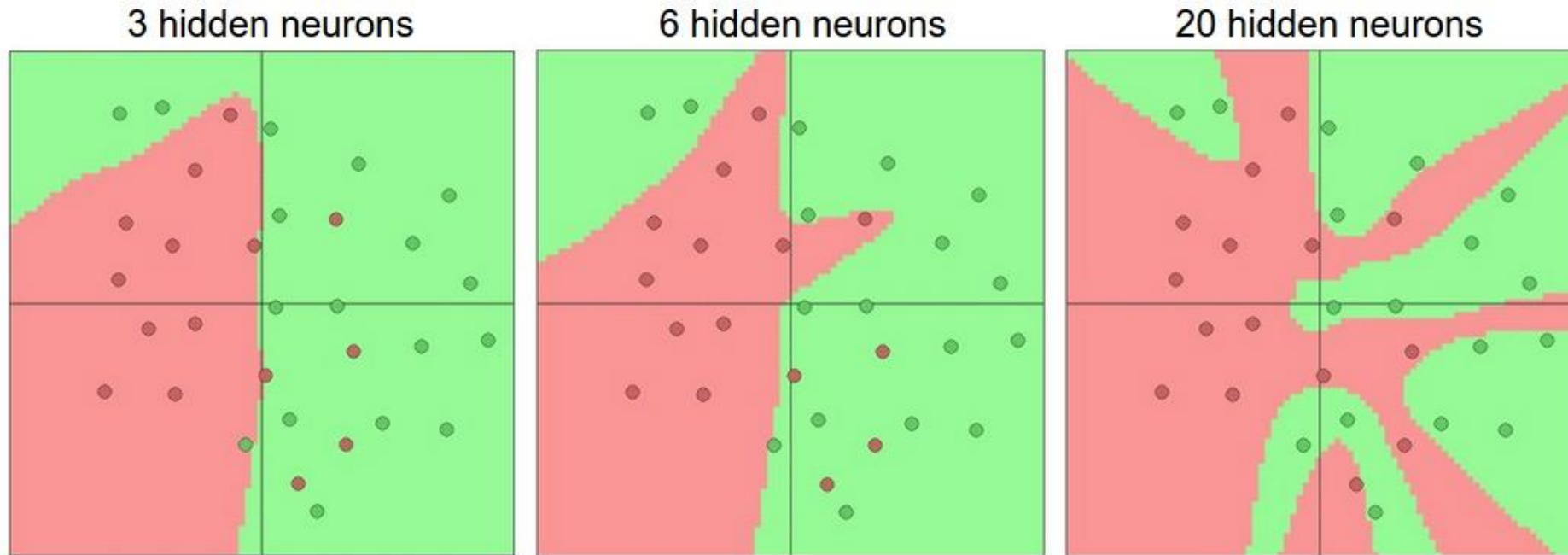
Any continuous function f

$$f : R^N \rightarrow R^M$$

can be realized by a network
with one hidden layer

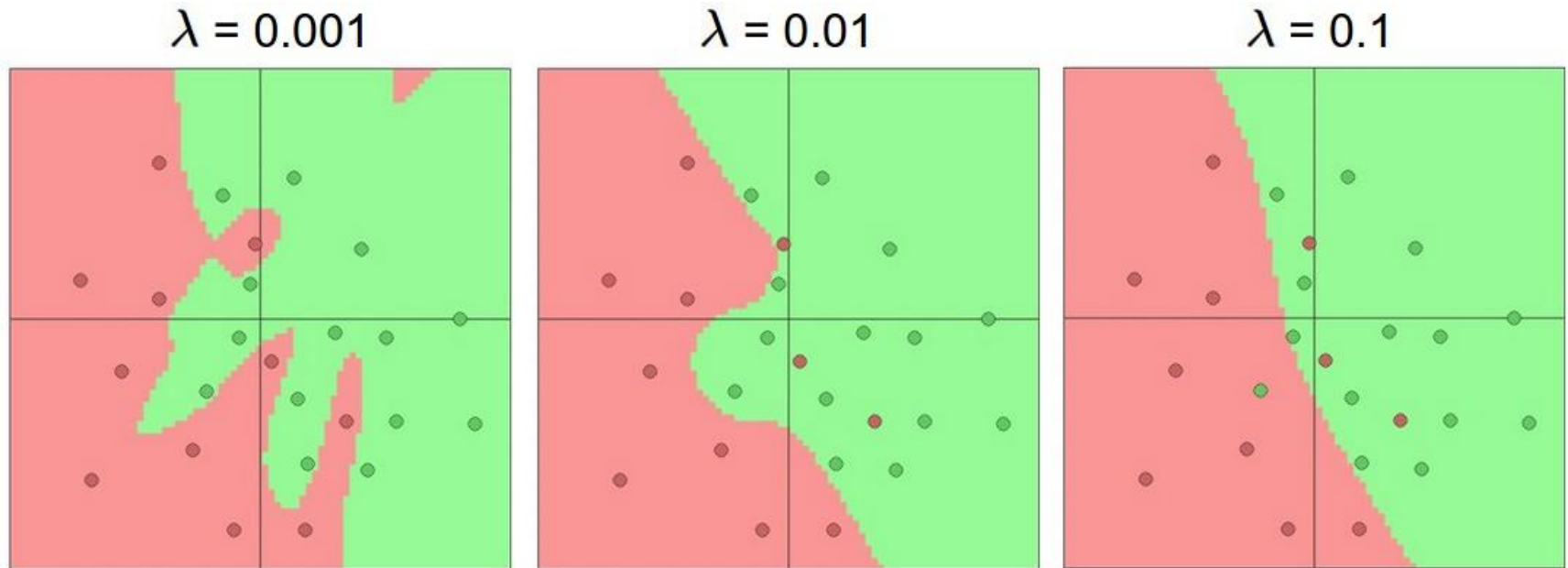(given **enough** hidden neurons)



Reference for the reason:
http://neuralnetworksandde
eplearning.com/chap4.html

# Setting number of layers and their sizes


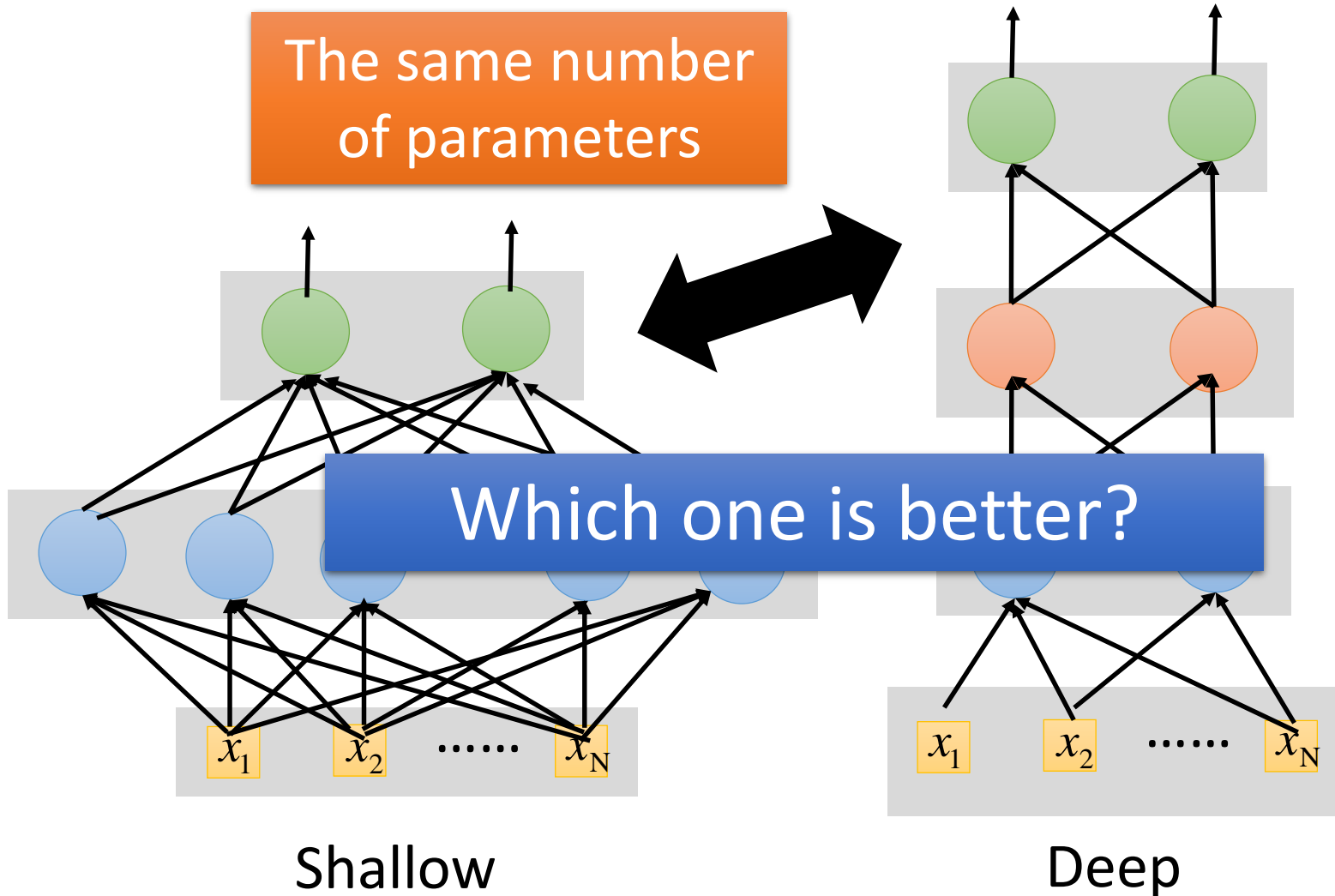
3 hidden neurons     6 hidden neurons     20 hidden neurons

❖ Neural Networks with more neurons can express more complicated functions.

❖ But a model with high capacity fits the noise in the data instead of the (assumed) underlying relationship : **overfitting**.

❖ Use as big of a neural network as your computational budget allows, and use other regularization techniques to control overfitting.



λ = 0.001          λ = 0.01          λ = 0.1

# Fat + Short vs. Thin + Tall



The same number of parameters

Which one is better?

Shallow

Deep

# Thin + Tall vs. Fat + Short

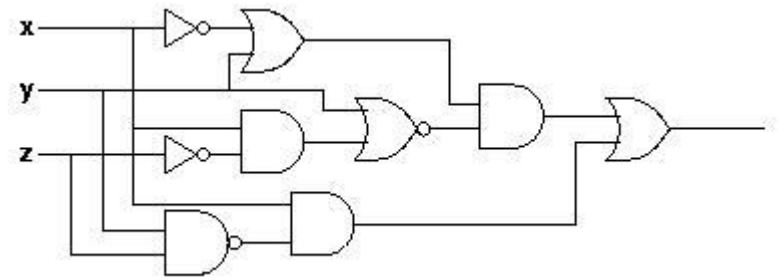| Layer X Size | Word Error Rate (%) | Layer X Size | Word Error Rate (%) |
|---|---|---|---|
| 1 X 2k | 24.2 | | |
| 2 X 2k | 20.4 | | Why? |
| 3 X 2k | 18.4 | | |
| 4 X 2k | 17.8 | | |
| 5 X 2k | 17.2 | 1 X 3772 | 22.5 |
| 7 X 2k | 17.1 | 1 X 4634 | 22.6 |
| | | 1 X 16k | 22.1 |

# Analogy



## Logic circuits

- Logic circuits consists of **gates**

- **A two layers of logic gates** can represent **any Boolean function.**

- Using multiple layers of logic gates to build some functions are much simpler
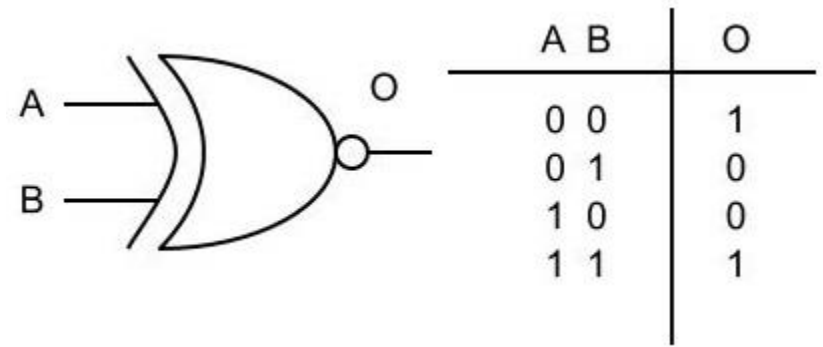
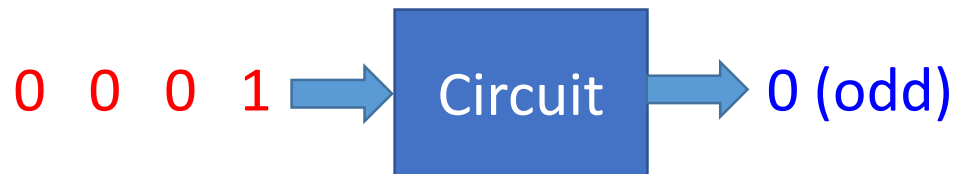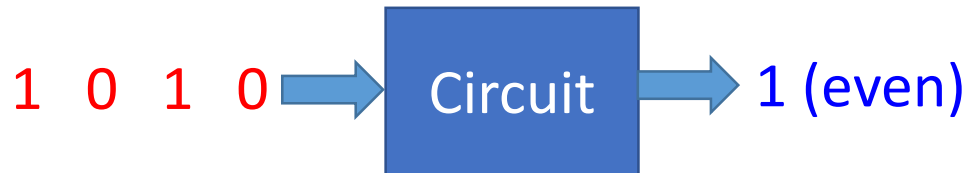→ less gates needed

## Neural network

- Neural network consists of **neurons**

- **A hidden layer network** can represent **any continuous function.**

- Using multiple layers of neurons to represent some functions are much simpler
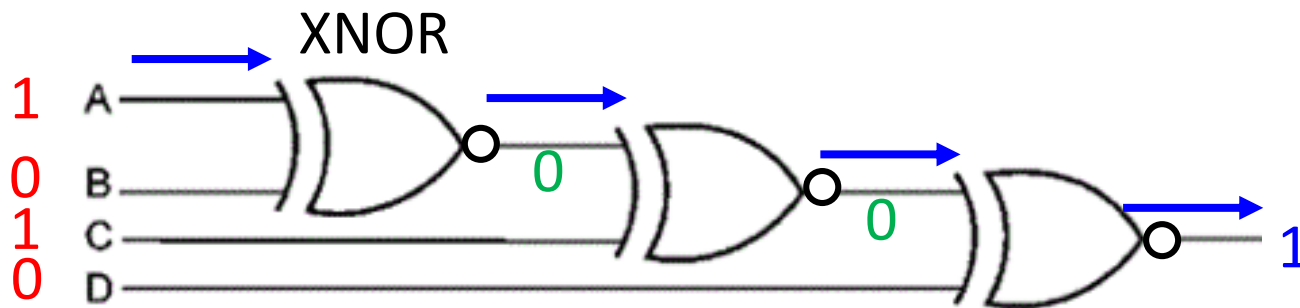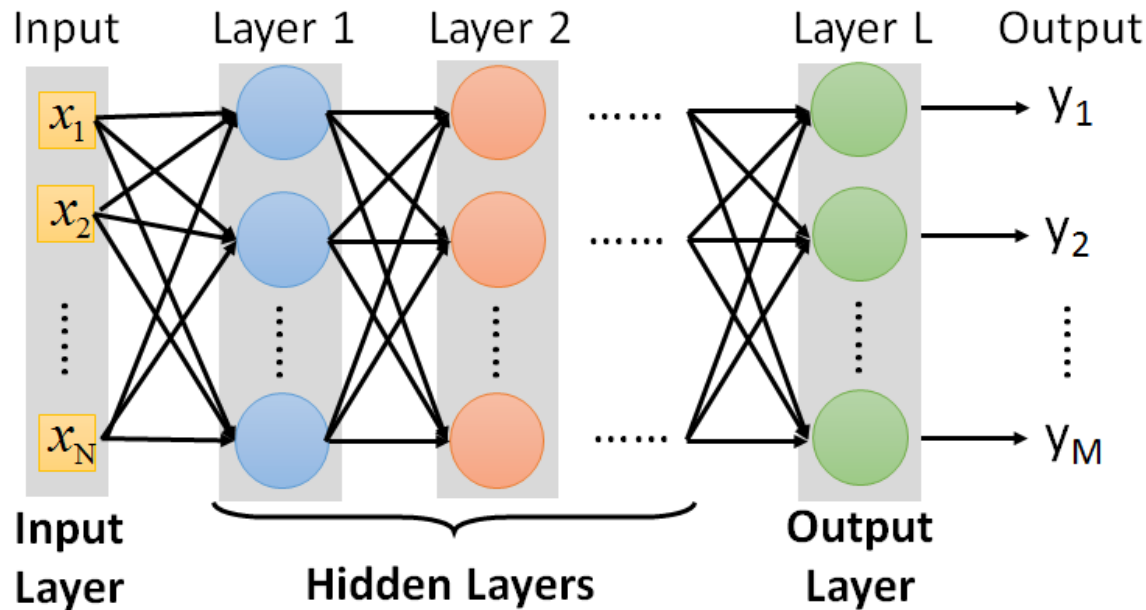
→ less parameters → less data?

# Analogy



| A B | O |
|-----|---|
| 0 0 | 1 |
| 0 1 | 0 |
| 1 0 | 0 |
| 1 1 | 1 |

- E.g. *parity check*

1  0  1  0 → Circuit → 1 (even)

0  0  0  1 → Circuit → 0 (odd)

For input sequence with d bits,

Two-layer circuit need $O(2^d)$ gates.

XNOR

1  A
0  B
1  C
0  D

0    0    1

With multiple layers, we need only $O(d)$ gates.

# Design the Network



- How many layers? How many neurons for each layer?

Trial and Error + Intuition

课程部分材料来自他人和网络，仅限教学使用，请勿传播，谢谢！