

CS2109S: Introduction to AI and Machine Learning

Lecture 8: Intro to Neural Networks

13 October 2023

Recap

- Overfitting
- Regularization
 - Linear and logistic regression
- Support Vector Machine (SVM)
 - Hard-margin SVM
 - Soft-margin SVM
- Kernel
 - SVM with Kernel Trick

Logistic Regression / SVM
With x as features

Logistic Regression / SVM
With $\phi(x)$ as features

SVM with Kernel Trick
With $\phi(x)$ mapping to
finite-dimensional
features

SVM with Kernel Trick
With $\phi(x)$ mapping to
infinite-dimensional
features



Linear Regression w/ Regz: Gradient Descent

$$J(\mathbf{w}) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\mathbf{w}}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{i=1}^n \mathbf{w}_i^2 \right]$$

Optimization goal: $\min_{\mathbf{w}} J(\mathbf{w})$

Repeat {

$$\mathbf{w}_0 := \mathbf{w}_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\mathbf{w}}(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

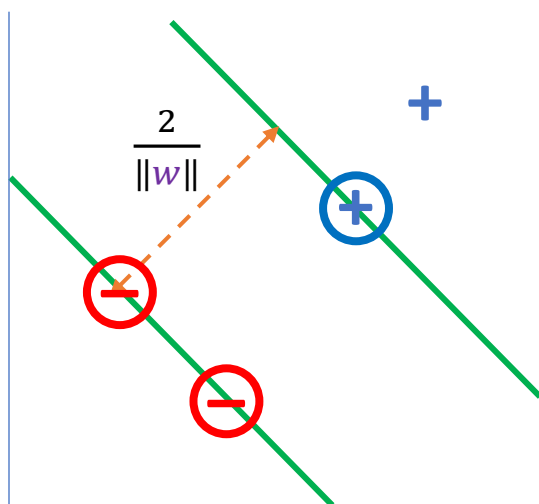
$$\mathbf{w}_1 := \mathbf{w}_1 - \alpha \frac{1}{m} \left[\sum_{i=1}^m (h_{\mathbf{w}}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)} + \lambda \mathbf{w}_1 \right]$$

$$\vdots$$

$$\mathbf{w}_n := \mathbf{w}_n - \alpha \frac{1}{m} \left[\sum_{i=1}^m (h_{\mathbf{w}}(x^{(i)}) - y^{(i)}) \cdot x_n^{(i)} + \lambda \mathbf{w}_n \right] \quad \}$$

Why does
this work?

Support Vector Machines (SVM)



$$\max_{\alpha} \min_{w, b} L(w, b, \alpha)$$

Minimized when gradient = 0

Objective:

$$\max \frac{2}{\|w\|} \rightarrow \max \frac{1}{\|w\|} \rightarrow \min \|w\| \rightarrow \min \frac{1}{2} \|w\|^2 \quad \text{"Maximize gap"}$$

$$\text{s.t. } \bar{y}^{(i)}(w \cdot x^{(i)} + b) - 1 \geq 0 \quad \text{"Classify correctly"}$$

$$b = \bar{y}^{(i)} - w \cdot x^{(i)}$$

Objective (Lagrange Dual):

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i \alpha^{(i)} [\bar{y}^{(i)}(w \cdot x^{(i)} + b) - 1], \forall_i \alpha^{(i)} \geq 0$$

$$\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_i \alpha^{(i)} \bar{y}^{(i)} x^{(i)} = 0$$

$$w = \sum_i \alpha^{(i)} \bar{y}^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_i \alpha^{(i)} \bar{y}^{(i)} = 0$$

Samples with non-zero $\alpha^{(i)}$ = support vectors

... a few math later ...

Maximize

$$L(\alpha) = \sum_i \alpha^{(i)} - \frac{1}{2} \sum_i \sum_j \alpha^{(i)} \alpha^{(j)} \bar{y}^{(i)} \bar{y}^{(j)} x^{(i)} \cdot x^{(j)}$$

SVM with Kernel Trick

Objective:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha^{(i)} [\bar{y}^{(i)} (\mathbf{w} \cdot \phi(x^{(i)}) + b) - 1]$$

... a few math later ...

$$L(\alpha) = \sum_i \alpha^{(i)} - \frac{1}{2} \sum_i \sum_j \alpha^{(i)} \alpha^{(j)} \bar{y}^{(i)} \bar{y}^{(j)} \phi(x^{(i)}) \cdot \phi(x^{(j)})$$

$$L(\alpha) = \sum_i \alpha^{(i)} - \frac{1}{2} \sum_i \sum_j \alpha^{(i)} \alpha^{(j)} \bar{y}^{(i)} \bar{y}^{(j)} K(x^{(i)}, x^{(j)})$$

From Before:

$$\mathbf{w} = \sum_i \alpha^{(i)} \hat{y}^{(i)} x^{(i)}$$

Decision Rule:

$\mathbf{w} \cdot \phi(x) + b \geq 0$ then +

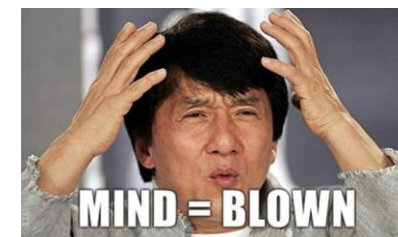
$\sum_i \alpha^{(i)} \hat{y}^{(i)} \phi(x^{(i)}) \cdot \phi(x) + b \geq 0$ then +

$\sum_i \alpha^{(i)} \hat{y}^{(i)} K(x^{(i)}, x) + b \geq 0$ then +



There is no need to compute the transformed features **explicitly**!

Can have SVM with **infinite-dimensional** features!



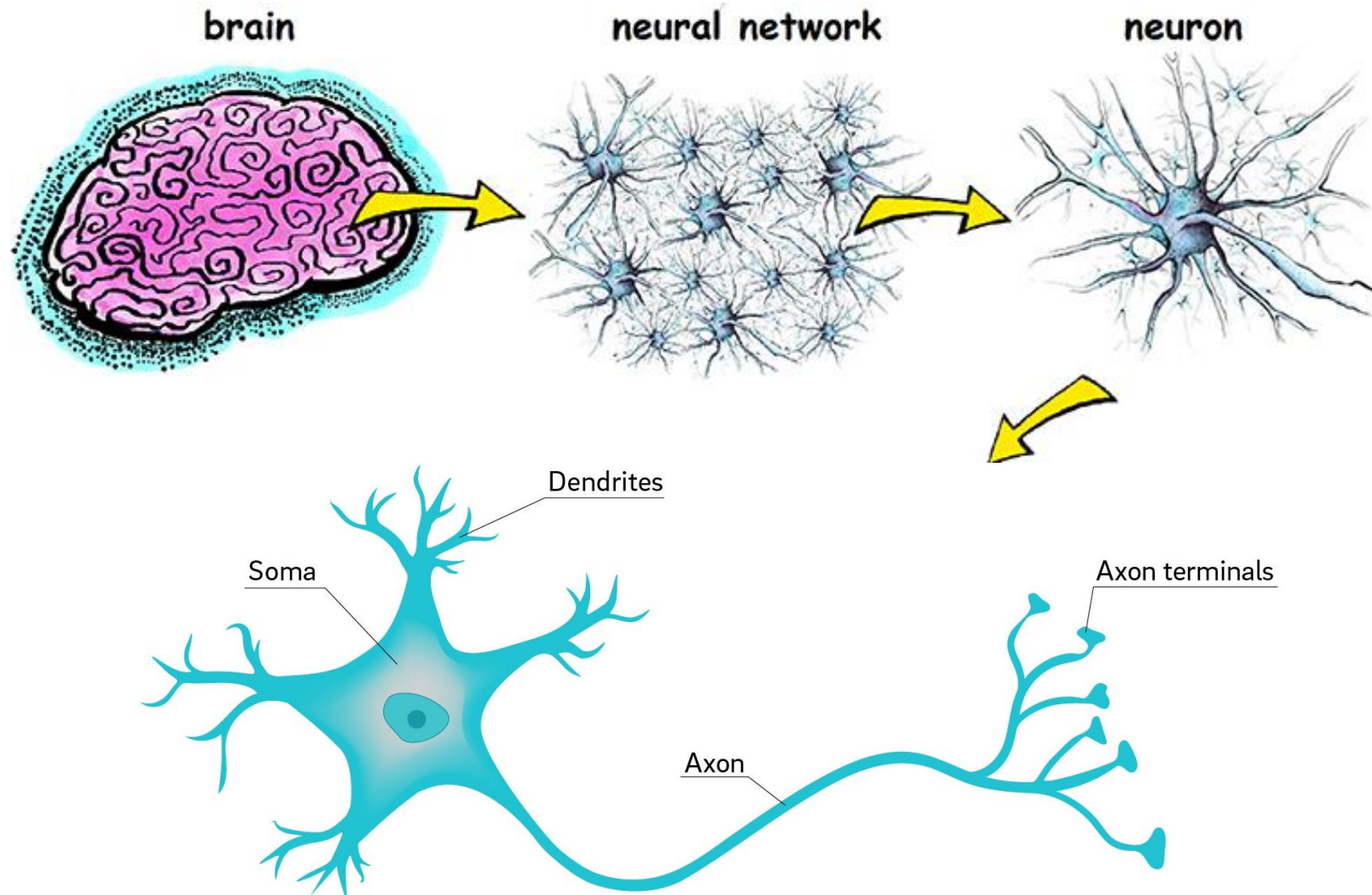
Outline

- Perceptron
 - Biological inspiration
 - Perceptron Learning Algorithm
- Neural Networks
 - Single-layer Neural Networks
 - Multi-layer Neural Networks
 - Regression and Classification
- Neural Networks with Gradient Descent
- Neural Networks vs Other Models

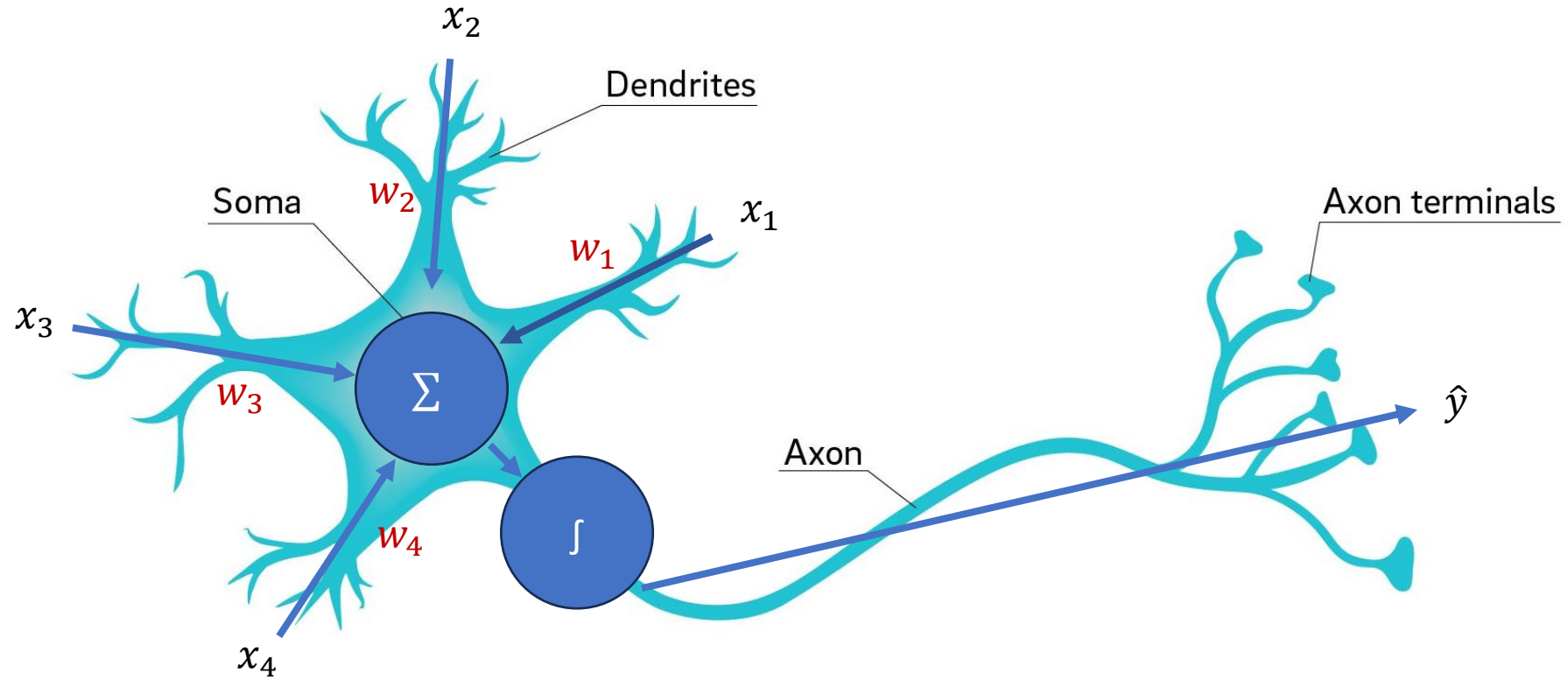
Outline

- **Perceptron**
 - Biological inspiration
 - Perceptron Learning Algorithm
- Neural Networks
 - Single-layer Neural Networks
 - Multi-layer Neural Networks
 - Regression and Classification
- Neural Networks with Gradient Descent
- Neural Networks vs Other Models

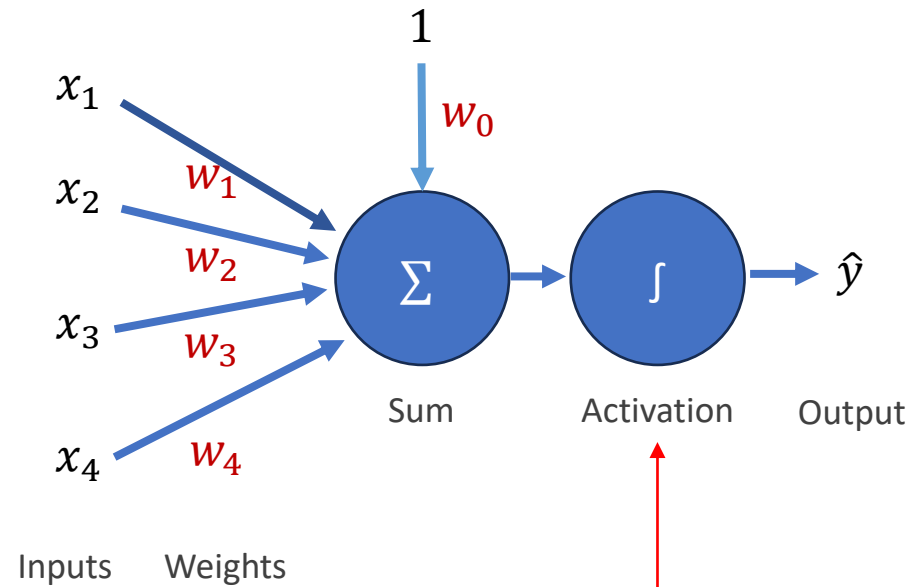
Brain and Neuron



Neuron

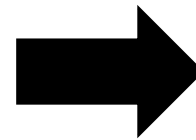


Perceptron



Sign function

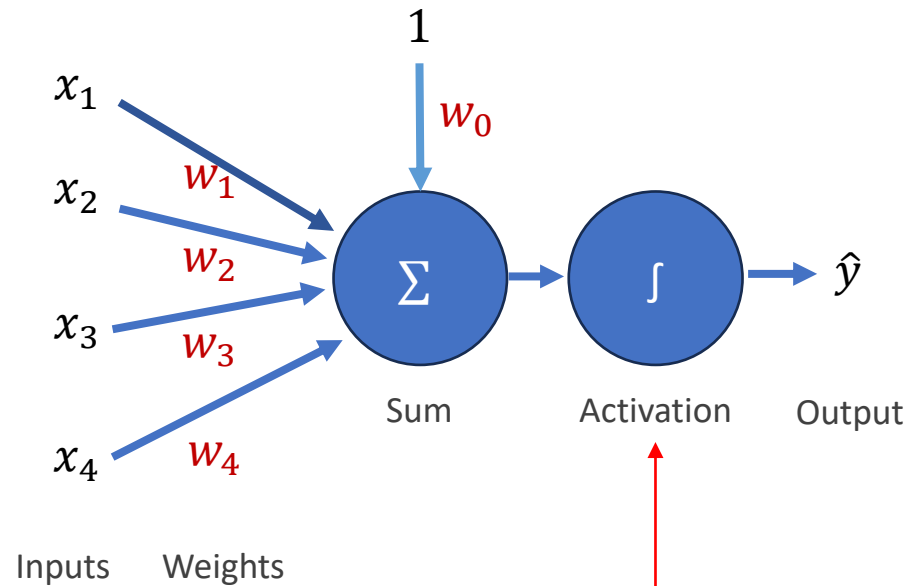
$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$



$$\hat{y} = h_w(x) = g\left(\sum_{i=0}^n w_i x_i\right)$$

Perceptron: An Example

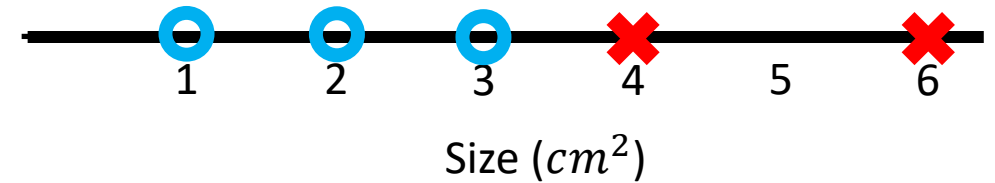
$$\hat{y} = h_{\mathbf{w}}(x) = g\left(\sum_{i=0}^n w_i x_i\right)$$



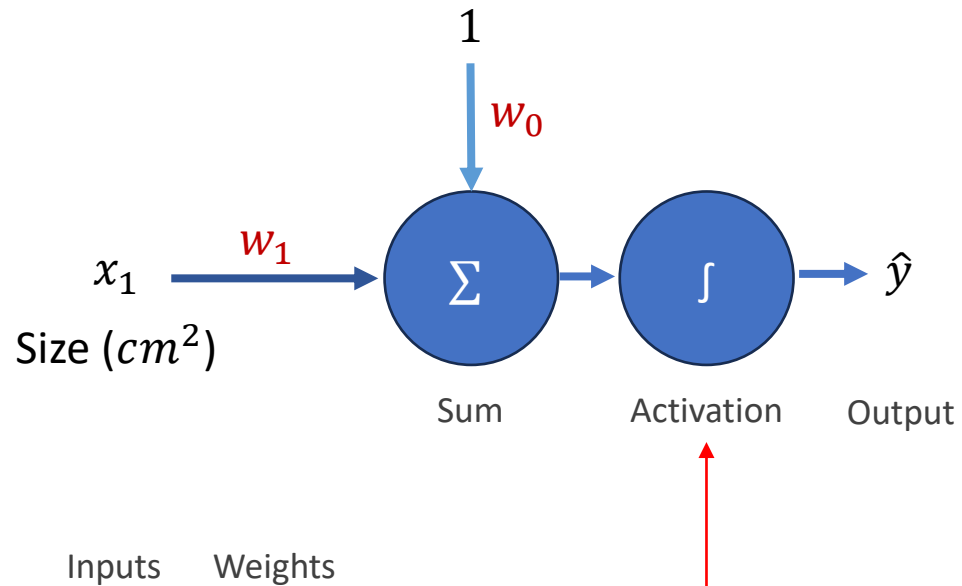
Sign function

$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

Cancer prediction: **benign** (-1), **malignant** (1)



Perceptron: An Example



Sign function

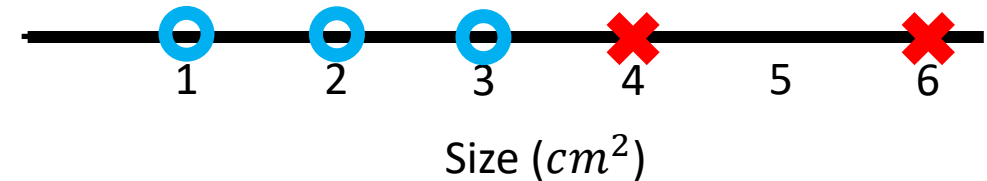
$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

$$\hat{y} = h_w(x) = g\left(\sum_{i=0}^n w_i x_i\right)$$

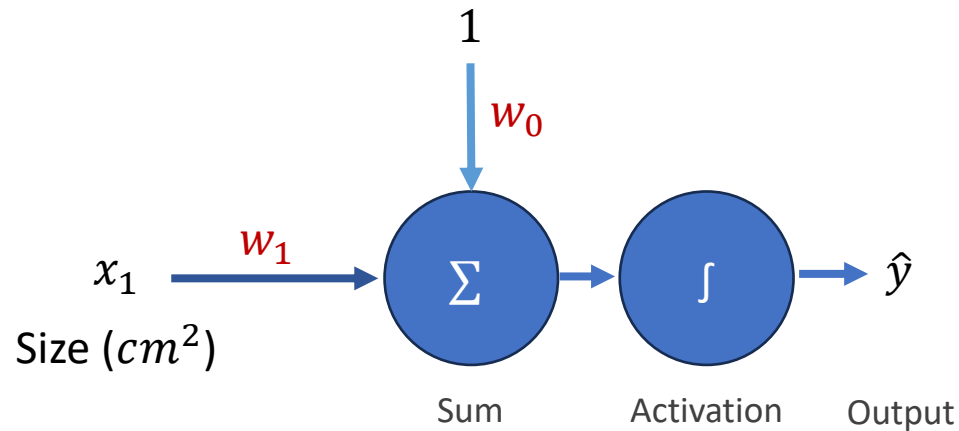
$$h_w(x) = \begin{cases} +1 & \text{if } w_0 + w_1 x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$h_w(x) = \begin{cases} +1 & \text{if } -3.5 + 1x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

Cancer prediction: **benign** (-1), **malignant** (1)



Perceptron: An Example



$$\hat{y} = h_{\mathbf{w}}(x) = g\left(\sum_{i=0}^n w_i x_i\right)$$

$$h_{\mathbf{w}}(x) = \begin{cases} +1 & \text{if } w_0 + w_1 x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$h_{\mathbf{w}}(x) = \begin{cases} +1 & \text{if } -3.5 + 1x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

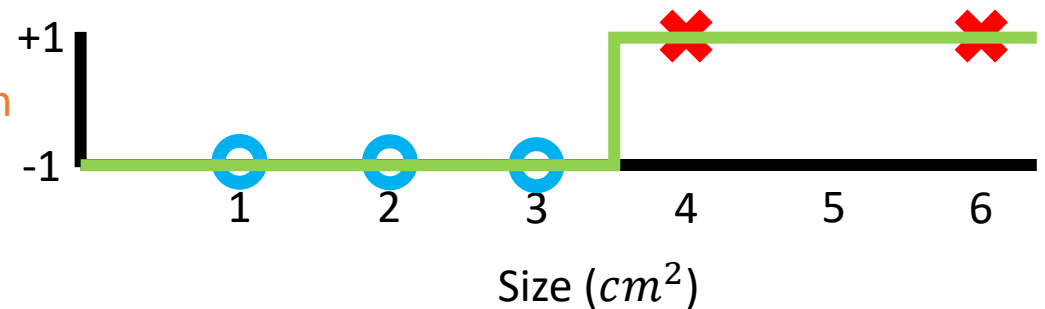
Cancer prediction: **benign** (-1), **malignant** (1)

Inputs Weights

Sign function

$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$


Prediction



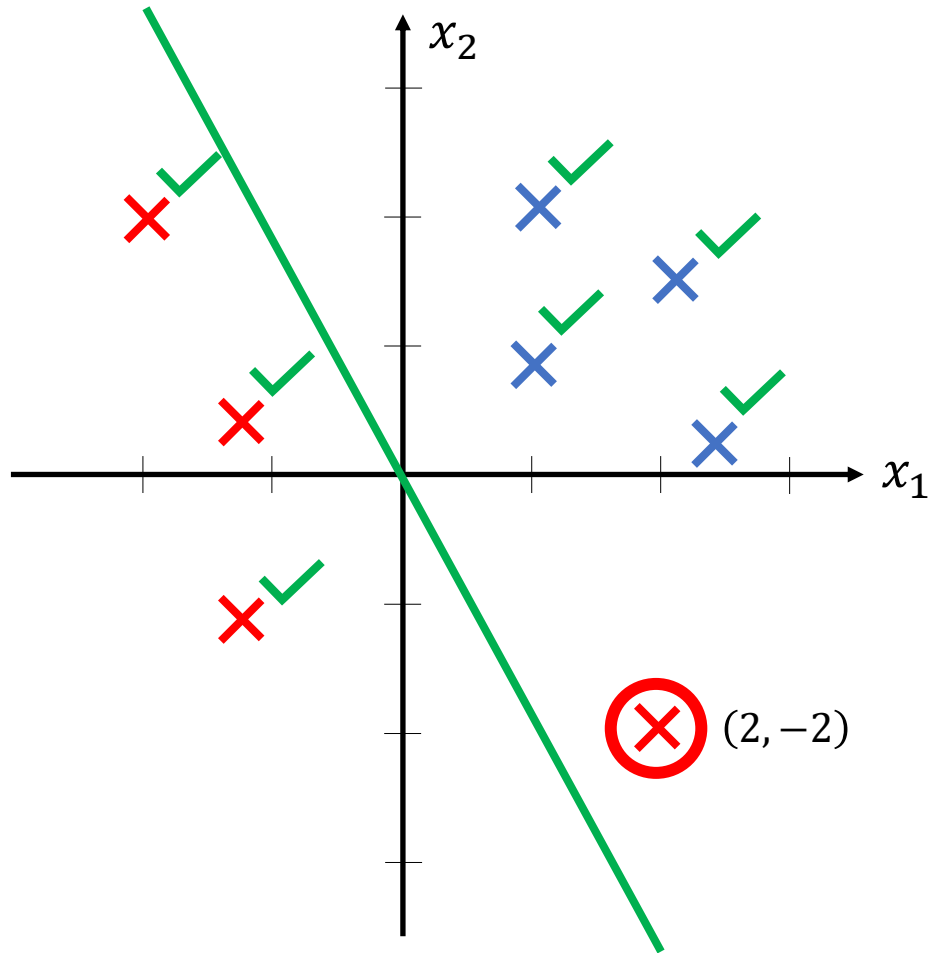
How do we learn w while g is not differentiable?

Perceptron Learning Algorithm

- Initialize $\forall_i \mathbf{w}_i = 0$
- Loop (until convergence or max steps reached)
 - For each instance $(x^{(i)}, y^{(i)})$, classify $\hat{y}^{(i)} = h_{\mathbf{w}}(x^{(i)})$
 - Select one misclassified instance $(x^{(j)}, y^{(j)})$
 - Update weights: $\mathbf{w} \leftarrow \mathbf{w} + \gamma(y^{(j)} - \hat{y}^{(j)})x^{(j)}$


Learning rate

Perceptron Learning Algorithm: An Example



Sign function

$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

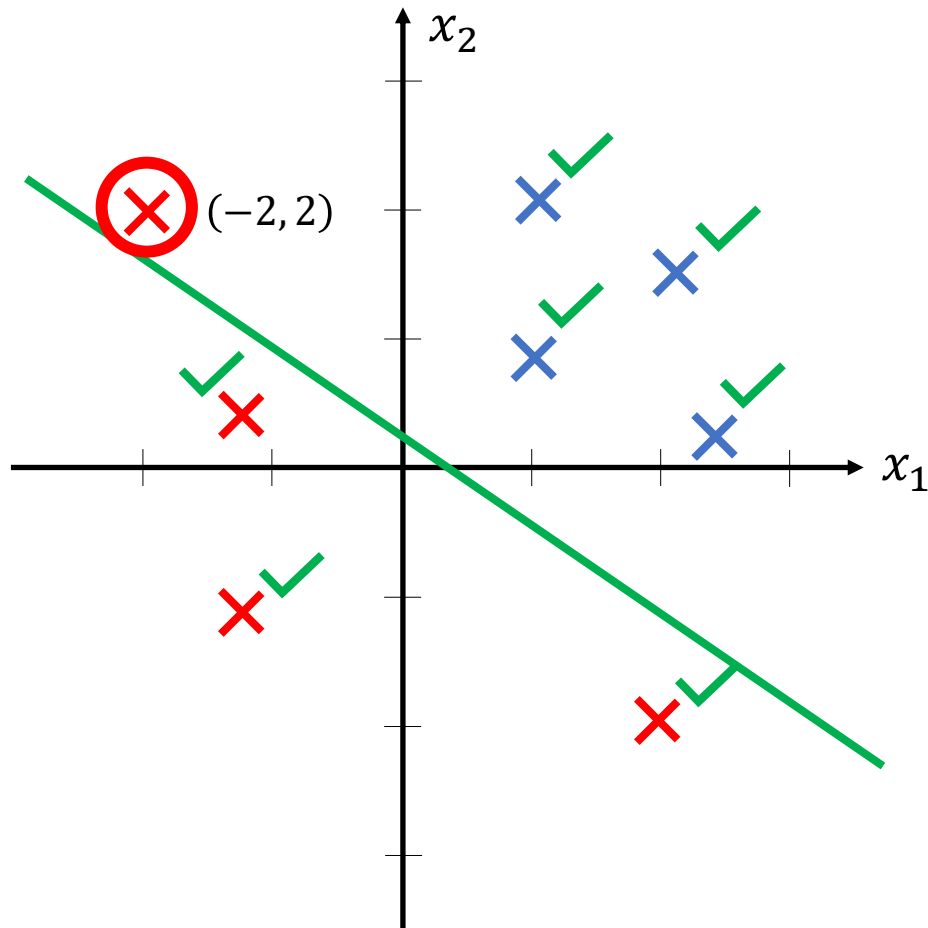
$$\begin{aligned} \hat{y} = h_w(x) &= g(w_0 + w_1 x_1 + w_2 x_2) \\ &= g(0 + 1x_1 + 0.5x_2) \end{aligned}$$

$$w \leftarrow w + \gamma(y^{(i)} - \hat{y}^{(i)}) x^{(i)}$$

$$\begin{bmatrix} 0 \\ 1 \\ 0.5 \end{bmatrix} + 0.1(-1 - 1) \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix} = \begin{bmatrix} -0.2 \\ 0.6 \\ 0.9 \end{bmatrix}$$

$$\text{New } h_w(x) = g(-0.2 + 0.6x_1 + 0.9x_2)$$

Perceptron Learning Algorithm: An Example



Sign function

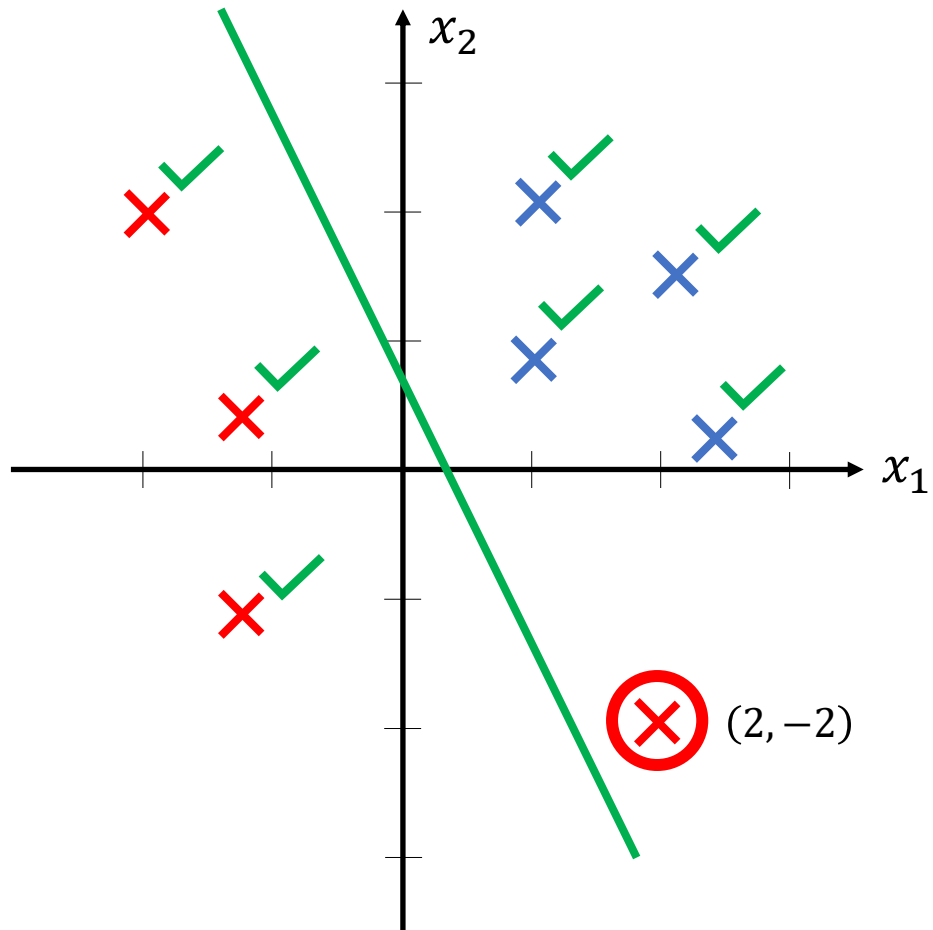
$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

$$\begin{aligned} \hat{y} = h_w(x) &= g(w_0 + w_1x_1 + w_2x_2) \\ &= g(-0.2 + 0.6x_1 + 0.9x_2) \end{aligned}$$

$$\begin{aligned} w &\leftarrow w + \gamma(y^{(i)} - \hat{y}^{(i)})x^{(i)} \\ \begin{bmatrix} -0.2 \\ 0.6 \\ 0.9 \end{bmatrix} &+ 0.1(-1 - 1) \begin{bmatrix} 1 \\ -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -0.4 \\ 1 \\ 0.5 \end{bmatrix} \end{aligned}$$

$$\text{New } h_w(x) = g(-0.4 + 1x_1 + 0.5x_2)$$

Perceptron Learning Algorithm: An Example



Sign function

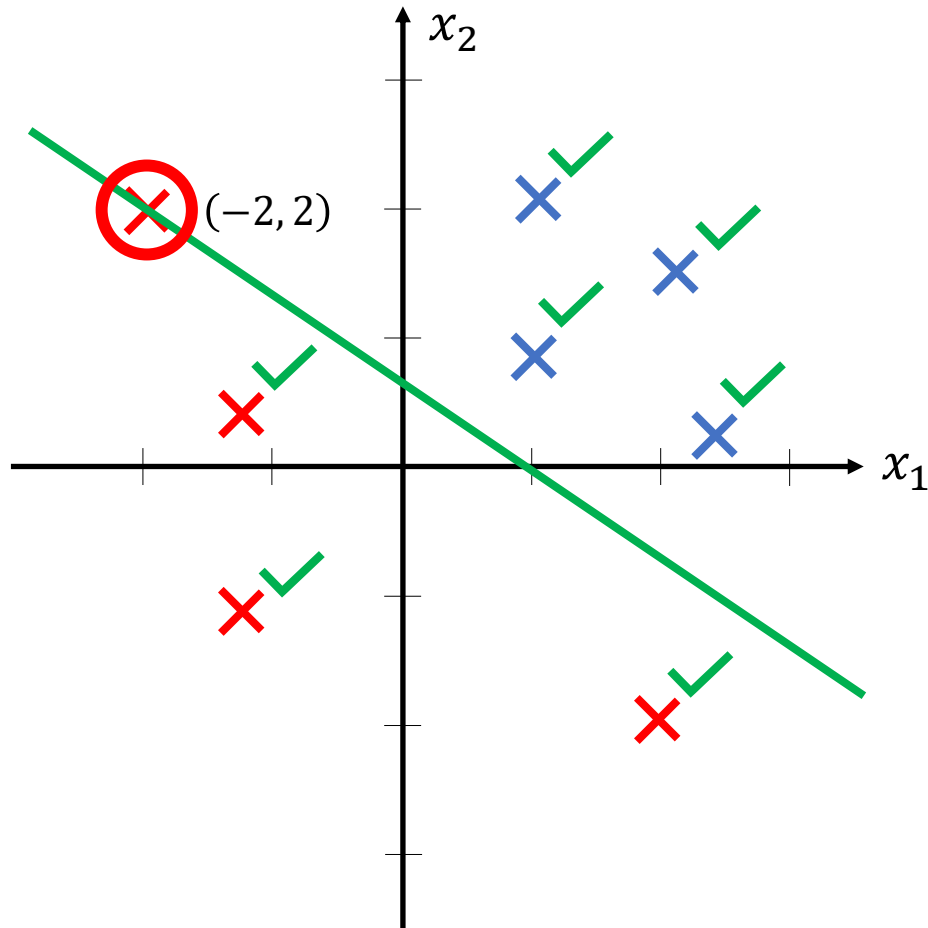
$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

$$\begin{aligned} \hat{y} = h_w(x) &= g(w_0 + w_1 x_1 + w_2 x_2) \\ &= g(-0.4 + 1x_1 + 0.5x_2) \end{aligned}$$

$$\begin{aligned} w &\leftarrow w + \gamma(y^{(i)} - \hat{y}^{(i)}) x^{(i)} \\ &\begin{bmatrix} -0.4 \\ 1 \\ 0.5 \end{bmatrix} + 0.1(-1 - 1) \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix} = \begin{bmatrix} -0.6 \\ 0.6 \\ 0.9 \end{bmatrix} \end{aligned}$$

$$\text{New } h_w(x) = g(-0.6 + 0.6x_1 + 0.9x_2)$$

Perceptron Learning Algorithm: An Example



Sign function

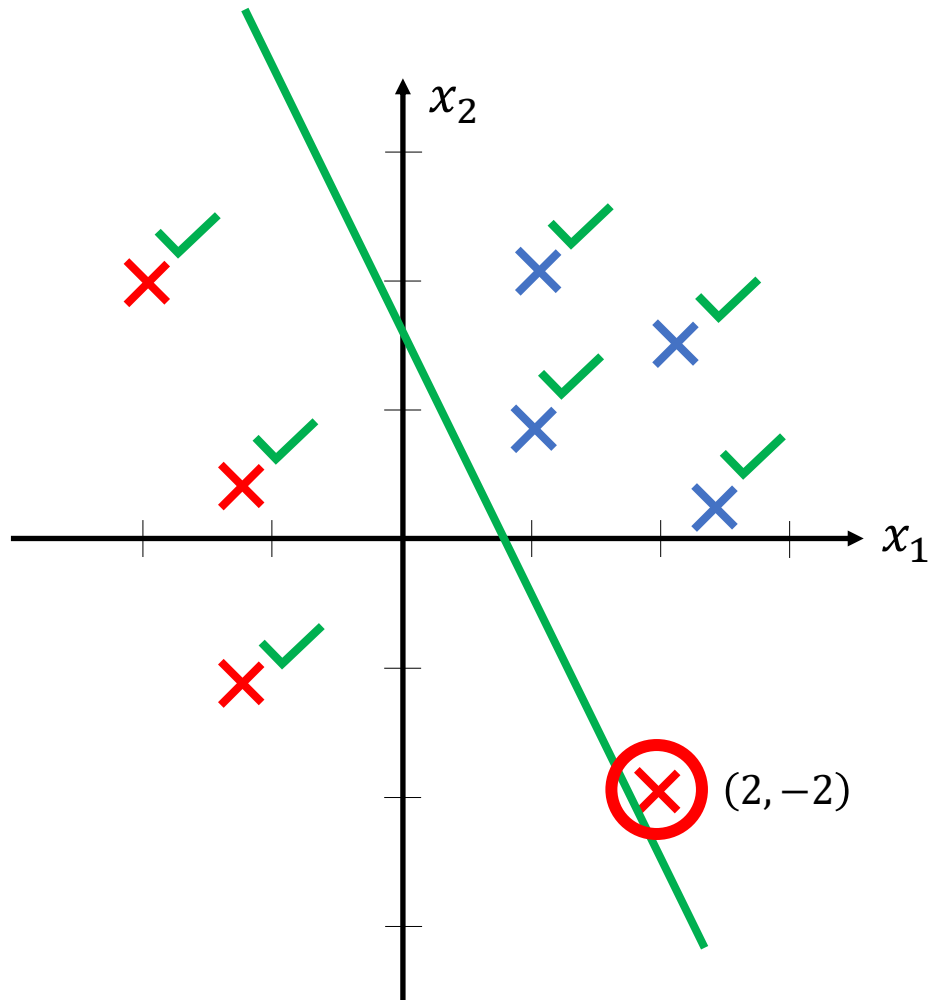
$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

$$\begin{aligned} \hat{y} = h_w(x) &= g(w_0 + w_1 x_1 + w_2 x_2) \\ &= g(-0.6 + 0.6x_1 + 0.9x_2) \end{aligned}$$

$$\begin{aligned} w &\leftarrow w + \gamma(y^{(i)} - \hat{y}^{(i)}) x^{(i)} \\ \begin{bmatrix} -0.6 \\ 0.6 \\ 0.9 \end{bmatrix} &+ 0.1(-1 - 1) \begin{bmatrix} 1 \\ -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -0.8 \\ 1 \\ 0.5 \end{bmatrix} \end{aligned}$$

$$\text{New } h_w(x) = g(-0.8 + 1x_1 + 0.5x_2)$$

Perceptron Learning Algorithm: An Example



Sign function

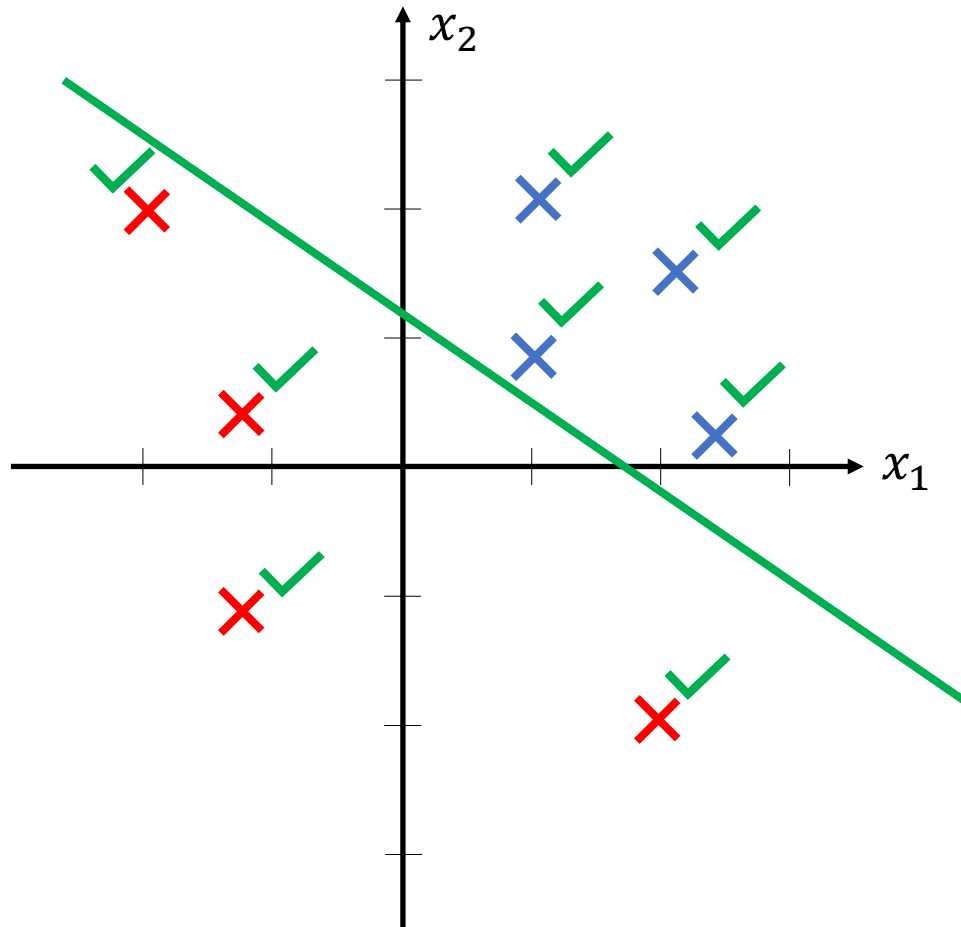
$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

$$\begin{aligned} \hat{y} = h_w(x) &= g(w_0 + w_1x_1 + w_2x_2) \\ &= g(-0.8 + 1x_1 + 0.5x_2) \end{aligned}$$

$$\begin{aligned} w &\leftarrow w + \gamma(y^{(i)} - \hat{y}^{(i)})x^{(i)} \\ \begin{bmatrix} -0.8 \\ 1 \\ 0.5 \end{bmatrix} &+ 0.1(-1 - 1) \begin{bmatrix} 1 \\ 2 \\ -2 \end{bmatrix} = \begin{bmatrix} -1 \\ 0.6 \\ 0.9 \end{bmatrix} \end{aligned}$$

$$\text{New } h_w(x) = g(-1 + 0.6x_1 + 0.9x_2)$$

Perceptron Learning Algorithm: An Example



Sign function

$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

$$\begin{aligned} \hat{y} = h_w(x) &= g(w_0 + w_1x_1 + w_2x_2) \\ &= g(-1 + 0.6x_1 + 0.9x_2) \end{aligned}$$

No misclassifications! Converged!

What if it's not linearly separable?
The algorithm will not converge

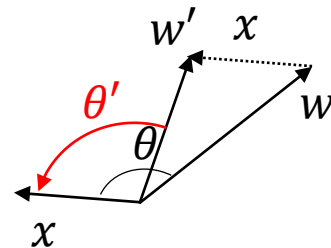
Perceptron Learning Algorithm: Why?

$$\hat{y} = g\left(\sum_{i=0}^n w_i x_i\right) \quad g(z) = \begin{cases} +1 & z \geq 0 \\ -1 & z < 0 \end{cases}$$

When there is a misclassification: $y = +1, \hat{y} = -1$, or vice versa

Doesn't affect the sign

$$\begin{aligned} \hat{y} &= -1 \\ w^T x &< 0 \\ w \cdot x &< 0 \\ ||w|| \cdot ||x|| \cos \theta &< 0 \\ \cos \theta &< 0 \\ \frac{\pi}{2} < \theta &< \pi \end{aligned}$$



$$w' \leftarrow w + x$$

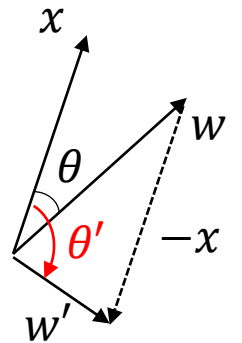
θ' will be smaller
(as required)

What we want:

$$\begin{aligned} y &= +1 \\ \cos \theta &> 0 \\ 0 < \theta &< \frac{\pi}{2} \end{aligned}$$

Doesn't affect the sign

$$\begin{aligned} \hat{y} &= +1 \\ w^T x &> 0 \\ w \cdot x &> 0 \\ ||w|| \cdot ||x|| \cos \theta &> 0 \\ \cos \theta &> 0 \\ 0 < \theta &< \frac{\pi}{2} \end{aligned}$$



$$w' \leftarrow w - x$$

θ' will be larger
(as required)

What we want:

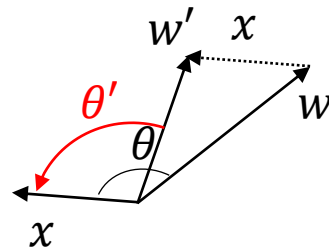
$$\begin{aligned} y &= -1 \\ \cos \theta &< 0 \\ \frac{\pi}{2} < \theta &< \pi \end{aligned}$$

Perceptron Learning Algorithm: Why?

$$\hat{y} = g\left(\sum_{i=0}^n w_i x_i\right) \quad g(z) = \begin{cases} +1 & z \geq 0 \\ -1 & z < 0 \end{cases}$$

When there is a misclassification: $y = +1, \hat{y} = -1$, or vice versa

$$\begin{aligned} \hat{y} &= -1 \\ y &= +1 \end{aligned}$$

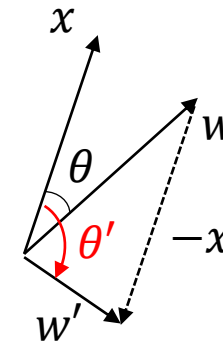


$$w' \leftarrow w + x$$

θ' will be smaller
(as required)

$$\begin{aligned} w &\leftarrow w + \gamma(y - \hat{y})x \\ &\leftarrow w + \gamma(+1 - (-1))x \\ &\leftarrow w + 2\gamma x \end{aligned}$$

$$\begin{aligned} \hat{y} &= +1 \\ y &= -1 \end{aligned}$$



$$w' \leftarrow w - x$$

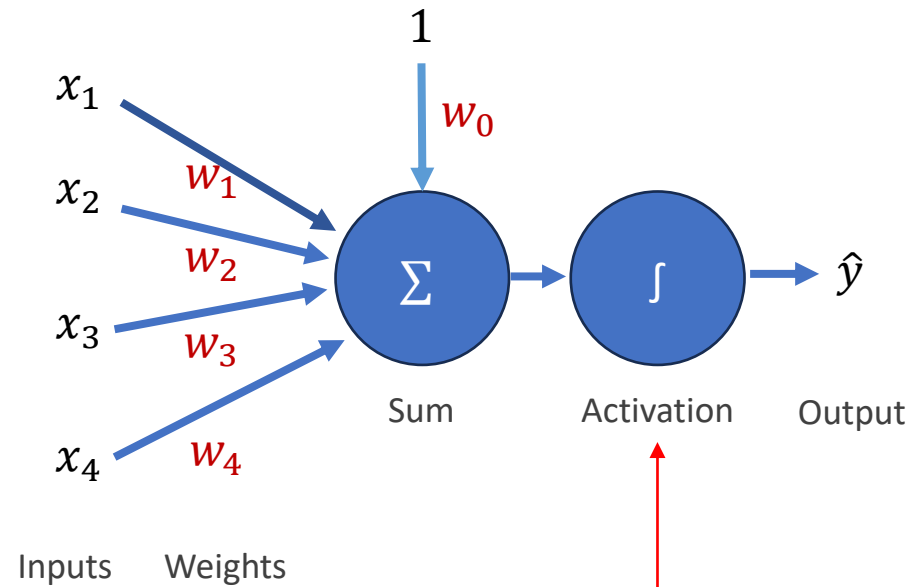
θ' will be larger
(as required)

$$\begin{aligned} w &\leftarrow w + \gamma(y - \hat{y})x \\ &\leftarrow w + \gamma(-1 - (+1))x \\ &\leftarrow w - 2\gamma x \end{aligned}$$

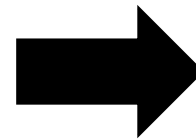
Outline

- Perceptron
 - Biological inspiration
 - Perceptron Learning Algorithm
- **Neural Networks**
 - Single-layer Neural Networks
 - Multi-layer Neural Networks
 - Regression and Classification
- Neural Networks with Gradient Descent
- Neural Networks vs Other Models

Perceptron

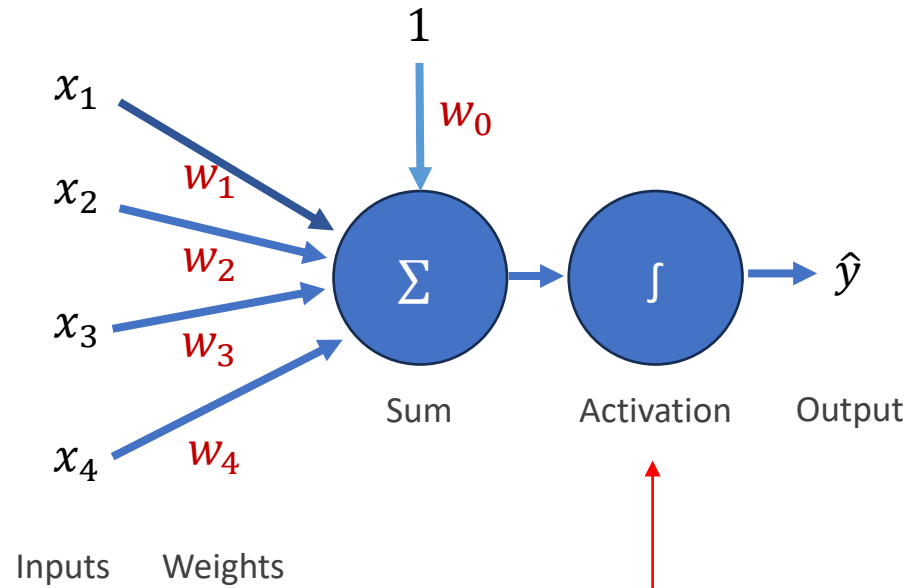


$$g(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

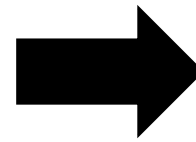


$$\hat{y} = h_w(x) = g\left(\sum_{i=0}^n w_i x_i\right)$$

Single-layer Neural Networks



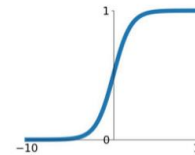
Activation Function
(usually non-linear)



$$\hat{y} = h_w(x) = g\left(\sum_{i=0}^n w_i x_i\right)$$

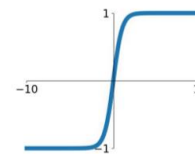
Sigmoid

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



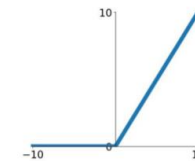
tanh

$$\tanh(x)$$



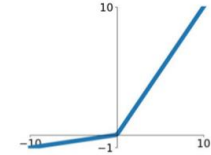
ReLU

$$\max(0, x)$$



Leaky ReLU

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

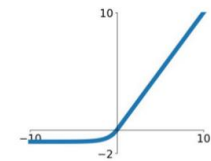


Maxout

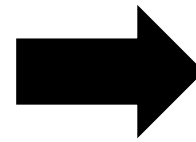
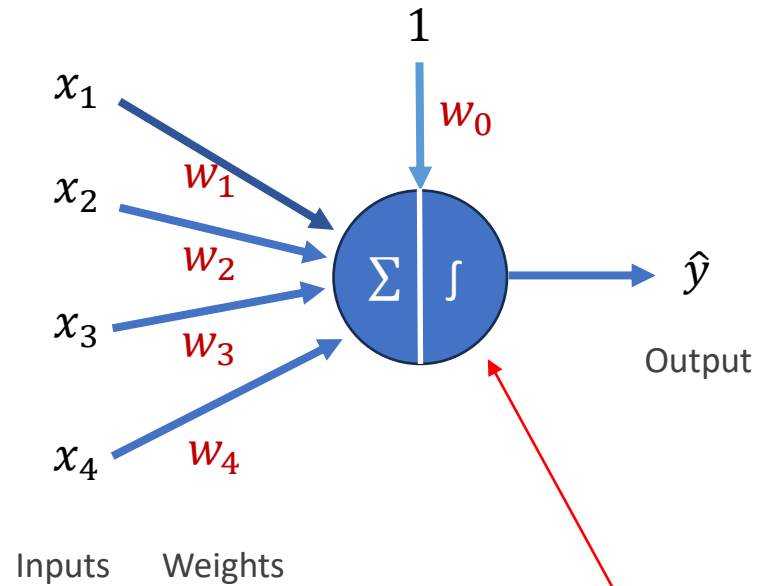
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

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



Single-layer Neural Networks

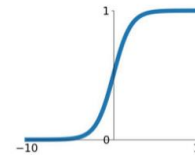


$$\hat{y} = h_w(x) = g\left(\sum_{i=0}^n w_i x_i\right)$$

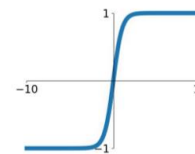
Activation Function
(usually non-linear)

Sigmoid

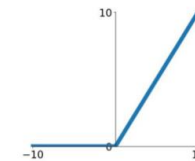
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



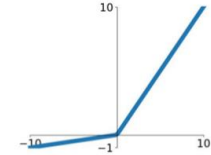
tanh
 $\tanh(x)$



ReLU
 $\max(0, x)$



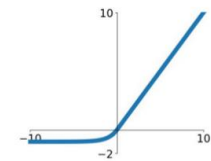
Leaky ReLU
 $\max(0.1x, x)$



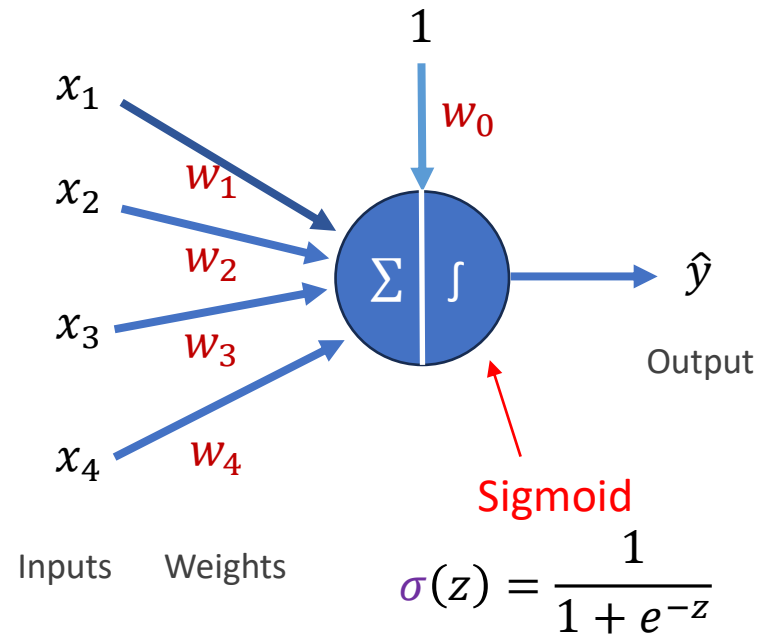
Maxout
 $\max(w_1^T x + b_1, w_2^T x + b_2)$

ELU

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

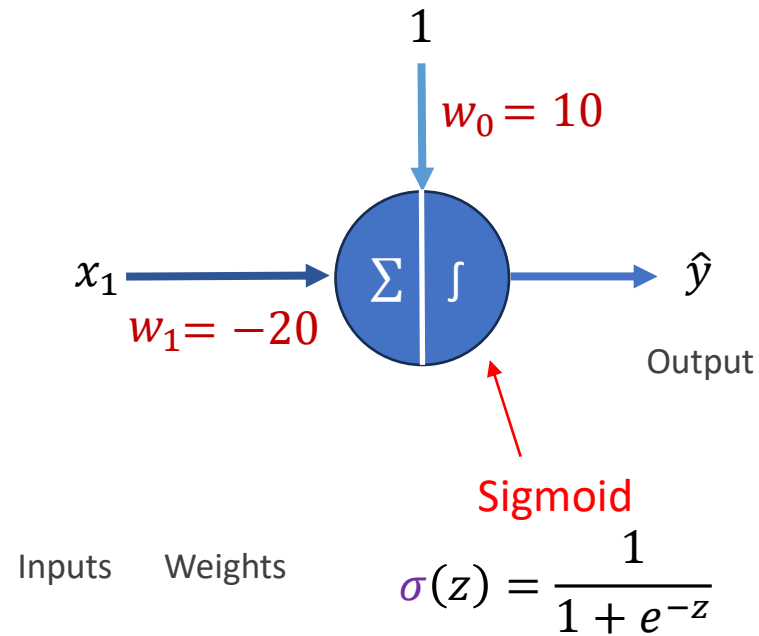


Single-layer Neural Networks



$$\hat{y} = h_{\mathbf{w}}(x) = \sigma(\sum_{i=0}^n w_i x_i)$$

Single-layer Neural Networks: NOT

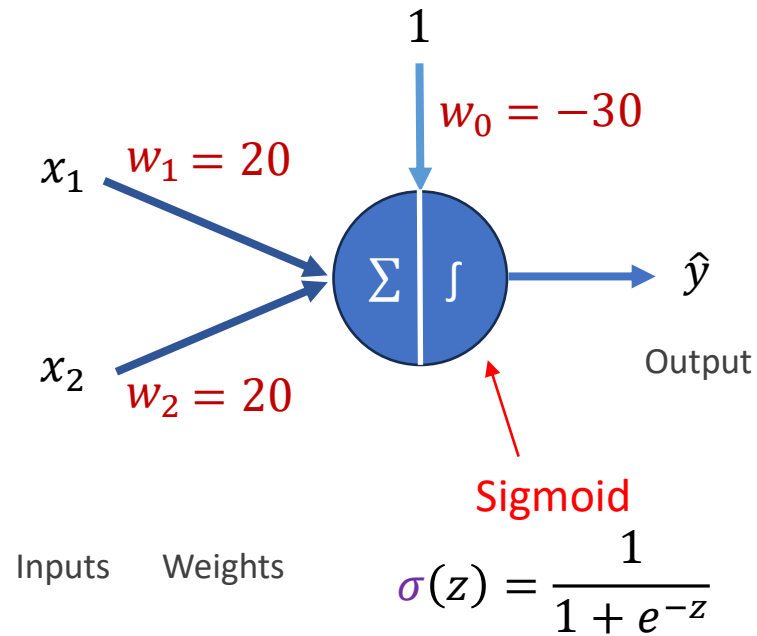


$$\hat{y} = h_{\mathbf{w}}(x) = \sigma(\sum_{i=0}^n w_i x_i)$$

Consider $x_1 \in \{1,0\}$

x_1	y	Σ	\hat{y}
0	1	10	0.999
1	0	-10	0.00004

Single-layer Neural Networks: AND

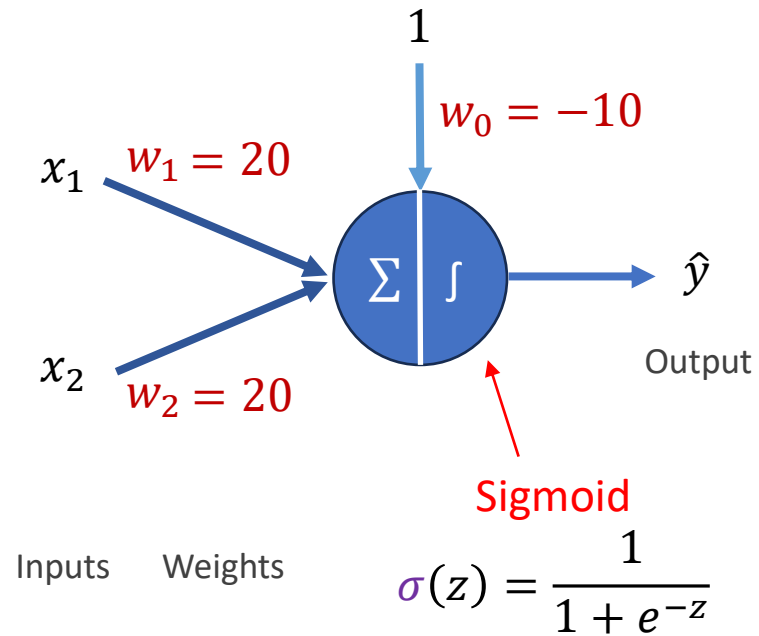


$$\hat{y} = h_{\mathbf{w}}(x) = \sigma\left(\sum_{i=0}^n w_i x_i\right)$$

Consider $x_1, x_2 \in \{1, 0\}$

x_1	x_2	y	Σ	\hat{y}
0	0	0	-30	0.000 ...
0	1	0	-10	0.00004
1	0	0	-10	0.00004
1	1	1	10	0.999

Single-layer Neural Networks: OR



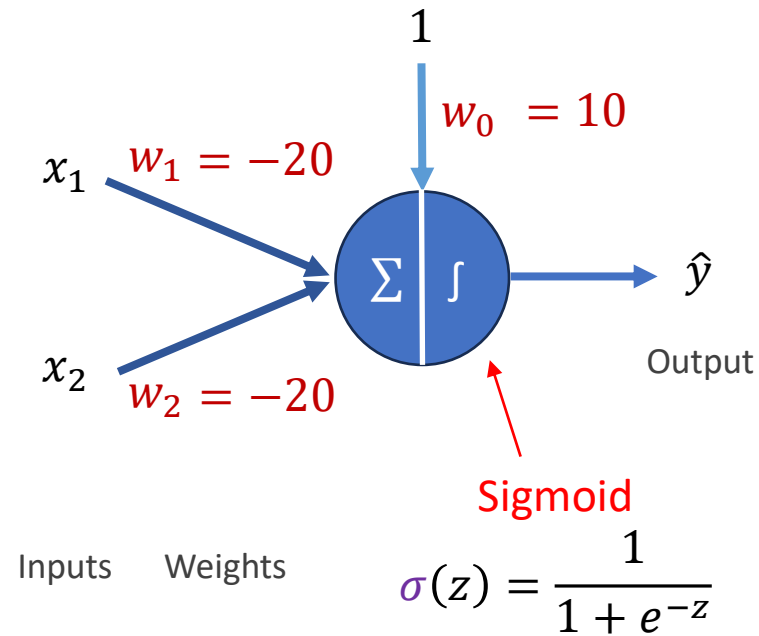
$$\hat{y} = h_{\mathbf{w}}(x) = \sigma\left(\sum_{i=0}^n w_i x_i\right)$$

Consider $x_1, x_2 \in \{1, 0\}$

x_1	x_2	y	Σ	\hat{y}
0	0	0	-10	0.00004
0	1	1	10	0.999
1	0	1	10	0.999
1	1	1	30	0.999

Single-layer Neural Networks: NOR

(NOT x_1) AND (NOT x_2)



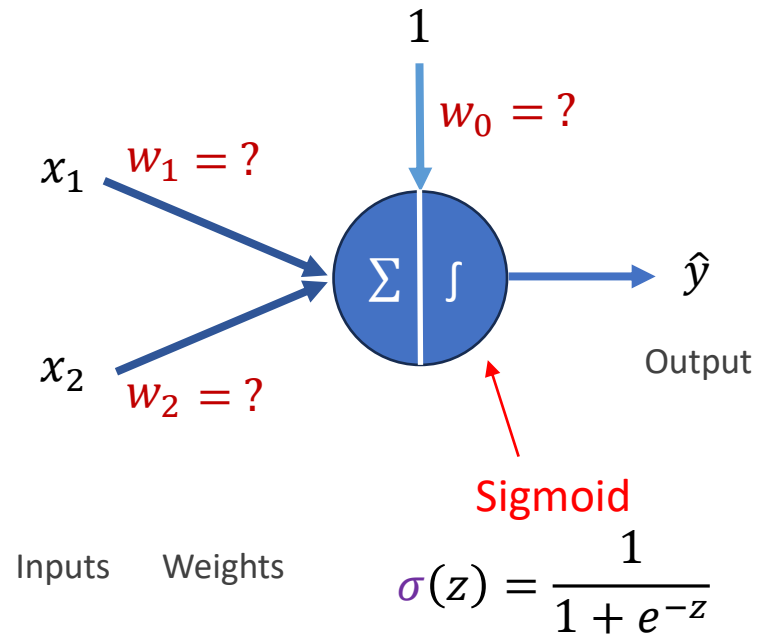
$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \sigma\left(\sum_{i=0}^n w_i x_i\right)$$

Consider $x_1, x_2 \in \{1, 0\}$

x_1	x_2	y	Σ	\hat{y}
0	0	1	10	0.999
0	1	0	-10	0.00004
1	0	0	-10	0.00004
1	1	0	-30	0.000 ...

Single-layer Neural Networks: XNOR

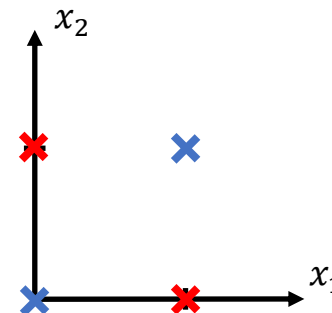
eXclusive Not OR



$$\hat{y} = h_{\mathbf{w}}(x) = \sigma(\sum_{i=0}^n w_i x_i)$$

Consider $x_1, x_2 \in \{1, 0\}$

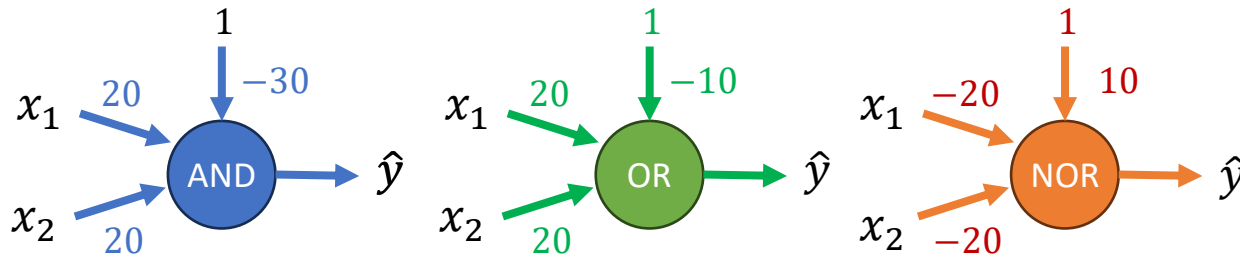
x_1	x_2	y
0	0	1
0	1	0
1	0	0
1	1	1



Not linearly separable!

Single-layer Neural Networks: XNOR

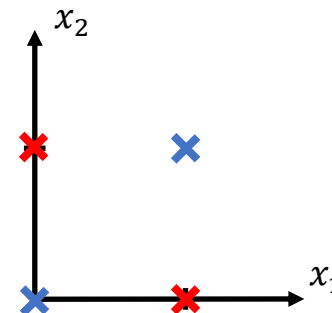
eXclusive Not OR



$$\hat{y} = h_w(x) = \sigma(\sum_{i=0}^n w_i x_i)$$

Consider $x_1, x_2 \in \{1, 0\}$

x_1	x_2	y
0	0	1
0	1	0
1	0	0
1	1	1

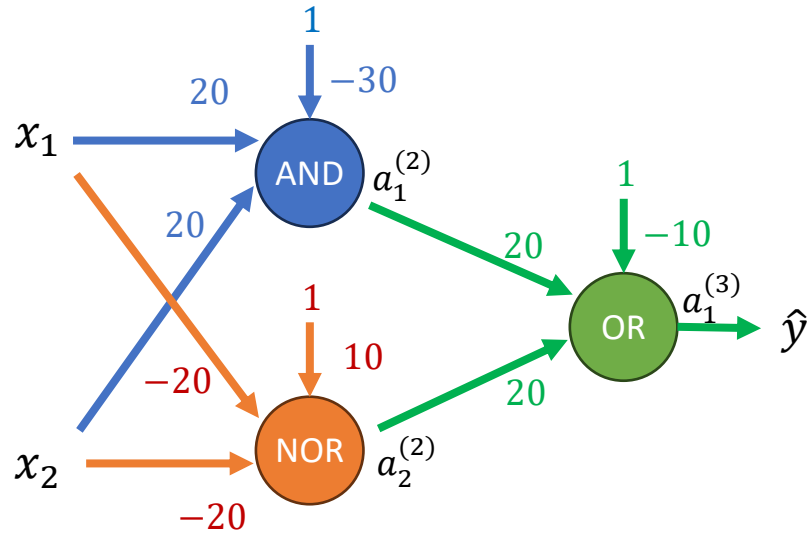


Not linearly separable!

Multi ~~Single~~-layer Neural Networks: XNOR

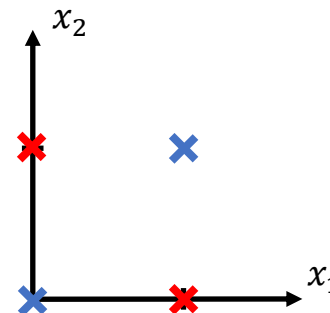
eXclusive Not OR

$$\hat{y} = h_{\text{out}}(r) = \sigma\left(\sum_{i=0}^n w_i x_i\right)$$



Consider $x_1, x_2 \in \{1, 0\}$

x_1	x_2	y
0	0	1
0	1	0
1	0	0
1	1	1



Not linearly separable!

Multi-layer Neural Networks: XNOR

eXclusive Not OR

$$\hat{y} = \sigma \left(\sum_{i=0}^n w_{1i}^{(3)} a_i^{(2)} \right)$$

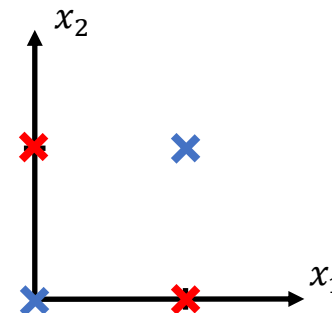
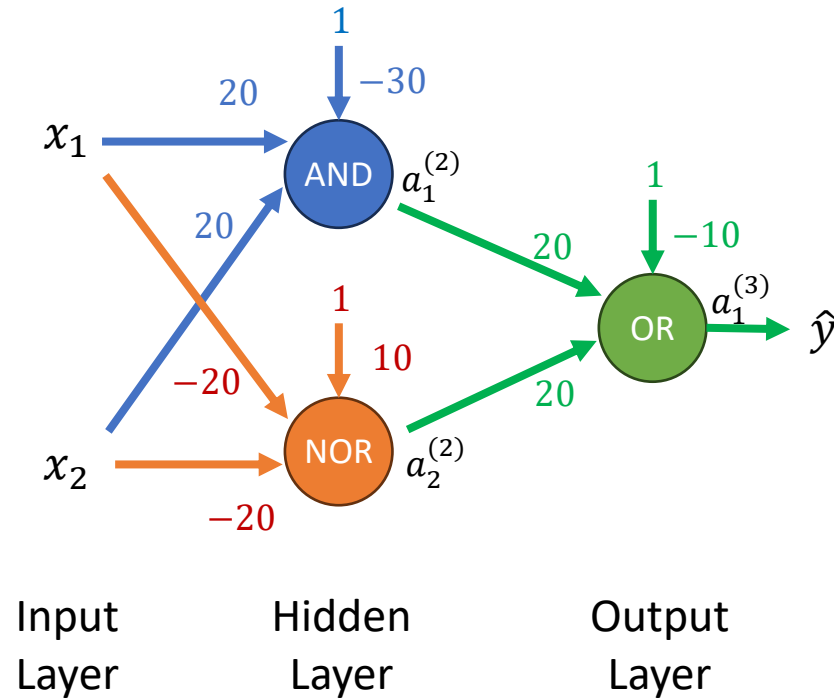
$$a_i^{(2)} = \sigma \left(\sum_{j=0}^n w_{ij}^{(2)} x_j \right)$$

Function composition

$$a_0^{(2)} = 1$$

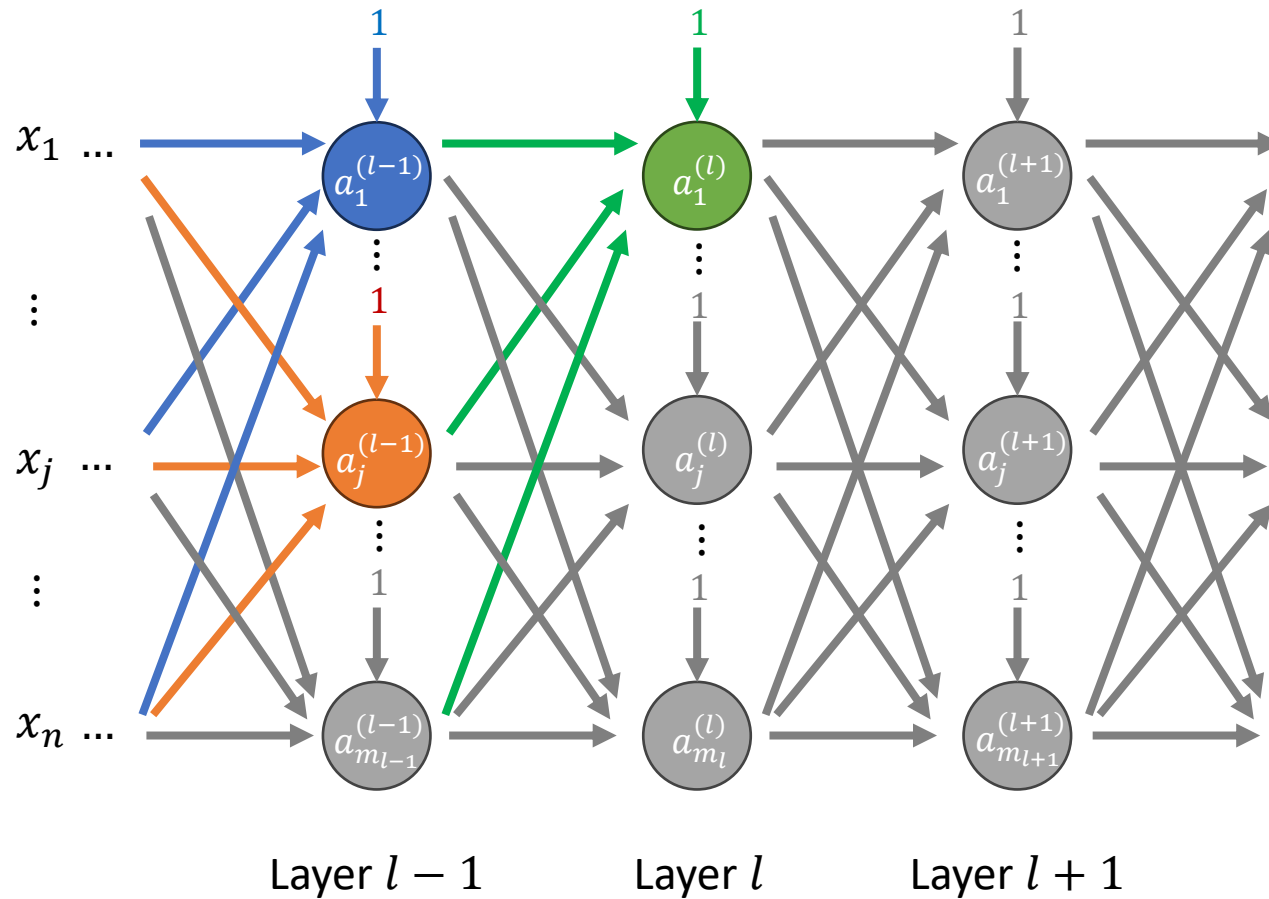
Consider $x_1, x_2 \in \{1, 0\}$

x_1	x_2	y	$a_1^{(2)}$	$a_2^{(2)}$	\hat{y}
0	0	1	0	1	1
0	1	0	0	0	0
1	0	0	0	0	0
1	1	1	1	0	1

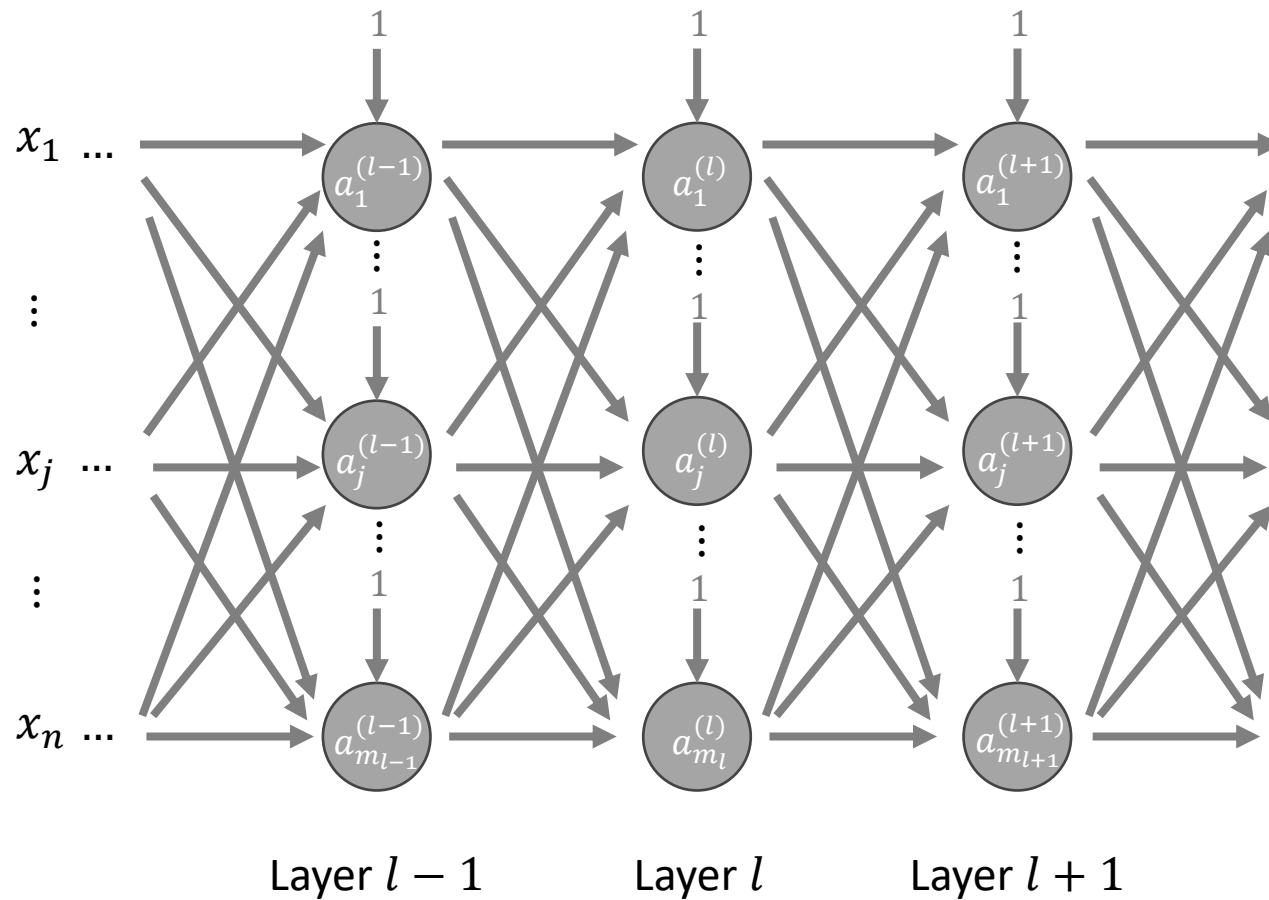


Not linearly separable!

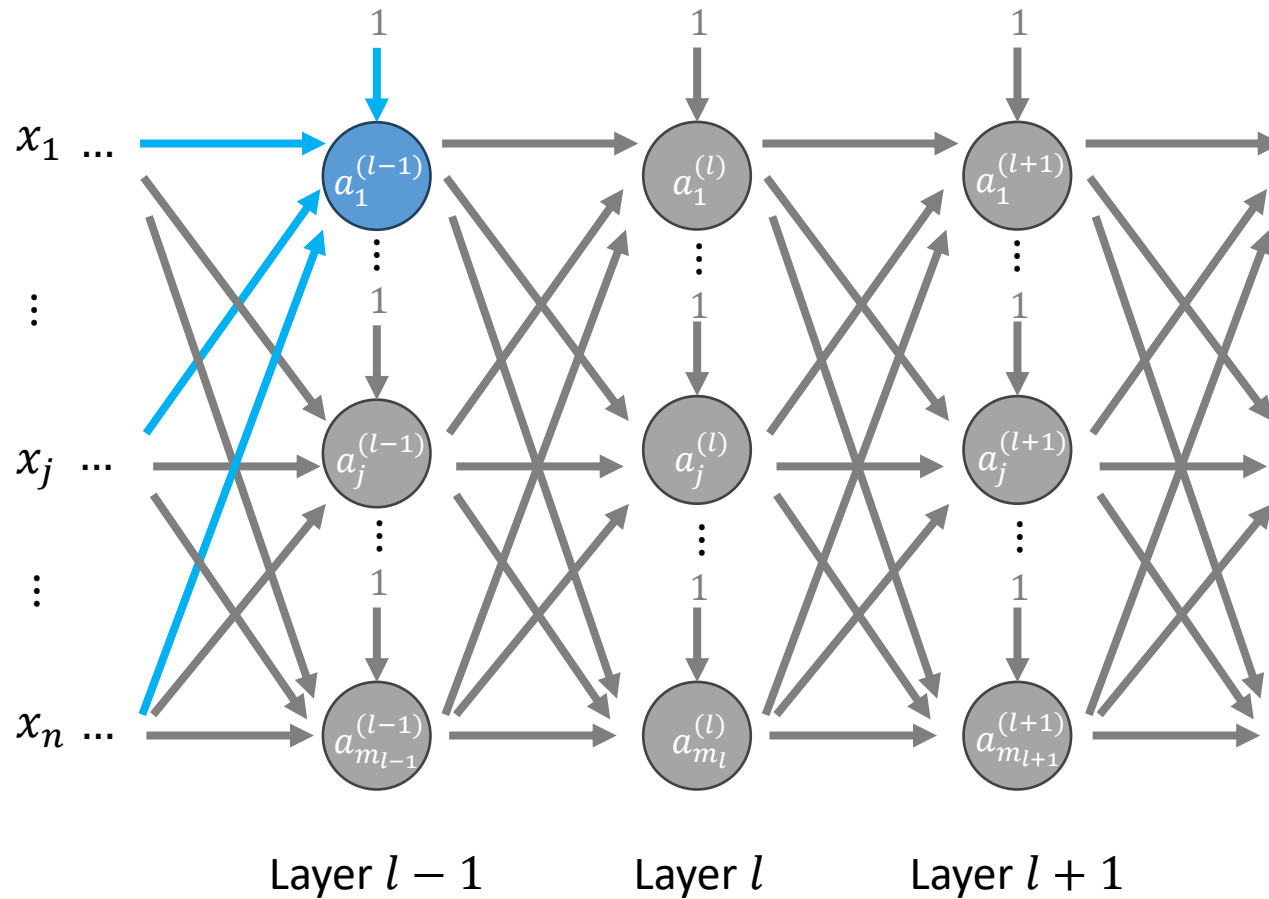
Multi-layer Neural Networks



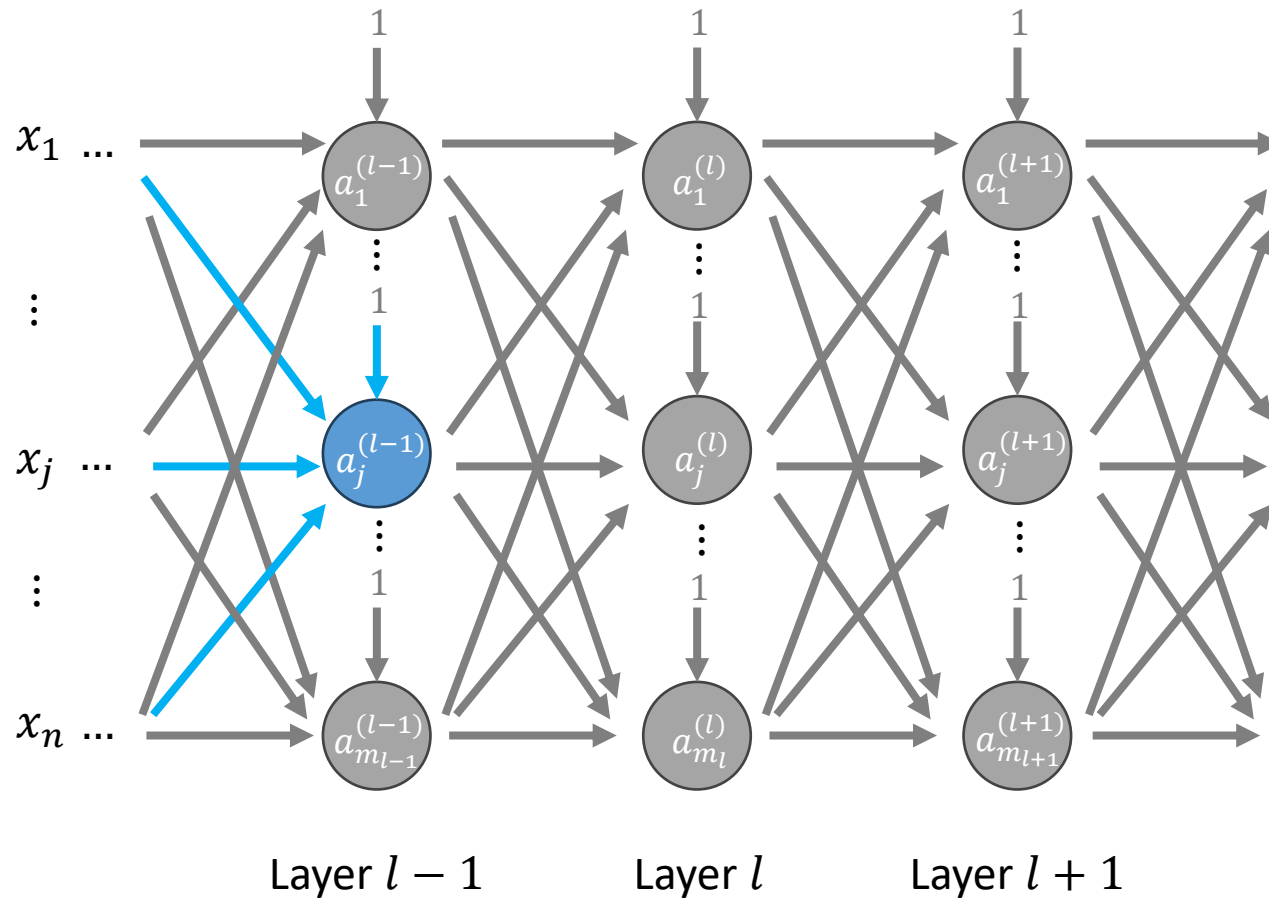
Multi-layer Neural Networks



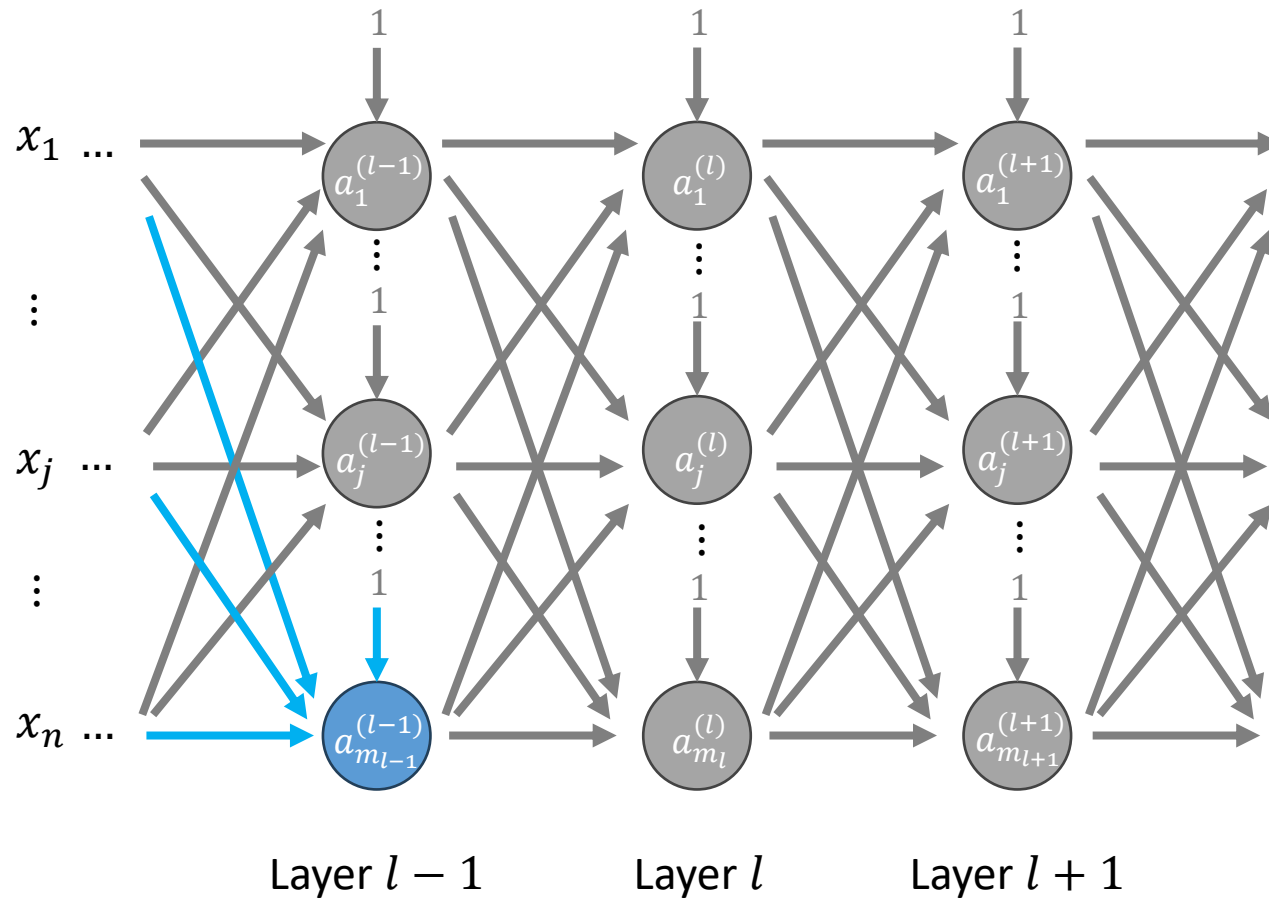
Multi-layer Neural Networks



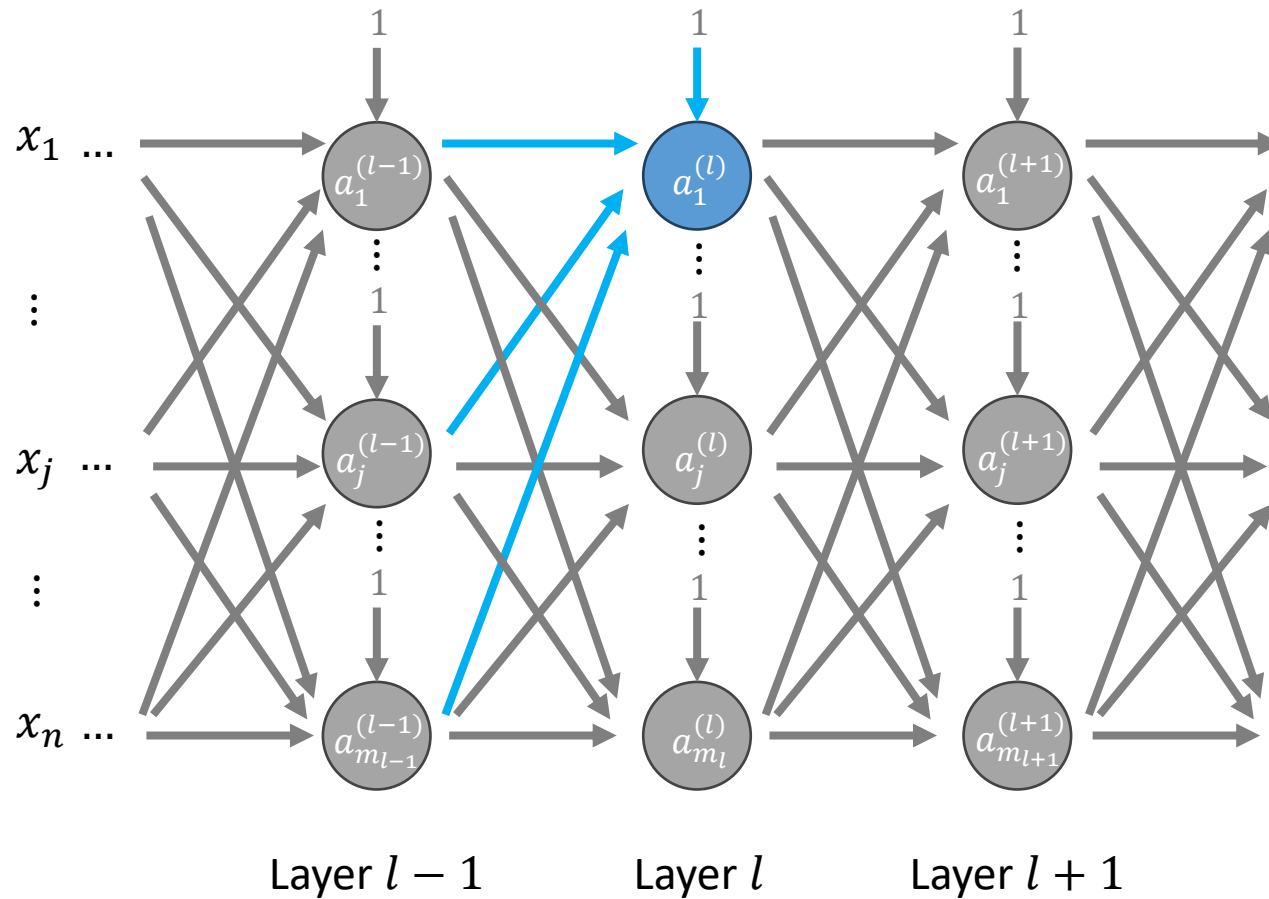
Multi-layer Neural Networks



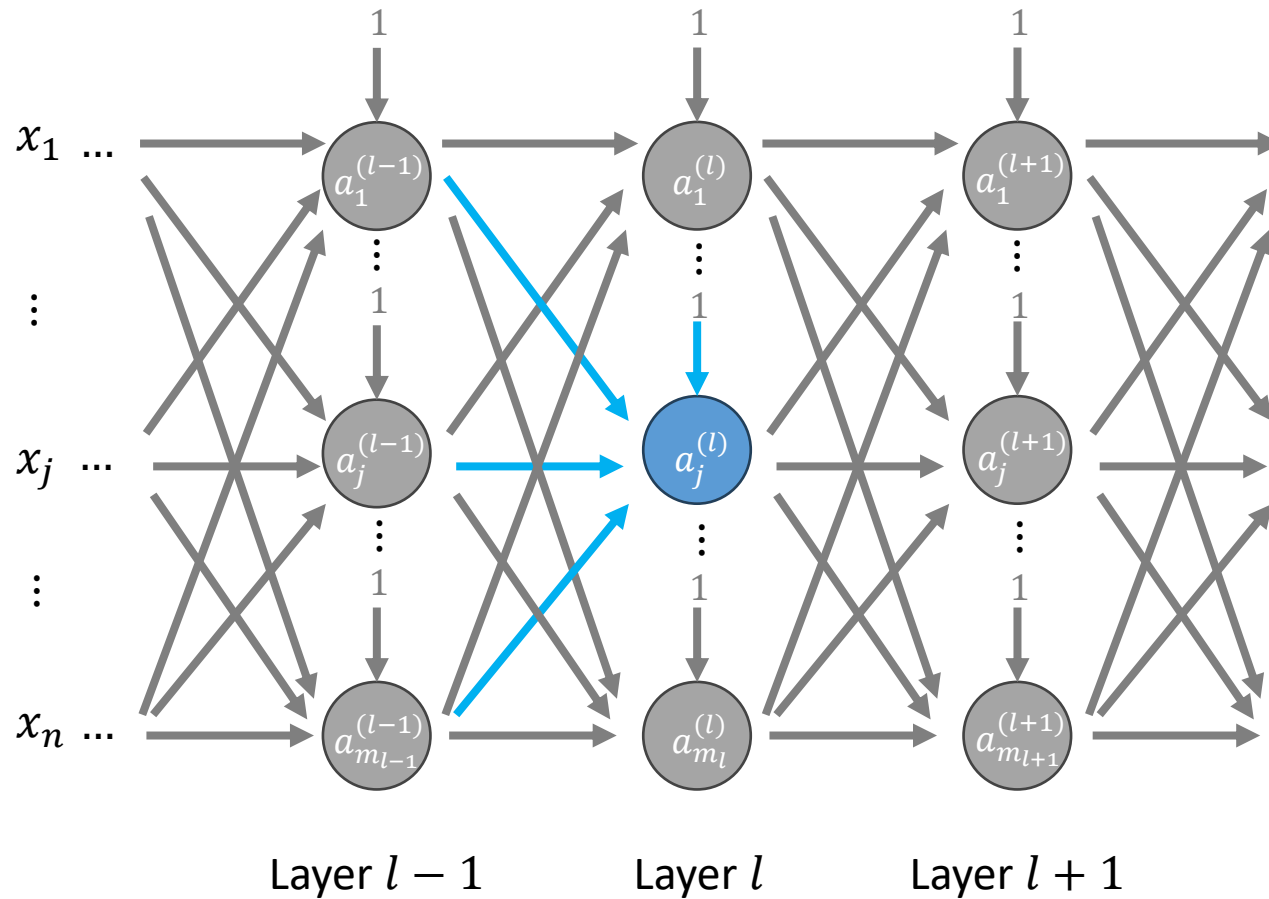
Multi-layer Neural Networks



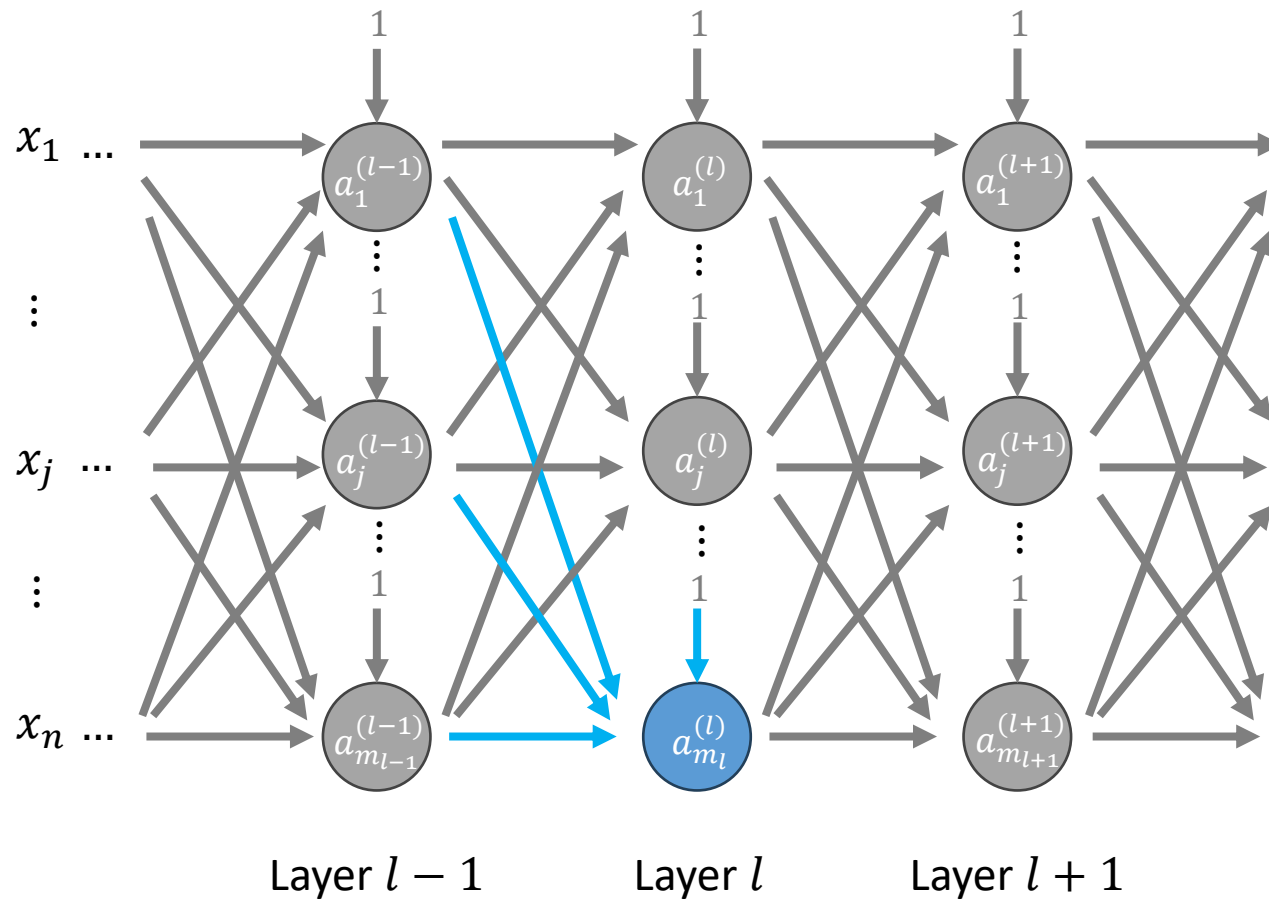
Multi-layer Neural Networks



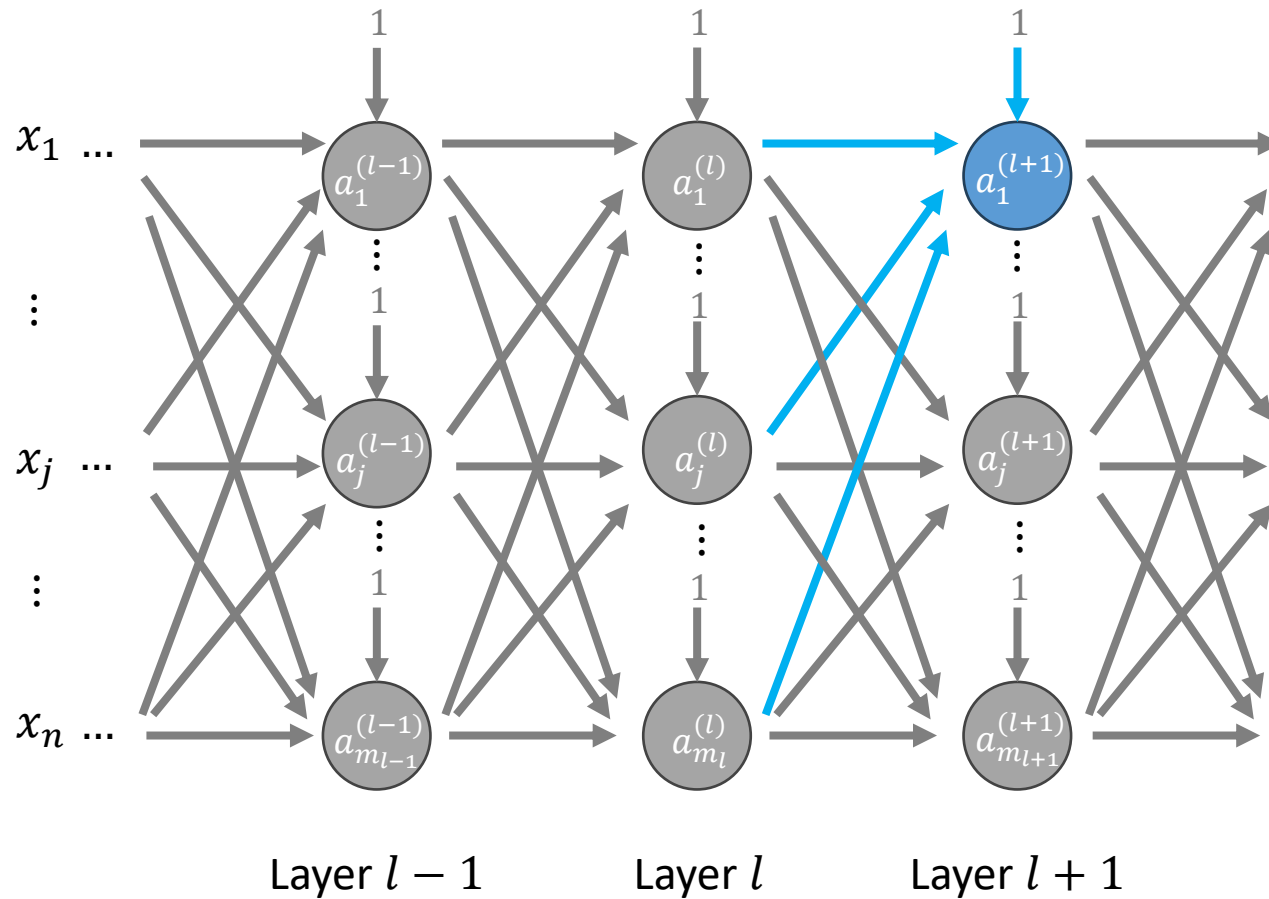
Multi-layer Neural Networks



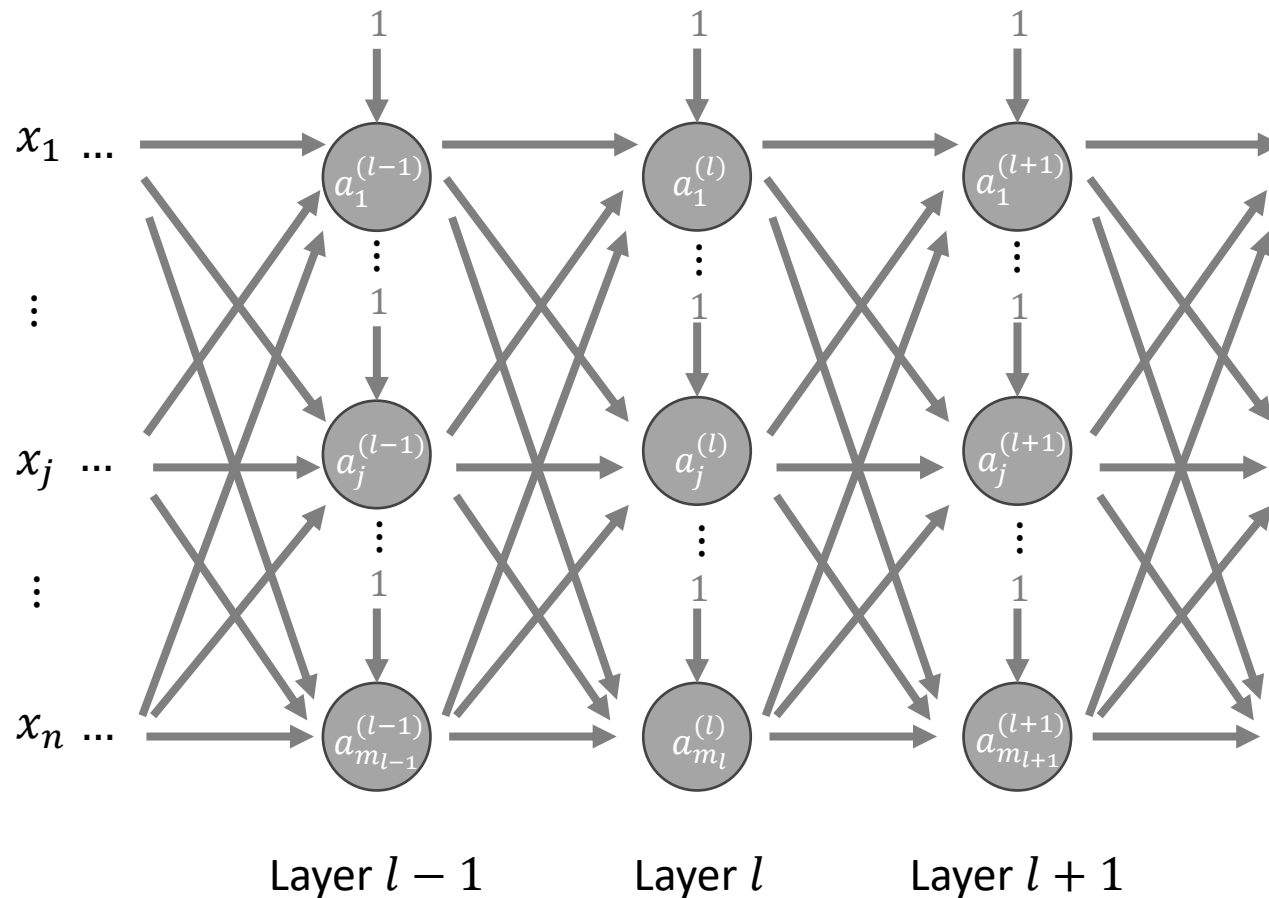
Multi-layer Neural Networks



Multi-layer Neural Networks



Multi-layer Neural Networks

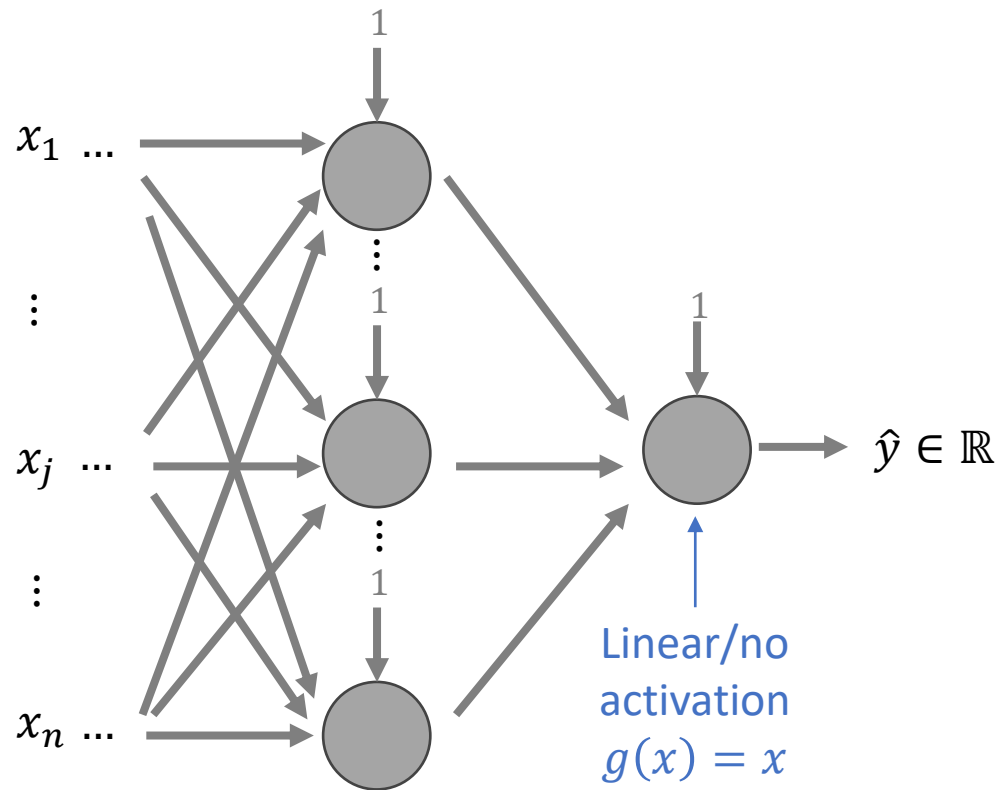


Universal Function Approximation Theorem

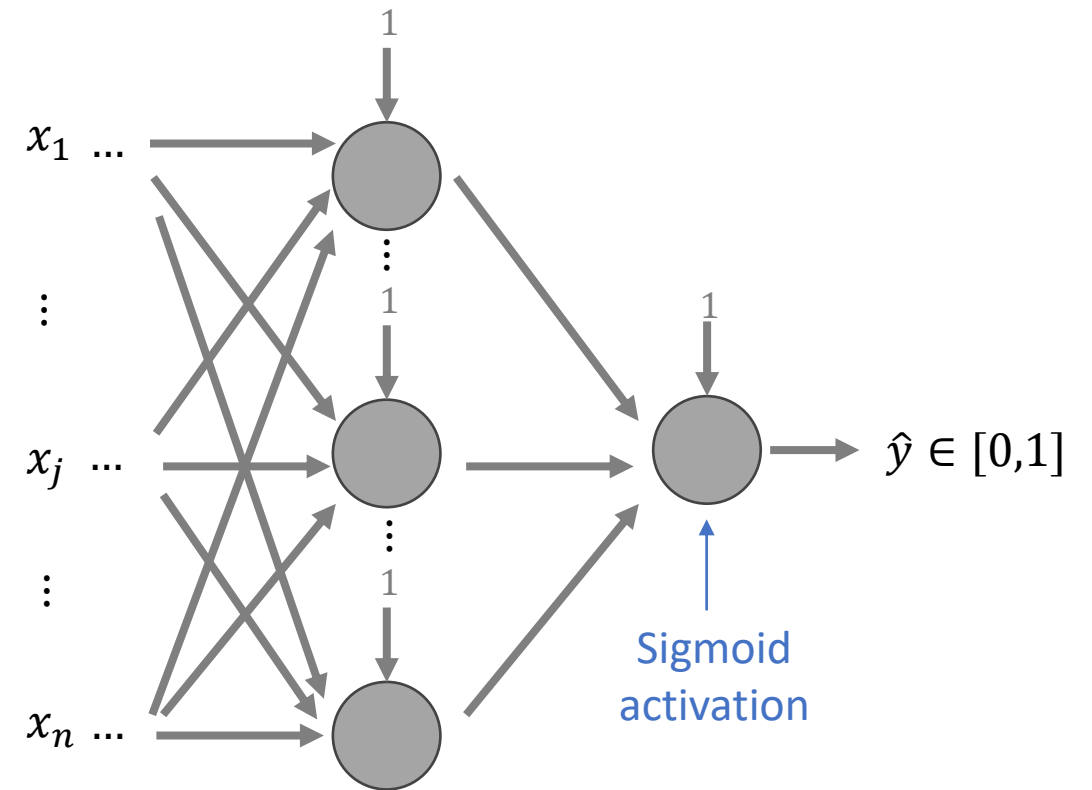
Neural networks can *represent* a wide variety of interesting functions with appropriate weights

- A single hidden layer network can approximate any continuous function within a specific range

Regression and Classification

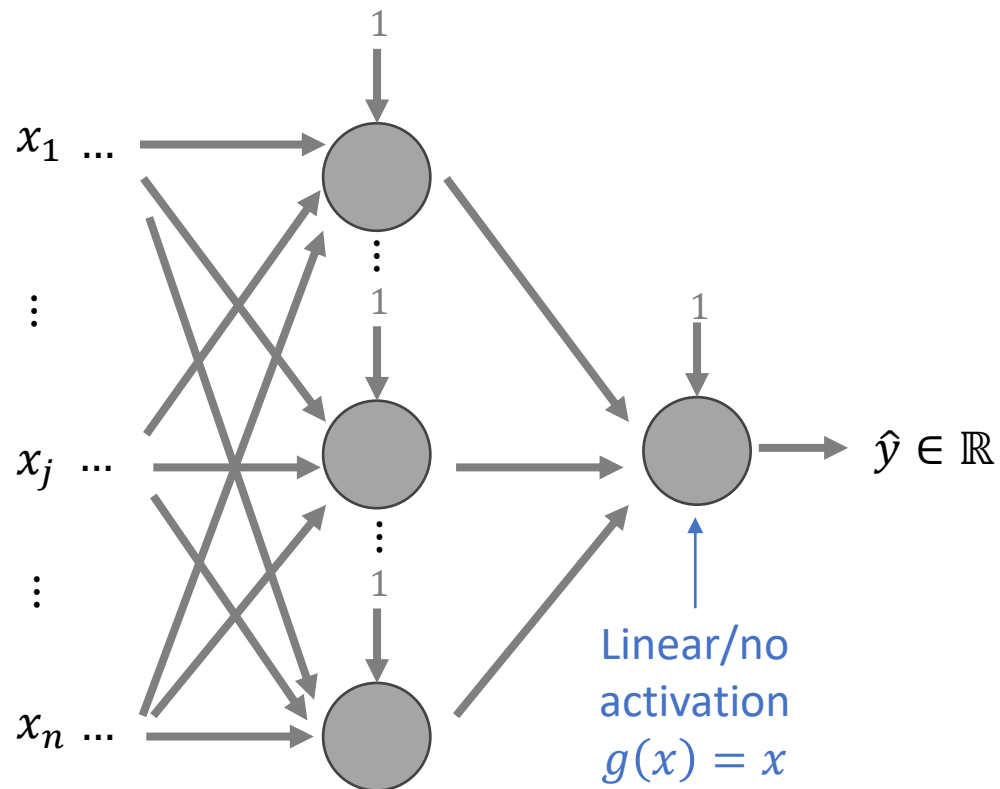


Regression

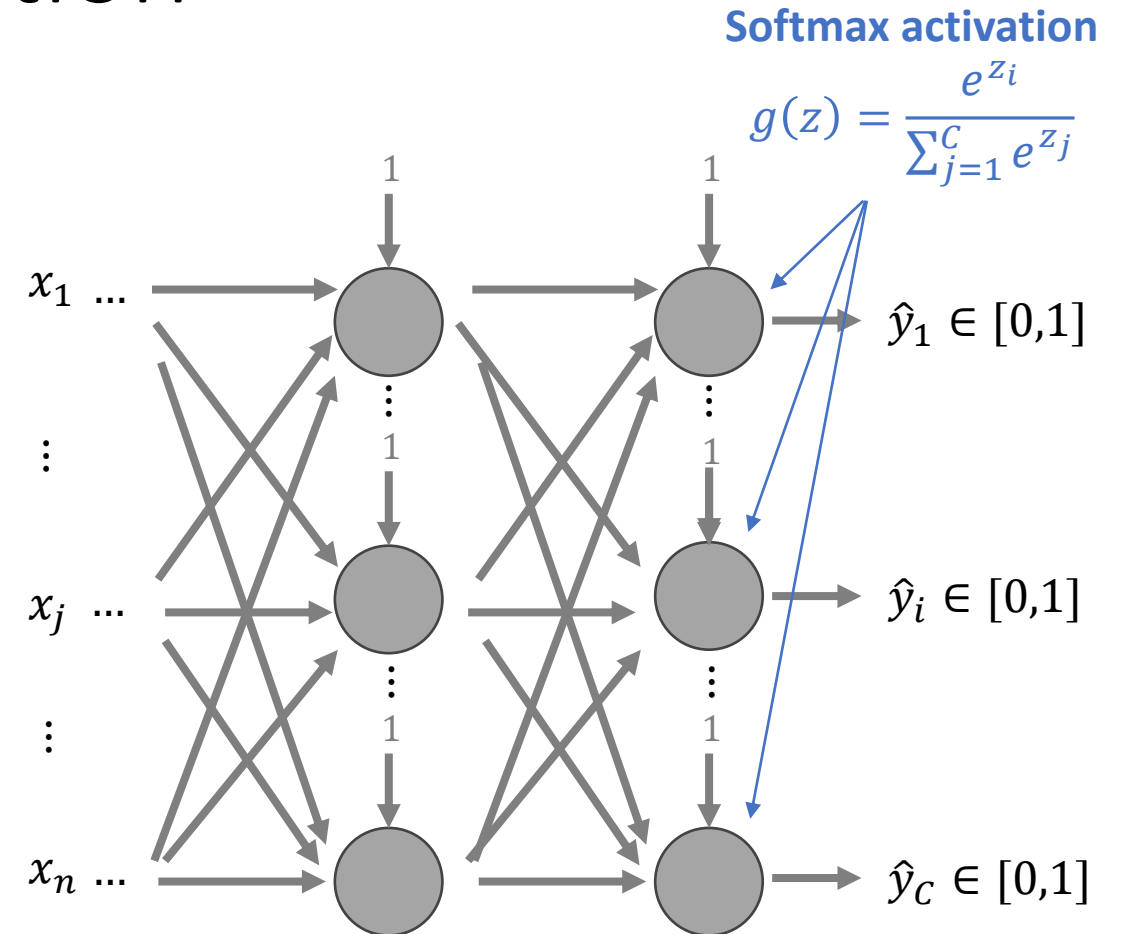


Binary Classification

Regression and Classification



Regression



Multi-class Classification
with C classes

$$\sum_{i=1}^C \hat{y}_i = 1$$

Outline

- Perceptron
 - Biological inspiration
 - Perceptron Learning Algorithm
- Neural Networks
 - Single-layer Neural Networks
 - Multi-layer Neural Networks
 - Regression and Classification
- **Gradient Descent with Neural Networks**
- Neural Networks vs Other Models

Background: Chain Rule

$$z(x) = h(g(f(x)))$$

$$\frac{dz}{dx} = \frac{dz}{dh} \frac{dh}{dg} \frac{dg}{df} \frac{df}{dx}$$

Aside: Derivative of Sigmoid Function

$$\begin{aligned}\sigma(x) &= \frac{1}{1 + e^{-x}} & \sigma'(x) &= \frac{d\left(\frac{1}{1 + e^{-x}}\right)}{dx} \\ & & &= \frac{d((1 + e^{-x})^{-1})}{dx} \\ & & &= -(1 + e^{-x})^{-2} (-e^{-x}) \\ & & &= \frac{e^{-x}}{(1 + e^{-x})^2} \\ & & &= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \\ & & &= \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}}\right) \\ & & &= \boxed{\sigma(x)(1 - \sigma(x))}\end{aligned}$$

Gradient Descent on Single-layer NN

$$\hat{y} = \sigma(f(x)) \quad f(x) = \sum_{i=0}^n \mathbf{w}_i x_i \quad L = \frac{1}{2}(\hat{y} - y)^2$$

Mean-squared Error

$$\mathbf{w}_i \leftarrow \mathbf{w}_i - \gamma \frac{dL}{d\mathbf{w}_i}$$

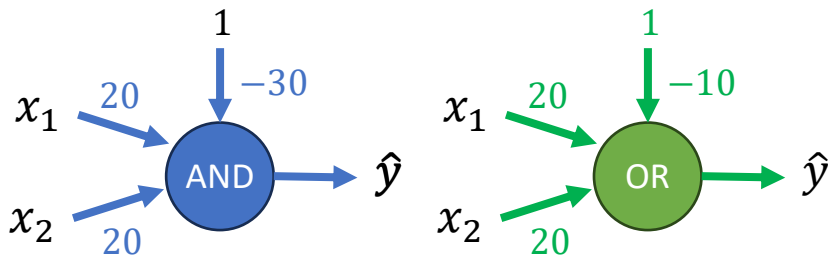
$$\mathbf{w}_i \leftarrow \mathbf{w}_i - \gamma(\hat{y} - y)\hat{y}(1 - \hat{y})x_i$$

$$\begin{aligned} \frac{dL}{d\mathbf{w}_i} &= \frac{dL}{d\hat{y}} \frac{d\hat{y}}{df} \frac{df}{d\mathbf{w}_i} \\ &= (\hat{y} - y)\hat{y}(1 - \hat{y})x_i \end{aligned}$$

$$\frac{dL}{d\hat{y}} = \frac{d(\frac{1}{2}(\hat{y} - y)^2)}{d\hat{y}} = (\hat{y} - y)$$

$$\frac{d\hat{y}}{df} = \frac{d(\sigma(f(x)))}{df} = \sigma(f(x))(1 - \sigma(f(x))) = \hat{y}(1 - \hat{y})$$

$$\frac{df}{d\mathbf{w}_i} = \frac{d(\sum_{i=0} \mathbf{w}_i x_i)}{d\mathbf{w}_i} = \frac{d(\mathbf{w}_i x_i)}{d\mathbf{w}_i} + \frac{d(\sum_{k \neq i} \mathbf{w}_k x_k)}{d\mathbf{w}_i} = x_i$$



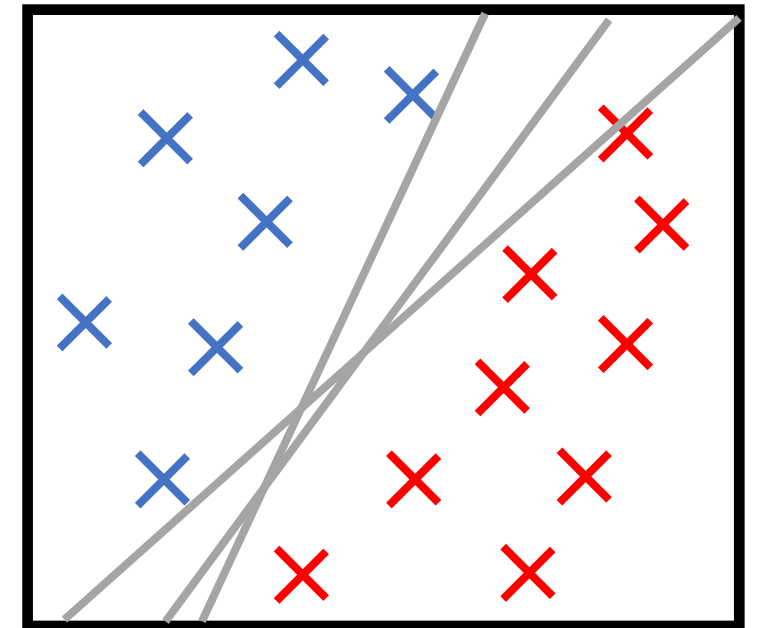
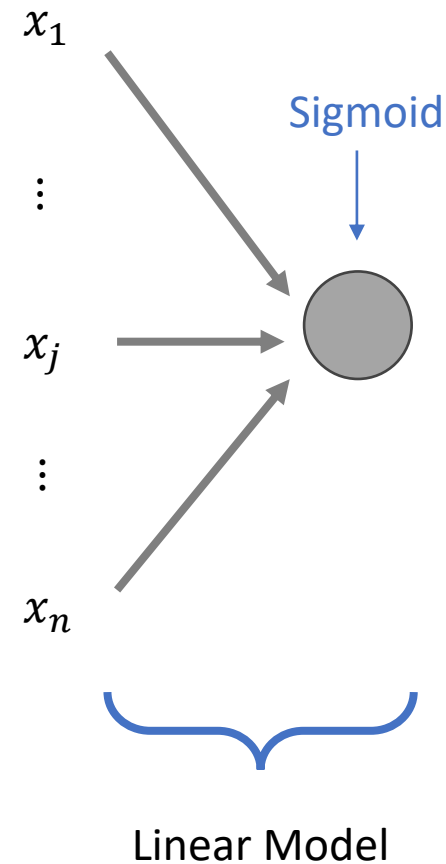
Gradient Descent on Multi-layer NN

Backpropagation is covered in the next Lecture!

Outline

- Perceptron
 - Biological inspiration
 - Perceptron Learning Algorithm
- Neural Networks
 - Single-layer Neural Networks
 - Multi-layer Neural Networks
 - Regression and Classification
- Neural Networks with Gradient Descent
- **Neural Networks vs Other Models**

Neural Networks vs Logistic Regression

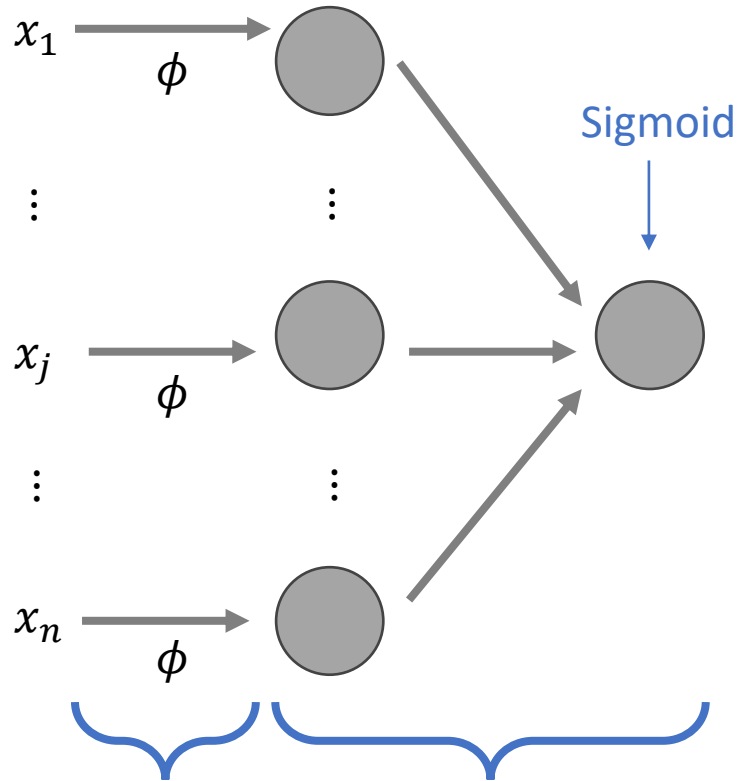


Linear, non-robust decision boundary

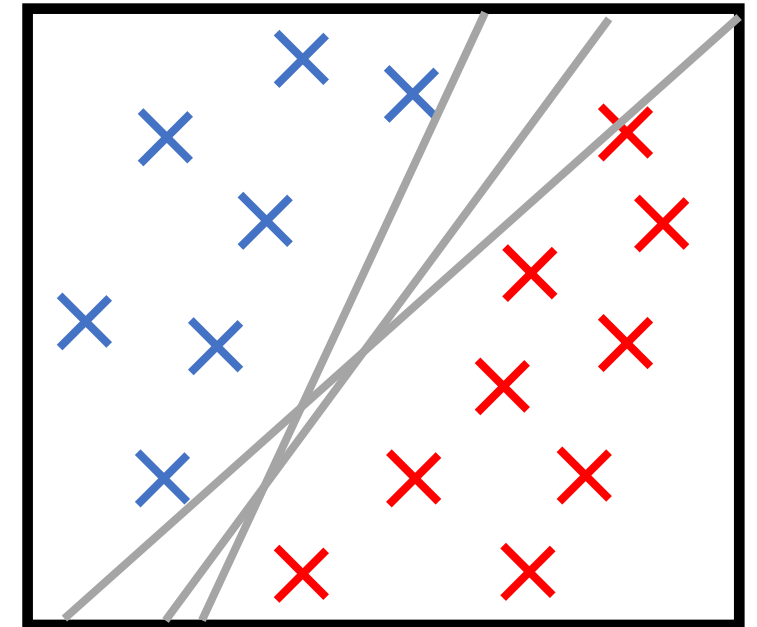
Prone to misclassification since the decision boundary can be too close to data points

Neural Networks vs Logistic Regression

With Feature Mapping



Handcrafted Feature Mapping Linear Model

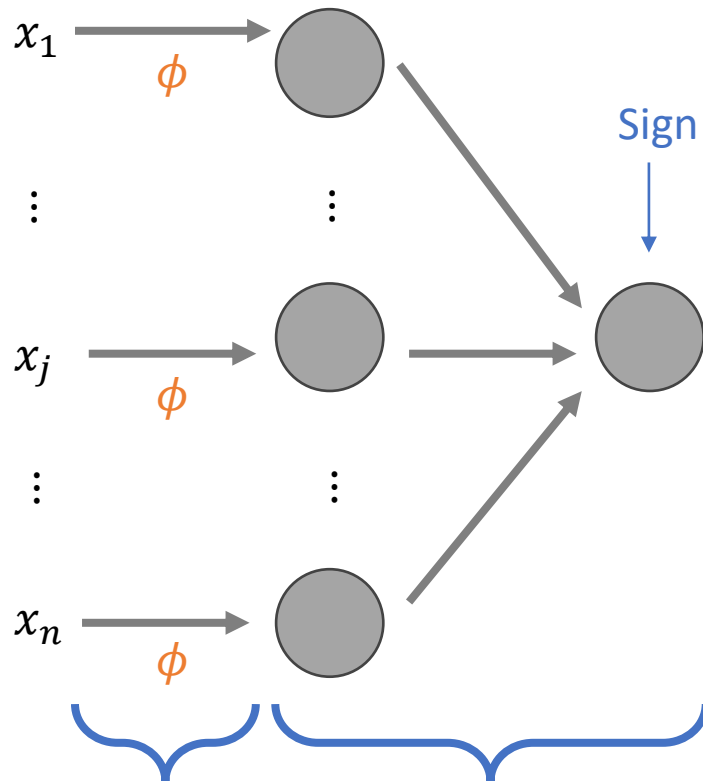


Non-linear, non-robust decision boundary

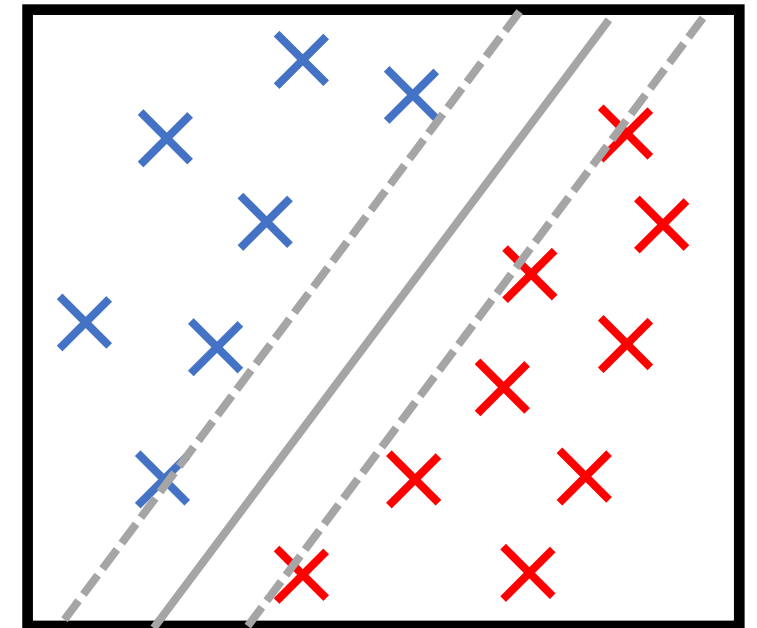
Prone to misclassification since the decision boundary can be too close to data points

Neural Networks vs Support Vector Machines

Implicit with **Kernel**: can have many, possibly infinite dimensional features



Handcrafted Feature Mapping Linear Model



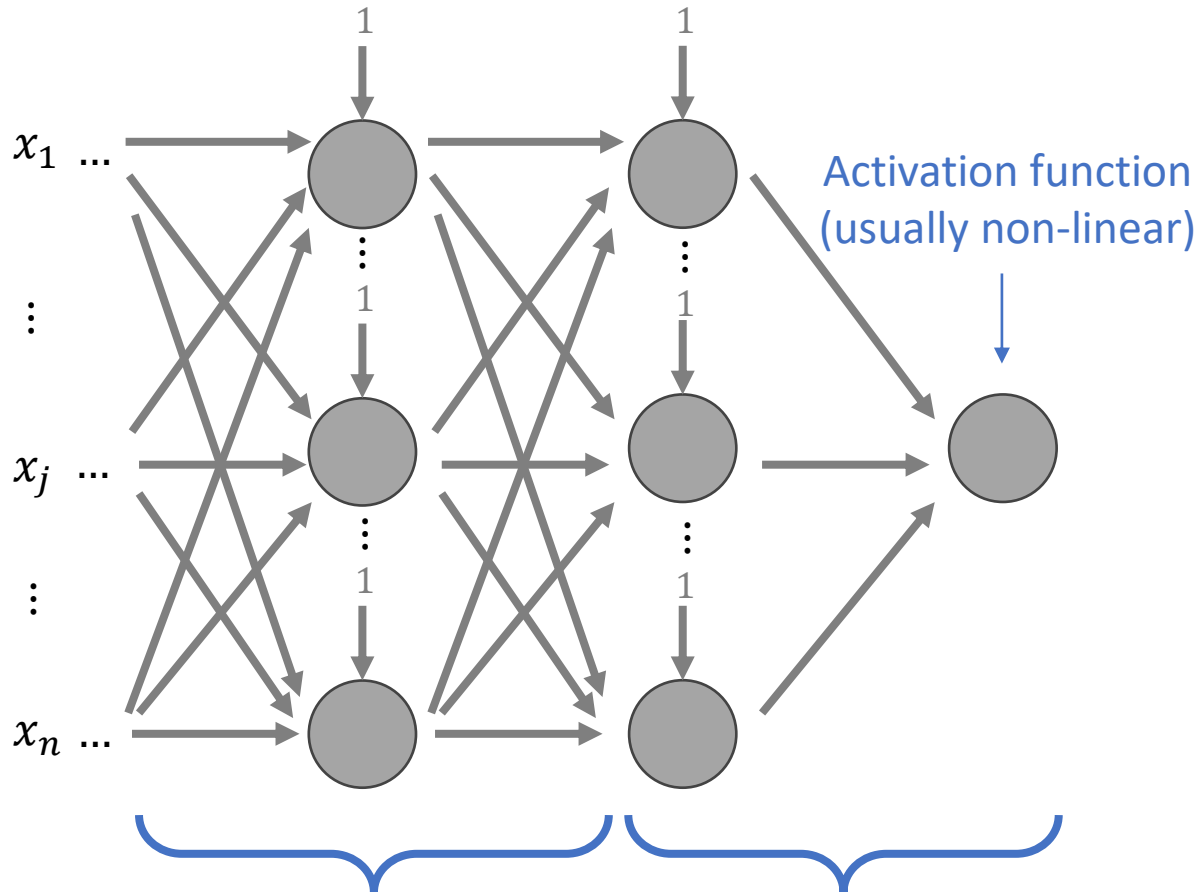
Non-linear, robust decision boundary

Decision boundary is guaranteed to be far from the data points

Neural Networks vs Other Methods

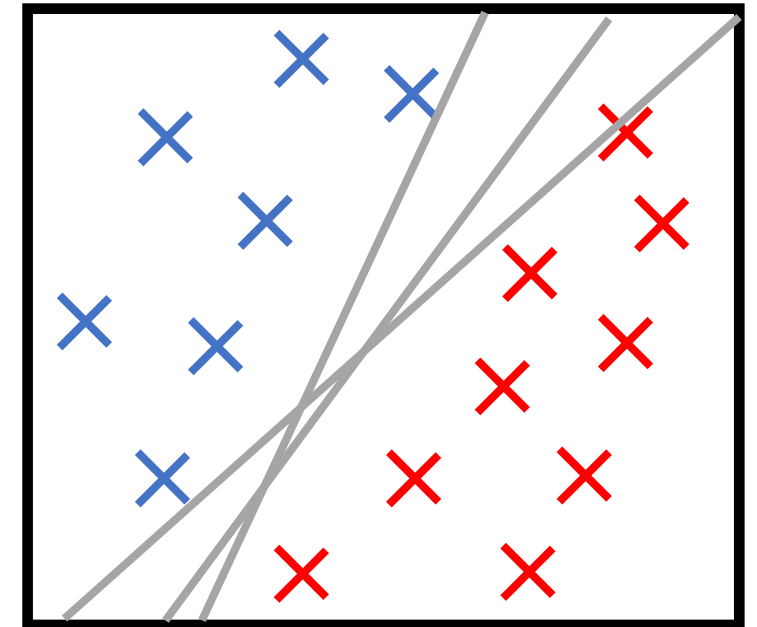


Credit: Internet meme, original source unknown



Learned Feature Mapping

Linear Model



Non-linear, non-robust decision boundary

Prone to misclassification since the decision boundary can be too close to data points

Summary

- Perceptron
 - Biological inspiration: brain, neural network, neuron
 - Perceptron Learning Algorithm:
 - $w \leftarrow w + \gamma(y^{(j)} - \hat{y}^{(j)})x^{(j)}$ on a misclassified instance
- Neural Networks
 - Single-layer Neural Networks: AND, OR, NOR
 - Multi-layer Neural Networks: XNOR, Universal Approximation Theorem
 - Regression and Classification
- Neural Networks with Gradient Descent
- Neural Networks vs Other Models: learned feature mapping!

Coming Up Next Week

- **Math revision**

- Linear algebra: scalar, vector, matrix, and their operations
- Calculus: partial derivative, matrix calculus, chain rule

- **Backpropagation**

- **Introduction to PyTorch**

- Training neural networks with gradient descent

To Do

- **Lecture Training 8**
 - +100 Free EXP
 - +50 Early bird bonus
- **PS5**
 - Due **tomorrow** (Saturday, 21st October) 23:59
- **PS6 is out today!**
 - May need some knowledge from the next lecture (Lecture 9)