

# Training Camp 2024

First Day – 4<sup>th</sup> September 2024

**Alessandro Nicolosi**



<https://github.com/alenic>



<https://www.linkedin.com/in/alessandro-nicolosi/>



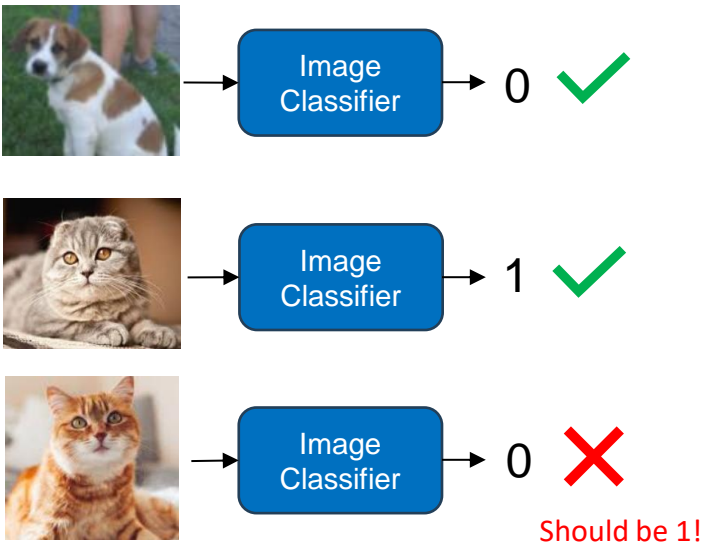
<https://www.leonardo.com/it/>

# Single Modality – Classification Problem

## Computer Vision

- **Task:** Image Classification
- **Input Modality:** Image

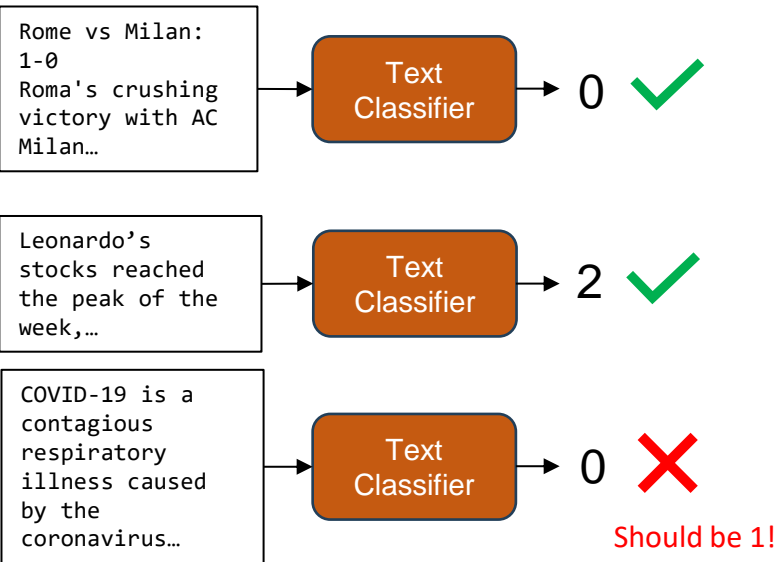
```
idx2label = {  
  0: "dog"  
  1: "cat"  
}
```



## Natural Language Processing

- **Task:** Text Classification
- **Input Modality:** Text

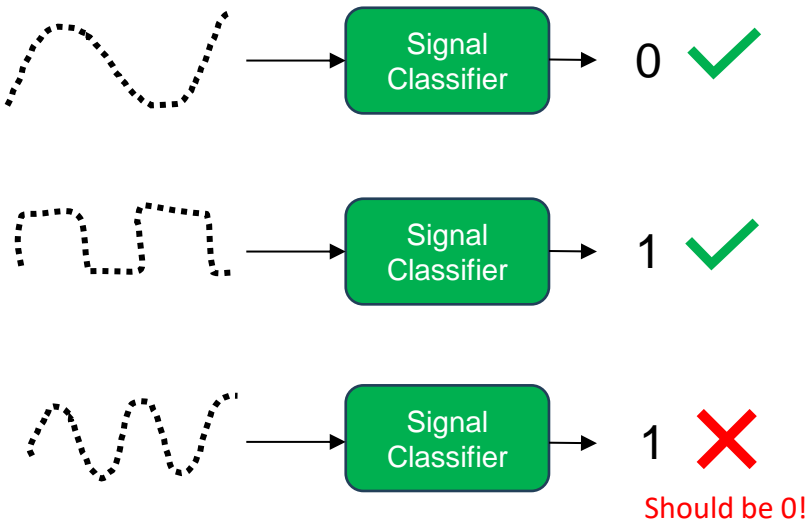
```
idx2label = {  
  0: "sport"  
  1: "health"  
  2: "economy"  
}
```



## Time Series Analysis

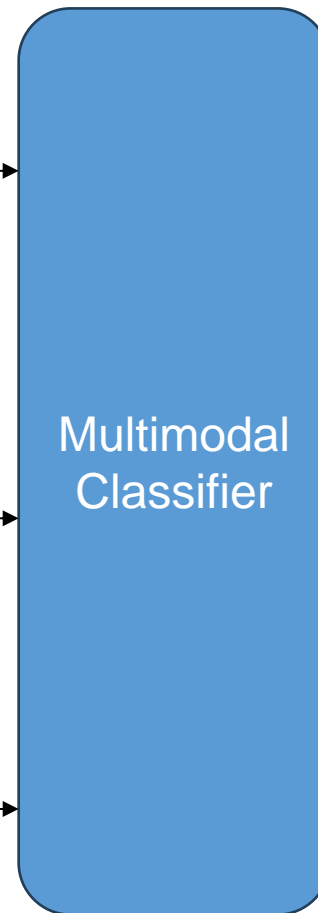
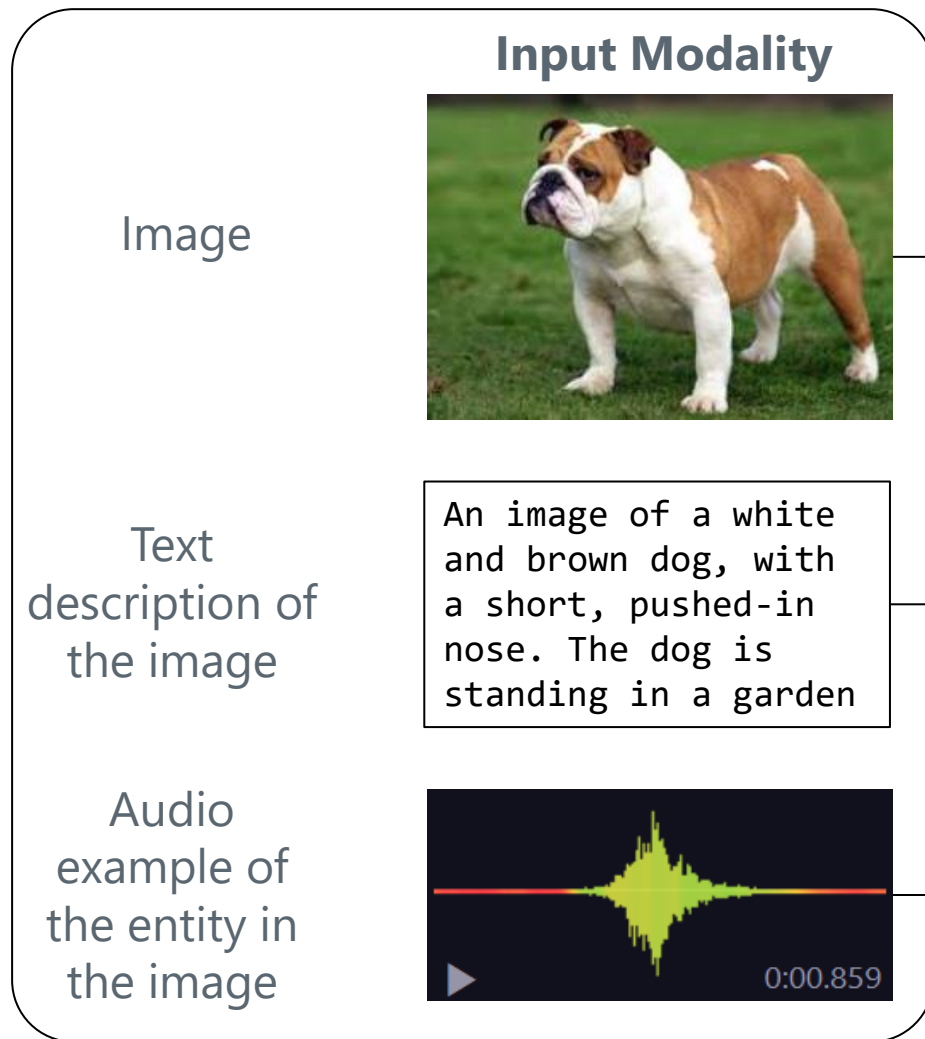
- **Task:** Time Series Classification
- **Input Modality:** Time series

```
idx2label = {  
  0: "sinusoid"  
  1: "pulse"  
}
```



# Multimodal Classification Problem

**A single example**  $x_i = (\text{Image}_i, \text{Text}_i, \text{Audio}_i)$



50 ✓

**Dataset**

$$D_{train} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

$$x_i = (\text{Image}_i, \text{Text}_i, \text{Audio}_i)$$

$$y_i \in \{0, 1, \dots, 99\}$$

```
idx2label = {  
  0: "cat - persian"  
  1: "cat - siamese"  
  ...  
  49: "dog - pug"  
  50: "dog - bulldog"  
  ...  
  99: "dog - rottweiler"  
}
```

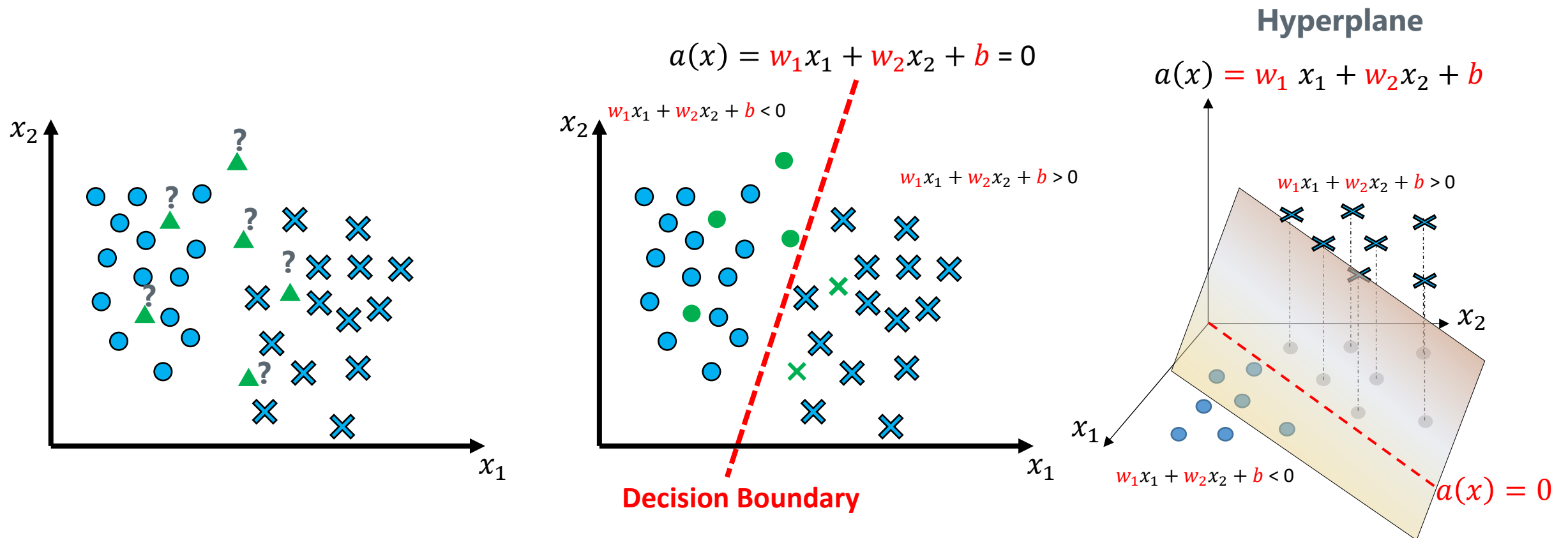
# Binary Classification Problem

Given a training dataset

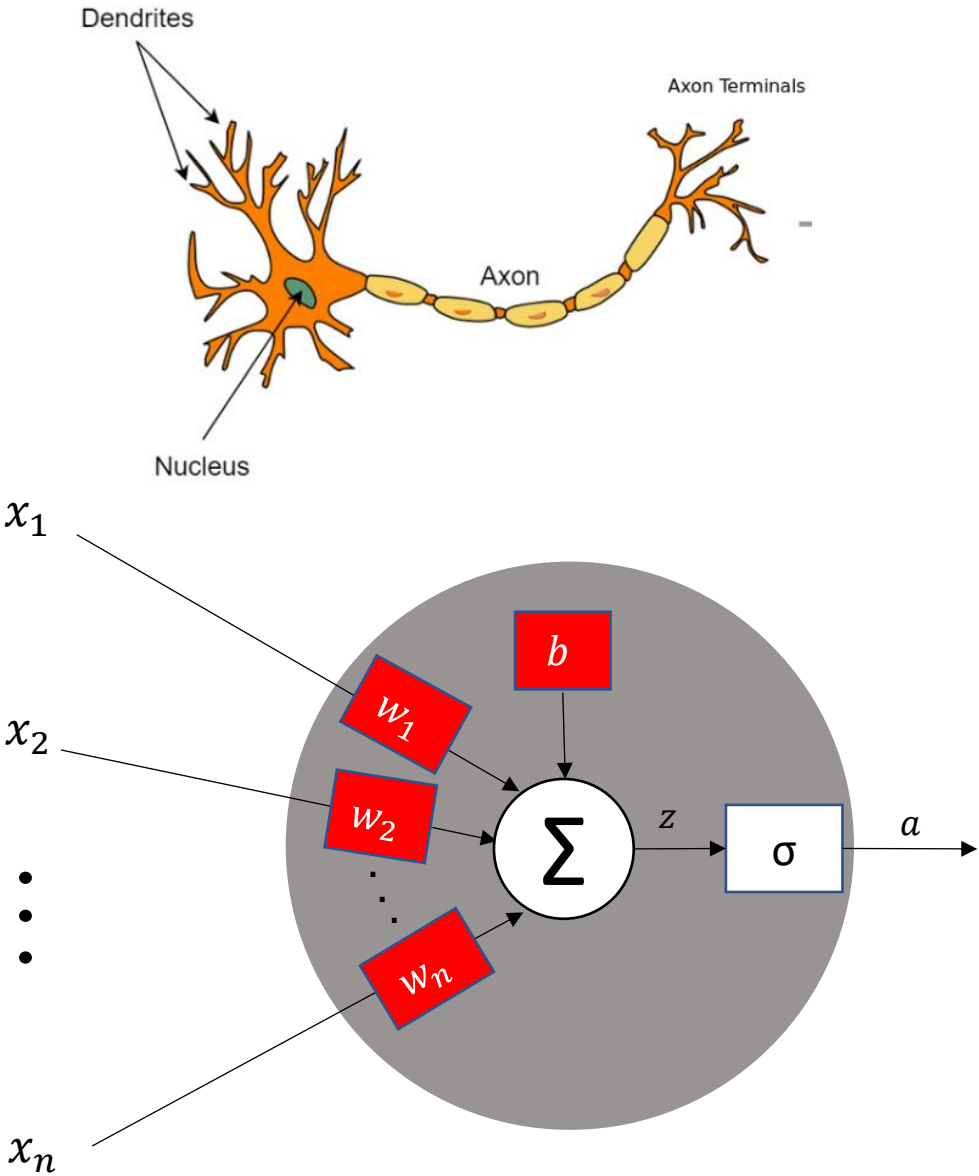
$$D_{train} = \{(x^{(i)}, y^{(i)})\}_{i=1}^{N_{train}}: \quad x^{(i)} \in \mathbb{R}^n, y^{(i)} \in \{0, 1\}\}$$

$$D_{test} = \{x_{test}^{(1)}, x_{test}^{(2)}, \dots, x_{test}^{(N_t)}\}$$

What is the probability that a point  $x_{test}^{(i)}$  belong to the class 1 ?  $P(Y = 1 | X = x_{test}^{(i)})$  ?

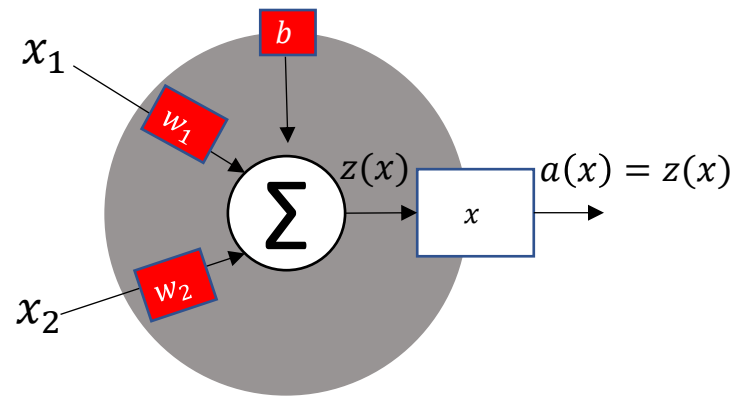


# Neural Networks – Single Neuron



Input	$x = [x_1 \ x_2 \ \dots \ x_n] \in \mathbb{R}^n$
Output	$z = \left( \sum_{i=1}^n w_i x_i \right) + b$ $a(x) = \sigma( z(x) ) = \sigma( w^T x + b )$
Weights or Parameters	$w = [w_1 \ w_2 \ \dots \ w_n] \in \mathbb{R}^n$ $b \in \mathbb{R}$
Some Activation Functions: $\sigma(x)$	
Linear	$x$
Hyperbolic Tangent	$\tanh(x)$
ReLU	$\begin{cases} x, & x \geq 0 \\ 0, & x < 0 \end{cases}$
LeakyReLU	$\begin{cases} x, & x \geq 0 \\ \alpha x, & x < 0 \end{cases}$

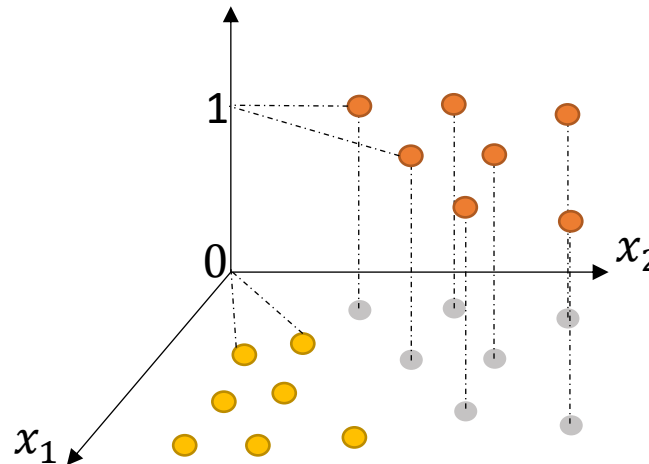
# Single Neuron for classification



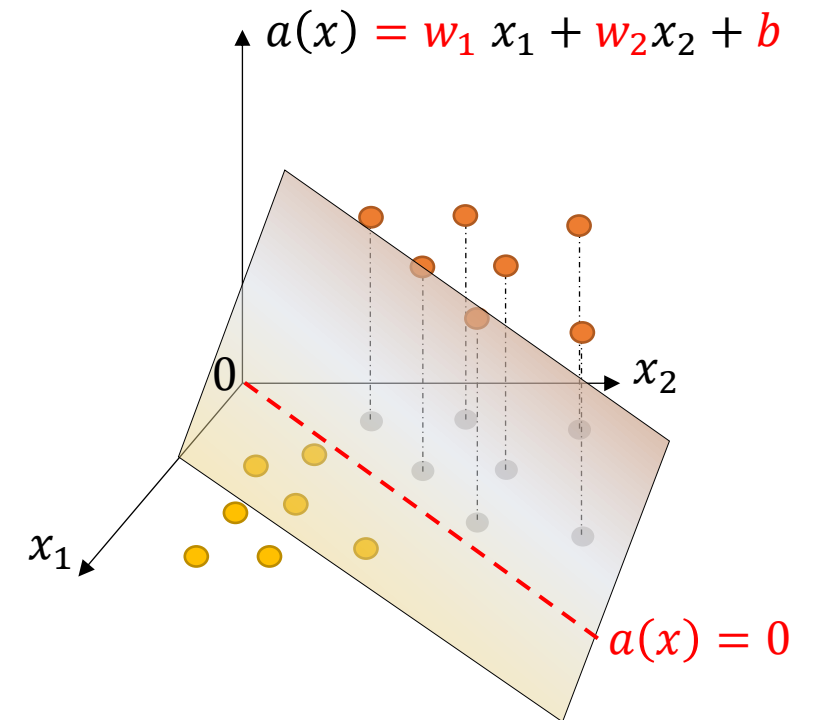
Activation Functions:  $\sigma(x) = 0$

Training Dataset

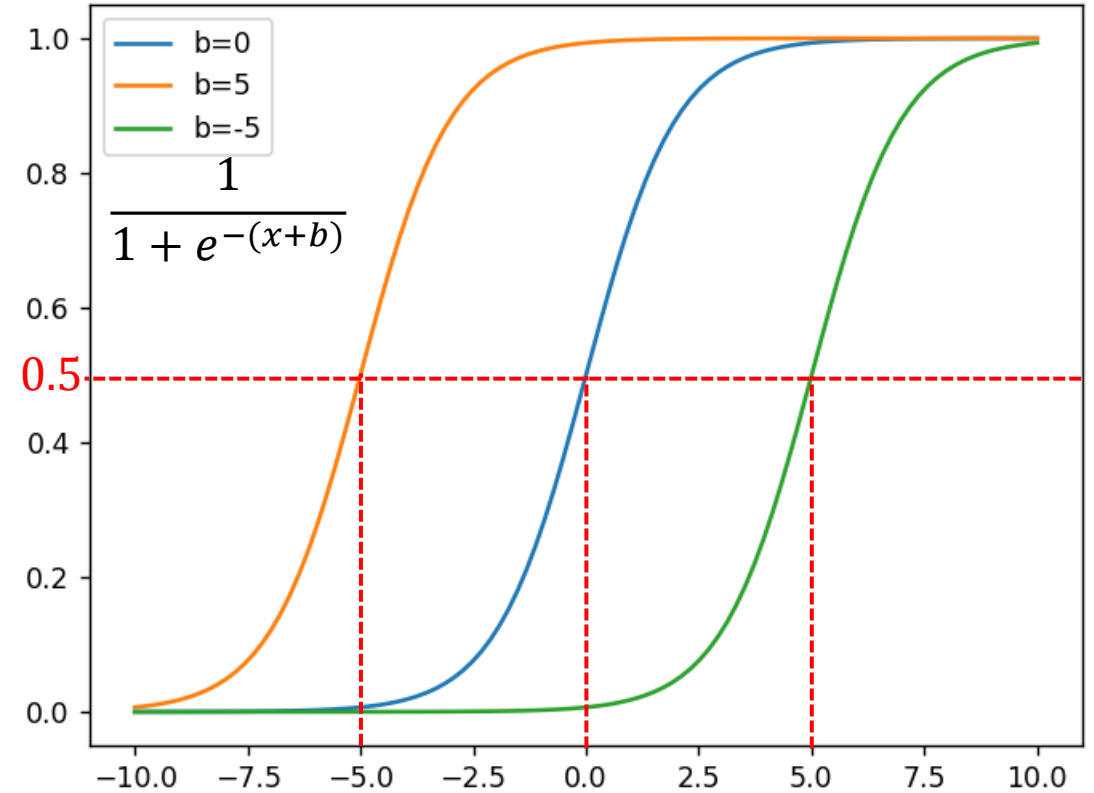
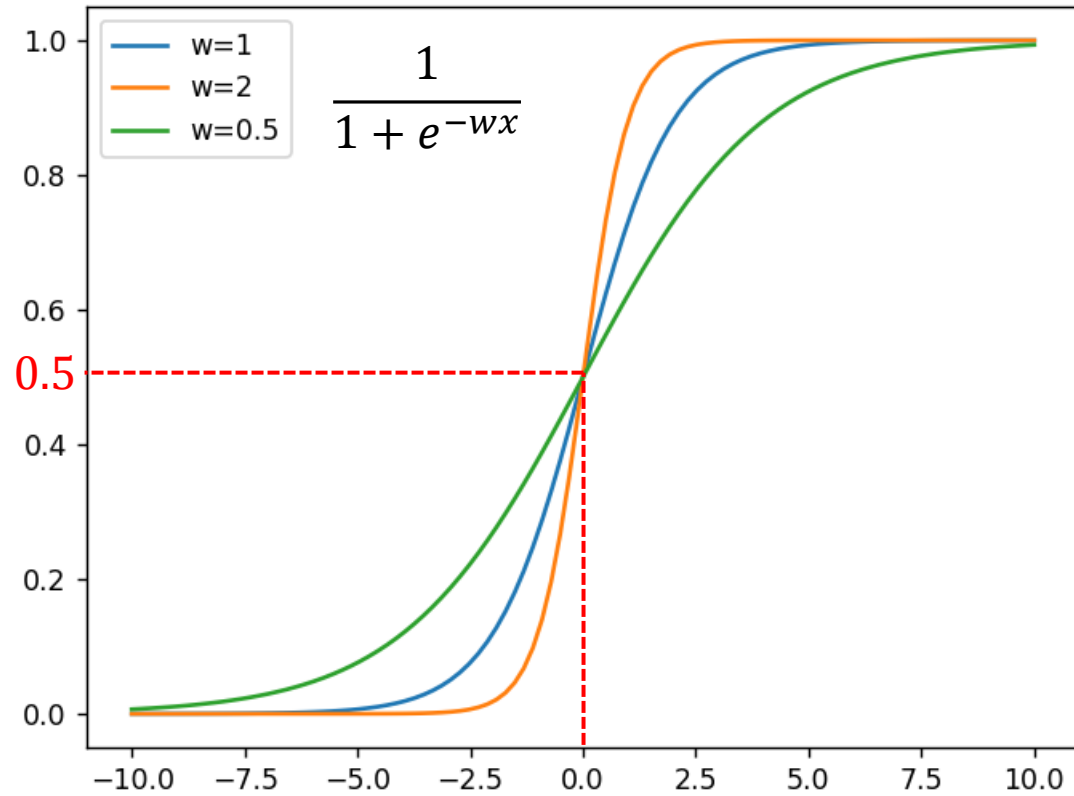
$$\{(x^{(i)}, y^{(i)})_{i=1}^N : x^{(i)} \in \mathbb{R}^2, y^{(i)} \in \{0, 1\}\}$$



Hyperplane



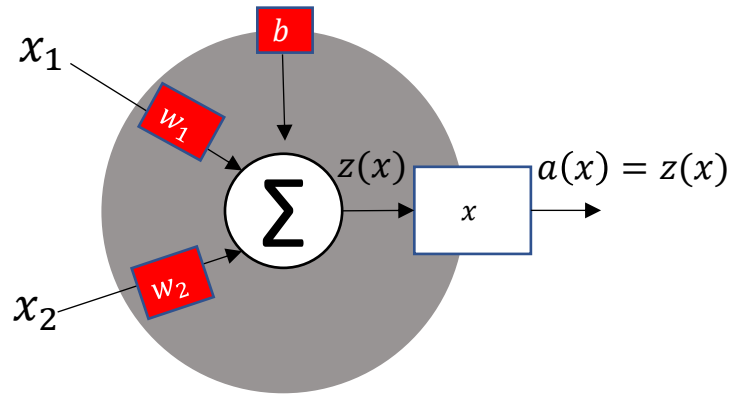
# Single Neuron for classification – Logistic function



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

$$y'(x) = y(x)(1 - y(x))$$

# Single Neuron for classification

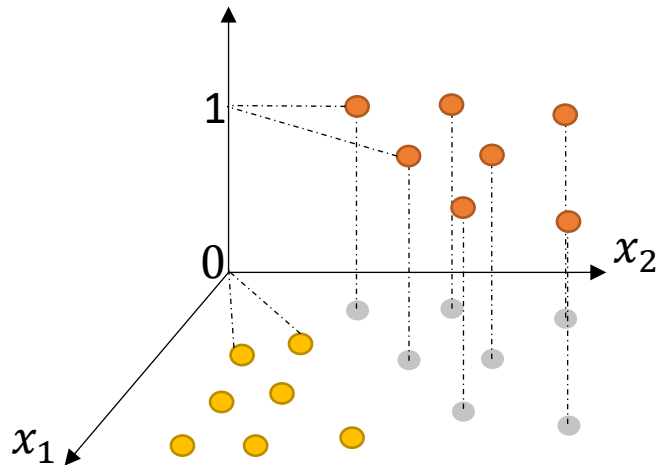


Apply Logistic Function to squash in 0-1

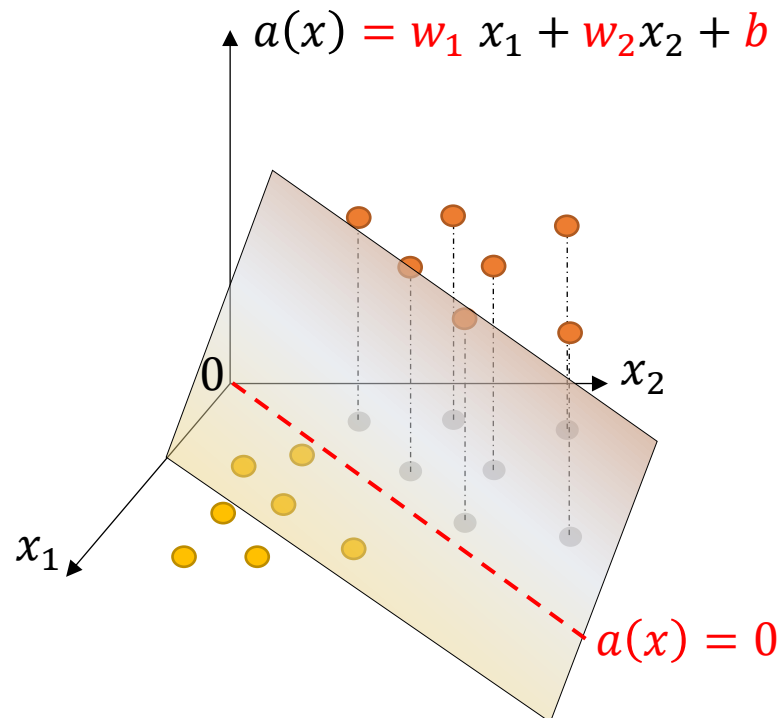
$$\hat{y}(x^{(i)}) = \frac{e^{a(x^{(i)})}}{1 + e^{a(x^{(i)})}} = P(Y = 1 | X = x^{(i)})$$

Training Dataset

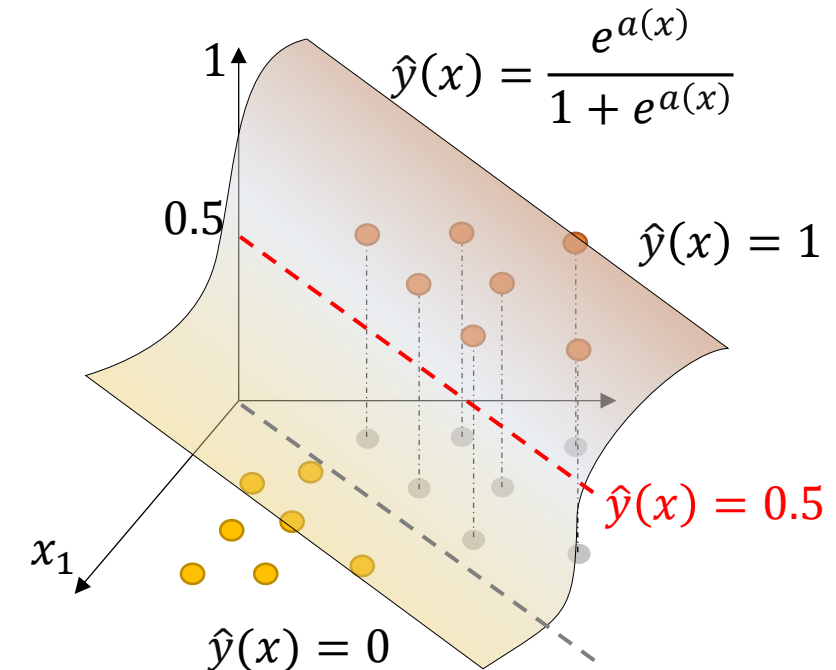
$$\{(x^{(i)}, y^{(i)})\}_{i=1}^N: x^{(i)} \in \mathbb{R}^2, y^{(i)} \in \{0, 1\}\}$$



Hyperplane



Squash in 0-1





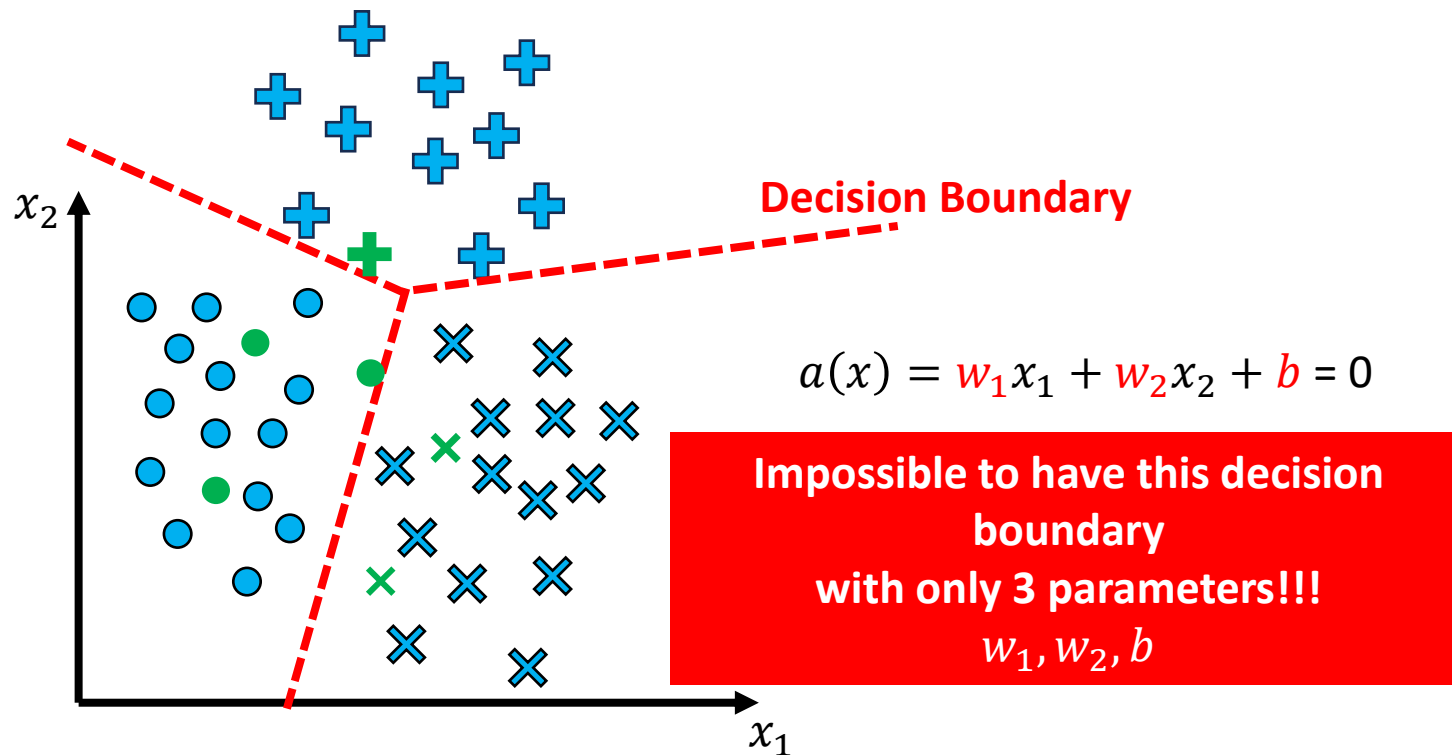
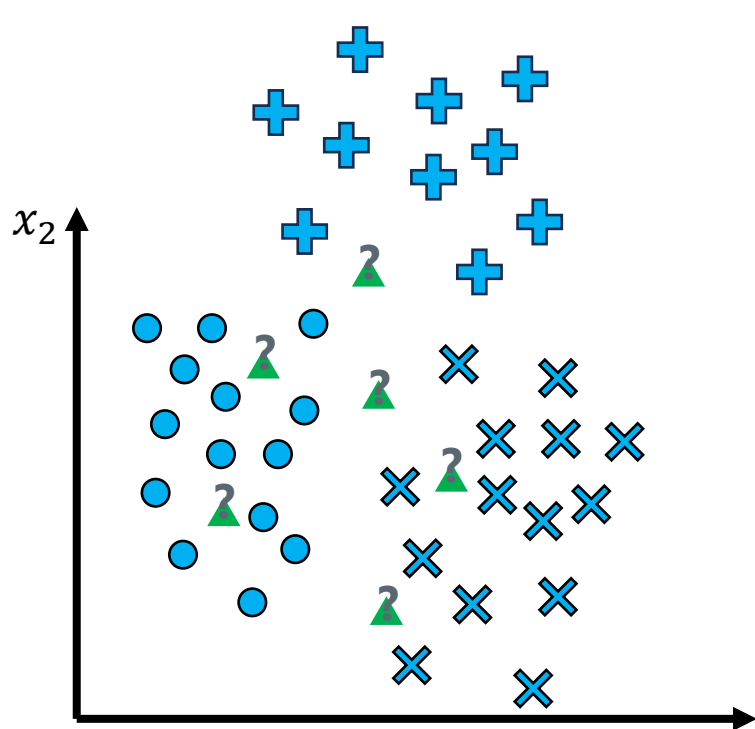
# Multiclass Classification Problem

Given a training dataset

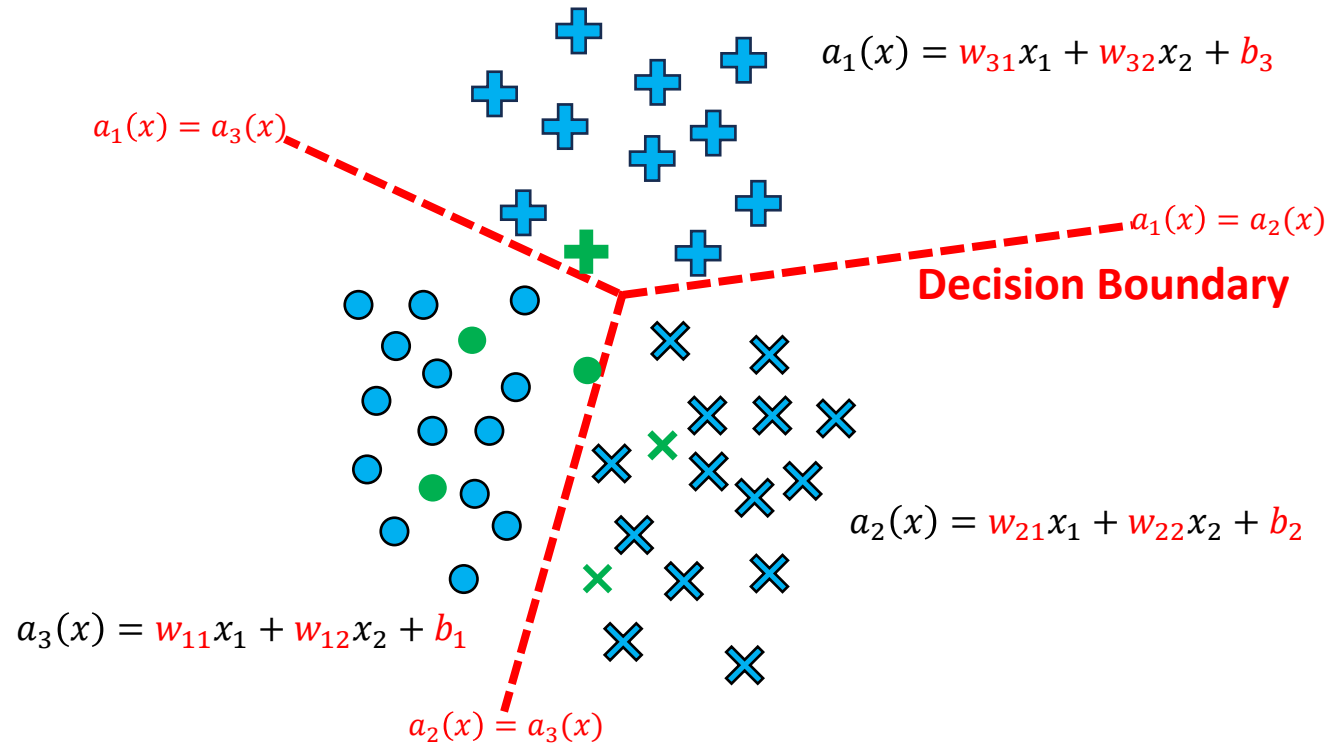
$$D_{train} = \{(x^{(i)}, y^{(i)})\}_{i=1}^{N_{train}}: \quad x^{(i)} \in \mathbb{R}^n, y^{(i)} \in \{0, 1, \dots, C-1\}\}$$

$$D_{test} = \{x_{test}^{(1)}, x_{test}^{(2)}, \dots, x_{test}^{(N_t)}\}$$

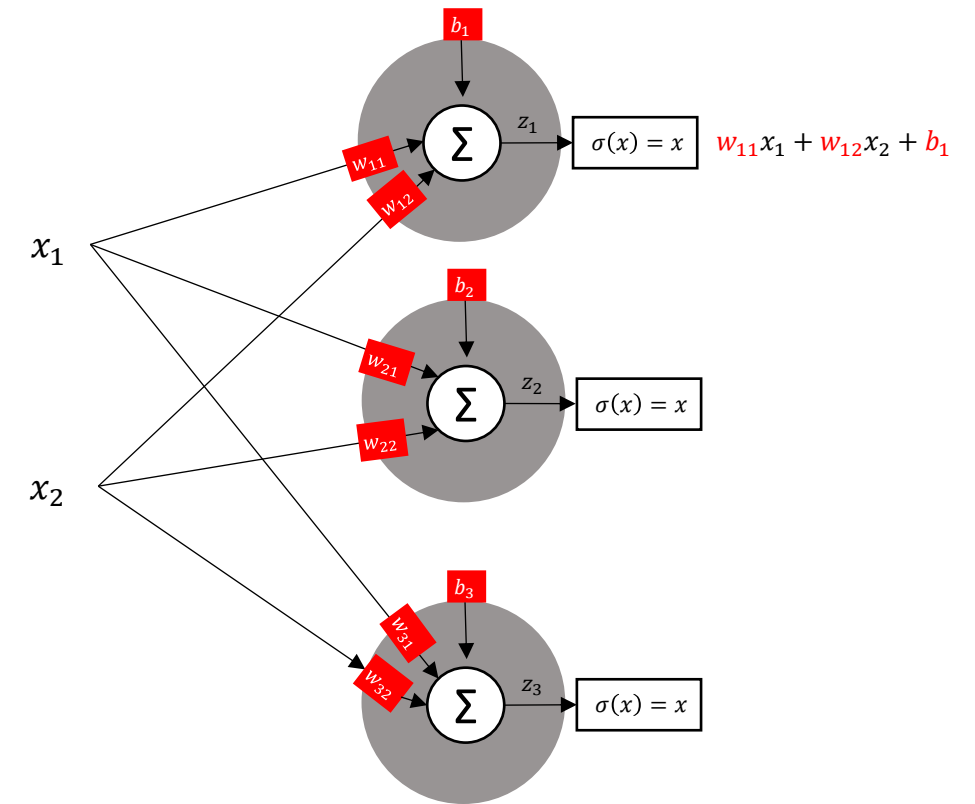
What is the probability that a point  $x_{test}^{(i)}$  belong to the class 1 ?  $P(Y = 1 | X = x_{test}^{(i)})$  ?



# Multiclass Classification Problem

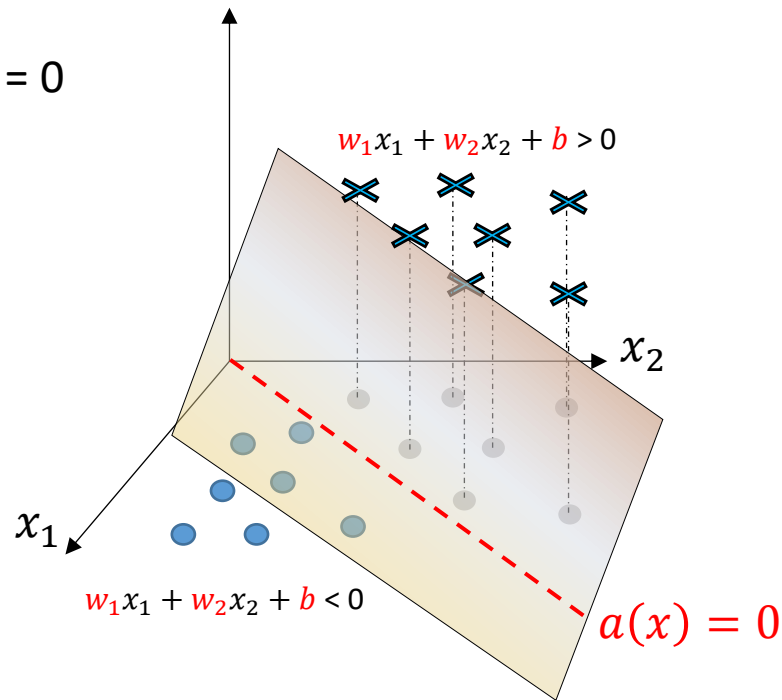
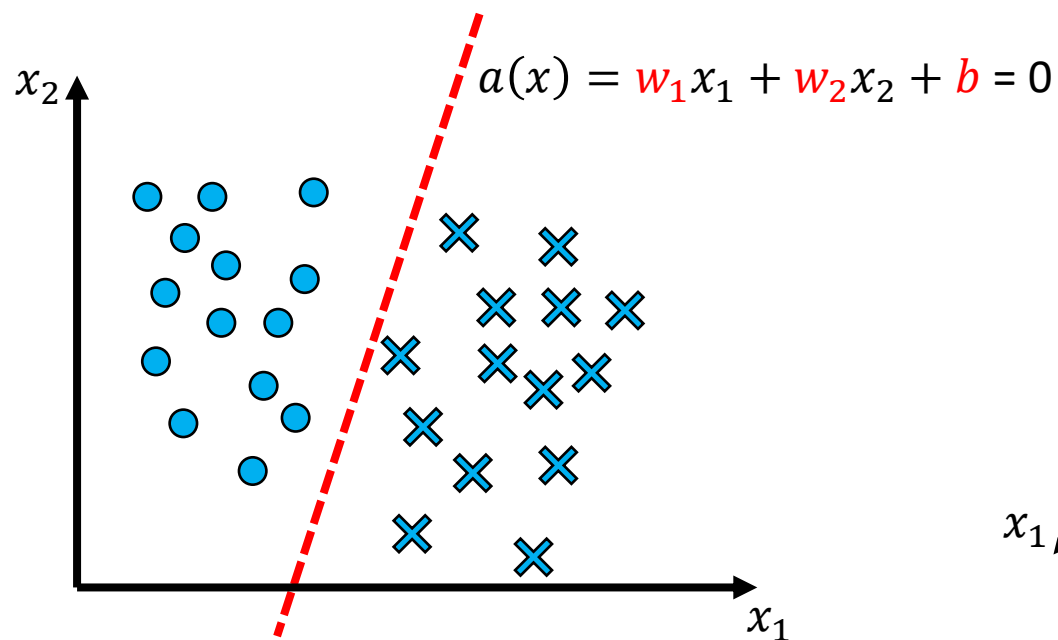


$$\hat{y}(x^{(i)}, \mathbf{W}, \mathbf{b}) = \text{softmax}(x^{(i)}) = \begin{bmatrix} \frac{e^{a_1(x^{(i)})}}{\sum_{j=1}^3 e^{a_j(x^{(i)})}} \\ \frac{e^{a_2(x^{(i)})}}{\sum_{j=1}^3 e^{a_j(x^{(i)})}} \\ \frac{e^{a_3(x^{(i)})}}{\sum_{j=1}^3 e^{a_j(x^{(i)})}} \end{bmatrix} = \begin{bmatrix} P(Y = 0 | X = x^{(i)}; \mathbf{W}, \mathbf{b}) \\ P(Y = 1 | X = x^{(i)}; \mathbf{W}, \mathbf{b}) \\ P(Y = 2 | X = x^{(i)}; \mathbf{W}, \mathbf{b}) \end{bmatrix}$$

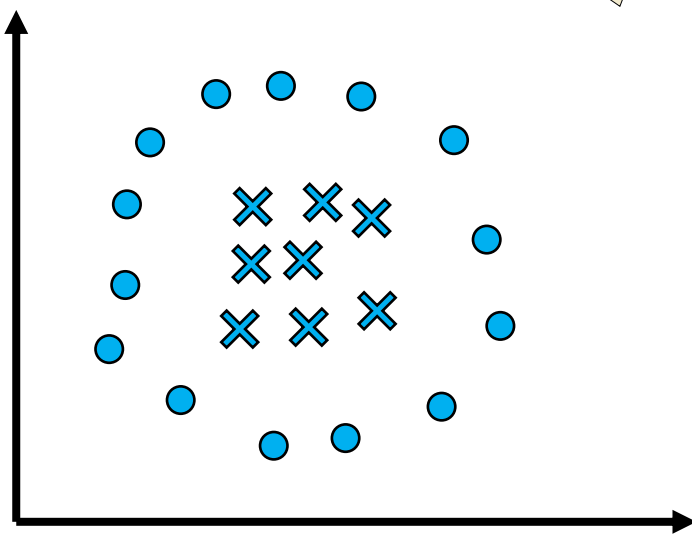
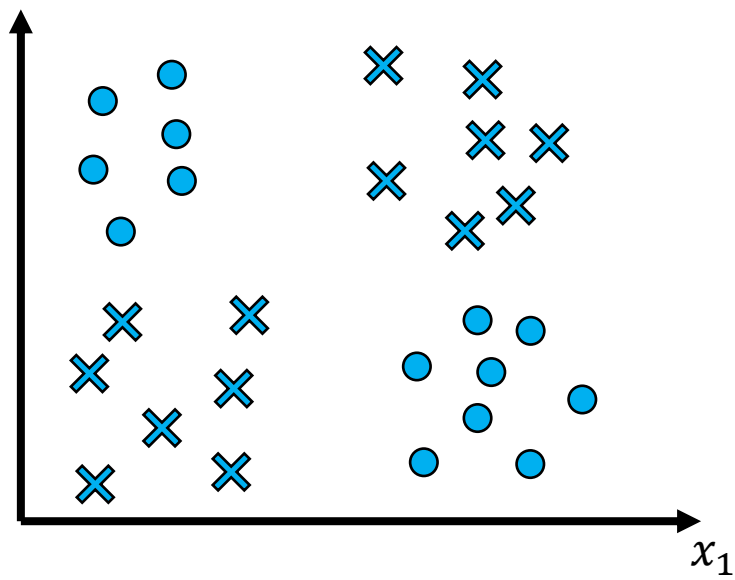


$$y_{pred}(x^{(i)}) = \text{argmax } \hat{y}(x^{(i)}, \mathbf{W}, \mathbf{b}) \in \{0, 1, 2\}$$

# Linear Separability

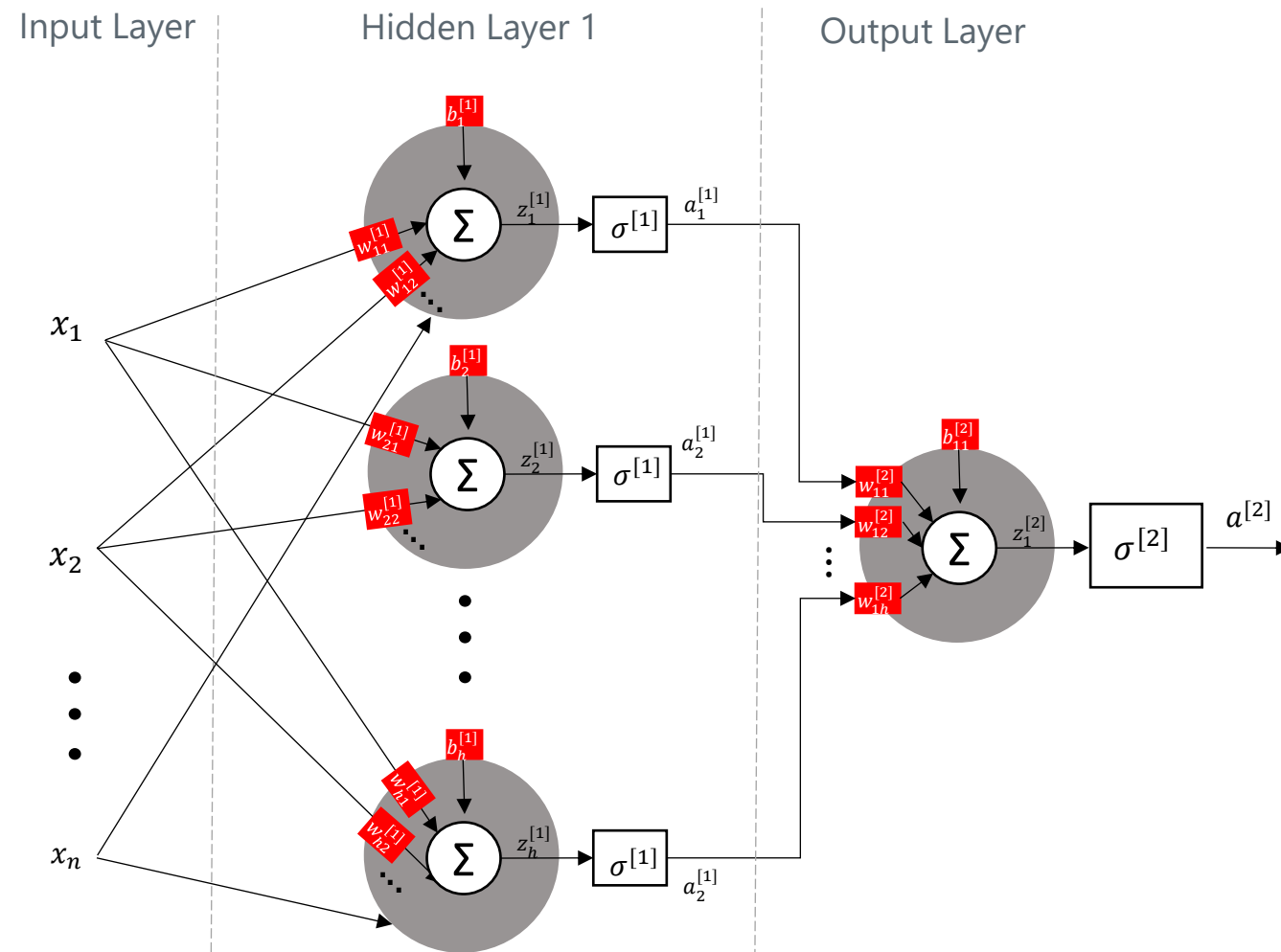


**Linearly Separable data**



**NOT linearly Separable data!**

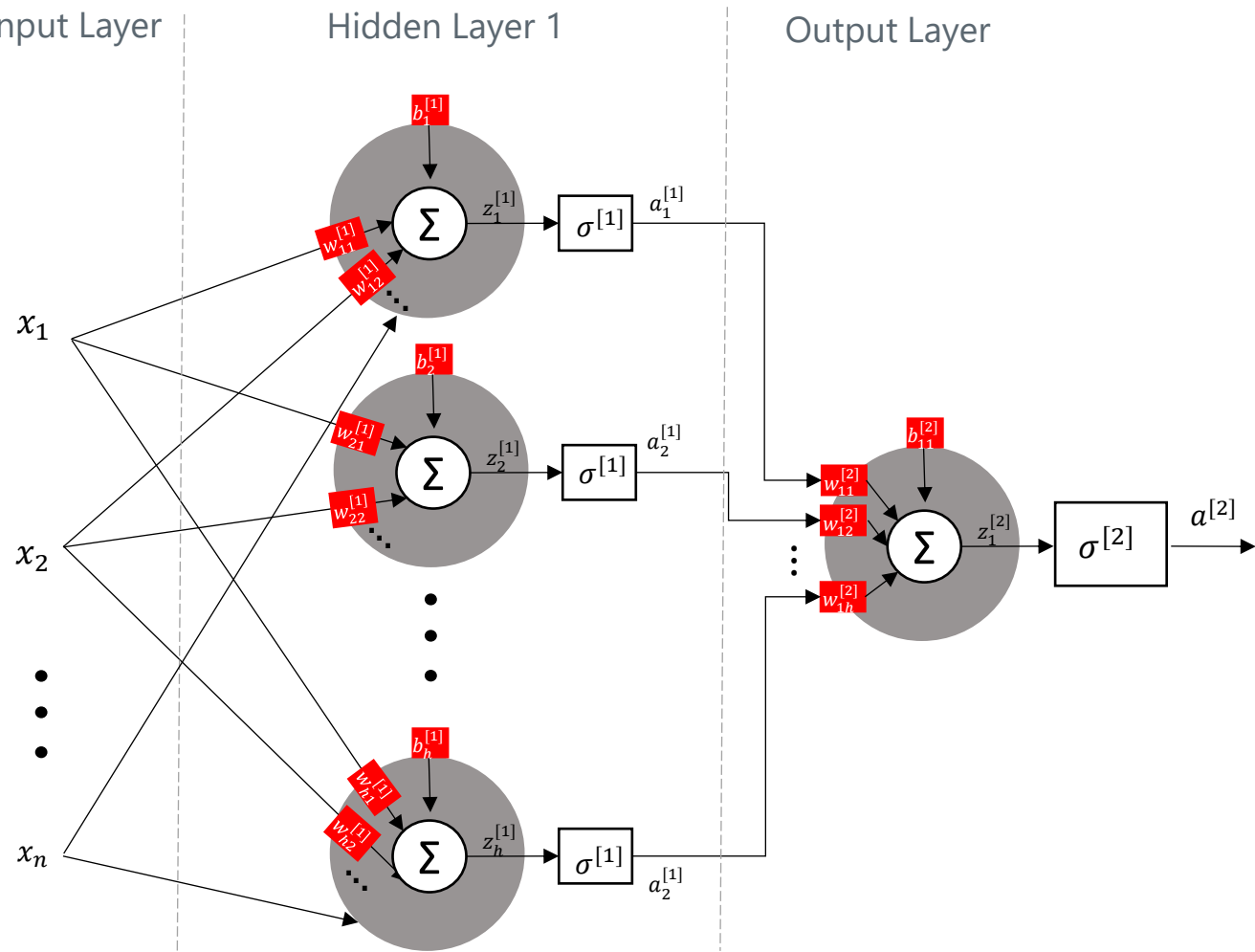
# Single-Layer Perceptron



$$z_1^{[2]} = \left( \sum_{i=1}^h w_{1i}^{[2]} a_i^{[1]} \right) + b_1^{[2]}$$

- ✓  $a_i^{[1]}$  are Non-linear functions!
- ✓  $z_1^{[2]}$  is a linear combination of non-linear functions
- ✓ If  $h$  increases, then  $z_1^{[2]}$  have more approximation power
- ✓ In the ideal case,  $a^{[1]}(x) = [a_1^{[1]}, a_2^{[1]}, \dots, a_h^{[1]}]$  is forced to be a new space, where the features are linearly separable!

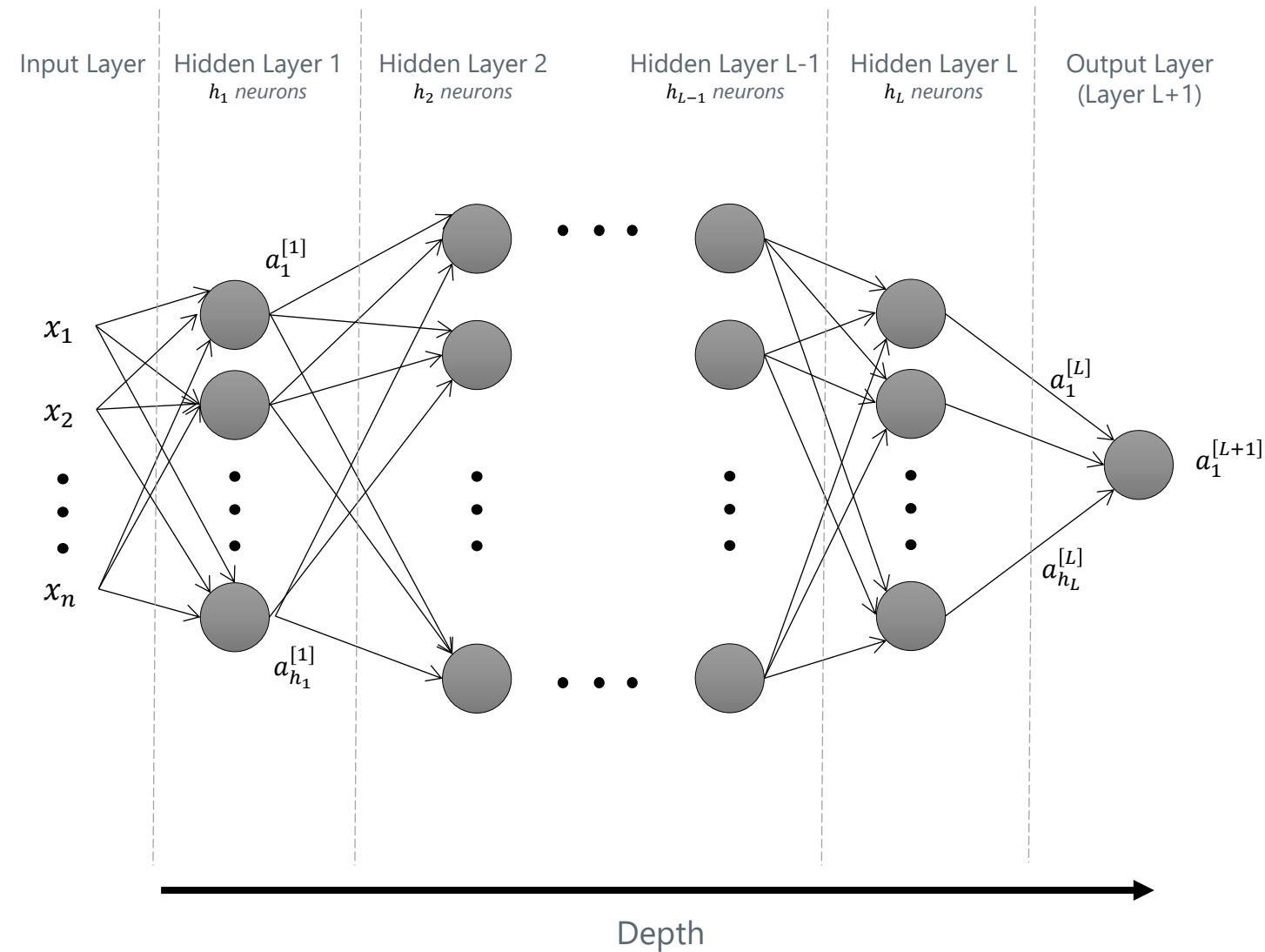
# Single-Layer Perceptron



Input	$x = [x_1 \ x_2 \ \dots \ x_n] \in \mathbb{R}^n$
Output	$z_1^{[2]} = \left( \sum_{j=1}^h w_{1j}^{[2]} a_j^{[1]} \right) + b_1^{[2]} \quad a^{[2]} = \sigma^{[2]}(z_1^{[2]})$ $z_i^{[1]} = \left( \sum_{j=1}^n w_{ij}^{[1]} x_j \right) + b_i^{[1]} \quad a_i^{[1]} = \sigma^{[1]}(z_i^{[1]})$
Weights or Parameters	$\begin{cases} W^{[1]} \in \mathbb{R}^{h \times n} & b^{[1]} \in \mathbb{R}^h \\ W^{[2]} \in \mathbb{R}^h & b^{[2]} \in \mathbb{R} \end{cases}$

Some Activation Functions: $\sigma(x)$	
Sigmoid	$x$
Hyperbolic Tangent	$\tanh(x)$
ReLU	$\begin{cases} x, & x \geq 0 \\ 0, & x < 0 \end{cases}$
LeakyReLU	$\begin{cases} x, & x \geq 0 \\ \alpha x, & x < 0 \end{cases}$

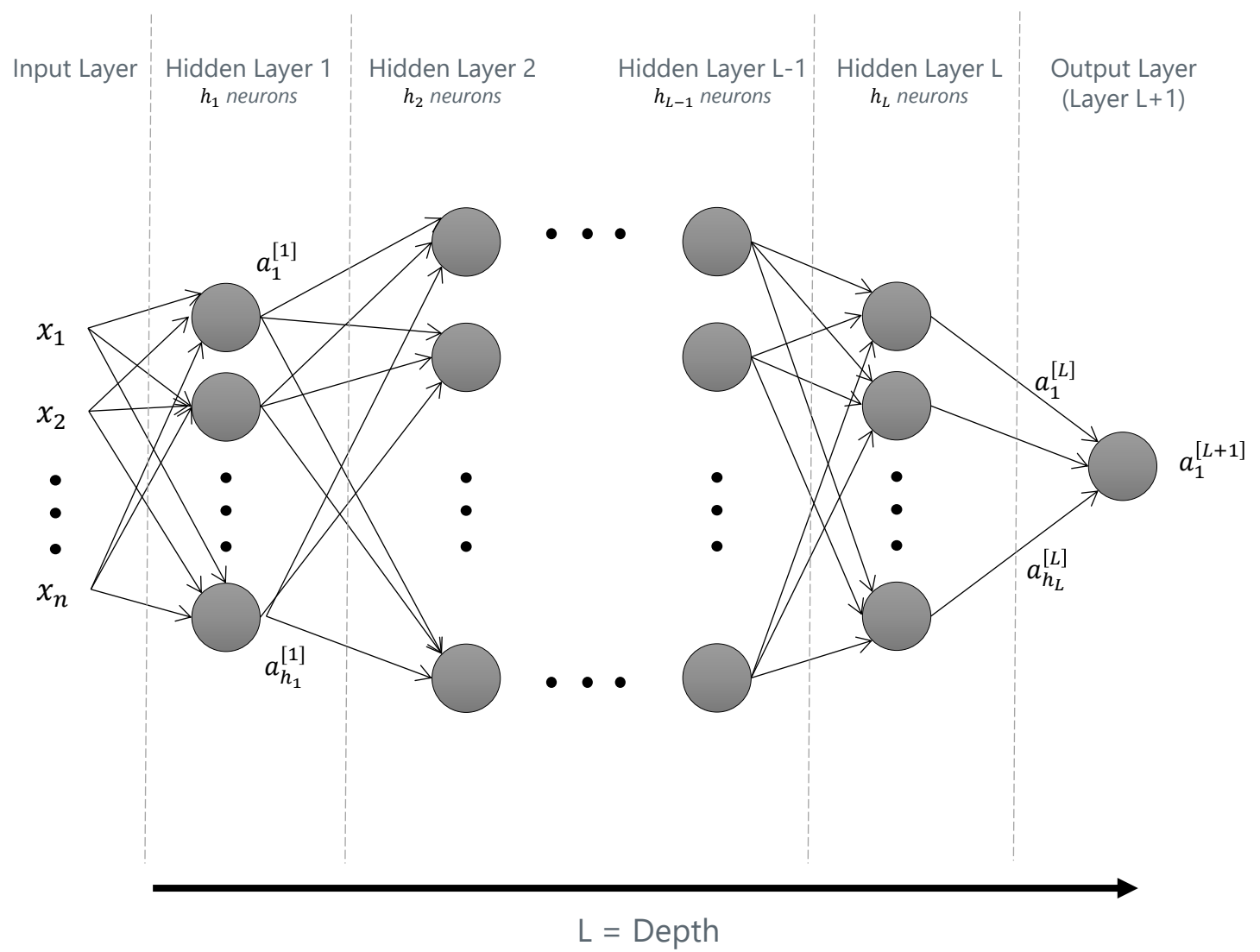
# Multi-layer Perceptron (MLP)



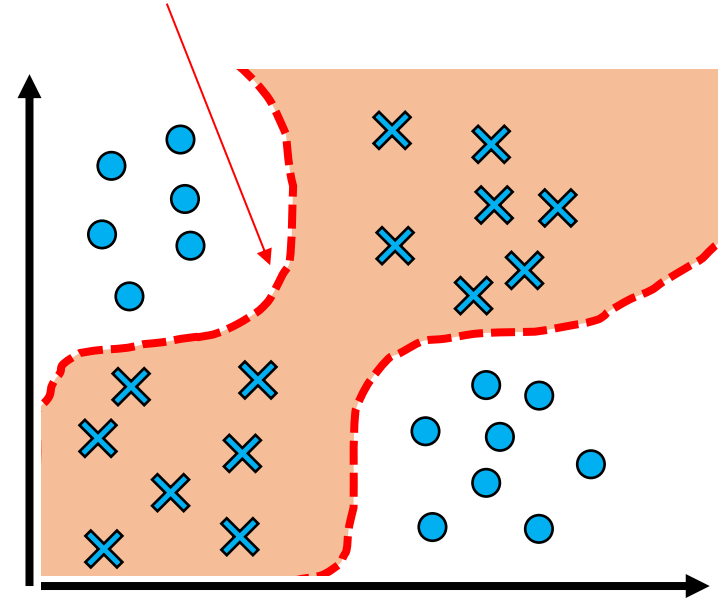
<b>Input</b>	$x = [x_1 \ x_2 \ \dots \ x_n] \in \mathbb{R}^n$
<b>Output</b>	$z_i^{[l]} = \left( \sum_{j=1}^{h_{l-1}} w_{ij}^{[l]} a_j^{[l-1]} \right) + b_i^{[l]} \quad a_i^{[l]} = \sigma^{[l]}(z_i^{[l]})$ $a_i^{[0]} = x_i$
<b>Weights or Parameters</b>	$W^{[l]} \in \mathbb{R}^{h_l \times h_{l-1}} \quad b^{[l]} \in \mathbb{R}^{h_l}$ $l = 1, \dots, L$

Some Activation Functions: $\sigma(x)$	
<b>Sigmoid</b>	$\frac{e^x}{1 + e^x}$
<b>Hyperbolic Tangent</b>	$\tanh(x)$
<b>ReLU</b>	$\begin{cases} x, & x \geq 0 \\ 0, & x < 0 \end{cases}$
<b>LeakyReLU</b>	$\begin{cases} x, & x \geq 0 \\ \alpha x, & x < 0 \end{cases}$

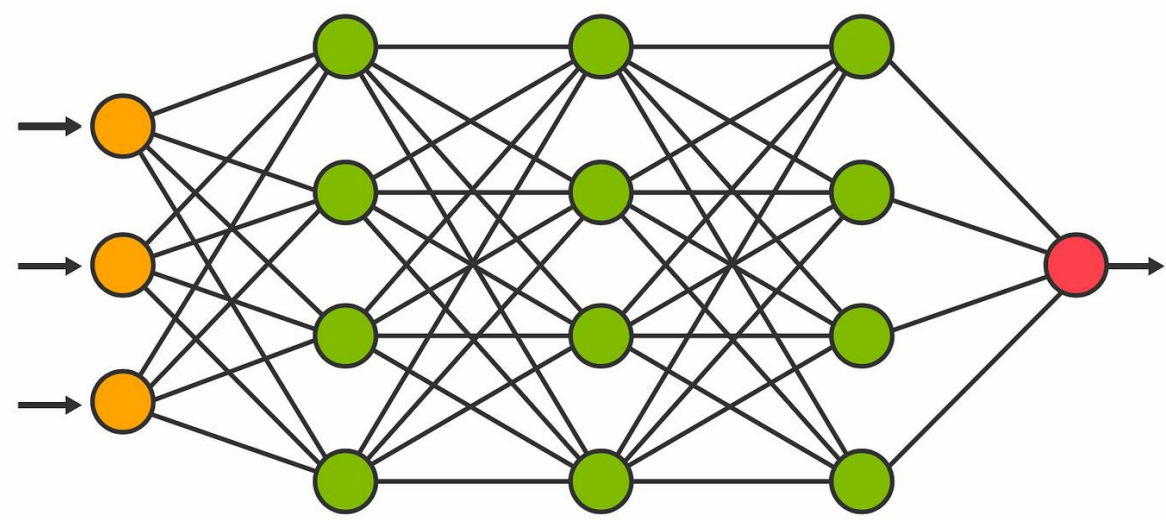
# Multi-layer Perceptron (MLP)



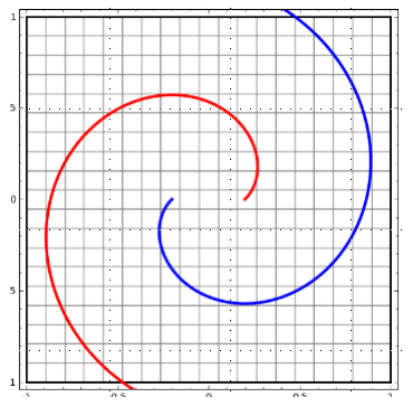
Decision Boundary



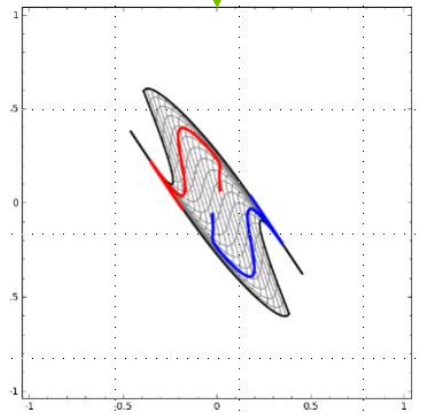
# Multi-layer Perceptron (MLP)



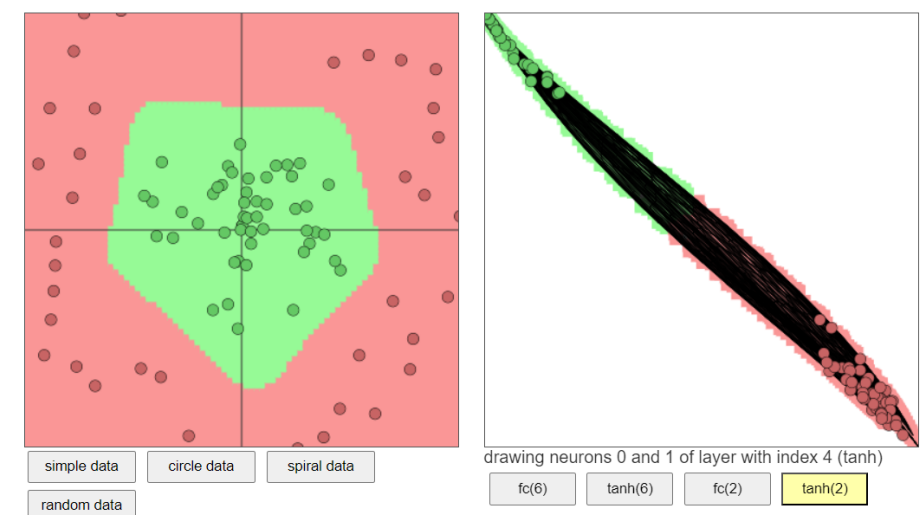
INPUT  
LAYER



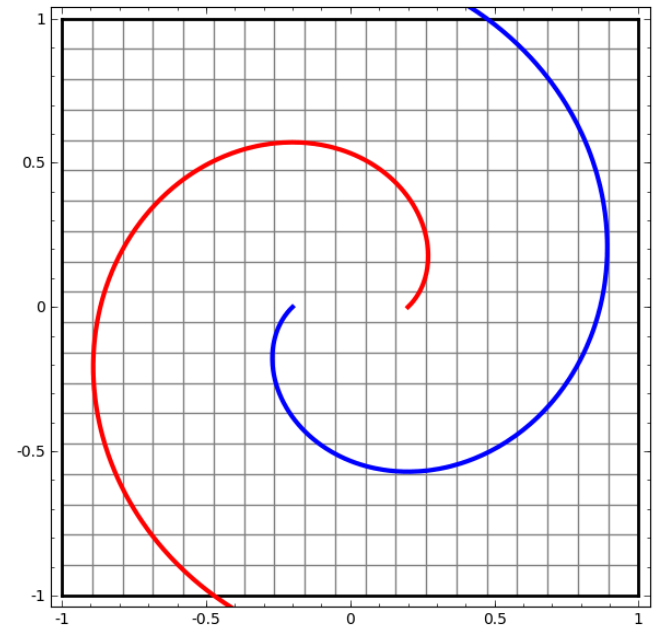
HIDDEN  
LAYERS



OUTPUT  
LAYER

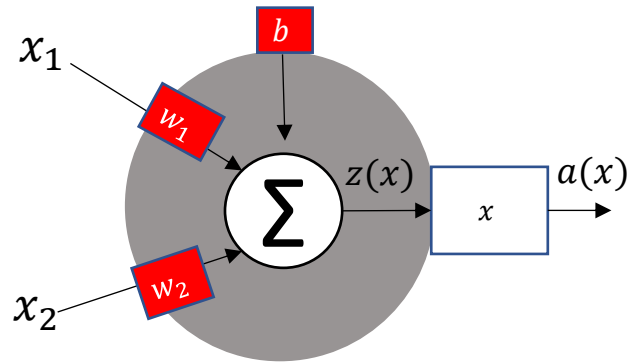


<https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>





# Single Neuron for classification – Loss



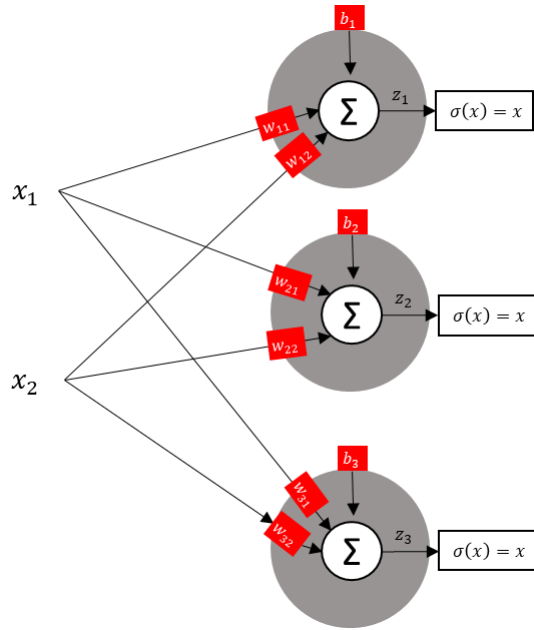
$$\hat{y}(x^{(i)}) = \frac{e^{a(x^{(i)})}}{1 + e^{a(x^{(i)})}} = \hat{y}(x^{(i)}, \mathbf{w}, b) = \frac{e^{(\mathbf{w}^T x^{(i)} + b)}}{1 + e^{(\mathbf{w}^T x^{(i)} + b)}} = P(Y = 1 | X = x^{(i)}; \mathbf{w}, b)$$

## Binary Cross Entropy minimization

$$(\mathbf{w}, b)^* = \underset{\mathbf{w}, b}{\operatorname{argmin}} \mathcal{C}(\mathbf{w}, b) = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\hat{y}(x^{(i)}, \mathbf{w}, b), y^{(i)})$$

$$\mathcal{L}(\hat{y}(x^{(i)}, \mathbf{w}, b), y^{(i)}) = -y^{(i)} \log(\hat{y}(x^{(i)}, \mathbf{w}, b)) - (1 - y^{(i)}) \log(1 - \hat{y}(x^{(i)}, \mathbf{w}, b))$$

# Multi-class classification – Loss



$$\hat{y}(x^{(i)}, \mathbf{W}, \mathbf{b}) = \text{softmax}(x^{(i)}) = \begin{bmatrix} \frac{e^{a_1(x^{(i)})}}{\sum_{j=1}^3 e^{a_j(x^{(i)})}} \\ \frac{e^{a_2(x^{(i)})}}{\sum_{j=1}^3 e^{a_j(x^{(i)})}} \\ \frac{e^{a_3(x^{(i)})}}{\sum_{j=1}^3 e^{a_j(x^{(i)})}} \end{bmatrix} = \begin{bmatrix} P(Y = 0 | X = x^{(i)}; \mathbf{W}, \mathbf{b}) \\ P(Y = 1 | X = x^{(i)}; \mathbf{W}, \mathbf{b}) \\ P(Y = 2 | X = x^{(i)}; \mathbf{W}, \mathbf{b}) \end{bmatrix}$$

## Cross Entropy minimization

If  $x^{(i)}$  belong to class 0:

$$\mathbf{y}^{(i)} = [1, 0, 0]$$

If  $x^{(i)}$  belong to class 1:

$$\mathbf{y}^{(i)} = [0, 1, 0]$$

If  $x^{(i)}$  belong to class 2:

$$\mathbf{y}^{(i)} = [0, 0, 1]$$

$$(\mathbf{W}, \mathbf{b})^* = \underset{\mathbf{w}, \mathbf{b}}{\operatorname{argmin}} \mathcal{C}(\mathbf{W}, \mathbf{b}) = \underset{\mathbf{w}, \mathbf{b}}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\hat{\mathbf{y}}(x^{(i)}, \mathbf{W}, \mathbf{b}), \mathbf{y}^{(i)})$$

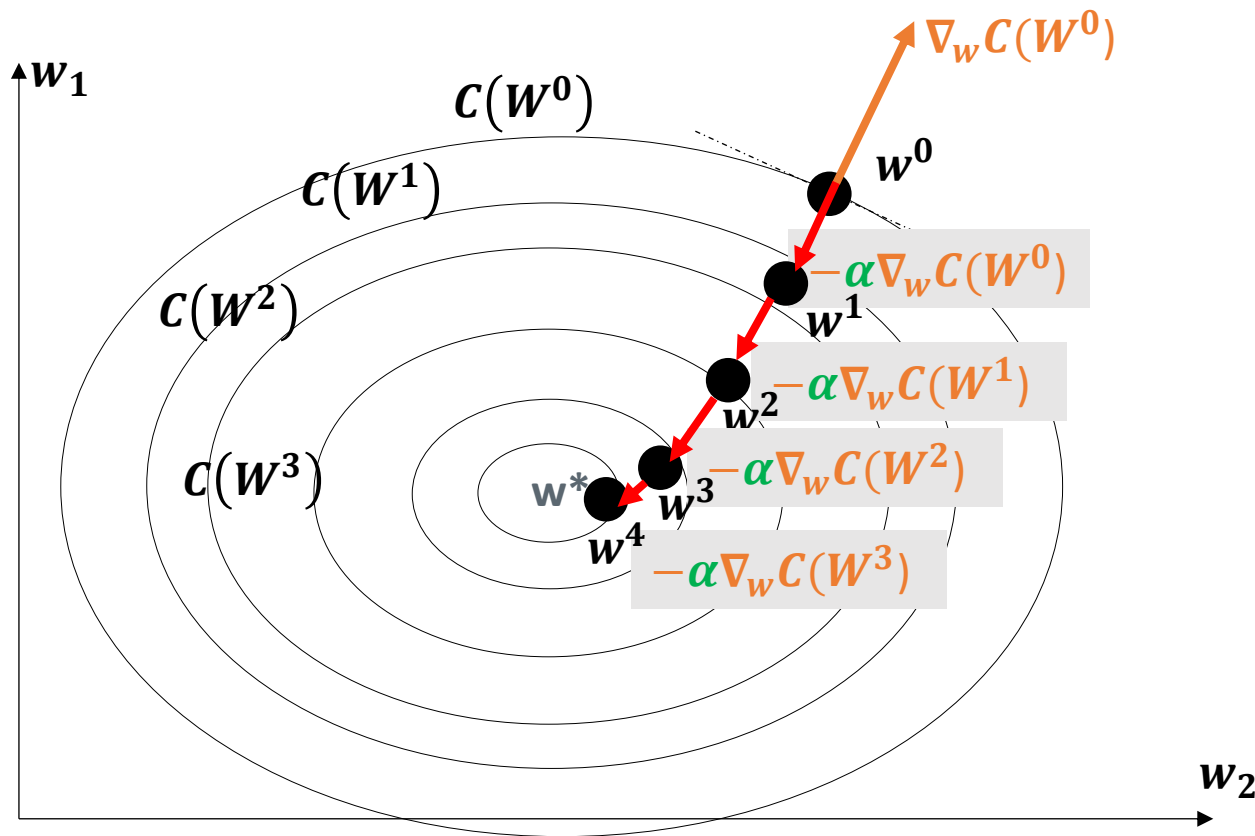
$$\mathcal{L}(\hat{\mathbf{y}}(x^{(i)}, \mathbf{W}, \mathbf{b}), \mathbf{y}^{(i)}) = - \sum_{j=1}^c y_j^{(i)} \log(\hat{y}_j(x^{(i)}, \mathbf{W}, \mathbf{b}))$$

# Optimizers

Without loss of generality suppose  $\mathbf{b}^* = \mathbf{0}$

$$\mathbf{W} = (\mathbf{w}, \mathbf{b})$$

$$\mathbf{W}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \mathcal{C}(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\hat{\mathbf{y}}(\mathbf{x}^{(i)}, \mathbf{W}), \mathbf{y}^{(i)})$$



## Gradient Descent Algorithm

- 0.1 Choose learning rate  $\alpha$
- 0.2 Initialize weights  $\mathbf{W} = \mathbf{W}^0$

For epoch=1..num\_epochs do

1. Compute  $\nabla_{\mathbf{W}} \mathcal{C}(\mathbf{W}^k)$
2. Update weights

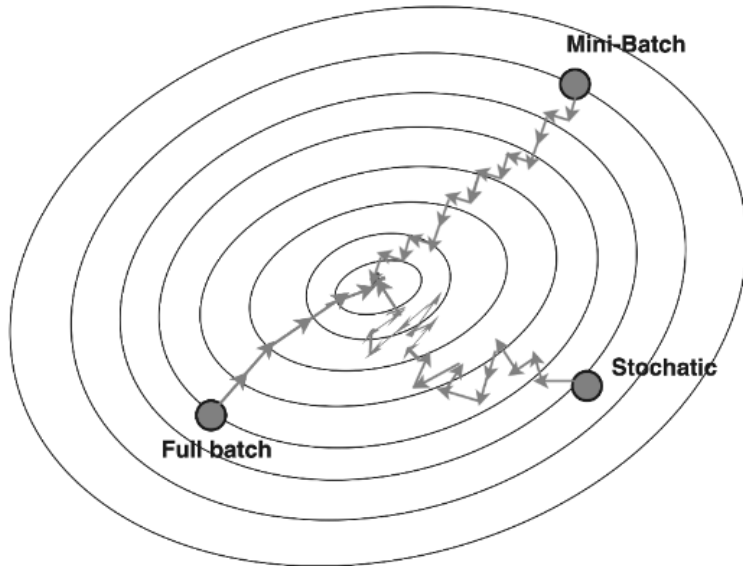
$$\mathbf{W} = \mathbf{W} - \alpha \nabla_{\mathbf{W}} \mathcal{C}(\mathbf{W}^k)$$

# Optimizers – Batch Size

batch\_size = N

## Gradient Descent Algorithm

- 0.1. Choose learning rate  $\alpha$
- 0.2. Initilize weights  $W=W^0$
- for epoch=1..num\_epochs do
  - 1. Compute  $\nabla_W \mathcal{C}(W)$  on the
  - 2. Update the weights  
 $W = W - \alpha \nabla_W \mathcal{C}(W)$



batch\_size = 1

## Stochastic Gradient Descent

- 0.1. Choose learning rate  $\alpha$
- 0.2. Initilize weights  $W=W^0$
- for epoch=1..num\_epochs do
  - 1. shuffle the dataset
  - for i=1..N do
    - 2. Compute  $\nabla_W \mathcal{L}_W(y^{(i)}, a(x^{(i)}, W))$
    - 3. Update the weights  
 $W = W - \alpha \nabla_W \mathcal{L}_W(y^{(i)}, a(x^{(i)}, W))$

batch\_size = b, 1<b<N

## Mini-batch Gradient Descent

- 0.1. Choose learning rate  $\alpha$
- 0.2. Choose the batch size b
- 0.3. Initilize weights  $W=W^0$

- for epoch=1..num\_epochs do
  - 1. shuffle the dataset
  - for k=0,..., floor(N/b) - 1
    - 2. Compute

$$\nabla_W B(W) = \frac{1}{b} \sum_{i=kb}^{(k+1)b} \nabla_W \mathcal{L}_W(y^{(i)}, a(x^{(i)}, W))$$

- 3. Update the weights  
 $W = W - \alpha \nabla_W B(W)$

# Optimizers – Algorithms

## Basic Mini-Batch algorithms

### SGD: Stochastic Gradient Descent

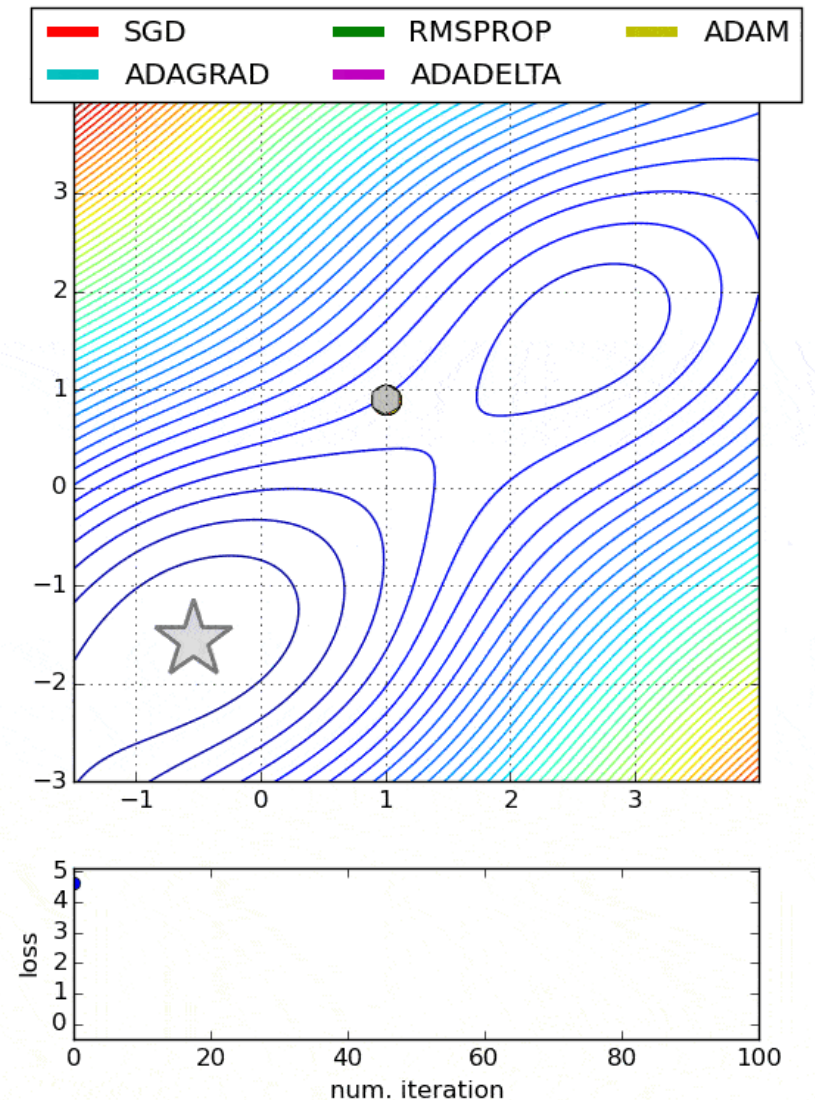
SGD is a foundational algorithm that can sometimes outperform others, especially in certain problem domains. However, it can be slow to converge and may struggle with local minima. The momentum parameter can help improve convergence by smoothing out the updates

### Adam

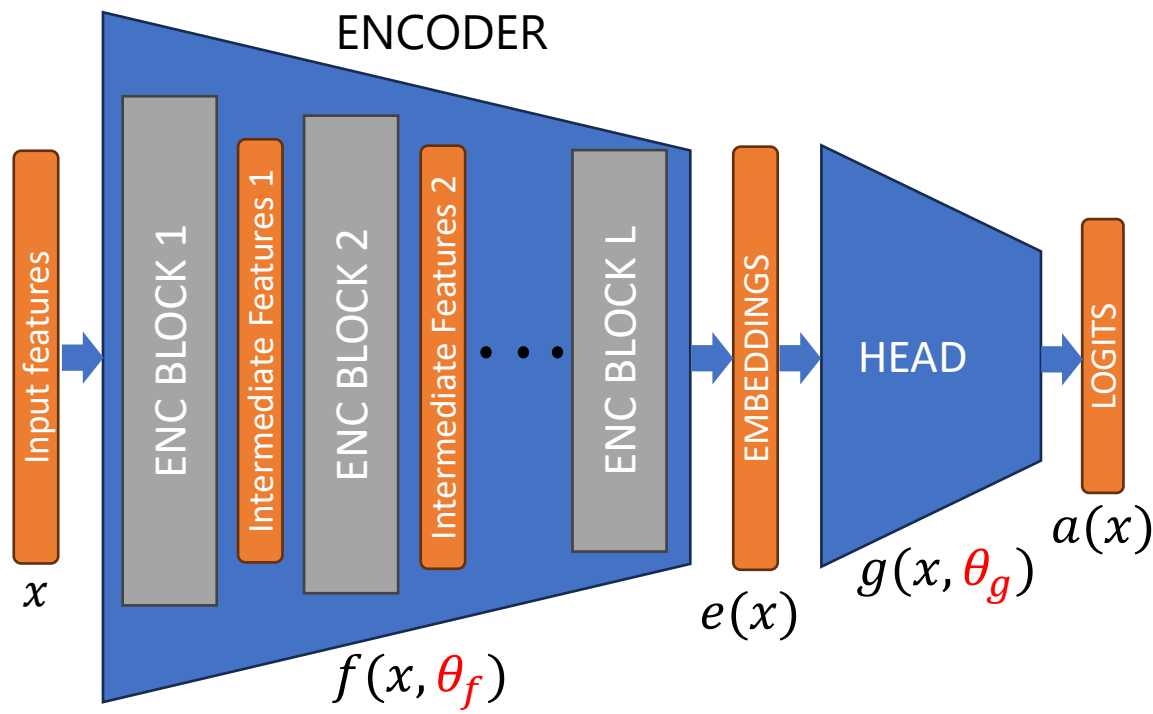
Adam is a popular optimization algorithm that combines the benefits of adaptive learning rates and momentum, making it an excellent starting choice for many applications

### AdamW

AdamW is similar to Adam but decouples weight decay from the gradient updates, providing more effective regularization. This makes AdamW more suitable when weight regularization is desired, helping to improve generalization

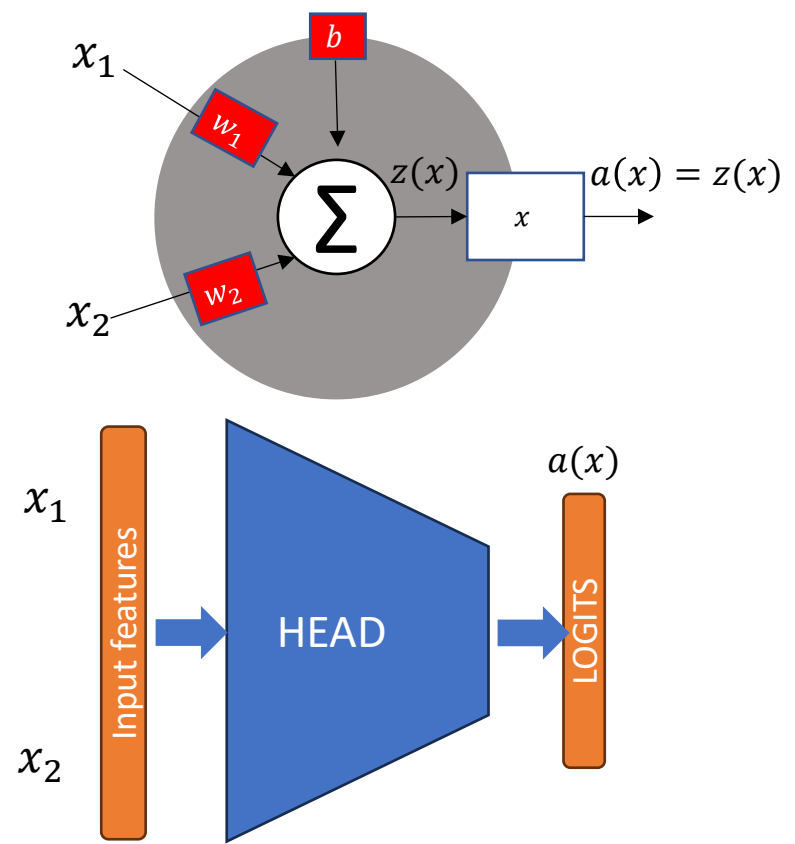
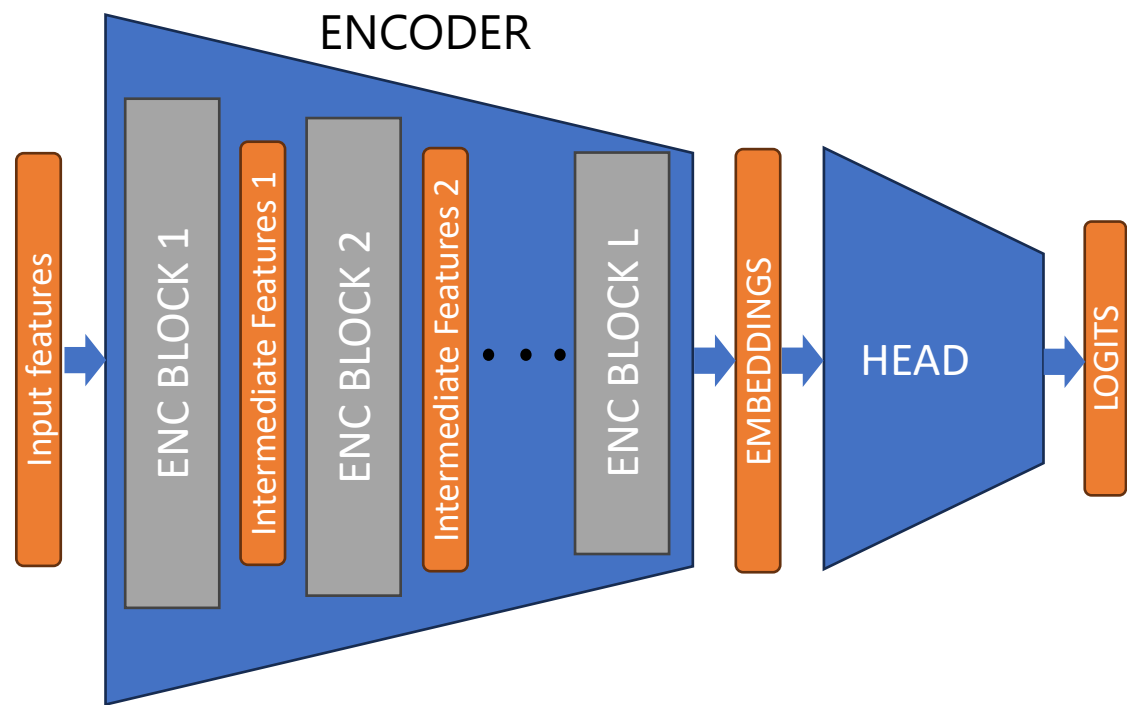


# Anatomy of a Neural Network

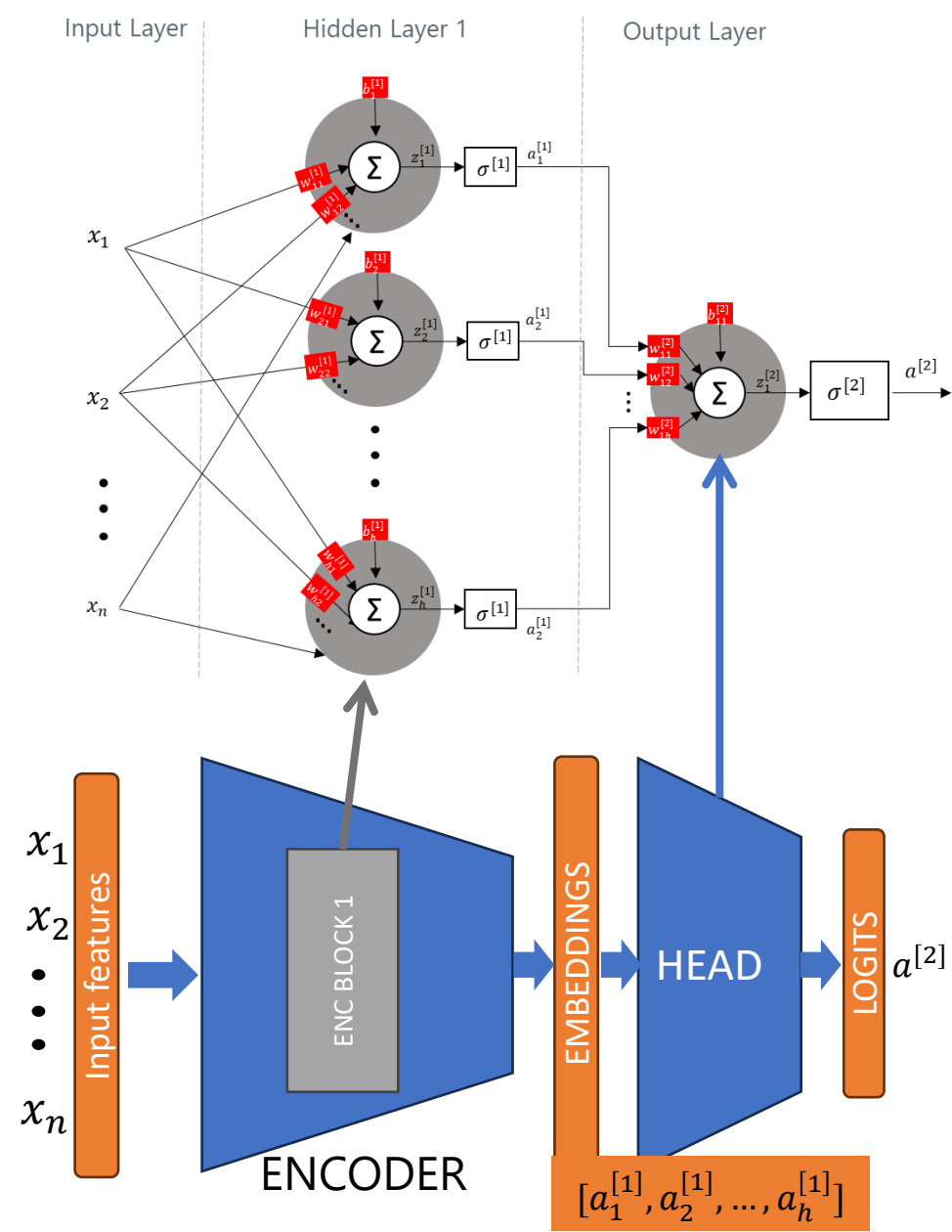
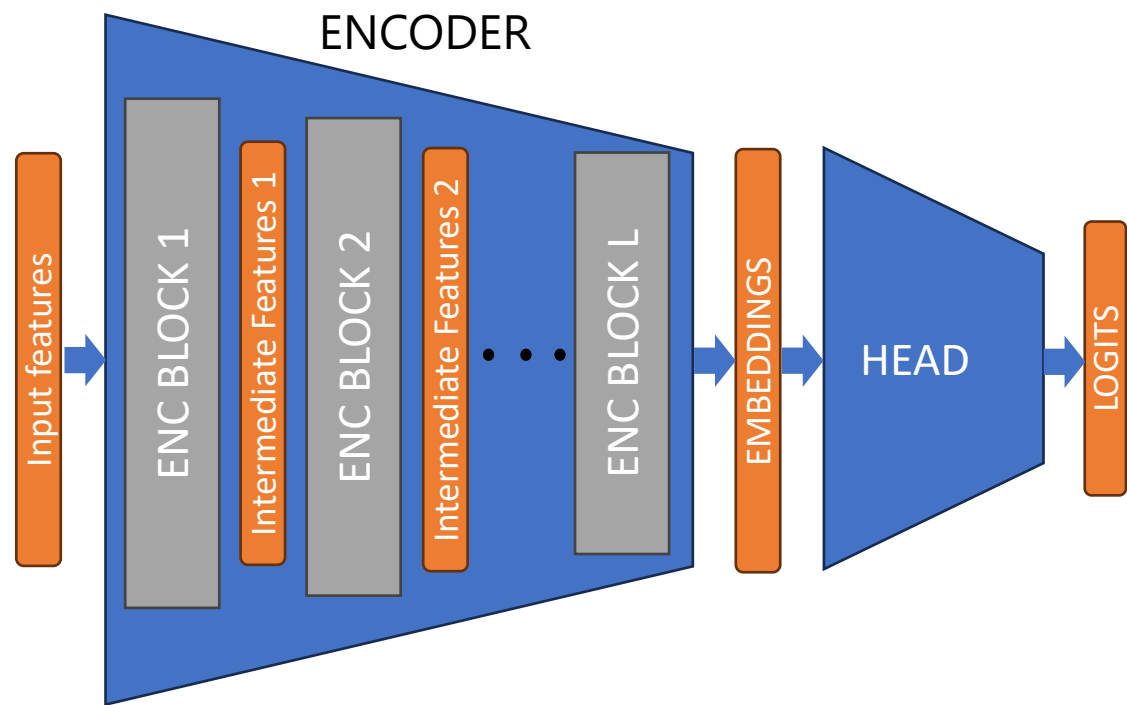


<b>Embeddings (or high level features)</b>	$e(x^{(i)}) = f(x^{(i)}, \theta_f)$
<b>Logits</b>	$a(x^{(i)}) = g(e(x^{(i)}), \theta_g)$
<b>Output</b>	$\hat{y}(x^{(i)}) = \sigma(a(x^{(i)}))$
<b>Parameters/Weights</b>	$\theta_f, \theta_g$ $\theta_f$ : All weights and biases of the encoder $\theta_g$ : All weights and biases of the head

# Anatomy of a Neural Network – A single neuron

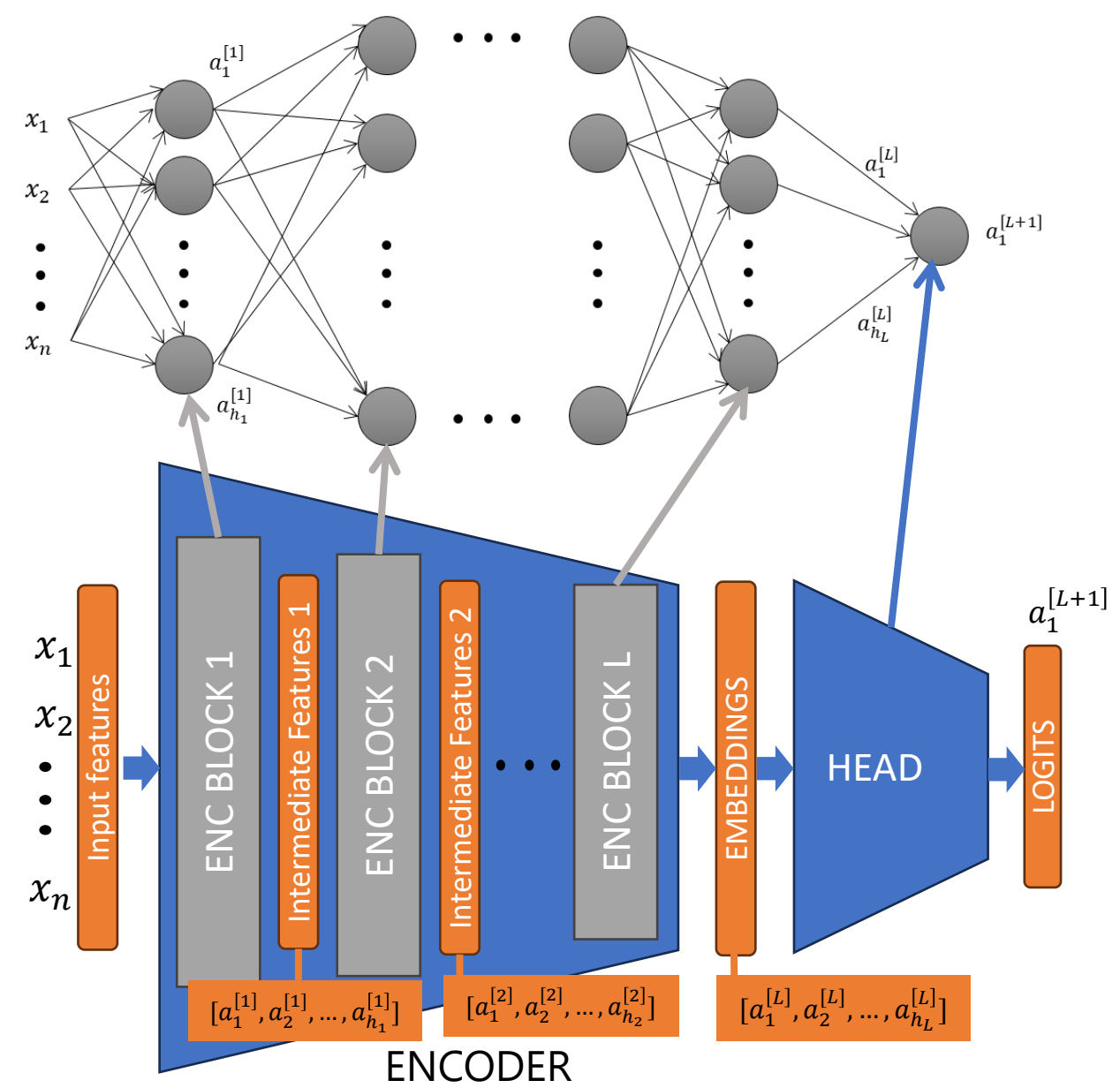


# Anatomy of a Neural Network – Single-Layer Perceptron

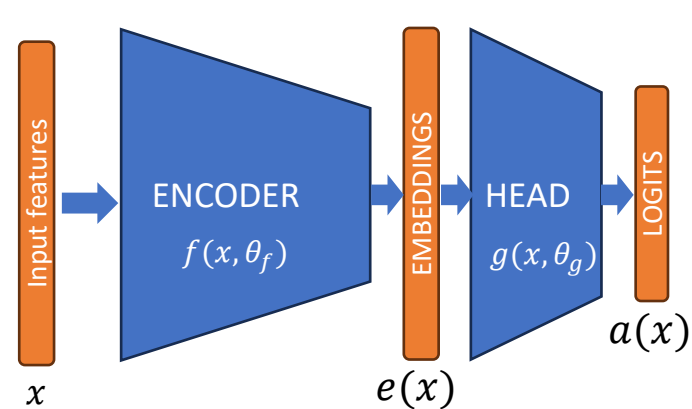




# Anatomy of a Neural Network – Multi-Layer Perceptron



# Classification Summary



Number of classes $C$	Dimensions	Output function $\hat{y}(x)$	Loss function
Binary Classification ( $C=2$ )	$a(x) \in \mathbb{R}$ $\hat{y}(x) \in [0,1]$ $y(x) \in \{0,1\}$	$\hat{y}(x) = \text{logistic}(x)$	Binary Cross Entropy
Multi-class classification ( $C>2$ )	$a(x) \in \mathbb{R}^C$ $\hat{y}(x) \in \mathbb{S}_1^C(0)$ (unit sphere) $y(x) \in \{\mathbf{1}^1, \mathbf{1}^2, \dots, \mathbf{1}^C\}$ $\mathbf{1}^i = 1$ in $i$ -th, 0 otherwise	$\hat{y}(x) = \text{softmax}(x)$	Cross Entropy

## Logistic function / sigmoid

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

## Binary Cross Entropy (BCE) Loss:

$$\mathcal{L}(\hat{y}(x^{(i)}), y^{(i)}) = -y^{(i)} \log(\hat{y}(x^{(i)})) - (1 - y^{(i)}) \log(1 - \hat{y}(x^{(i)}))$$

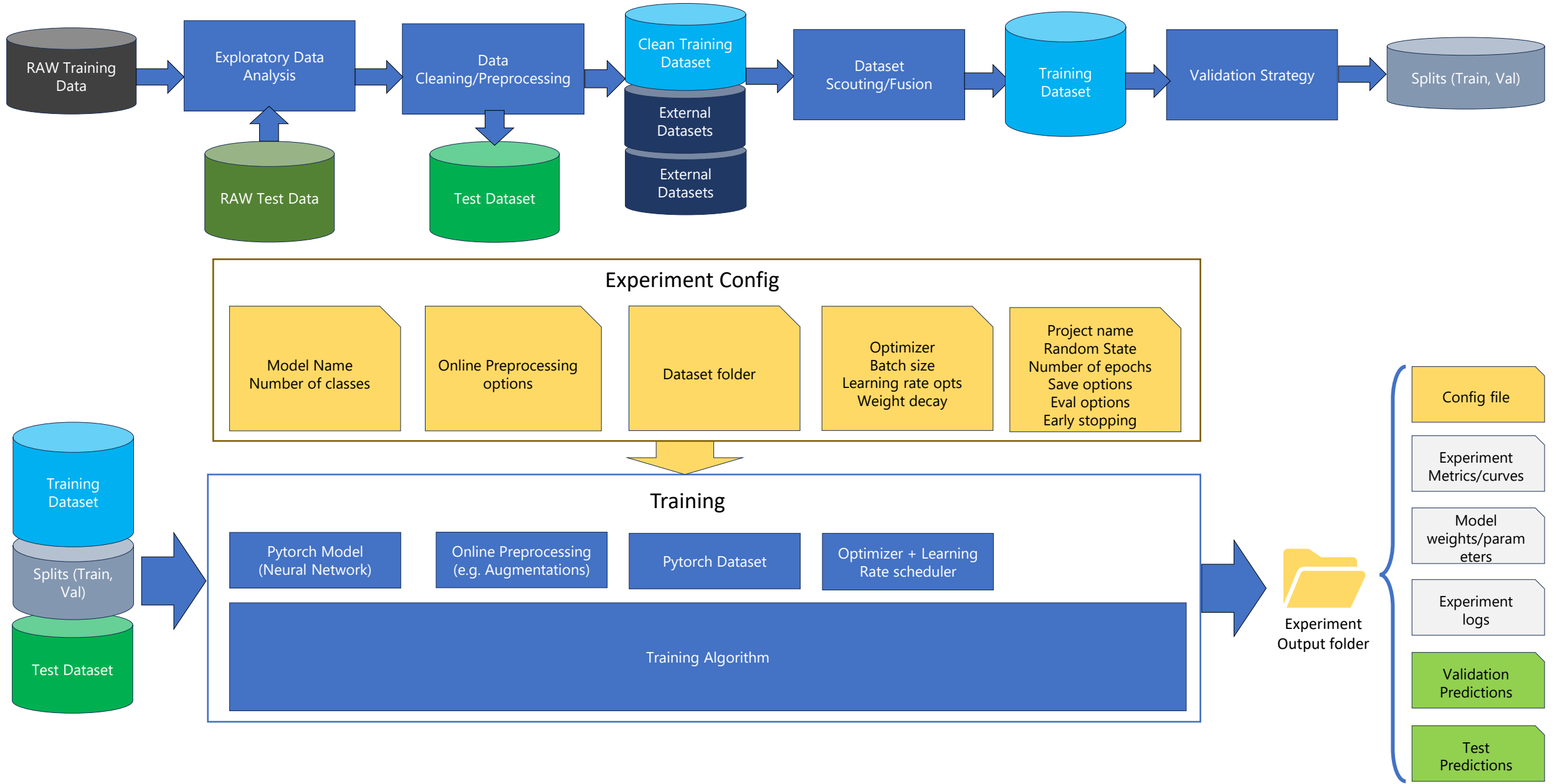
## Softmax

$$y(x^{(i)}) = \text{softmax}(x^{(i)}) = \left[ \frac{e^{a_1(x^{(i)})}}{\sum_{j=1}^C e^{a_j(x^{(i)})}}, \dots, \frac{e^{a_C(x^{(i)})}}{\sum_{j=1}^C e^{a_j(x^{(i)})}} \right]$$

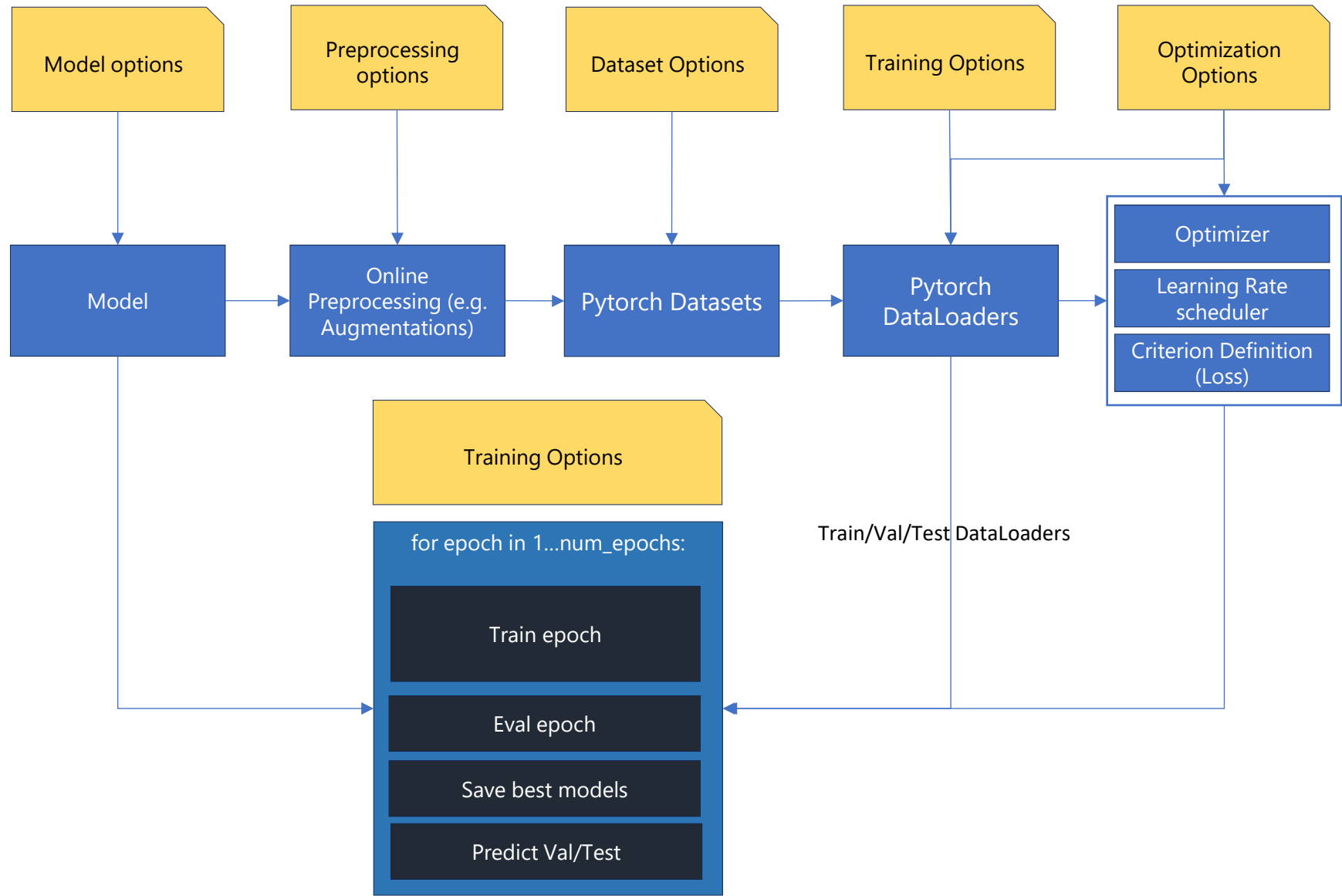
## Cross Entropy (CE) Loss:

$$\mathcal{L}(\hat{y}(x^{(i)}), y^{(i)}) = - \sum_{j=1}^C y_j^{(i)} \log(\hat{y}(x^{(i)}))$$

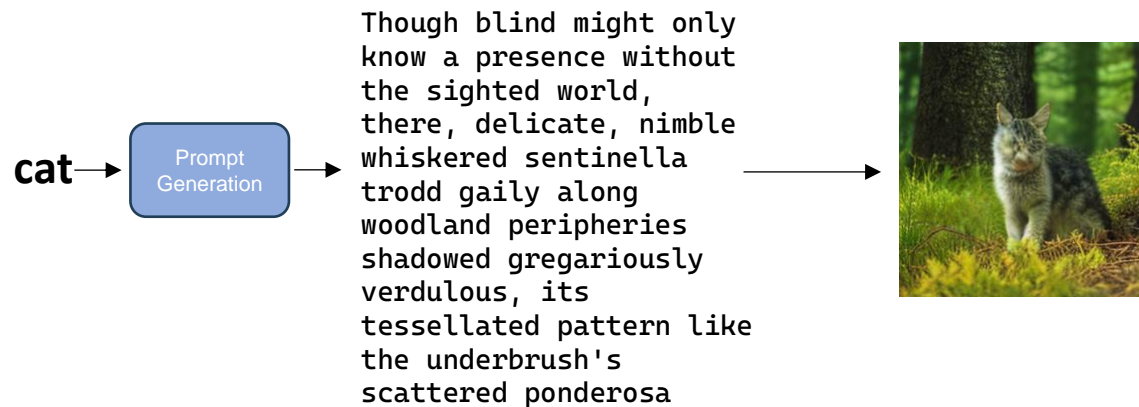
# Training recipe



# Classification Summary

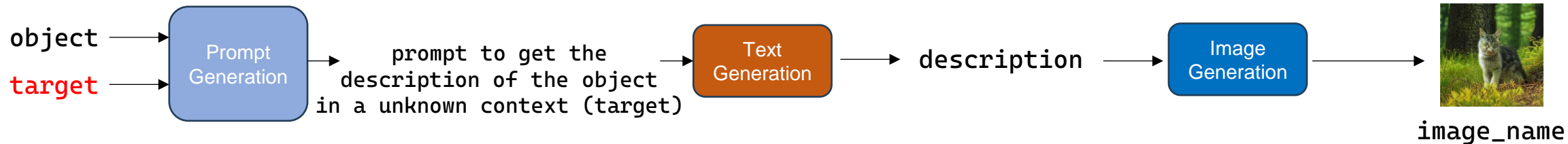


# The challenge!



## Description of the training features

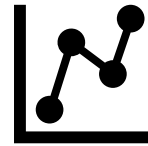
- **image\_name** : the file name of the image generated from an unknown text-2-image generator from the relative **description**
- **object**: categorical value in {**cat**, **dog**, **bus**, **car**, **person**, **house**}
- **description**: the “scene description”, generated from an unknown LLM
- **target**: the unknown context to be predicted, there are only 4 different contexts (from 0 to 3)



# The challenge!



- **End Date:** Friday 06<sup>th</sup> at 15:00
- **Private Leaderboard:** Friday 06<sup>th</sup> at 15:00!
- **Award and Presentations:** Friday 06<sup>th</sup> at 16:00



- **Maximum Daily submission limit:** 12
- **Scored private submissions:** 2



- **Team:** Groups should be of 4 people. Groups of 3 or 5 are also accepted, but no exception will be granted for other sizes



Thank you for your  
attention!