

Lecture 11  
Deep Learning 01  
*Shallow Neural Networks*

2024-11-06

Sébastien Valade



1. Introduction

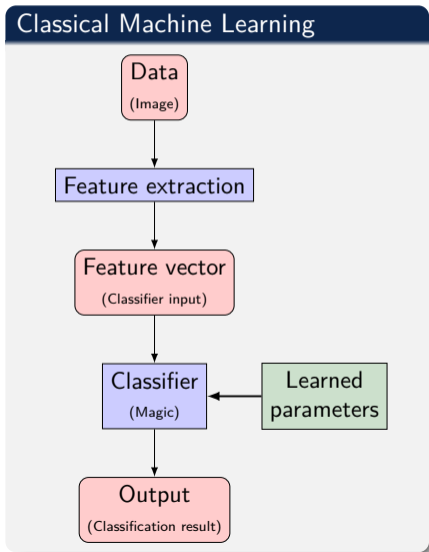
2. Perceptron

3. Multilayer perceptron (MLP)

4. Application

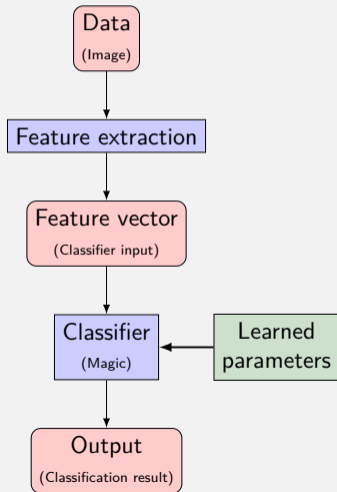
5. Glossary

## From classical learning to modern learning

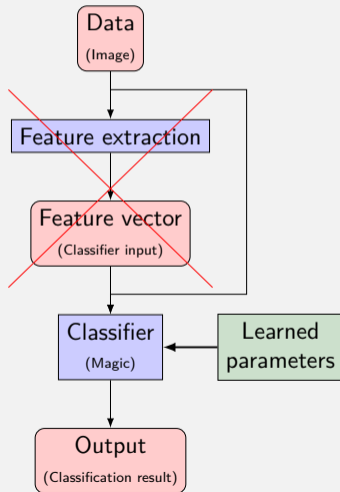


## From classical learning to modern learning

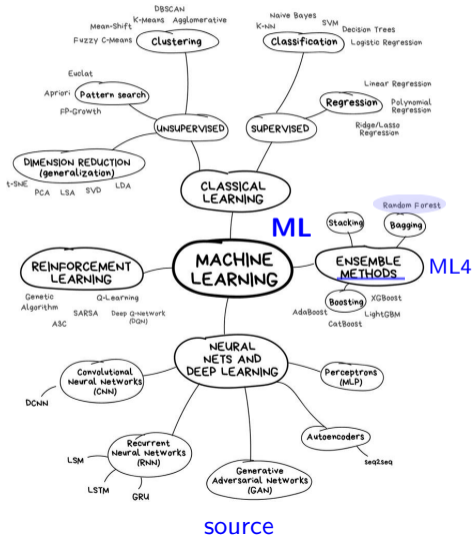
### Classical Machine Learning



### Modern Machine Learning

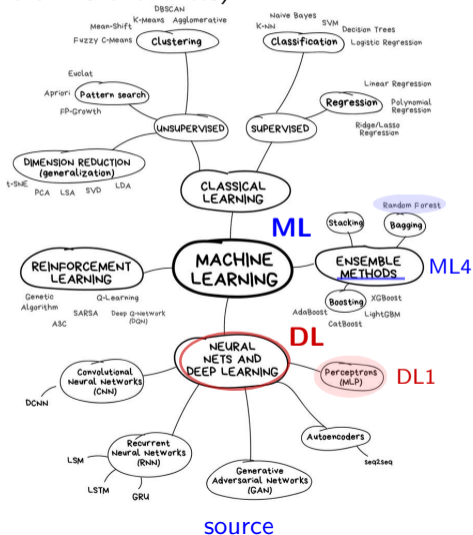


Last week: **Random Forests** (*Ensemble Method*) - ML4

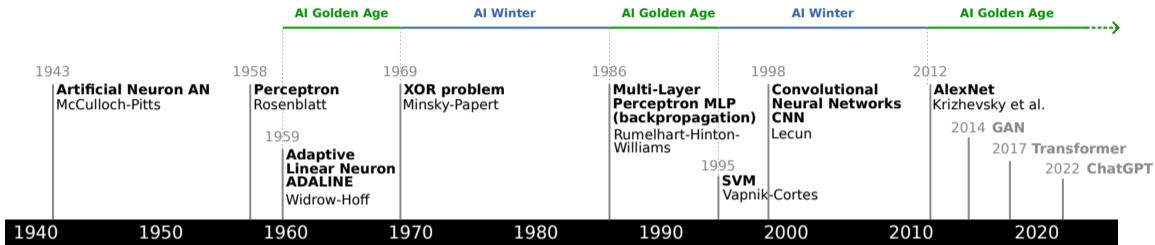


Last week: **Random Forests** (*Ensemble Method*) - ML4

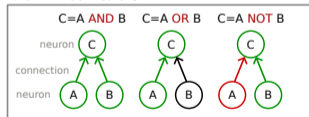
This week: **Neural Networks** (*Part-1: shallow nets*) - DL1



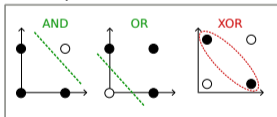
# Brief history of Neural Networks



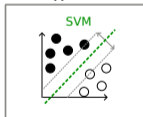
1943 - Artificial Neurons



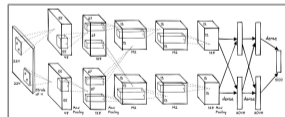
1969 - XOR problem



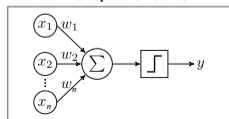
1995 - Support Vector Machines (SVM)



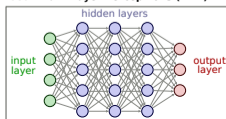
2012 - AlexNet



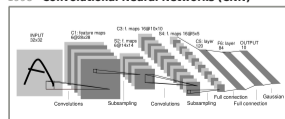
1958-1959 - Perceptron & Adaline



1986 - Multi-Layer Perceptrons (MLP)



1998 - Convolutional Neural Networks (CNN)



1. Introduction

**2. Perceptron**

1. definition
2. learning algorithm
3. gradient descent
4. limitations

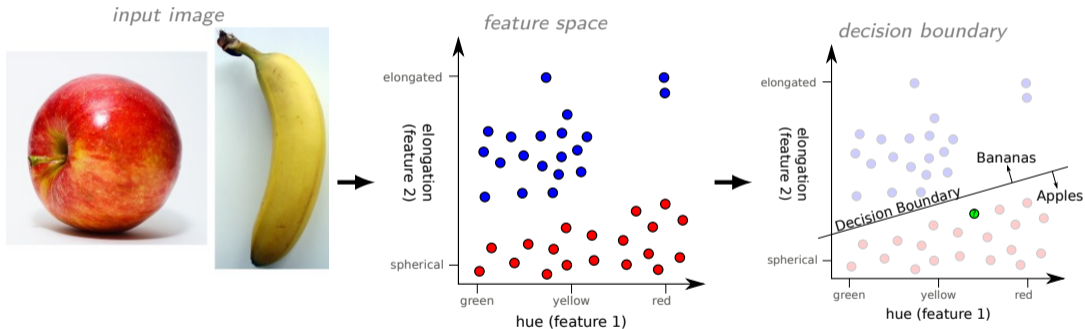
3. Multilayer perceptron (MLP)

4. Application

5. Glossary



Recall our toy example from lecture 9: classify fruit images into either bananas or apples



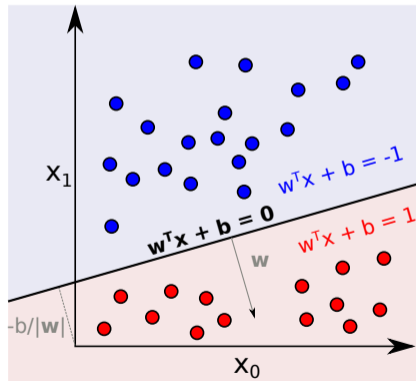
## Perceptron classifier

⇒ algorithm which classifies data based on linear decision boundary

⇒ perceptron:

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

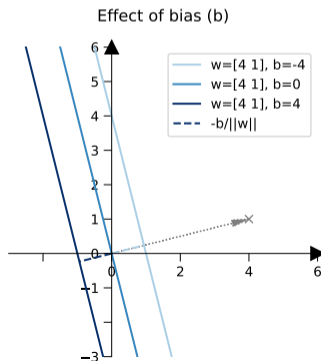
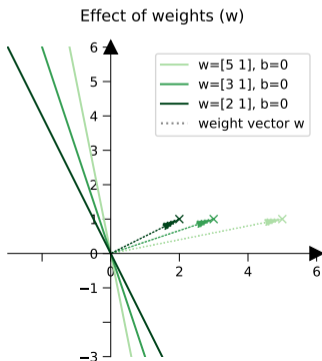
- $\hat{y} \in \{-1, 1\}$ : predicted class → *banana or apple*
- $\mathbf{x} \in \mathbb{R}^2$ : feature matrix (n, 2) → [hue, elongation]
- $\mathbf{w} \in \mathbb{R}^2$ : weight vector (2,) → needs to be learned
- $b \in \mathbb{R}$ : bias → needs to be learned
- *sign*: sign function returning the sign of a real number



## Perceptron classifier

⇒ influence of **weight vector  $w$**  & **bias  $b$**  on decision boundary:

- **weight vector  $w$** : determines the *normal vector* to the decision boundary
- **bias  $b$** : shifts the decision boundary in the direction of the weight vector ( $b > 0$  shifts the boundary in the negative direction of the weight vector, and vice versa)



## Perceptron

⇒ representation as an artificial neuron called the threshold logic unit (TLU), which:

1. computes a linear function of the input vectors  $\mathbf{x}_i$  and associated weights  $\mathbf{w}_i$ , plus a bias term  $b$ :

$$z = \mathbf{w}_0 \mathbf{x}_0 + \mathbf{w}_1 \mathbf{x}_1 + b = \mathbf{w}^T \mathbf{x} + b$$

2. applies a **step function**  $\sigma$ , typically the modified sign function

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

## Perceptron

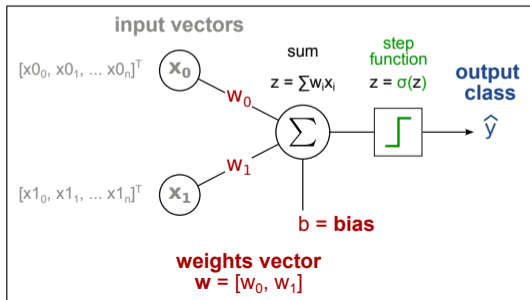
⇒ representation as an artificial neuron called the threshold logic unit (TLU), which:

1. computes a linear function of the input vectors  $\mathbf{x}_i$  and associated **weights**  $w_i$ , plus a **bias term**  $b$ :

$$z = \mathbf{w}_0 \mathbf{x}_0 + \mathbf{w}_1 \mathbf{x}_1 + b = \mathbf{w}^T \mathbf{x} + b$$

2. applies a **step function**  $\sigma$ , typically the modified sign function

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



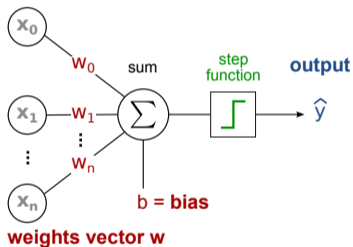
$$\begin{aligned} \hat{y} &= \sigma(\mathbf{w}^T \mathbf{x} + b) \\ &= \sigma \left( \sum_{i=1}^n \mathbf{w}_i x_i + b \right) \\ &= \sigma \left( \begin{pmatrix} w_0 & w_1 \end{pmatrix} \cdot \begin{pmatrix} x_{0_0} & x_{1_0} \\ x_{0_1} & x_{1_1} \\ \vdots & \vdots \\ x_{0_n} & x_{1_n} \end{pmatrix} + b \right) \end{aligned}$$

## Perceptron

⇒ another representation is with the weight vector  $\tilde{\mathbf{w}}$  and augmented input vector  $\tilde{\mathbf{x}}$ :

→ the bias term  $b$  is learned as a weight  $w_n = b$ , and the input vector  $x$  is augmented with a vector of values  $x_n = 1$

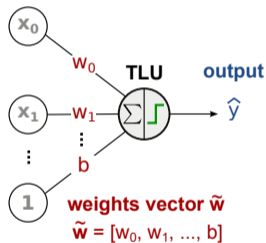
inputs  $x$



$$\hat{y} = \sigma(\mathbf{w}^T x + b)$$

$$= \sigma\left(\sum_{i=1}^n w_i x_i + b\right)$$

inputs  $\tilde{\mathbf{x}} = [x_0, x_1, \dots, 1]$



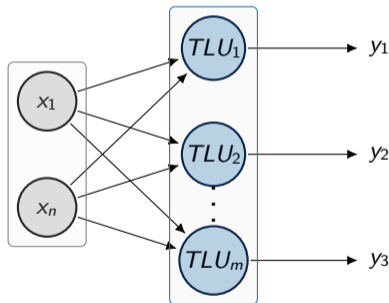
$$\hat{y} = \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})$$

$$= \sigma\left(\sum_{i=1}^n \tilde{w}_i \tilde{x}_i\right)$$

where:  $\tilde{\mathbf{w}}_n = b, \tilde{\mathbf{x}}_n = \text{vector of } 1$

## Perceptron

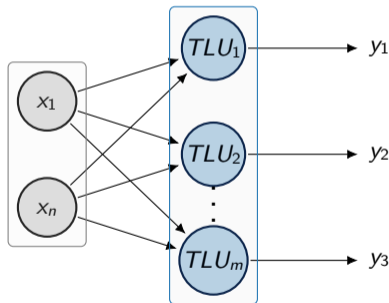
- ⇒ a perceptron can be composed of one or more TLUs organized in a single layer, where every TLU is connected to every input: such a layer is called a fully connected layer, a.k.a. a dense layer
- ⇒ when the output layer contains several TLUs, the perceptron becomes a multilabel classifier



- ⇒ the perceptron is one of the simplest Artificial Neural Network (ANN) architecture

## Perceptron

- ⇒ a perceptron can be composed of one or more TLUs organized in a single layer, where every TLU is connected to every input: such a layer is called a fully connected layer, a.k.a. a dense layer
- ⇒ when the output layer contains several TLUs, the perceptron becomes a multilabel classifier



- ⇒ the perceptron is one of the simplest Artificial Neural Network (ANN) architecture



Nota Bene:

- when TLUs are organized in multiple layers, the network is called a **multilayer perceptron (MLP)**
- in MLPs, the step function is often called the **activation function**, which can take various forms (e.g., step, sigmoid, ReLU, Tanh, etc.)

⇒ how is the perceptron trained?

*i.e., how are the weights  $w$  and bias  $b$  learned?*

Nota Bene:

- when TLUs are organized in multiple layers, the network is called a **multilayer perceptron (MLP)**
- in MLPs, the step function is often called the **activation function**, which can take various forms (e.g., step, sigmoid, ReLU, Tanh, etc.)

⇒ how is the perceptron trained?

*i.e., how are the weights  $\mathbf{w}$  and bias  $b$  learned?*

## Perceptron learning algorithm

- ⇒ there is no analytical solution to the perceptron learning problem
- ⇒ **iterative algorithm** to learn the weight vector  $\mathbf{w}$  and bias  $b$  that minimize the classification error:

initialize  $\tilde{\mathbf{w}}$  with random values

for each training samples  $\tilde{\mathbf{x}}_i$ :

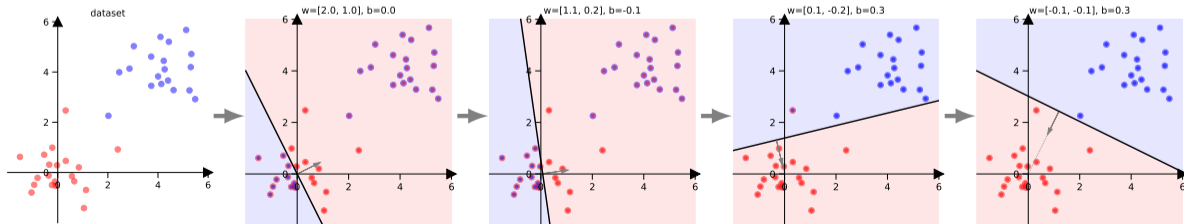
predict  $\hat{y}_i = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i)$

update  $\tilde{\mathbf{w}}$  if  $\hat{y}_i \neq y_i$ :

$\tilde{\mathbf{w}} = \tilde{\mathbf{w}} + \eta \cdot (y_i - \hat{y}_i) \cdot \tilde{\mathbf{x}}_i$  where  $\eta$  is the learning rate

where  $\mathbf{w} = \tilde{\mathbf{w}}[: -1]$ ,  $b = \tilde{\mathbf{w}}[-1]$

- ⇒ Illustration of the convergence of the perceptron learning algorithm:



## Perceptron learning algorithm

- ⇒ there is no analytical solution to the perceptron learning problem
- ⇒ **iterative algorithm** to learn the weight vector  $\mathbf{w}$  and bias  $b$  that *minimize the classification error*:

initialize  $\tilde{\mathbf{w}}$  with random values

for each training samples  $\tilde{\mathbf{x}}_i$ :

predict  $\hat{y}_i = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i)$

update  $\tilde{\mathbf{w}}$  if  $\hat{y}_i \neq y_i$ :

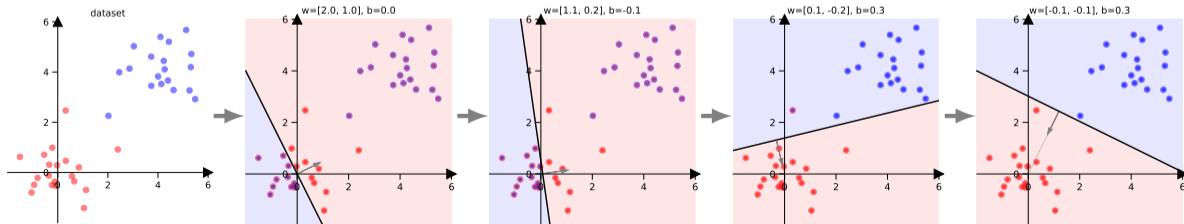
$$\tilde{\mathbf{w}} = \tilde{\mathbf{w}} + \eta \cdot (y_i - \hat{y}_i) \cdot \tilde{\mathbf{x}}_i$$

where  $\mathbf{w} = \mathbf{w}[: -1]$ ,  $b = \mathbf{w}[-1]$

where  $\eta$  is the learning rate

where does this come from?

- ⇒ Illustration of the convergence of the perceptron learning algorithm:



**Gradient descent of the Loss function** (for the perceptron)

⇒ learning the weights  $\tilde{\mathbf{w}}$  means modifying them such that predicted labels  $\hat{y}$  get closer to true labels  $y$

**Gradient descent of the Loss function** (for the perceptron)

⇒ learning the weights  $\tilde{\mathbf{w}}$  means modifying them such that predicted labels  $\hat{y}$  get closer to true labels  $y$

⇒ **loss function  $\mathcal{L}$**  = measure of the difference between predicted and true labels

$$\text{L2 loss} = \text{mean squared error (MSE)} = \boxed{\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2}$$

## Gradient descent of the Loss function (for the perceptron)

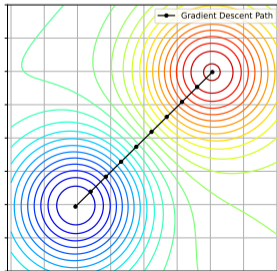
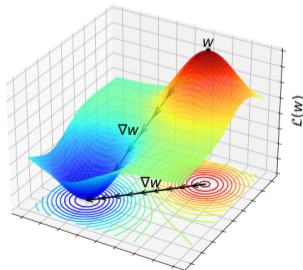
⇒ learning the weights  $\tilde{w}$  means modifying them such that predicted labels  $\hat{y}$  get closer to true labels  $y$

⇒ **loss function  $\mathcal{L}$**  = measure of the difference between predicted and true labels

$$\text{L2 loss} = \text{mean squared error (MSE)} = \boxed{\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2}$$

⇒ **gradient descent**: starting from a random point on the function  $\mathcal{L}(w)$ , move in the direction of the steepest descending gradient with a step size  $\eta$  (learning rate)

update rule:  $\boxed{w = w - \eta \nabla w}$  where  $\nabla w = \frac{d\mathcal{L}}{dw}$  (for simplicity  $\tilde{w}$  is here annotated  $w$ )



## Gradient descent of the Loss function

⇒ computing the derivative  $\frac{d\mathcal{L}}{dw}$ :

$$\begin{aligned} \mathcal{L} &= \frac{1}{2}(y - \hat{y})^2 \\ \frac{d\mathcal{L}}{dw} &= \frac{d}{dw} \left\{ \frac{1}{2}(y - \hat{y})^2 \right\} \\ &= \frac{1}{2} \cdot 2 \cdot (y - \hat{y}) \cdot \frac{d}{dw}(y - \hat{y})^{2-1} \quad \text{applying the power rule, where } \frac{d}{dw}(u^n) = n \cdot u^{n-1} \cdot \frac{du}{dw} \\ &= (y - \hat{y}) \cdot \frac{d}{dw}(y - w \cdot x) \quad \text{substituting } \hat{y} = w \cdot x \\ &= (y - \hat{y}) \cdot -\frac{d}{dw}(w \cdot x) \quad \text{considering } \frac{dy}{dw} = 0 \text{ since } y \text{ is constant} \\ \nabla w &= -(y - \hat{y}) \cdot x \quad \text{considering that } \frac{d(w \cdot x)}{dw} = x, \text{ since } x \text{ is treated as a constant with respect to } w \end{aligned}$$

⇒ the weight update rule for gradient descent becomes:  $w = w - \eta \nabla w = w + \eta \cdot (y - \hat{y}) \cdot x$



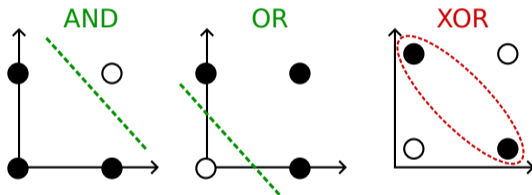
## Reminder: common derivation rules

Rule	Function	Derivative
Constant Rule	$f(x) = c$	$f'(x) = 0$
Power Rule	$f(x) = x^n$	$f'(x) = n \cdot x^{n-1}$
Generalized Power Rule	$f(x) = u(x)^n$	$f'(x) = n \cdot u(x)^{n-1} \cdot u'(x)$
Sum Rule	$f(x) = u(x) + v(x)$	$f'(x) = u'(x) + v'(x)$
Difference Rule	$f(x) = u(x) - v(x)$	$f'(x) = u'(x) - v'(x)$
Product Rule	$f(x) = u(x) \cdot v(x)$	$f'(x) = u'(x) \cdot v(x) + u(x) \cdot v'(x)$
Quotient Rule	$f(x) = \frac{u(x)}{v(x)}$	$f'(x) = \frac{u'(x) \cdot v(x) - u(x) \cdot v'(x)}{[v(x)]^2}$
Chain Rule	$f(x) = u(v(x))$	$f'(x) = u'(v(x)) \cdot v'(x)$
Exponential Rule	$f(x) = e^x$	$f'(x) = e^x$
Exponential with Constant Base	$f(x) = a^x$	$f'(x) = a^x \cdot \ln(a)$
Logarithmic Rule	$f(x) = \ln(x)$	$f'(x) = \frac{1}{x}$
Sine Function	$f(x) = \sin(x)$	$f'(x) = \cos(x)$
Cosine Function	$f(x) = \cos(x)$	$f'(x) = -\sin(x)$
Tangent Function	$f(x) = \tan(x)$	$f'(x) = \sec^2(x) = \frac{1}{\cos^2(x)}$

## Limitation of the perceptron classifier

⇒ the perceptron is a *linear classifier*, so it cannot deal with even trivial *non-linear classifications*

EX: the XOR problem (*exclusive OR*)



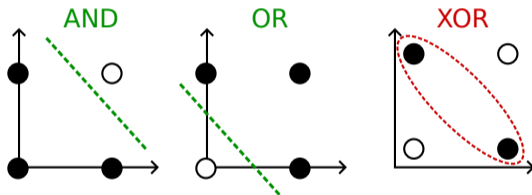
⇒ the limitations of perceptrons can be eliminated by stacking multiple perceptrons into several layers:

⇒ the resulting ANN is called a multilayer perceptron (MLP)

## Limitation of the perceptron classifier

⇒ the perceptron is a *linear classifier*, so it cannot deal with even trivial *non-linear classifications*

EX: the XOR problem (*exclusive OR*)



⇒ the limitations of perceptrons can be eliminated by stacking multiple perceptrons into several layers:

⇒ the resulting ANN is called a multilayer perceptron (MLP)

1. Introduction

2. Perceptron

**3. Multilayer perceptron (MLP)**

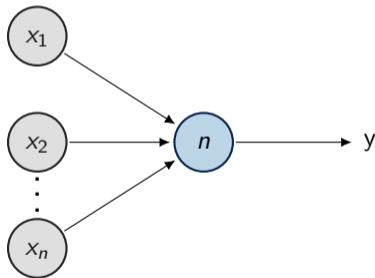
1. from single to multilayer perceptron
2. activation functions
3. backpropagation

4. Application

5. Glossary

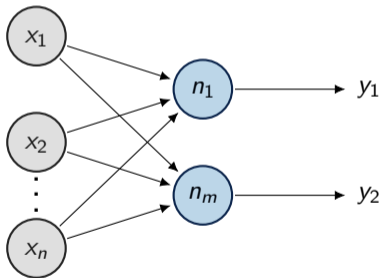
## From single to multilayer perceptron

output layer with 1 neuron



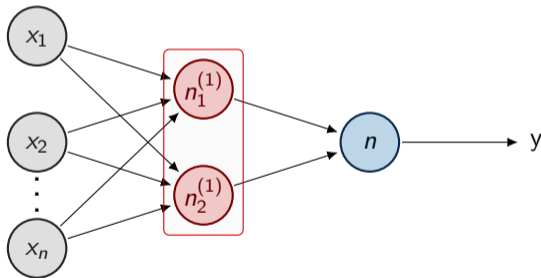
## From single to multilayer perceptron

output layer with 2 neurons



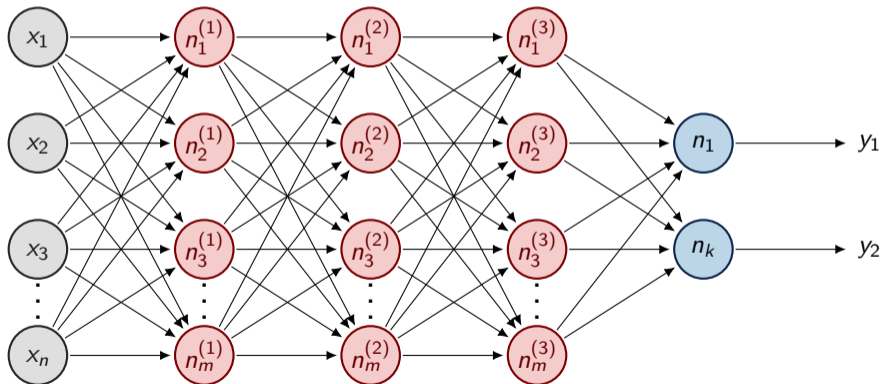
## From single to multilayer perceptron

1 hidden layer with 2 neurons + output layer with 1 neuron (= 2 fully connected layers)



## From single to multilayer perceptron

3 hidden layers (m neurons) + output layer (2 neurons) (= 4 fully connected layers)

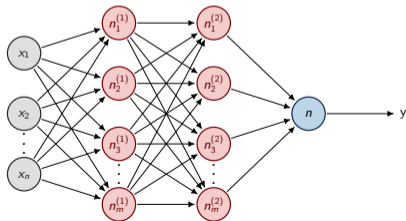




## From single to multilayer perceptron

- ⇒ a **multilayer perceptron** (MLP) is also called a **feedforward neural network**, because information flows from input to output
- ⇒ a MLP consists of a concatenation of multiple **fully connected (dense) layers**, i.e. each neuron in one layer is connected to every neuron in the next layer
- ⇒ a MLP is a concatenation of multiple functions, the entire network can be written out as a long equation:

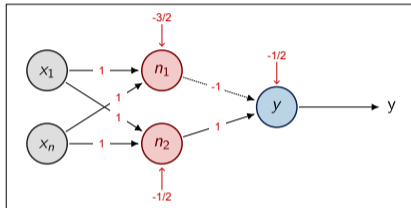
$$\hat{y} = f^{(out)}(w^{out} \cdot f^2(w^2 \cdot f^1(w^1 \cdot x)))$$



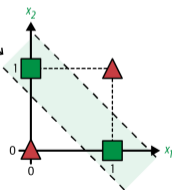
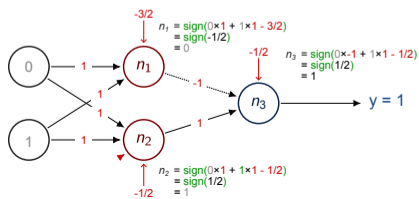
```
# Forward-pass calculations of a 3-layer neural network
f = lambda x: 1.0 / (1.0 + np.exp(-x)) # activation function (sigmoid)
x = np.random.randn(3, 1) # input vector (3x1)
h1 = f(np.dot(w1, x) + b1) # first hidden layer activations
h2 = f(np.dot(w2, h1) + b2) # second hidden layer activations
out = np.dot(w3, h2) + b3 # output neuron
```

## From single to multilayer perceptron

⇒ 1 hidden layer to solve the XOR problem (example taken from book "Geron 2022")

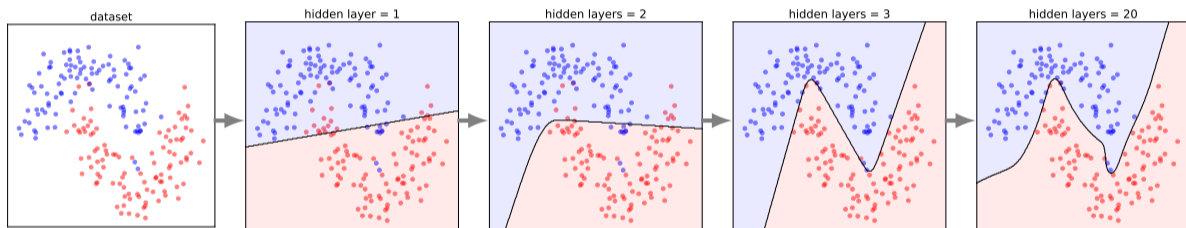


EX: classification of datapoint ( $x_1=0$ ,  $x_2=1$ )



## Effects of the number of hidden layers?

- ⇒ increasing the number of hidden layers allows the network to *learn more complex functions*
- ⇒ no need to worry about feature engineering, the *hidden layer of the network learn the features!*



## Effects of the number of nodes in layers?

⇒ number of nodes in each layer:

- **input layer:** number determined by the *data dimensionality*
  - EX: *previous examples  $x_1, x_2 \Rightarrow 2$  input nodes*
  - EX: *MNIST handwritten digits dataset = 28x28 pixel images  $\Rightarrow 784$  input nodes*
- **output layer:** number determined by the *number of classes* to classify
  - EX: *binary classification  $\Rightarrow 1$  output node*
  - EX: *MNIST handwritten digits dataset  $\Rightarrow 10$  output nodes*
- **hidden layer:** number determined by the *complexity of the function to learn*
  - EX: *XOR problem  $\Rightarrow 2$  hidden nodes is enough*
  - EX: *MNIST handwritten digits dataset  $\Rightarrow$  more nodes & hidden layers can capture more subtleties and improve performance*

Nota Bene: *higher network dimensionality (more nodes & layers)  $\Leftrightarrow$  more computation and prone to overfitting!*

## Effects of the number of nodes in layers?

⇒ number of nodes in each layer:

- **input layer**: number determined by the *data dimensionality*
  - EX: previous examples  $x_1, x_2 \Rightarrow 2$  input nodes
  - EX: MNIST handwritten digits dataset = 28x28 pixel images  $\Rightarrow 784$  input nodes
- **output layer**: number determined by the *number of classes* to classify
  - EX: binary classification  $\Rightarrow 1$  output node
  - EX: MNIST handwritten digits dataset  $\Rightarrow 10$  output nodes
- **hidden layer**: number determined by the *complexity of the function to learn*
  - EX: XOR problem  $\Rightarrow 2$  hidden nodes is enough
  - EX: MNIST handwritten digits dataset  $\Rightarrow$  more nodes & hidden layers can capture more subtleties and improve performance

*Nota Bene*: higher network dimensionality (more nodes & layers)  $\Leftrightarrow$  more computation and prone to overfitting!

## Effects of the number of nodes in layers?

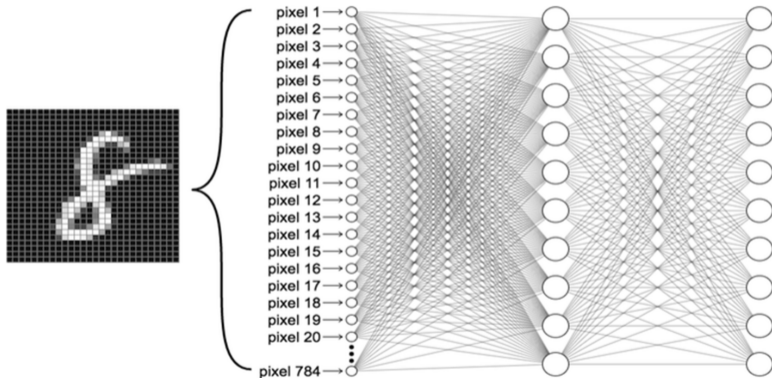
⇒ number of nodes in each layer:

- **input layer**: number determined by the *data dimensionality*
  - EX: previous examples  $x_1, x_2 \Rightarrow 2$  input nodes
  - EX: MNIST handwritten digits dataset = 28x28 pixel images  $\Rightarrow 784$  input nodes
- **output layer**: number determined by the *number of classes* to classify
  - EX: binary classification  $\Rightarrow 1$  output node
  - EX: MNIST handwritten digits dataset  $\Rightarrow 10$  output nodes
- **hidden layer**: number determined by the *complexity of the function to learn*
  - EX: XOR problem  $\Rightarrow 2$  hidden nodes is enough
  - EX: MNIST handwritten digits dataset  $\Rightarrow$  more nodes & hidden layers can capture more subtleties and improve performance

Nota Bene: higher network dimensionality (more nodes & layers)  $\Leftrightarrow$  more computation and prone to overfitting!

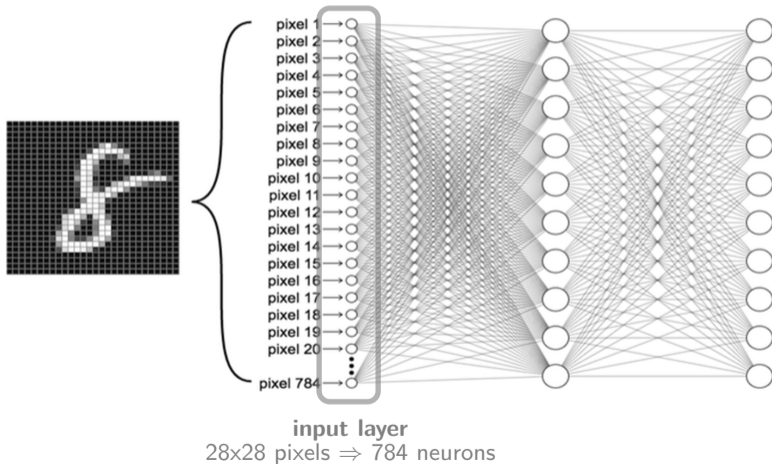
## Effects of the number of nodes in layers?

⇒ classification of the [MNIST dataset](#) with a MLP



## Effects of the number of nodes in layers?

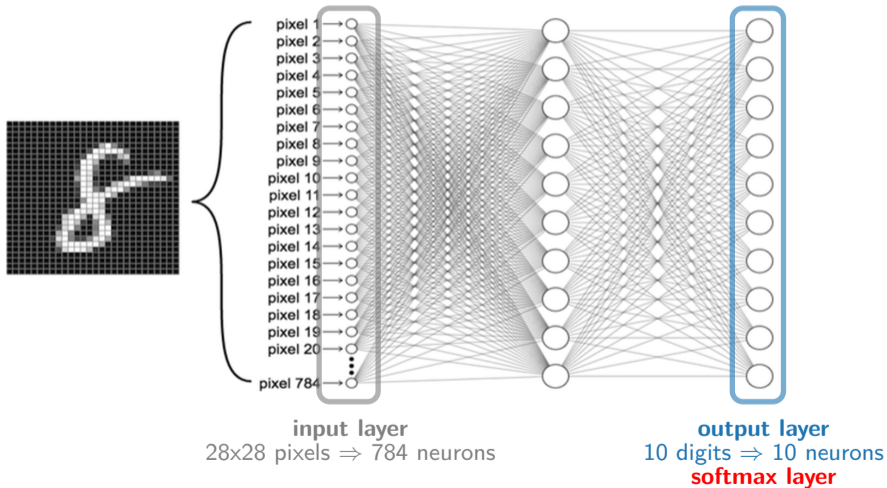
⇒ classification of the **MNIST dataset** with a MLP





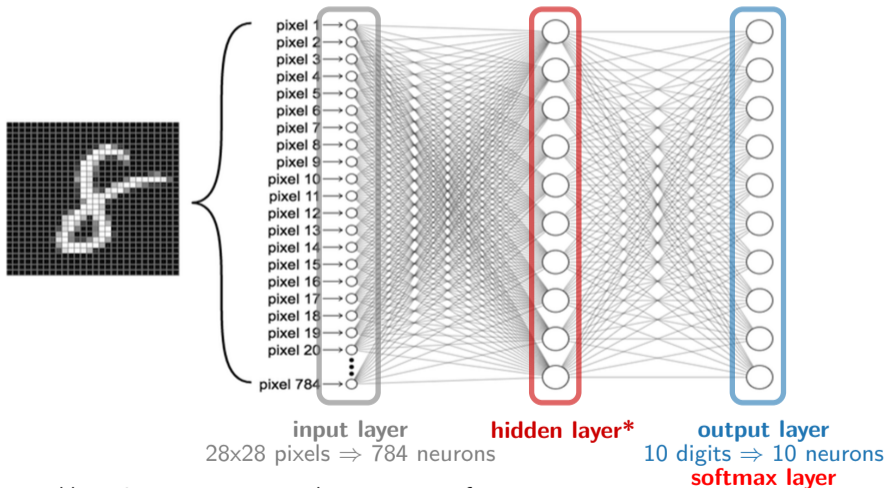
## Effects of the number of nodes in layers?

⇒ classification of the **MNIST dataset** with a MLP



## Effects of the number of nodes in layers?

⇒ classification of the **MNIST dataset** with a MLP



\* *hidden layer*: would require more neurons to have proper performance

## Types of activation functions

- **Sigmoid / Logistic**

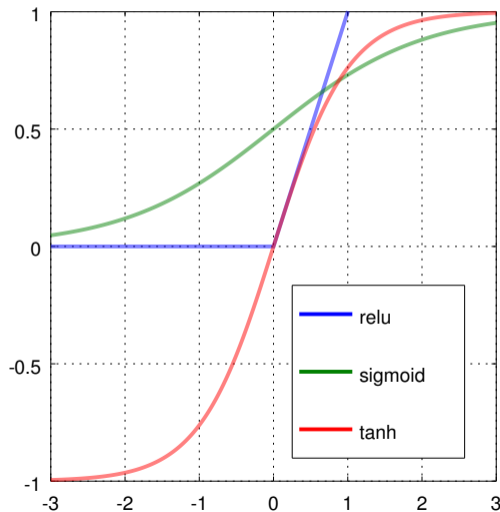
$$f(\mathbf{x})_i = \frac{1}{1+e^{-x_i}}$$

- **ReLU** (Rectified Linear Unit)

$$f(\mathbf{x})_i = \max(x_i, 0)$$

- **TanH**

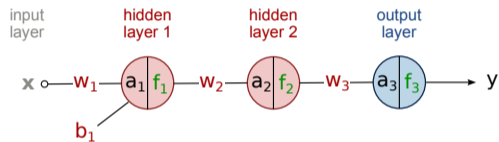
$$f(\mathbf{x})_i = \tanh(x_i) = \frac{e^{x_i} - e^{-x_i}}{e^{x_i} + e^{-x_i}}$$



## Using **backpropagation** to learn MLP

⇒ recall that the *feedforward* operations of a MLP is a concatenation of multiple functions:

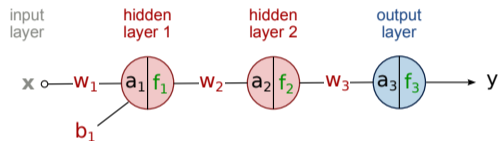
$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1) ) )$$



## Using **backpropagation** to learn MLP

⇒ recall that the *feedforward* operations of a MLP is a concatenation of multiple functions:

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1) ) )$$



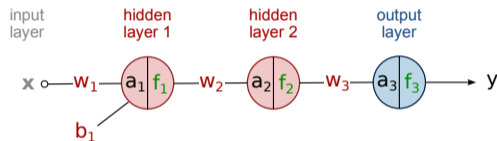
⇒ to learn the MLP  $y = f_{MLP}(x; \theta)$  means that the MLP parameters  $\theta = \{f, w, b\}$  need to be optimized to best fit the training dataset  $\{x, y\}$ , where:

- $f$  = node activation functions (*some activation functions have parameters, e.g. the Leaky ReLU  $f(x) = \max(\alpha x, x)$* )
- $w$  = node weights
- $b$  = node biases

## Using **backpropagation** to learn MLP

⇒ recall that the *feedforward* operations of a MLP is a concatenation of multiple functions:

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1) ) )$$



⇒ to learn the MLP  $y = f_{MLP}(x; \theta)$  means that the MLP parameters  $\theta = \{f, w, b\}$  need to be optimized to best fit the training dataset  $\{x, y\}$ , where:

- $f$  = node activation functions (*some activation functions have parameters, e.g. the Leaky ReLU  $f(x) = \max(\alpha x, x)$* )
- $w$  = node weights
- $b$  = node biases

⇒ to do so we need to compute the partial derivatives of the loss function  $\mathcal{L}$  with respect to  $\theta$ :

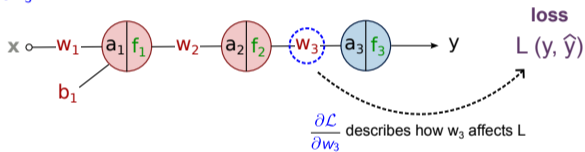
$$\frac{\partial \mathcal{L}}{\partial \theta} = \left[ \frac{\partial \mathcal{L}}{\partial w_3}, \frac{\partial \mathcal{L}}{\partial w_2}, \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial b} \right]$$

## Using **backpropagation** to learn MLP

Let's compute the partial derivatives  $\left[ \frac{\partial \mathcal{L}}{\partial w_3}, \frac{\partial \mathcal{L}}{\partial w_2}, \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial b} \right]$ :

1. compute  $\frac{\partial \mathcal{L}}{\partial w_3}$ :

$\Rightarrow$  the partial derivative  $\frac{\partial \mathcal{L}}{\partial w_3}$  describes how  $w_3$  will affect the Loss function  $L$

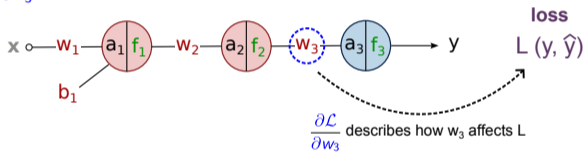


## Using **backpropagation** to learn MLP

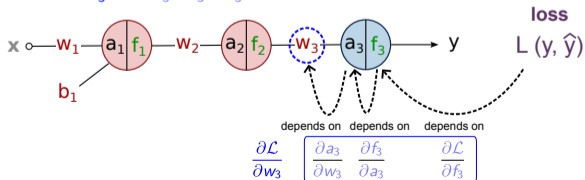
Let's compute the partial derivatives  $\left[ \frac{\partial \mathcal{L}}{\partial w_3}, \frac{\partial \mathcal{L}}{\partial w_2}, \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial b} \right]$ :

1. compute  $\frac{\partial \mathcal{L}}{\partial w_3}$ :

$\Rightarrow$  the partial derivative  $\frac{\partial \mathcal{L}}{\partial w_3}$  describes how  $w_3$  will affect the Loss function  $L$



$\Rightarrow$  according to the chain rule:  $\frac{\partial \mathcal{L}}{\partial w_3} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$





## Using **backpropagation** to learn MLP

Let's compute the partial derivatives  $\left[ \frac{\partial \mathcal{L}}{\partial w_3}, \frac{\partial \mathcal{L}}{\partial w_2}, \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial b} \right]$ :

1. compute  $\frac{\partial \mathcal{L}}{\partial w_3}$  (*continued*):

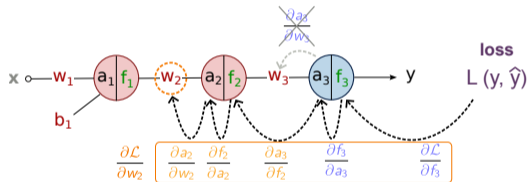
⇒ we can compute  $\frac{\partial \mathcal{L}}{\partial w_3} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$  as follows:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_3} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\ &= (y - \hat{y}) \frac{\partial f}{\partial a_3} \frac{\partial a_3}{\partial w_3} \quad \text{considering } \mathcal{L} = \text{L2 loss} = \frac{1}{2}(y - \hat{y})^2 \Rightarrow \frac{\partial \mathcal{L}}{\partial y} = (y - \hat{y}) \\ &= (y - \hat{y}) f_3 (1 - f_3) \frac{\partial a_3}{\partial w_3} \quad \text{considering } f_3 \text{ is a sigmoid function whose derivative is } f' = f(x)(1 - f(x)) \\ &= (y - \hat{y}) f_3 (1 - f_3) f_2 \quad \text{considering } a_3 = w_3 \cdot f_2 \end{aligned}$$

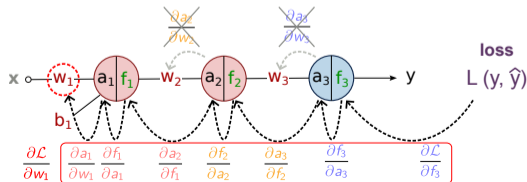
## Using **backpropagation** to learn MLP

Let's compute the partial derivatives  $\left[ \frac{\partial \mathcal{L}}{\partial w_3}, \frac{\partial \mathcal{L}}{\partial w_2}, \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial b} \right]$ :

1. already computed:  $\frac{\partial \mathcal{L}}{\partial w_3} = (y - \hat{y})f_3(1 - f_3)f_2$
2. compute:  $\frac{\partial \mathcal{L}}{\partial w_2}$ ,  $\frac{\partial \mathcal{L}}{\partial w_1}$ ,  $\frac{\partial \mathcal{L}}{\partial b}$  by re-using already computed derivatives (backpropagate!)



$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_3} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\ \frac{\partial \mathcal{L}}{\partial w_2} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} \\ \frac{\partial \mathcal{L}}{\partial w_1} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} \\ \frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b} \end{aligned}$$



## Summary of the procedure to train the MLP

⇒ for each training sample  $\{x_i, y_i\}$ :

### 1. Predict

- forward pass: compute the output of the network

$$\hat{y}_i = f_{MLP}(x_i; \theta)$$

- compute Loss  $\mathcal{L}$ : compare the predicted output with the true output

$$\mathcal{L} = \text{loss}(\hat{y}_i, y_i)$$

### 2. Update weights

- backpropagation: compute the gradients of the loss function with respect to the network parameters  $\theta$

$$\frac{\partial \mathcal{L}}{\partial \theta} = \left[ \frac{\partial \mathcal{L}}{\partial w_3}, \frac{\partial \mathcal{L}}{\partial w_2}, \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial b} \right]$$

- update weights: use the gradients to update the weights of the network

$$w_3 = w_3 - \eta \nabla w_3 \quad \text{where } \nabla w_3 = \frac{\partial \mathcal{L}}{\partial w_3}$$

$$w_2 = w_2 - \eta \nabla w_2 \quad \text{where } \nabla w_2 = \frac{\partial \mathcal{L}}{\partial w_2}$$

$$w_1 = w_1 - \eta \nabla w_1 \quad \text{where } \nabla w_1 = \frac{\partial \mathcal{L}}{\partial w_1}$$

$$b = b - \eta \nabla b \quad \text{where } \nabla b = \frac{\partial \mathcal{L}}{\partial b}$$

1. Introduction

2. Perceptron

3. Multilayer perceptron (MLP)

**4. Application**

1. MLP playground

2. Frameworks for Deep Learning

3. Installing Tensor Flow

4. From ML (sklearn) to DL (tensorflow)

5. Glossary

## MLP playground

- ⇒ build your own MLP using the interactive web platform developed by Google:  
[playground.tensorflow.org](https://playground.tensorflow.org)
- ⇒ see the effects of training in real time!

## Several frameworks exist:

- [Tensor Flow](#)
  - developed by Google



## Several frameworks exist:

- [Tensor Flow](#)
  - developed by Google
  - includes the high-level API [Keras](#) library (TF version  $\geq 2$ )



## Several frameworks exist:

- [Tensor Flow](#)
  - developed by Google
  - includes the high-level API [Keras](#) library (TF version  $\geq 2$ )
- [PyTorch](#)
  - developed by Facebook





## Several frameworks exist:

- [Tensor Flow](#)
  - developed by Google
  - includes the high-level API [Keras](#) library (TF version  $\geq 2$ )
- [PyTorch](#)
  - developed by Facebook
  - based on the [Torch](#) framework (Lua)

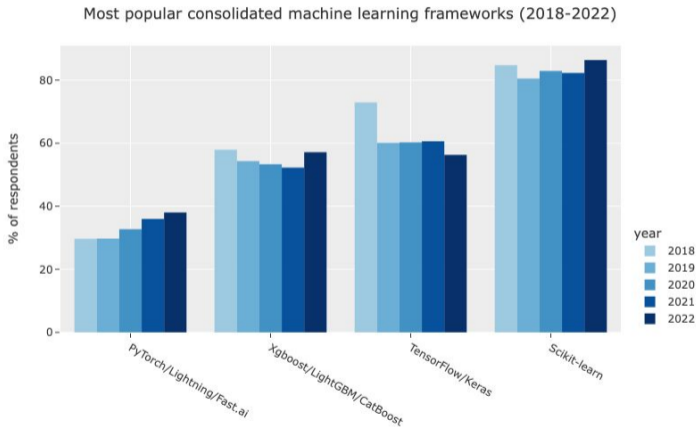


## Several frameworks exist:

- [Tensor Flow](#)
  - developed by Google
  - includes the high-level API [Keras](#) library (TF version  $\geq 2$ )
- [PyTorch](#)
  - developed by Facebook
  - based on the [Torch](#) framework (Lua)
- [MindSpore](#), [Caffe](#), [MXNet](#), [Theano](#), etc.



## Popularity of the main frameworks (from Kaggle 2022<sup>1</sup>)



<sup>1</sup>2022 - Kaggle Data Science & ML Survey ([link](#))

## Installing [Tensor Flow](#) with Anaconda ([instructions](#)):

⇒ we will install Tensor Flow in a "conda environment"

1. Create environment and install Tensor Flow package & dependencies inside

```
$ conda env list           # optional: list existing environments
$ conda create -n tf tensorflow # create environment called "tf" & install CPU-only TensorFlow
```

2. Activate the created environment

```
$ conda activate tf
```

3. Install additional packages in the active environment

```
$ conda install jupyter matplotlib pandas scikit-learn
```

4. Launch Jupyter from the active environment, import Tensor Flow, and you're good to go!

```
$ jupyter notebook
```

```
# Create a new notebook with Python 3 kernel
import tensorflow as tf
```

## Installing [Tensor Flow](#) with Anaconda ([instructions](#)):

⇒ we will install Tensor Flow in a "conda environment"

1. Create environment and install Tensor Flow package & dependencies inside

```
$ conda env list           # optional: list existing environments
$ conda create -n tf tensorflow # create environment called "tf" & install CPU-only TensorFlow
```

2. Activate the created environment

```
$ conda activate tf
```

3. Install additional packages in the active environment

```
$ conda install jupyter matplotlib pandas scikit-learn
```

4. Launch Jupyter from the active environment, import Tensor Flow, and you're good to go!

```
$ jupyter notebook
```

```
# Create a new notebook with Python 3 kernel
import tensorflow as tf
```

## Installing [Tensor Flow](#) with Anaconda ([instructions](#)):

⇒ we will install Tensor Flow in a "conda environment"

1. Create environment and install Tensor Flow package & dependencies inside

```
$ conda env list           # optional: list existing environments
$ conda create -n tf tensorflow # create environment called "tf" & install CPU-only TensorFlow
```

2. Activate the created environment

```
$ conda activate tf
```

3. Install additional packages in the active environment

```
$ conda install jupyter matplotlib pandas scikit-learn
```

4. Launch Jupyter from the active environment, import Tensor Flow, and you're good to go!

```
$ jupyter notebook
# Create a new notebook with Python 3 kernel
import tensorflow as tf
```

## Installing [Tensor Flow](#) with Anaconda ([instructions](#)):

⇒ we will install Tensor Flow in a "conda environment"

1. Create environment and install Tensor Flow package & dependencies inside

```
$ conda env list                # optional: list existing environments
$ conda create -n tf tensorflow # create environment called "tf" & install CPU-only TensorFlow
```

2. Activate the created environment

```
$ conda activate tf
```

3. Install additional packages in the active environment

```
$ conda install jupyter matplotlib pandas scikit-learn
```

4. Launch Jupyter from the active environment, import Tensor Flow, and you're good to go!

```
$ jupyter notebook
```

```
# Create a new notebook with Python 3 kernel
import tensorflow as tf
```

## Nota Bene

Two distinct versions of TF exist, depending on whether it should run on CPU (Central Processing Unit), or GPU (Graphics Processing Unit)

⇒ CPU-only TensorFlow (recommended for beginners)

```
$ conda create -n tf tensorflow
```

⇒ GPU TensorFlow

```
$ conda create -n tf-gpu tensorflow-gpu
```

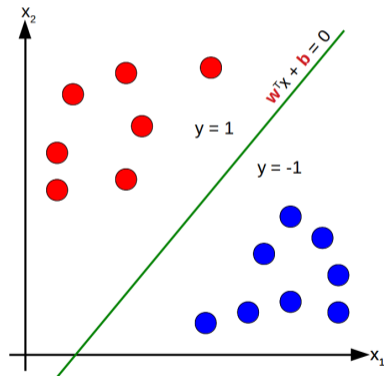
⇒ GPU will be much faster, but more expensive, and trickier to setup (requires CUDA)



## Toy example: linear classification task using scikit-learn and tensor flow

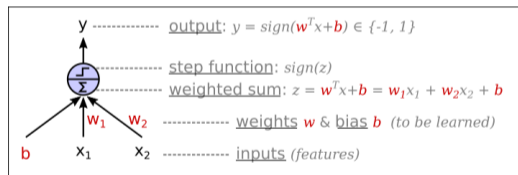
perceptron:  $y = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$

- $y \in \{-1, 1\}$ : predicted class  $\rightarrow$  *banana or apple*
- $\mathbf{x} \in \mathbb{R}^2$ : feature matrix (n, 2)  $\rightarrow$  [*hue, elongation*]
- $\mathbf{w} \in \mathbb{R}^2$ : "weight vector" (2,)  $\rightarrow$  *needs to be learned*
- $b \in \mathbb{R}$ : "bias"  $\rightarrow$  *needs to be learned*
- *sign*: [sign function](#) returning the sign of a real number



**Toy example:** linear classification task using scikit-learn and tensor flow**perceptron:**  $y = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$ 

- $y \in \{-1, 1\}$ : predicted class  $\rightarrow$  *banana or apple*
- $\mathbf{x} \in \mathbb{R}^2$ : feature matrix (n, 2)  $\rightarrow$  *[hue, elongation]*
- $\mathbf{w} \in \mathbb{R}^2$ : "weight vector" (2,)  $\rightarrow$  *needs to be learned*
- $b \in \mathbb{R}$ : "bias"  $\rightarrow$  *needs to be learned*
- *sign*: **sign function** returning the sign of a real number



**Solution with Scikit-Learn:** Perceptron classifier

```

from sklearn import datasets
from sklearn import linear_model
from sklearn.utils import shuffle
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score

# Load data
iris = datasets.load_iris()
X = iris.data[:, (2, 3)] # petal length, petal width
y = (iris.target == 0).astype('int') # Iris setosa?

# Preprocess data
X, y = shuffle(X, y, random_state=0)
scaler = StandardScaler()
X = scaler.fit_transform(X)

X_train = X[:75]
y_train = y[:75]
X_test = X[75:]
y_test = y[75:]

# Select model
clf = linear_model.Perceptron()

# Train model
clf.fit(X_train, y_train)

print('weights:', clf.coef_)
print('bias:', clf.intercept_)

# Evaluate
y_pred = clf.predict(X_train)
accuracy_score(y_train, y_pred)

# Predict from model
y_pred = clf.predict([[2, 0.5]])

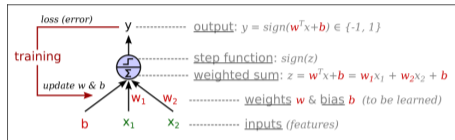
# Plot data + linear classifier
plt.scatter(X[:,0], X[:,1], c=y)
plt.scatter(X_train[:,0], X_train[:,1], c=y_train)
plt.scatter(X_test[:,0], X_test[:,1], c=y_test, alpha=.25)

weights = clf.coef_[0]
bias = clf.intercept_
slope = -weights[0] / weights[1]
yintercept = -bias / weights[1]
_x = np.linspace(-2,2)
_y = slope*_x + yintercept
plt.plot(_x, _y, '-r')

```

**1.1 Load data****1.2 Preprocess data**

- shuffle
- scale
- split into train/test

**2. Select model****3. Train model****4. Evaluate model****5. Predict from model**

Solution with **Scikit-Learn**: Perceptron classifier

```

from sklearn import datasets
from sklearn import linear_model
from sklearn.utils import shuffle
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score

# Load data
iris = datasets.load_iris()
X = iris.data[:, (2, 3)] # petal length, petal width
y = (iris.target == 0).astype('int') # Iris setosa?

# Preprocess data
X, y = shuffle(X, y, random_state=0)
scaler = StandardScaler()
X = scaler.fit_transform(X)

X_train = X[:75]
y_train = y[:75]
X_test = X[75:]
y_test = y[75:]

# Select model
clf = linear_model.Perceptron()

# Train model
clf.fit(X_train, y_train)

print('weights:', clf.coef_)
print('bias:', clf.intercept_)

# Evaluate
y_pred = clf.predict(X_train)
accuracy_score(y_train, y_pred)

# Predict from model
y_pred = clf.predict([[2, 0.5]])

# Plot data + linear classifier
plt.scatter(X[:,0], X[:,1], c=y)
plt.scatter(X_train[:,0], X_train[:,1], c=y_train)
plt.scatter(X_test[:,0], X_test[:,1], c=y_test, alpha=.25)

weights = clf.coef_[0]
bias = clf.intercept_
slope = -weights[0] / weights[1]
yintercept = -bias / weights[1]
_x = np.linspace(-2,2)
_y = slope*_x + yintercept
plt.plot(_x, _y, '-r')

```

## 1.1 Load data

## 1.2 Preprocess data

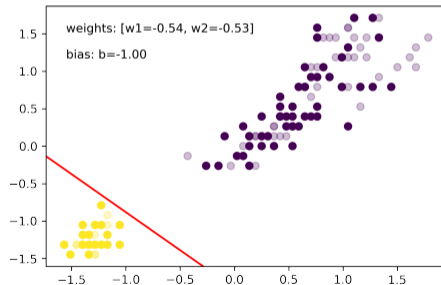
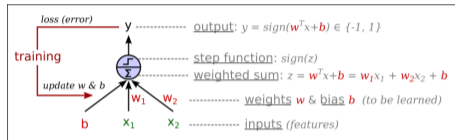
- shuffle
- scale
- split into train/test

## 2. Select model

## 3. Train model

## 4. Evaluate model

## 5. Predict from model



**Solution with Tensor Flow - Keras: 1 neuron network**

```
import tensorflow as tf
from sklearn import datasets
from sklearn import linear_model
from sklearn.utils import shuffle
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score
```

```
# Load data
iris = datasets.load_iris()
X = iris.data[:, (2, 3)] # petal length, petal width
y = (iris.target == 0).astype('int') # Iris setosa?
```

```
# Preprocess data
X, y = shuffle(X, y, random_state=0)
scaler = StandardScaler()
X = scaler.fit_transform(X)
```

```
X_train = X[:75]
y_train = y[:75]
X_test = X[75:]
y_test = y[75:]
```

```
model = tf.keras.Sequential([
    tf.keras.layers.Flatten(input_shape=(2,)),
    tf.keras.layers.Dense(1, activation='sigmoid')
])
model.summary()
```

```
# Compile model
model.compile(optimizer='sgd',
              loss='BinaryCrossentropy',
              metrics=['accuracy'])
```

```
# Train model
history = model.fit(X_train, y_train, epochs=50, batch_size=10)
```

```
# Evaluate model
test_loss, test_acc = model.evaluate(X_test, y_test, verbose=2)
print('Test accuracy:', test_acc)
```

```
# Predict (data should be preprocessed just like training data)
probability_model = tf.keras.Sequential([model, tf.keras.layers.Softmax()])
pred = probability_model.predict([[ -2, -2]])
print(pred)
```

**1.1 Load data****1.2 Preprocess data**

- shuffle
- scale
- split into train/test

**2.1 Build model**

- set layer type/order

**2.2 Compile model**

- set loss function
- set optimizer
- set metrics

**3. Train model**

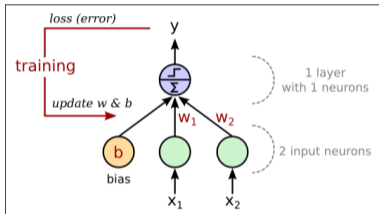
- learn layer parameters (weights/biases)
- plot training history (check for overfitting)

**4. Evaluate model**

- evaluate accuracy on test dataset

**5. Predict from model**

- predict image class using learned model



Model: "sequential"

Layer (type)	Output Shape	Param #
flatten (Flatten)	(None, 2)	0
dense (Dense)	(None, 1)	3
Total params: 3		
Trainable params: 3		
Non-trainable params: 0		

So if we can do the same thing, why switch from sklearn to tensor flow ?

Tensor Flow is a framework for Deep Learning

- ⇒ can design multi-layered networks, and train them in a very flexible/optimized manner
- ⇒ can solve much more complex problems, by optimizing several thousands/millions of weights during training!

So if we can do the same thing, why switch from sklearn to tensor flow ?

Tensor Flow is a framework for Deep Learning

- ⇒ can design multi-layered networks, and train them in a very flexible/optimized manner
- ⇒ can solve much more complex problems, by optimizing several thousands/millions of weights during training!

## “Hello World” example in Keras TensorFlow: MNIST fashion dataset classification task with MLP

```
import tensorflow as tf
```

```
# Load data
fashion_mnist = tf.keras.datasets.fashion_mnist

(X_train_full, y_train_full), (X_test, y_test) = fashion_mnist.load_data()
X_valid, X_train = X_train_full[:5000], X_train_full[5000:]
y_valid, y_train = y_train_full[:5000], y_train_full[5000:]
```

```
# Preprocess data
X_train, X_test, X_valid = X_train/255.0, X_test/255.0, X_valid/255.0
```

```
# Build model (using the Sequential API)
model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28)),
    tf.keras.layers.Dense(300, activation='relu'),
    tf.keras.layers.Dense(100, activation='relu'),
    tf.keras.layers.Dense(10, activation='softmax')
])
model.summary()
```

```
# Compile model
model.compile(loss='sparse_categorical_crossentropy',
              optimizer='sgd',
              metrics=['accuracy'])
```

```
# Train model
history = model.fit(X_train, y_train, validation_data=(X_valid, y_valid),
                   epochs=30, # nb of times X_train is seen
                   batch_size=32) # nb of images per training instance
print('training instances per epoch = {}'.format(X_train.shape[0] / 32))
```

```
# Plot training history
import pandas as pd
pd.DataFrame(history.history).plot()
```

```
# Evaluate model
test_loss, test_acc = model.evaluate(X_test, y_test)
print('Test accuracy:', test_acc)
```

```
# Predict
img = X_test[0,:]
img = (np.expand_dims(img,0)) # add image to a batch
y_proba = model.predict(img).round(2)
y_pred = np.argmax(model.predict(img), axis=-1)
```

```
plt.bar(range(10), y_proba[0])
plt.imshow(img[0,:], cmap='binary')
plt.title('class {} = {}'.format(y_pred, class_names[np.argmax(y_proba)]))
```

### 1.1 Load data

- training dataset
- validation dataset
- test dataset

### 1.2 Preprocess data

- scale pixel intensities to 0-1

### 2.1 Build model

- set layer type/order

### 2.2 Compile model

- set loss function
- set optimizer
- set metrics

### 3. Train model

- learn layer parameters (weights/biases)
- plot training history (check for overfitting)

### 4. Evaluate model

- evaluate accuracy on test dataset

### 5. Predict from model

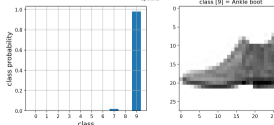
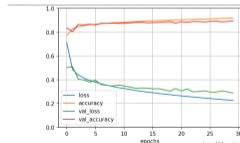
- predict image class using learned model



Model: "sequential"

Layer (type)	Output Shape	Param #
flatten (Flatten)	(None, 784)	0
dense (Dense)	(None, 300)	235500
dense_1 (Dense)	(None, 100)	30100
dense_2 (Dense)	(None, 10)	1010

Total params: 266,610  
Trainable params: 266,610  
Non-trainable params: 0





## “Hello World” example in Keras TensorFlow: MNIST fashion dataset classification task with MLP

```
import tensorflow as tf
```

```
# Load data
```

```
fashion_mnist = tf.keras.datasets.fashion_mnist
```

```
(X_train_full, y_train_full), (X_test, y_test) = fashion_mnist.load_data()
X_valid, X_train = X_train_full[:5000], X_train_full[5000:]
y_valid, y_train = y_train_full[:5000], y_train_full[5000:]
```

```
# Preprocess data
```

```
X_train, X_test, X_valid = X_train/255.0, X_test/255.0, X_valid/255.0
```

```
# Build model (using the Sequential API)
```

```
model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28)),
    tf.keras.layers.Dense(300, activation='relu'),
    tf.keras.layers.Dense(100, activation='relu'),
    tf.keras.layers.Dense(10, activation='softmax')
])
```

```
model.summary()
```

```
# Compile model
```

```
model.compile(loss='sparse_categorical_crossentropy',
              optimizer='sgd',
              metrics=['accuracy'])
```

```
# Train model
```

```
history = model.fit(X_train, y_train, validation_data=(X_valid, y_valid),
                   epochs=30, # nb of times X_train is seen
                   batch_size=32) # nb of images per training instance
print('training instances per epoch = {}'.format(X_train.shape[0] / 32))
```

```
# Plot training history
```

```
import pandas as pd
pd.DataFrame(history.history).plot()
```

```
# Evaluate model
```

```
test_loss, test_acc = model.evaluate(X_test, y_test)
print('Test accuracy:', test_acc)
```

```
# Predict
```

```
img = X_test[0,:]
img = (np.expand_dims(img,0)) # add image to a batch
y_proba = model.predict(img).round(2)
y_pred = np.argmax(model.predict(img), axis=-1)
```

```
plt.bar(range(10), y_proba[0])
plt.imshow(img[0,:], cmap='binary')
plt.title('class {} = {}'.format(y_pred, class_names[np.argmax(y_proba)]))
```

### 1.1 Load data

- training dataset
- validation dataset
- test dataset

### 1.2 Preprocess data

- scale pixel intensities to 0-1

### 2.1 Build model

- set layer type/order

### 2.2 Compile model

- set loss function
- set optimizer
- set metrics

### 3. Train model

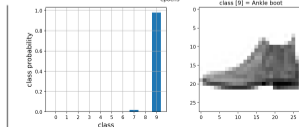
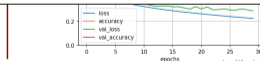
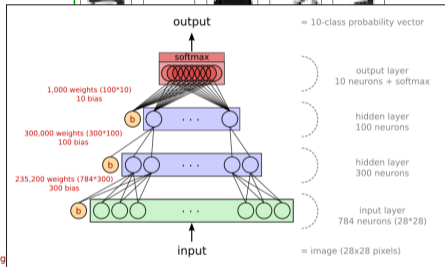
- learn layer parameters (weights/biases)
- plot training history (check for overfitting)

### 4. Evaluate model

- evaluate accuracy on test dataset

### 5. Predict from model

- predict image class using learned model



1. Introduction

2. Perceptron

3. Multilayer perceptron (MLP)

4. Application

5. Glossary

## Key parameters and definitions (from Google's [ML glossary](#), Chollet 2017, etc.)

- **loss function** (objective function)

The quantity that will be minimized during training. It represents a measure of success for the task at hand.

- **optimizer**

Determines how the network will be updated based on the loss function. It implements a specific variant of stochastic gradient descent (SGD).

- **accuracy**

The fraction of predictions that a classification model got right.

- **epoch**

Each iteration over all the training data.

- **batch\_size**

Number of samples per gradient update.

- **activation function**

A function (for example, ReLU or sigmoid) that takes in the weighted sum of all of the inputs from the previous layer and then generates and passes an output value (typically nonlinear) to the next layer.

- **softmax**

A function that provides probabilities for each possible class in a multi-class classification model. The probabilities add up to exactly 1.0. For example, softmax might determine that the probability of a particular image being a dog at 0.9, a cat at 0.08, and a horse at 0.02.