

# Deep Learning from Scratch

## Session #1: Introduction



by: Ali Tourani – Summer 2021

# Agenda

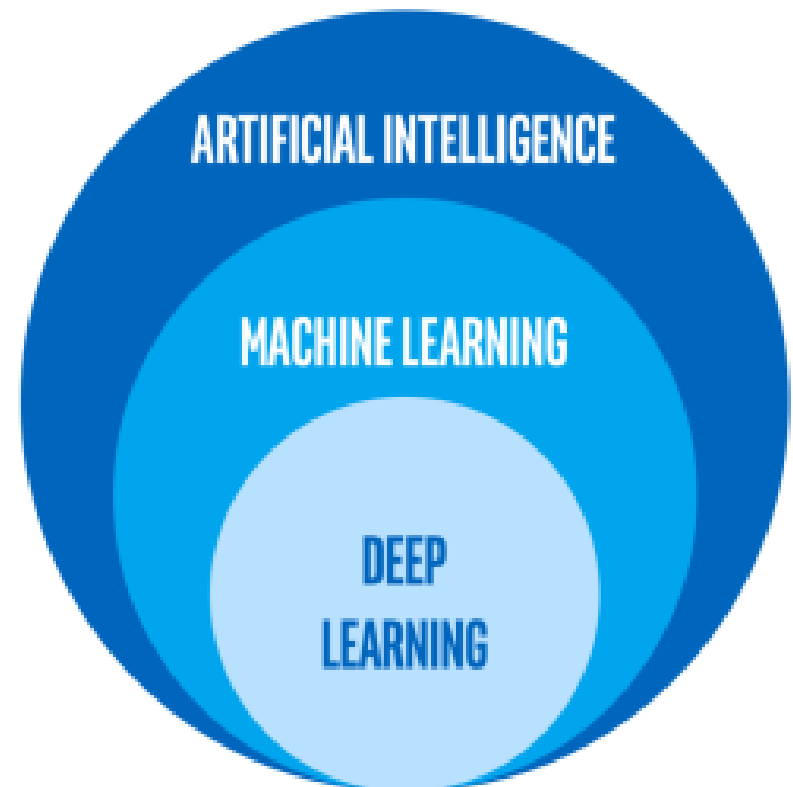
- ▶ What is Deep Learning?
- ▶ Fundamental Concepts
- ▶ Neural Networks
- ▶ Training the Network
- ▶ Use cases of ANNs
- ▶ Assignments

# What is Deep Learning?

How to enable computers to mimic advanced human capabilities?

Learning AI rules by analyzing collections of known examples

A type of machine learning needs to be told about the essential features



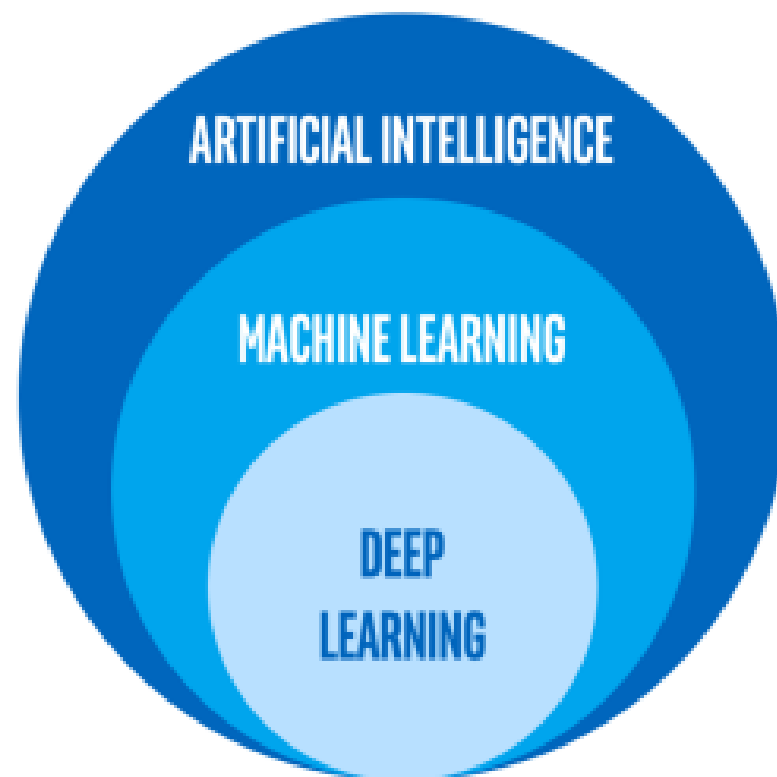
# What is Deep Learning?

## Machine Learning (ML)

- ▶ Using **parameters** from known and important features of data
- ▶ **Predict** the outcomes on similar data
- ▶ Final output: a **Model**

## Deep Learning (DL)

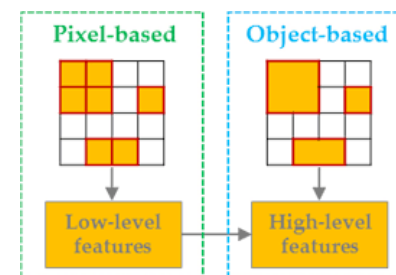
- ▶ Using **Artificial Neural Networks (ANNs)**
- ▶ Learning tasks directly from **raw data**
- ▶ Needs better **hardware** and **data**



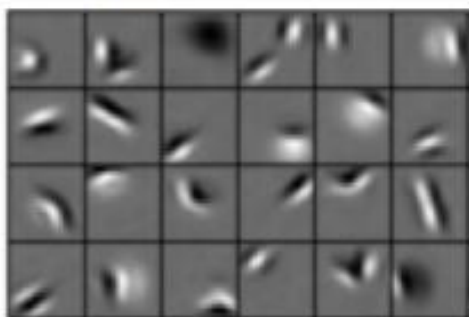
# What is Deep Learning?

## How to learn data?

- ▶ Using the **features** existing in data (e.g., images, sounds, etc.)
- ▶ Trying to find the **rules and patterns**
- ▶ **Decomposing** existing features



Low level features



Edges, dark spots

Mid level features



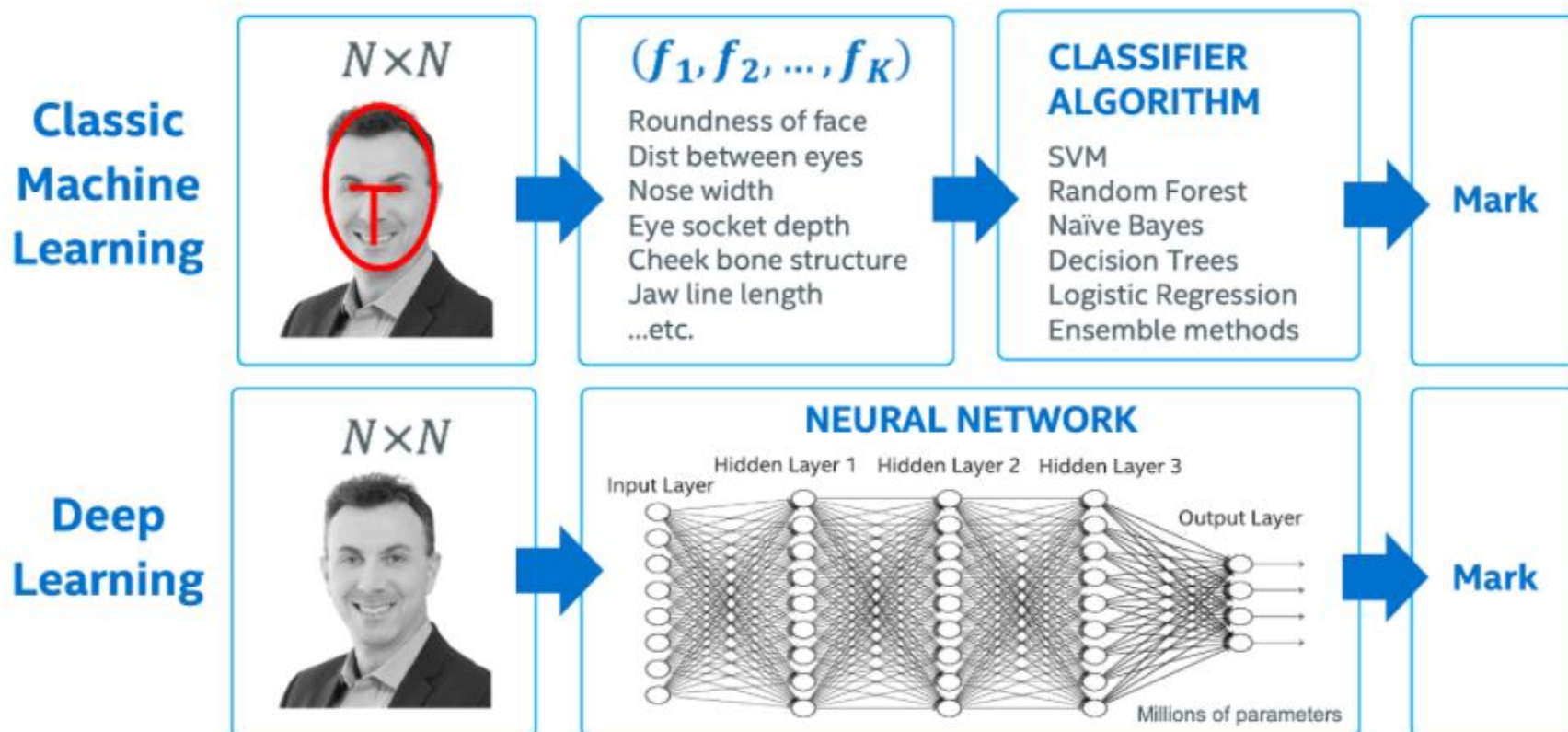
Eyes, ears, nose

High level features



Facial structure

# What is Deep Learning?



# What is Deep Learning?

## What are the requirements of Deep Learning?

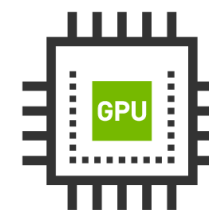
### ▶ Bigger Datasets

- ▶ To extract the **features**, we need to have all variations of data
- ▶ Large **datasets** with thousands and millions of samples



### ▶ Better Hardware

- ▶ Highly **parallelizable hardware** is needed to process big data
- ▶ **Graphics Processing Units (GPUs)** are the best option



### ▶ Efficient Software and Smart Algorithms

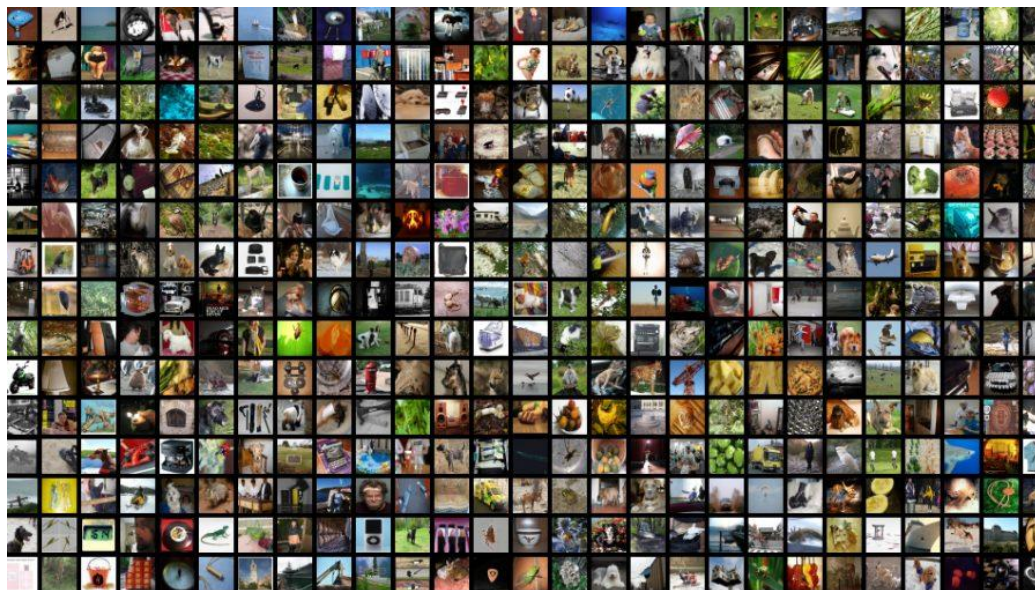
- ▶ Different **toolboxes** and **platforms** that enable us to process big data using GPUs



# What is Deep Learning?

What are the requirements of Deep Learning?

**Big Data** + Powerful Hardware + Efficient Software





# What is Deep Learning?

What are the requirements of Deep Learning?

Big Data + **Powerful Hardware** + Efficient Software

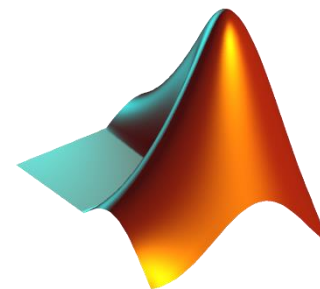


colab

# What is Deep Learning?

What are the requirements of Deep Learning?

Big Data + Powerful Hardware + **Efficient Software**



# What is Deep Learning?

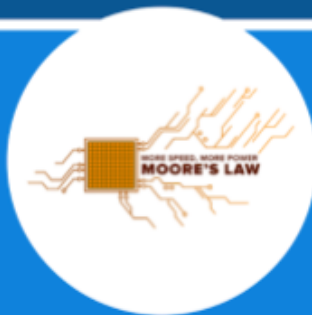
## Bigger Datasets



**IMAGENET** 10M labeled images

**YouTube** 8M categorized videos

## Better Hardware



### Moore's Law

Cost / GB in 1995: \$1000

Cost / GB in 2020: \$0.02

## Smarter Algorithms



**Recurrent Neural Networks**

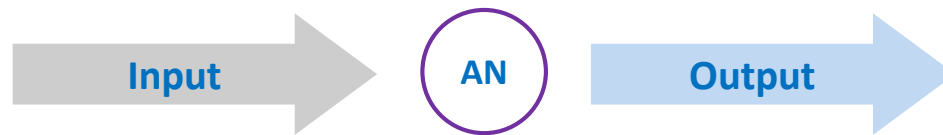
**Convolutional Neural Nets**

**LSTM**

# Fundamental Concepts

## Artificial Neural Networks (ANNs, or simply, NNs)

- ▶ Inspired from [neurons](#) in the human brain
- ▶ Contains at least one [Artificial Neuron \(AN\)](#)
  - ▶ Receives an input signal, processes it, and generates an output signal



- ▶ Connecting several ANs?
  - ▶ The output of neuron **A** will be the input of neuron **B**
  - ▶ They will shape an ANN that is organized into multiple [layers](#)

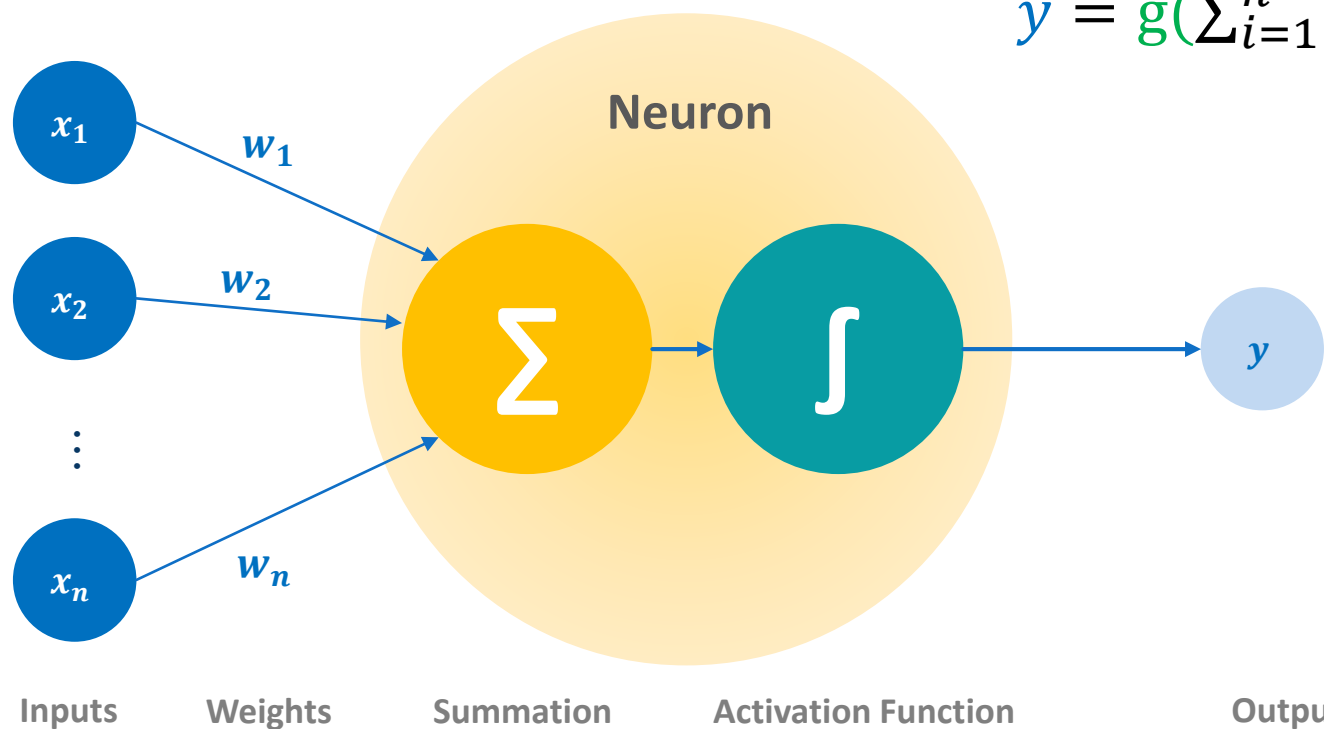
# Fundamental Concepts

## The Perceptron

- ▶ The most straightforward architecture of an ANN
  - ▶ At least **one input signal  $X$**  with a corresponding **weight  $W$** 
    - ▶ The weight shows the importance/priority of the input signal
  - ▶ Calculating the **summation** of each input signal **multiplied** by its weight
    - ▶ The output will be a single number
  - ▶ Passing the sum value to a non-linear function, called **Activation Function (AF)**
  - ▶ Now, the **final output signal  $Y$**  is ready!

# Fundamental Concepts

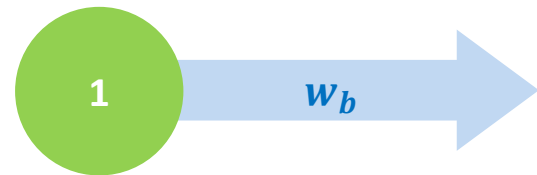
## The Perceptron



# Fundamental Concepts

## The Perceptron

- ▶ But how to make the AF independent from the input signals?
  - ▶ **Solution:** by adding a **Bias** node



## What is a Bias?

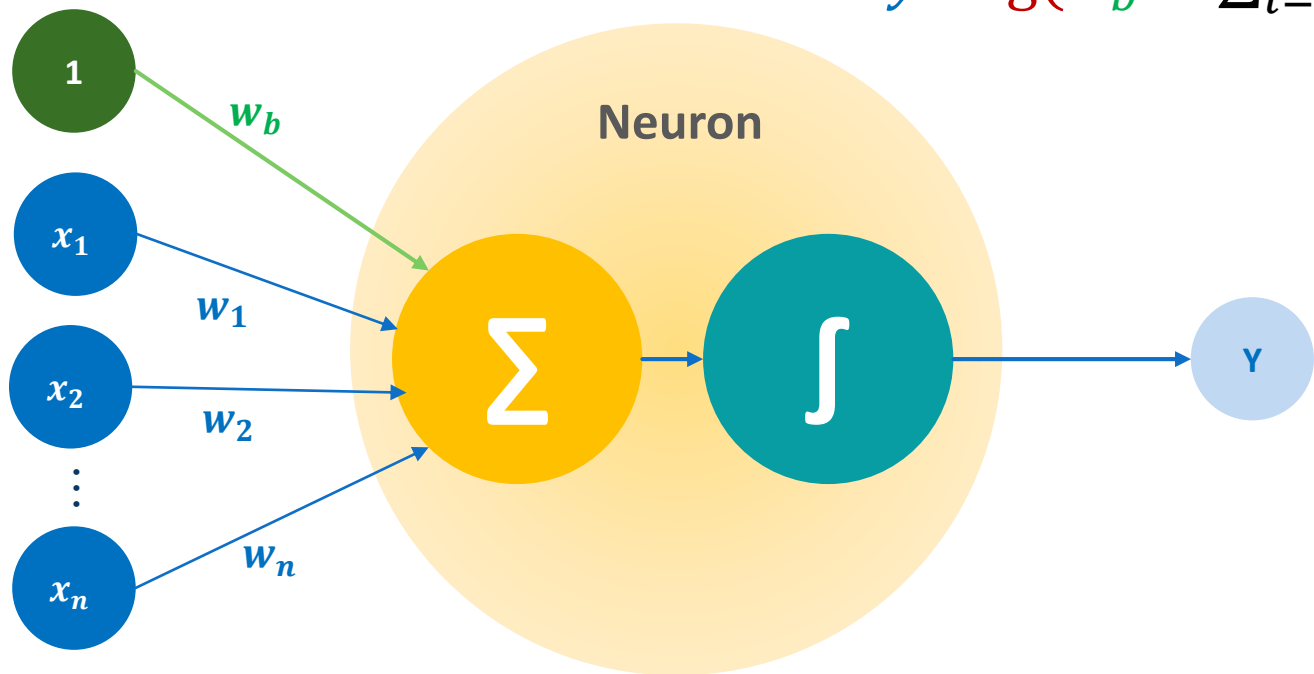
- ▶ A weighted input signal with the **value of 1**
- ▶ Using this, and by providing proper weight, we can **shift** our AFs to the left and right



# Fundamental Concepts

## The Perceptron

$$y = g(w_b + \sum_{i=1}^n x_i w_i)$$



Inputs

Weights

Summation

Activation Function

Output

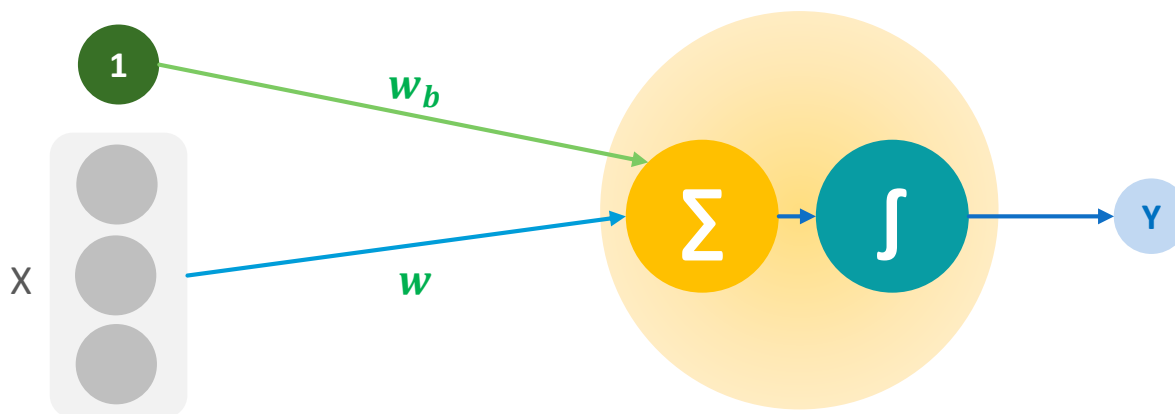
# Fundamental Concepts

## The Perceptron

- ▶ Now, we can mathematically define the output as:

$$y = g(w_b + X^T W)$$

Where  $X$  and  $W$  are the vector of inputs and weights, respectively



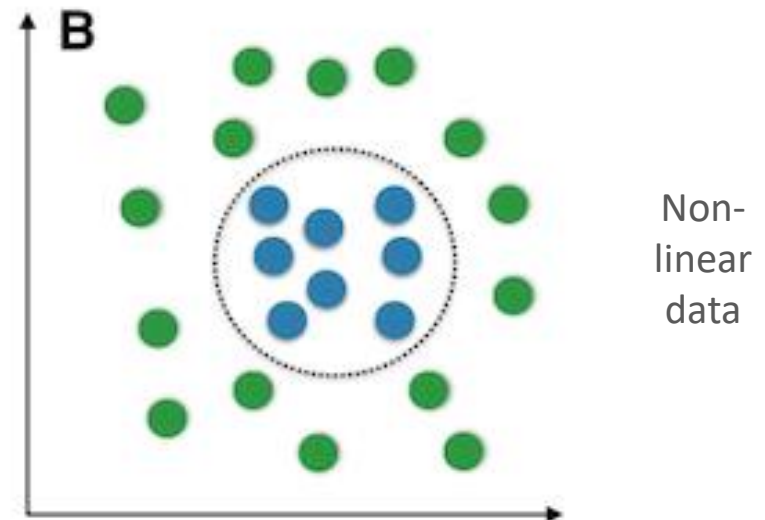
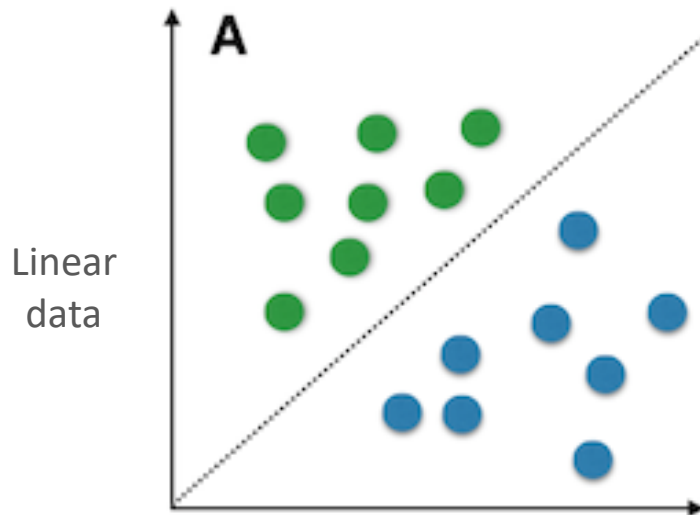
# Fundamental Concepts

## Activation Functions

- ▶ Why do we need AFs?
  - ▶ To provide **non linearity** and ...

But Why?

$$y = g(w_b + X^T W)$$



# Fundamental Concepts

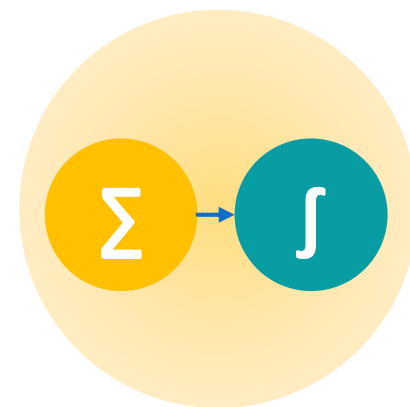
## Activation Functions

$$y = g(w_b + X^T W)$$

- ▶ Why do we need AFs?
  - ▶ To provide **nonlinearity** for using ANNs in real-world scenarios
  - ▶ To define how the **weighted sum** is transformed into an output
  - ▶ To get access to much richer **hypothesis** space (especially deep ones)

### Notes:

- ▶ We need to choose them carefully in our ANNs!
  - ▶ There are many of them!
- ▶ Technically, the same for all nodes in a layer of an ANN



# Fundamental Concepts

## Some of the popular Activation Functions

- ▶ Linear

$$f(x) = x$$

- ▶ Unit step

$$f(x) = \begin{cases} -1, & x < 0 \\ 0, & x = 0 \\ 1, & x > 0 \end{cases}$$

$$\begin{matrix} x < 0 \\ x = 0 \\ x > 0 \end{matrix}$$

- ▶ Sign (signum)

$$f(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

- ▶ ReLU (Rectified Linear Activation)

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

- ▶ Hyperbolic Tangent (Tanh)

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

- ▶ Sigmoid/Logistic (maps real numbers to [0, 1])

- ▶ Gaussian

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

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

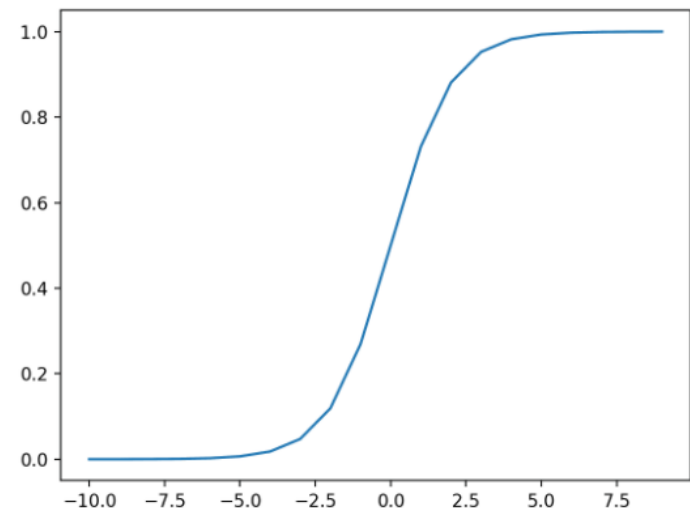
# Fundamental Concepts

## Activation Functions - Sigmoid

- ▶ **Input:** a real number
- ▶ **Output:** a number in the range of 0 to 1
  - ▶ Larger input → closer to 1
- ▶ A common use case: **probabilities**

```
# Applying the sigmoid function and  
# storing the result in 'b'  
b = tf.nn.sigmoid(a, name = 'sigmoid')
```

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



# Fundamental Concepts

## An example of Perceptron with Sigmoid AF

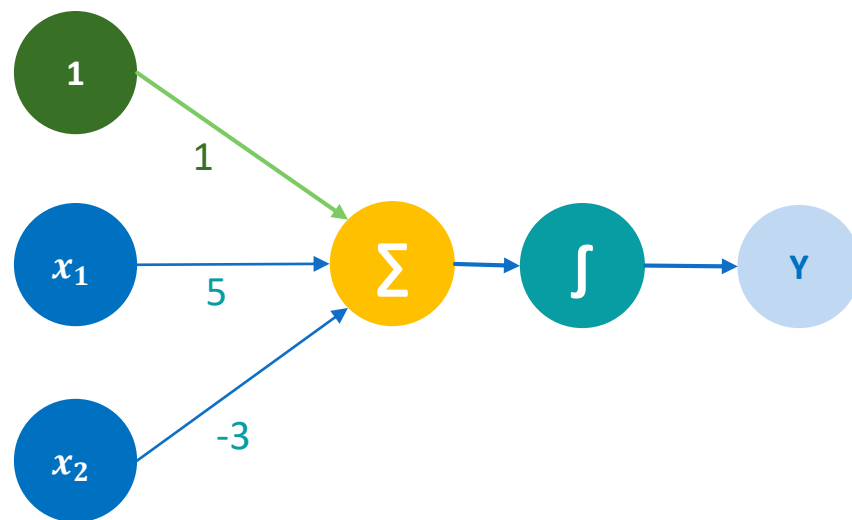
- ▶ Suppose we have this AN:

- ▶ The weights vector:

$$W_0 = 1, \quad W^T = [5 \quad -3]$$

- ▶ The summation function provides:

$$\begin{aligned} y &= g(w_b + X^T W) \\ &= g(1 + 5x_1 - 3x_2) \end{aligned}$$





# Fundamental Concepts

## An example of Perceptron with Sigmoid AF

- ▶ Thus, we will have:

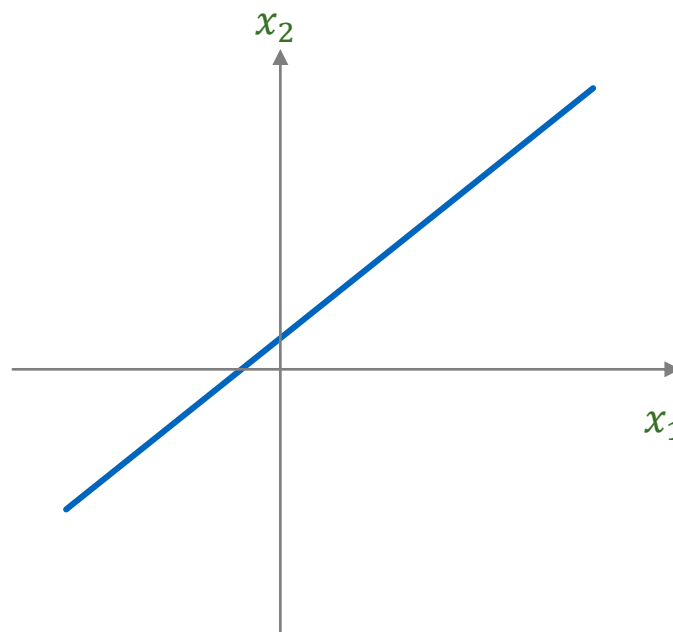
$$y = g(1 + 5x_1 - 3x_2)$$

- ▶ Some points:

$$x = \begin{bmatrix} -2 \\ 2 \end{bmatrix} \Rightarrow y = g(-15) = 0.03$$

$$x = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \Rightarrow y = g(0) = 0.5$$

$$x = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \Rightarrow y = g(5) = 0.99$$

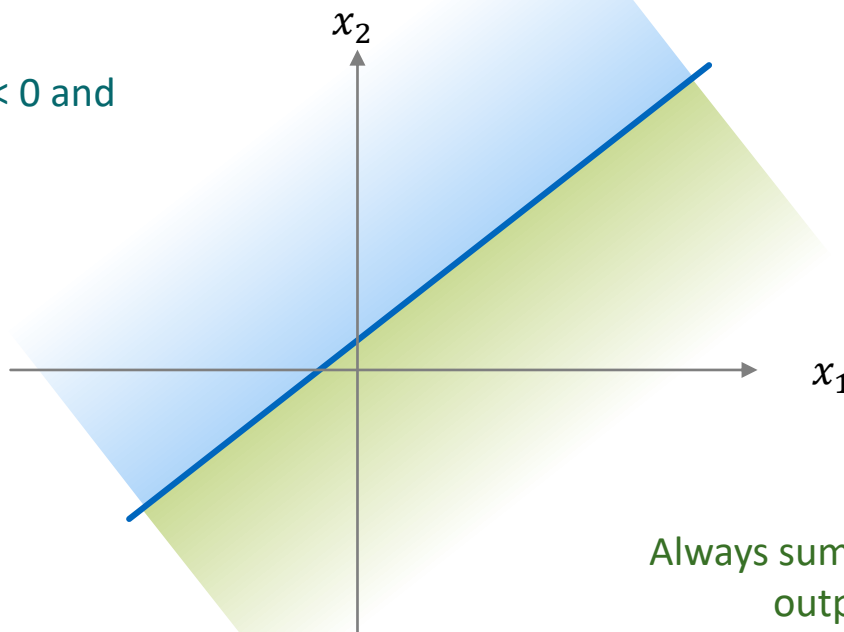


# Fundamental Concepts

## An example of Perceptron with Sigmoid AF

- Generally, we can conclude that:

Always summation  $z < 0$  and  
output  $y < 0.5$



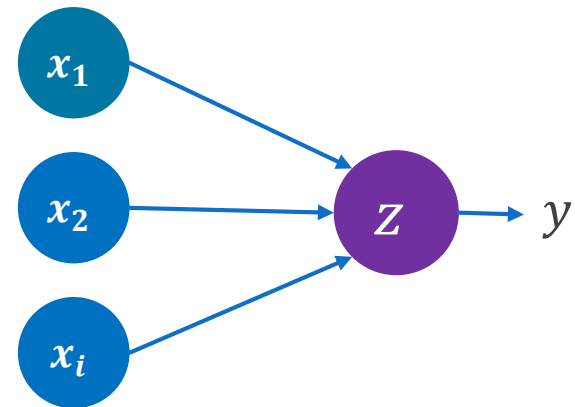
Always summation  $z > 0$  and  
output  $y > 0.5$

# Neural Networks

- ▶ Let's simplify the diagram:

$$z = w_b + \sum_{i=1}^m x_i w_i$$

$$y = g(z)$$



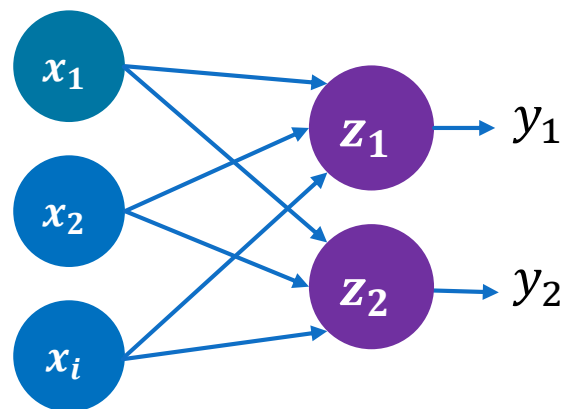
- ▶ What if the Perceptron provides **multiple outputs**?
  - ▶ Can we connect two Perceptron networks to have two different results?

# Neural Networks

- ▶ We can extend it to **Multi-Output Perceptron (MOP)**:

$$z_j = w_{b,j} + \sum_{i=1}^m x_i w_{i,j}$$

$$y_i = g(z_i)$$

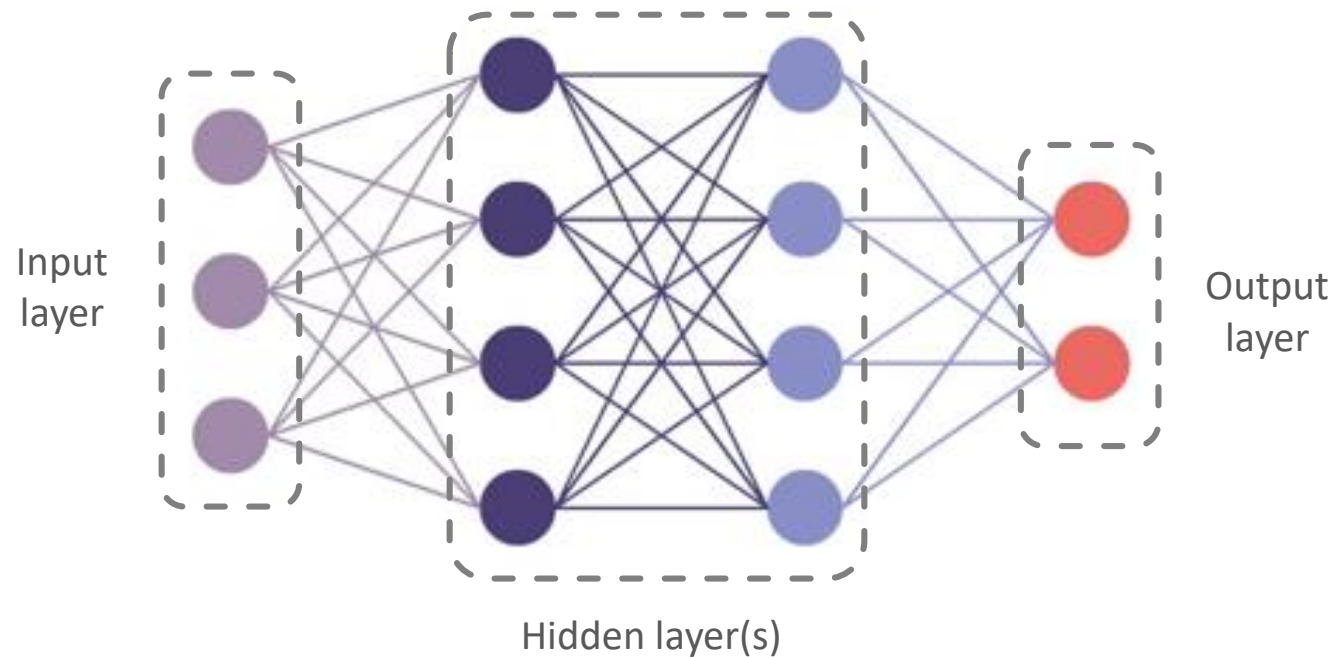


## Dense Layers

- ▶ All inputs are connected to all outputs

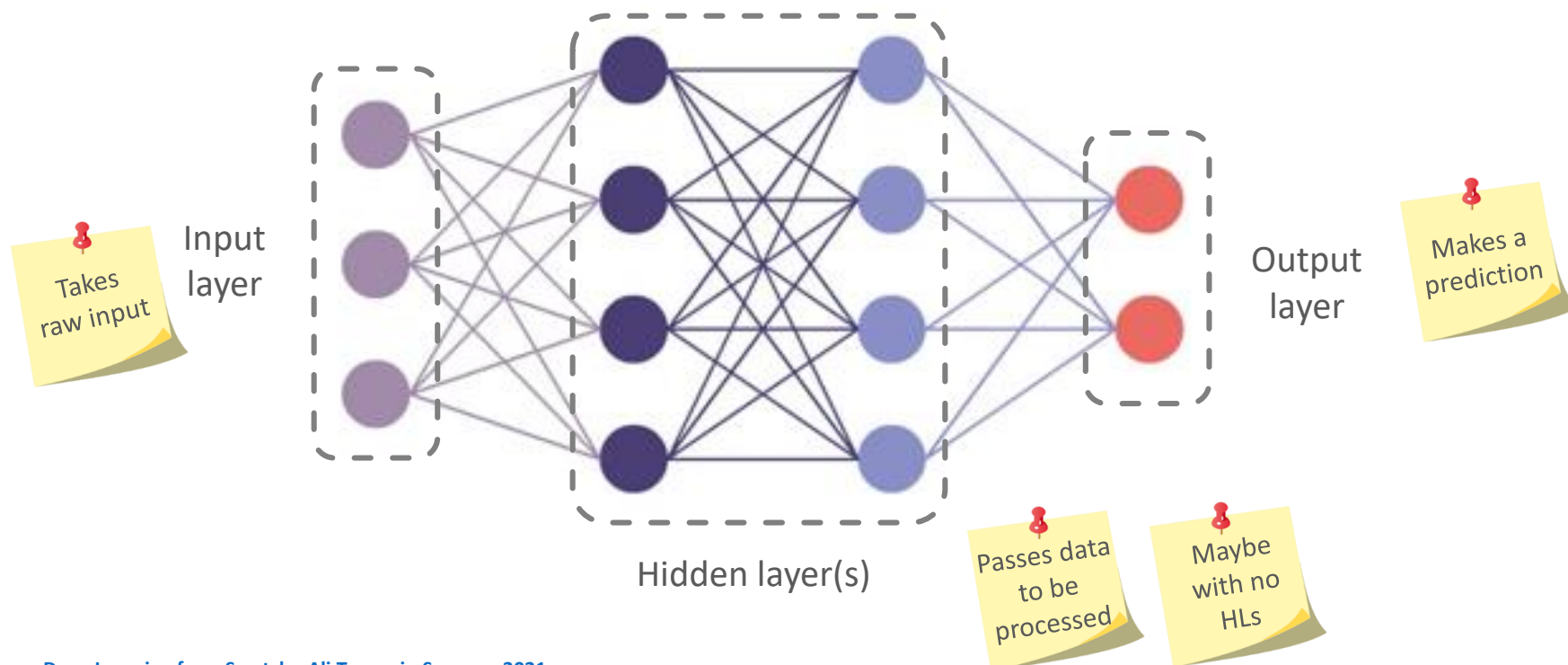
# Neural Networks

Let's connect several single nodes (Perceptron) together:



# Neural Networks

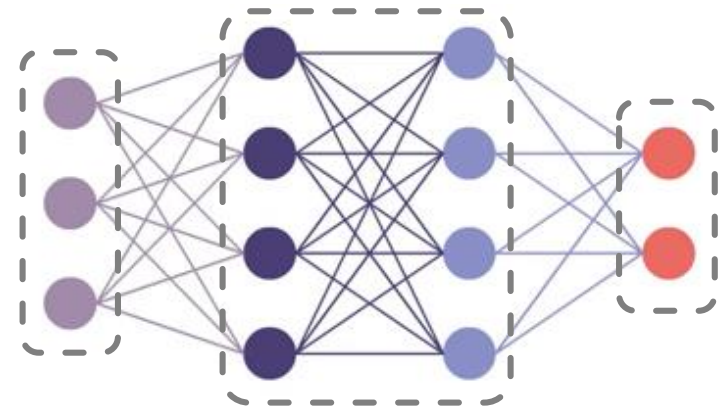
Let's connect several single nodes (Perceptron) together:



# Neural Networks

## Useful Notes

- ▶ If there is only one layer in HL, it is called a **Single Layer NN**.
- ▶ If the number of layers in HL is more than five, (usually) it is called it a **Deep NN**.



- ▶ We can always find sequential dense layers in NNs, like the image above with 3x4, 4x4, and 4x2 dense layers



# Neural Networks

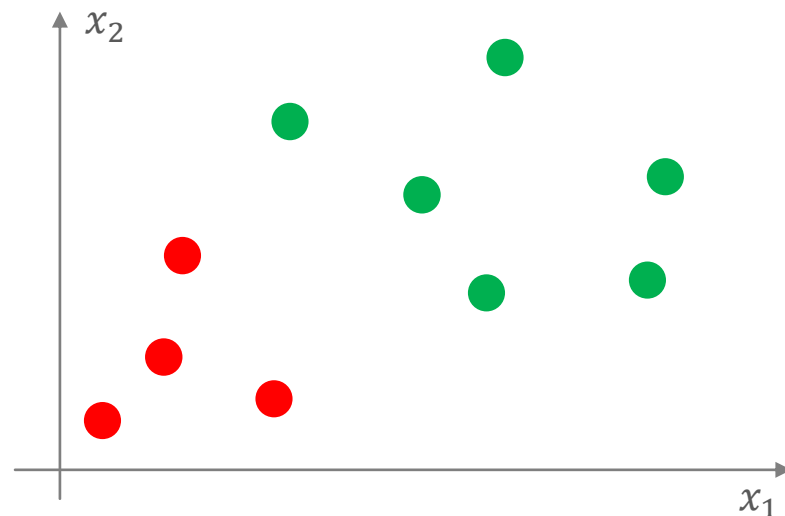
## Sample#1: passing a driving test

- ▶ Assume we want to pass both theory and practical tests.

$$x_1 = \{\text{topics achieved in theory classes}\}$$

$$x_2 = \{\text{hours spent in practical tests}\}$$

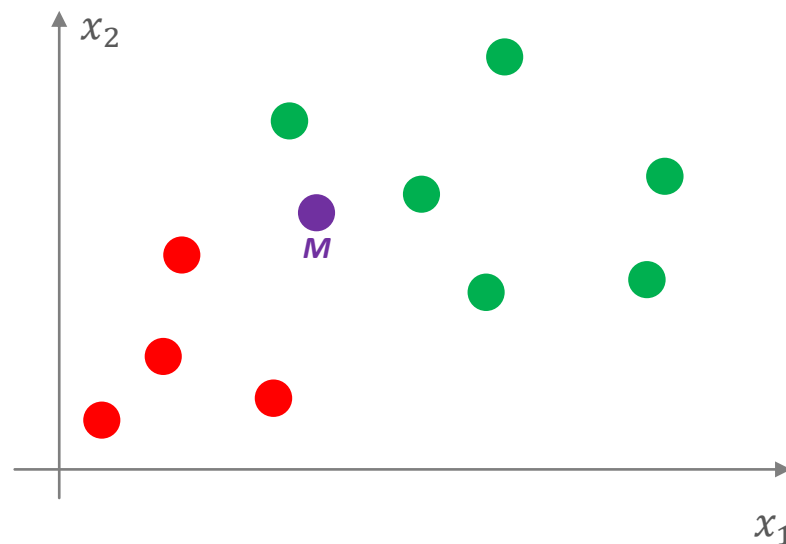
- ▶ Now, we can have:
  - ▶ The distribution of people who have passed or failed the test
  - ▶ The chart is 2D, as there are **two features**



# Neural Networks

## Sample#1: passing a driving test

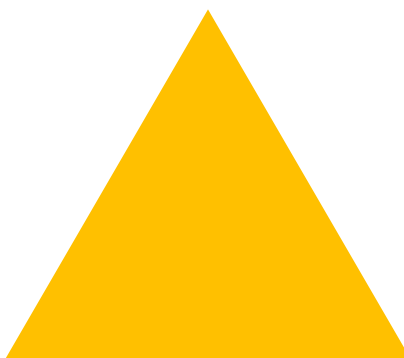
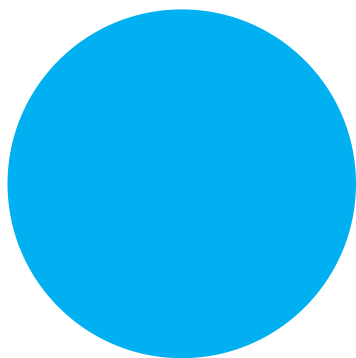
- ▶ Question: will a given node  $M$  pass the test?
- ▶ Assume  $M = [10, 12]$ , meaning that
  - ▶ 10: number of achieved topics
  - ▶ 12: hours spent training
- ▶ *Answer:* we need to **train** the ANN first
  - ▶ Training = adjusting the weights



# Neural Networks

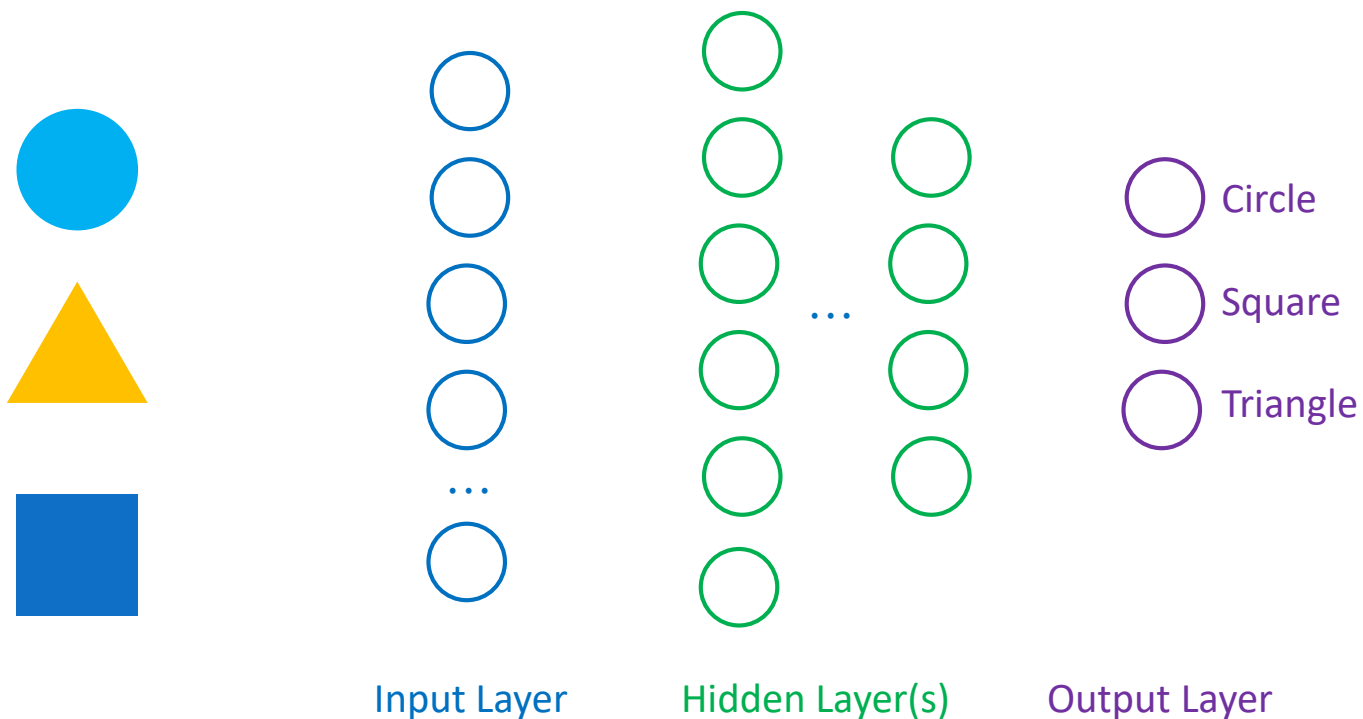
## Sample#2: Classification of objects into predefined classes

- ▶ **Inputs:** Rectangle, Circle, Triangle
- ▶ **Goal:** Correct classification



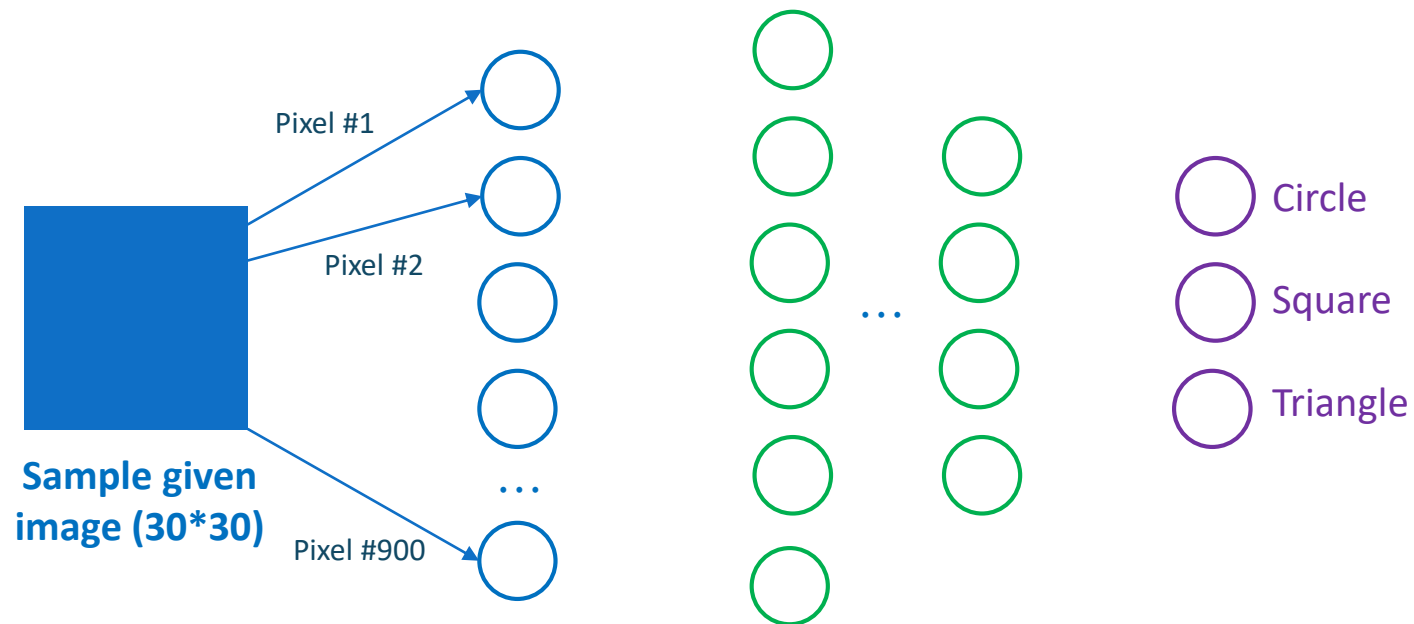
# Neural Networks

**Sample#2:** Classification of objects into predefined classes



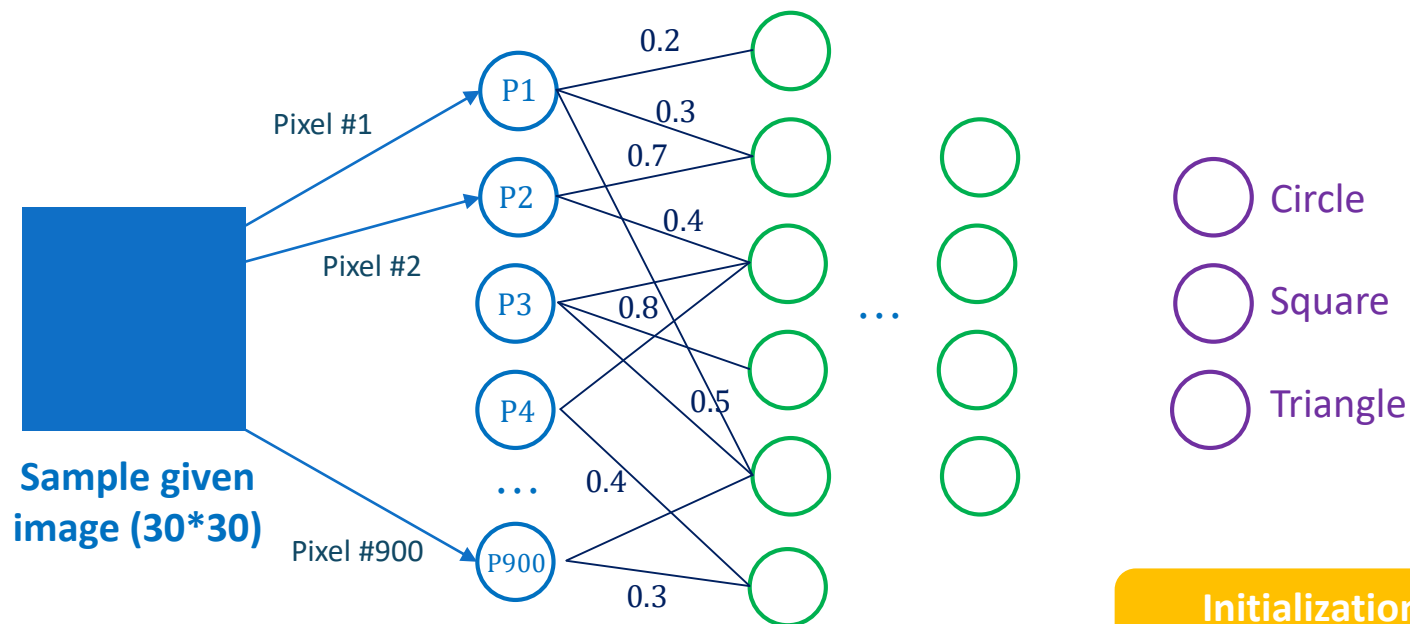
# Neural Networks

Each pixel is mapped to a **single neuron** from the input layer



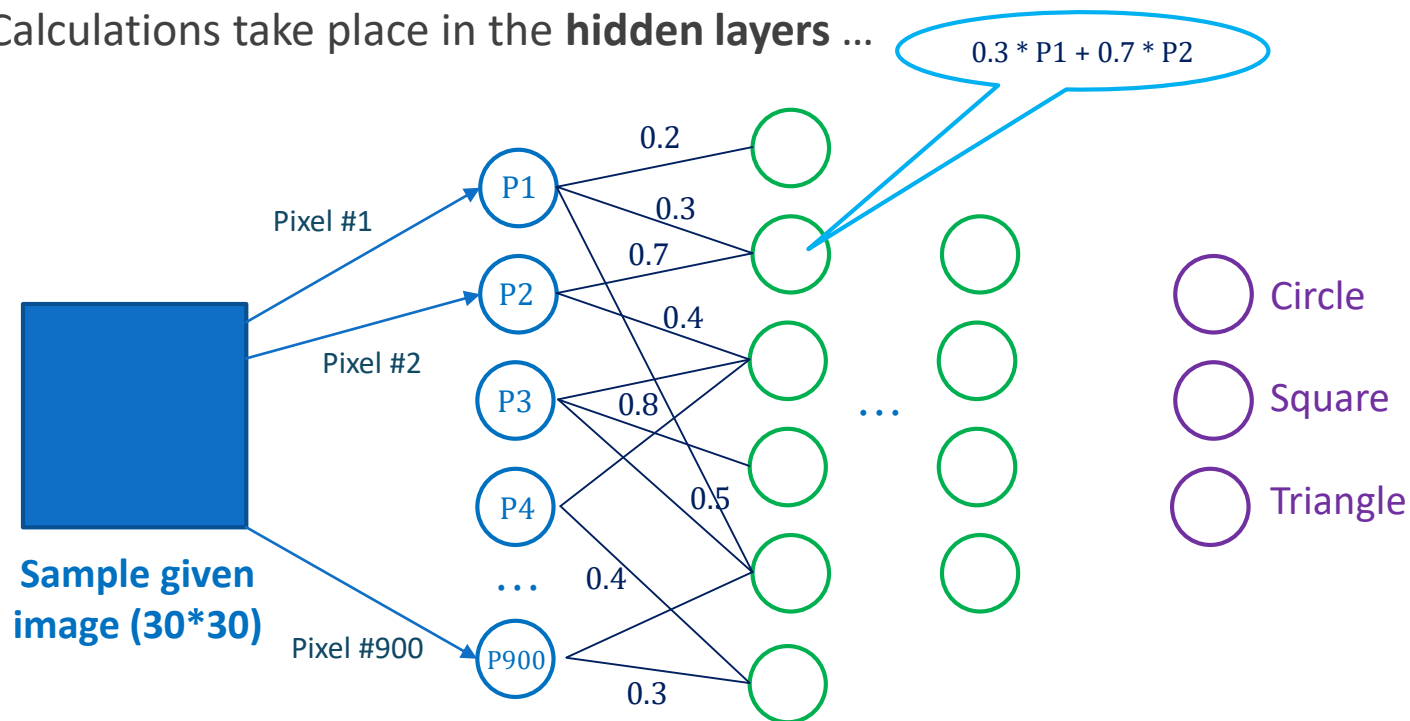
# Neural Networks

Calculations take place in the **hidden layers** ...



# Neural Networks

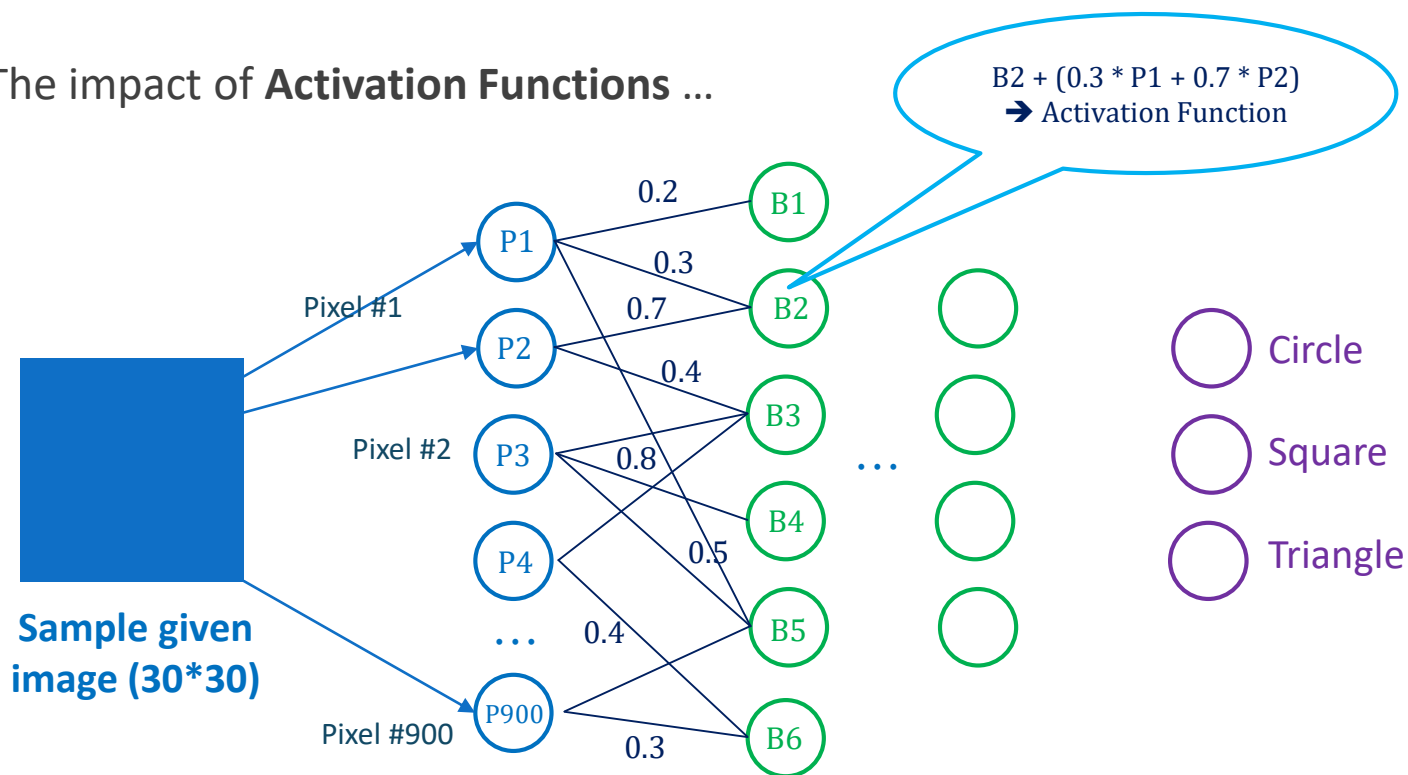
Calculations take place in the **hidden layers** ...





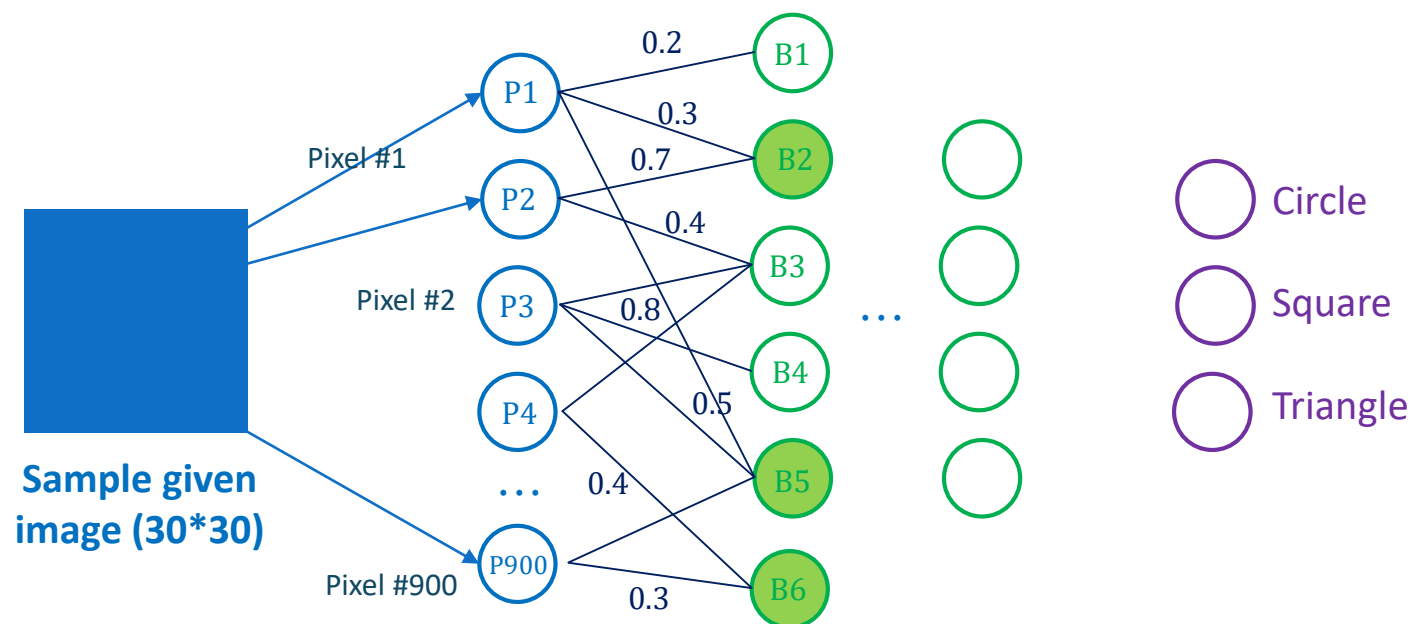
# Neural Networks

The impact of **Activation Functions** ...



# Neural Networks

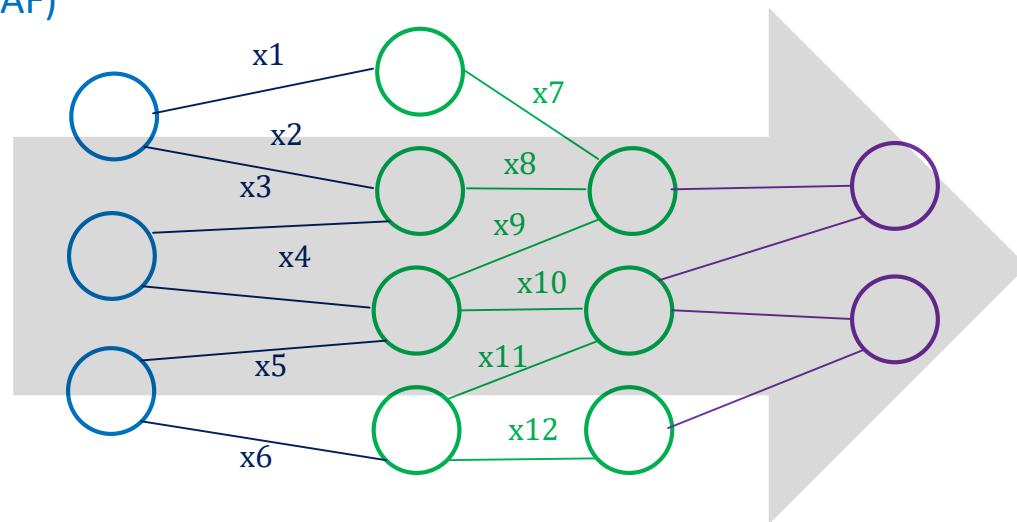
The impact of **Activation Functions** (Active nodes with **green** backgrounds) ...



# Neural Networks

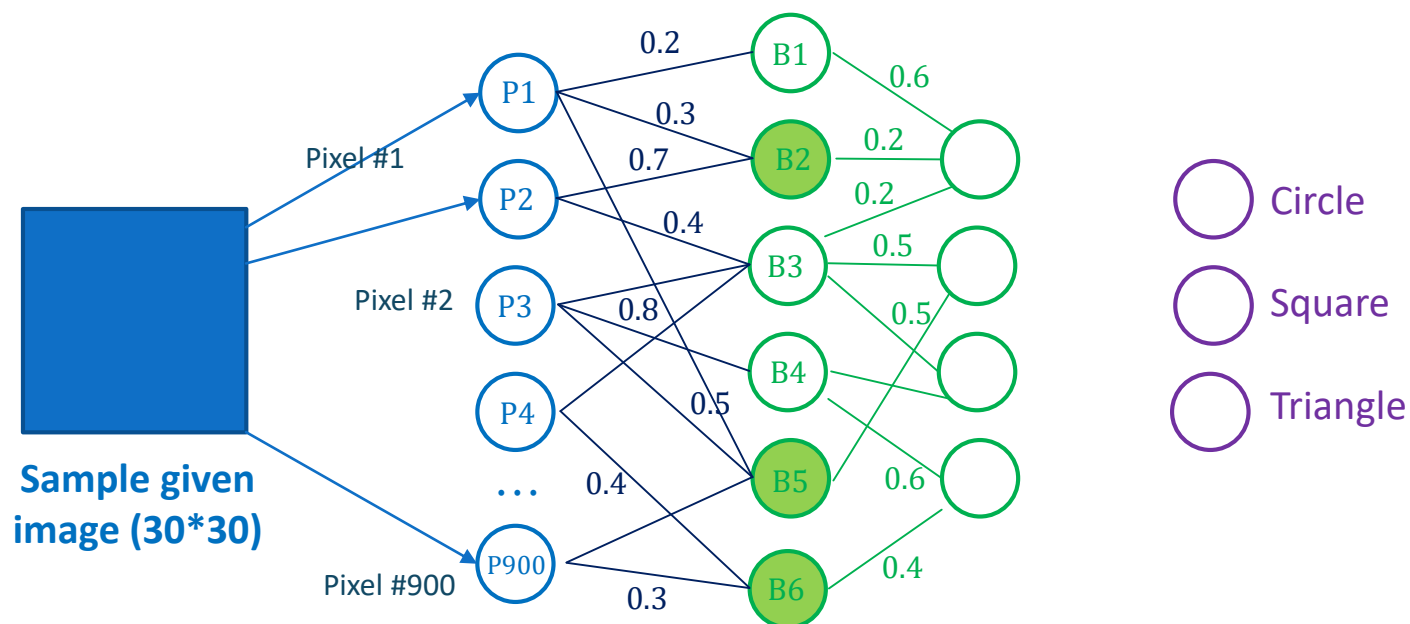
## Forward pass

- ▶ The main flow of calculations, where each layer:
  - ▶ Takes inputs from the previous layer
  - ▶ Processes them (Summation + AF)
  - ▶ Generates outputs
  - ▶ Passes them to the next layer
- ▶ Do while generating the *F.O*



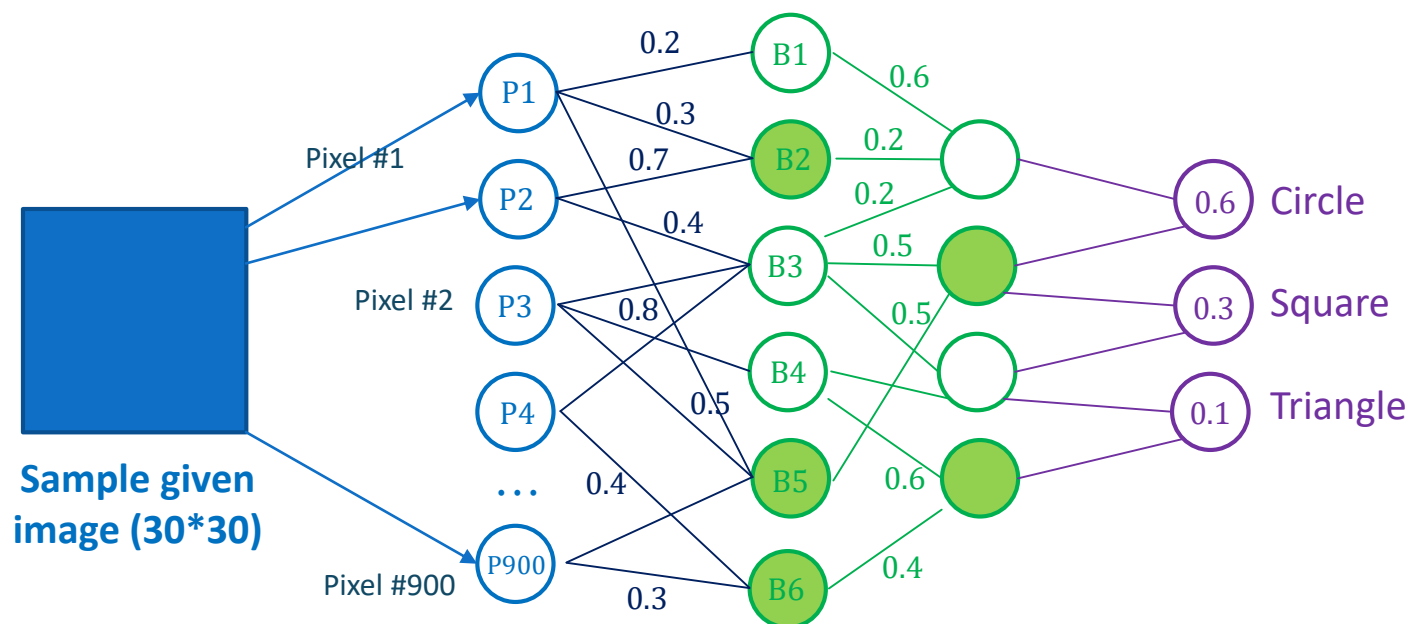
# Neural Networks

## Forward pass



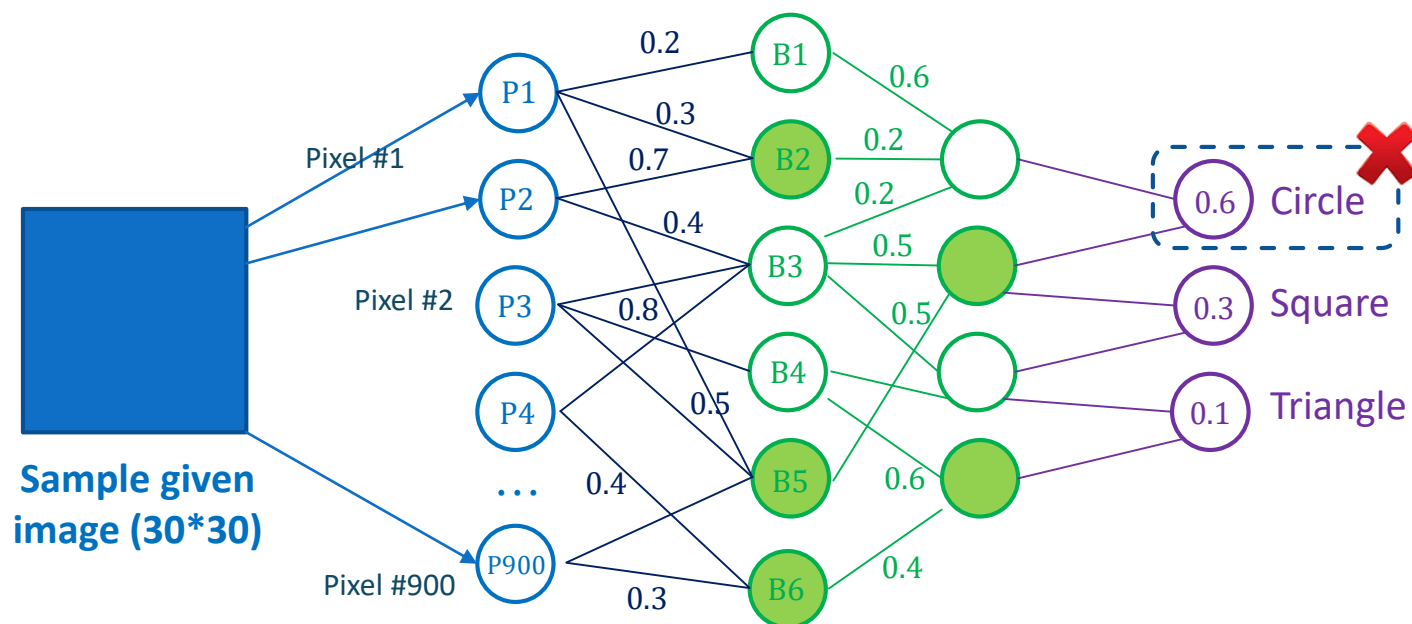
# Neural Networks

## Forward pass



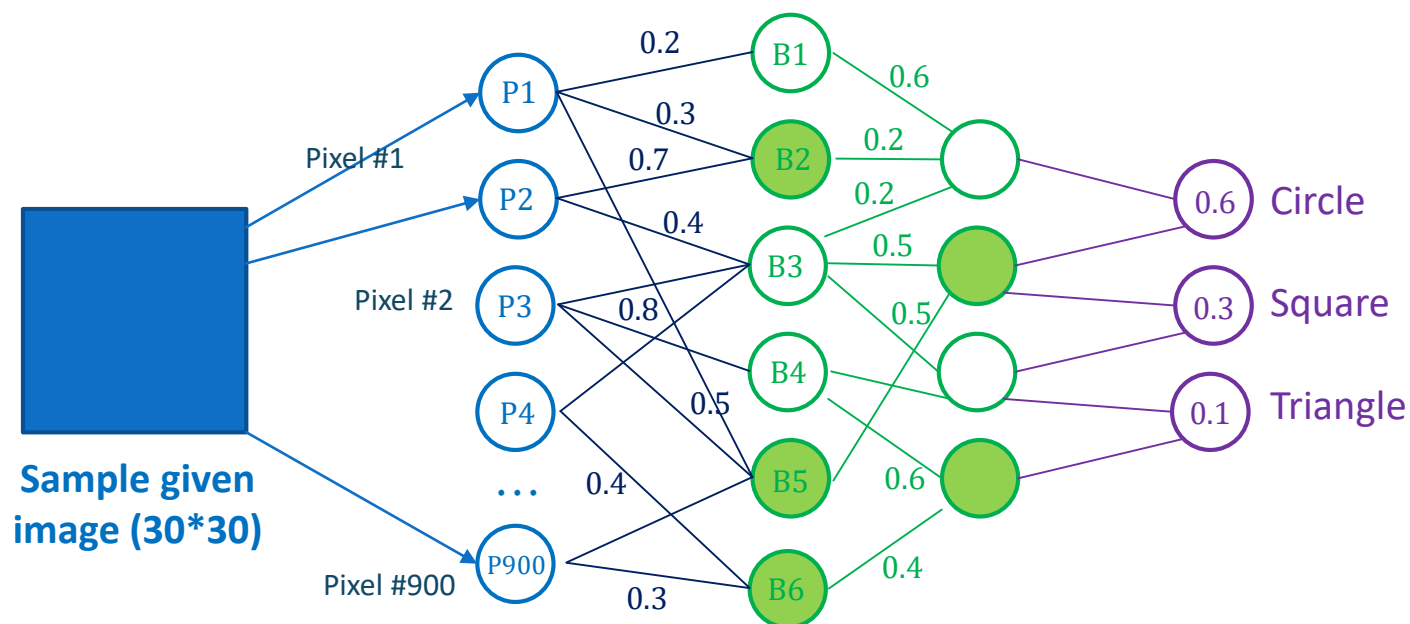
# Neural Networks

**Wrong prediction!** (The highest probability is circle with confidence score 0.6)



# Neural Networks

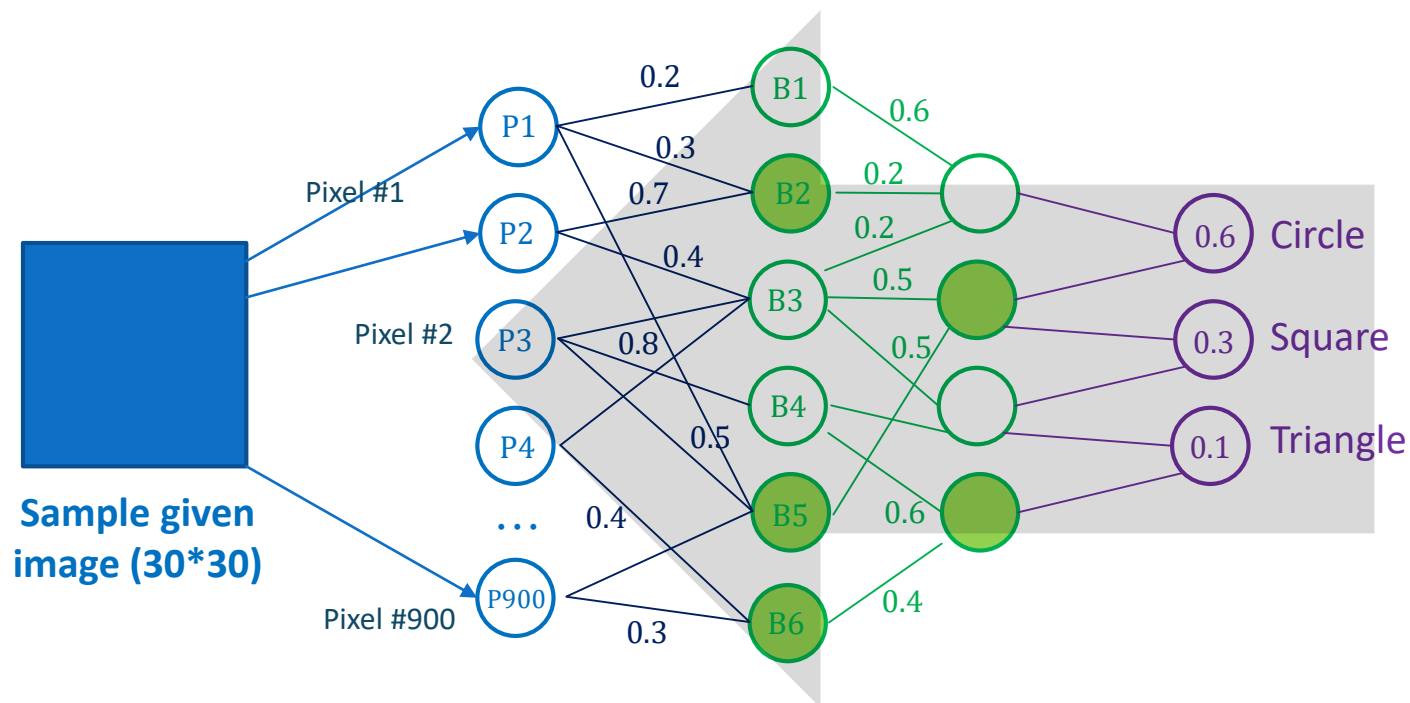
**How to know?** By comparing the results with the **Actual (expected) Outputs**



Actual	Error
0	-0.6
1	0.7
0	-0.1

# Neural Networks

How to resolve? **Backpropagation**



Actual	Error
0	-0.6
1	0.7
0	-0.1

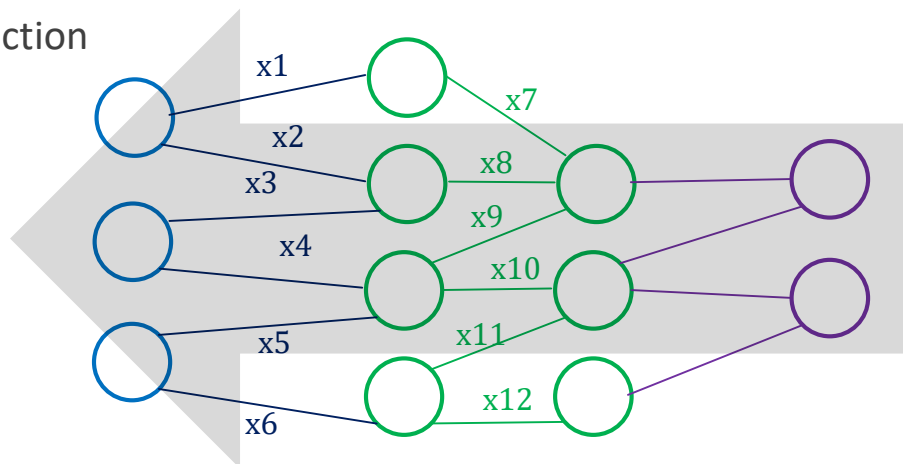


# Neural Networks

## Backpropagation (the backward pass)

- ▶ Why do we need this?
  - ▶ To discover the **optimal weights** for neurons
  - ▶ To minimize the **Loss Function (\*)**
  - ▶ To generate the best possible prediction
- ▶ How?
  - ▶ **Gradient Descent Algorithm (\*)**

**\* We'll talk about them soon!**



# Training the Network

- ▶ In a **supervised** training process, data is fed to the neural network, and only then the network learns how to adjust the weights based on it
  - ▶ We already know the **Correct Answers (Actual)**
- ▶ By comparing the **Actual and Predicted** results, we tell the network:
  - ▶ **Correct results?** No need to change the weights!
  - ▶ **Incorrect results?** Change the weights and get closer to the correct answer!

**Loss** is defined as:

- ▶ The cost of incurred from incorrect predictions



# Training the Network

## Empirical Loss (*Empirical Risk Minimization*)

- ▶ The average of all individual losses
  - ▶ Comparing the predicted and actual values
- ▶ Total loss applies over the whole dataset samples

$$J(w) = 1/n \sum_{i=1}^n \text{Loss}(h(x_i), y_i)$$

## Binary Cross-Entropy Loss (*Log Loss for Binary Classification*)

- ▶ Where the predicted values are probability rates, ranging between 0 and 1
- ▶ The final goal is to solve a binary problem (pass or fail)
- ▶ Comparing the real and estimated distributions

$$J(w) = 1/n \sum_{i=1}^n (y_i \cdot \log(h(x_i)) + (1 - y_i) \cdot \log(1 - h(x_i)))$$

# Training the Network

## Mean Squared Error Loss

- ▶ Used in **Regression** models (i.e., continuous real numbers outputs)
  - ▶ E.g., what will be someone's final grade? (non-binary)
- ▶ Subtracting the **actual** and **predicted** values
- ▶ Finally, calculating the average of its squared errors

$$J(w) = 1/n \sum_{i=1}^n (h(x_i) - y_i)^2$$

So, we can use different types of **Loss functions** for different problems. Using these quantified errors, we can train our ANN to **find optimal weights**.

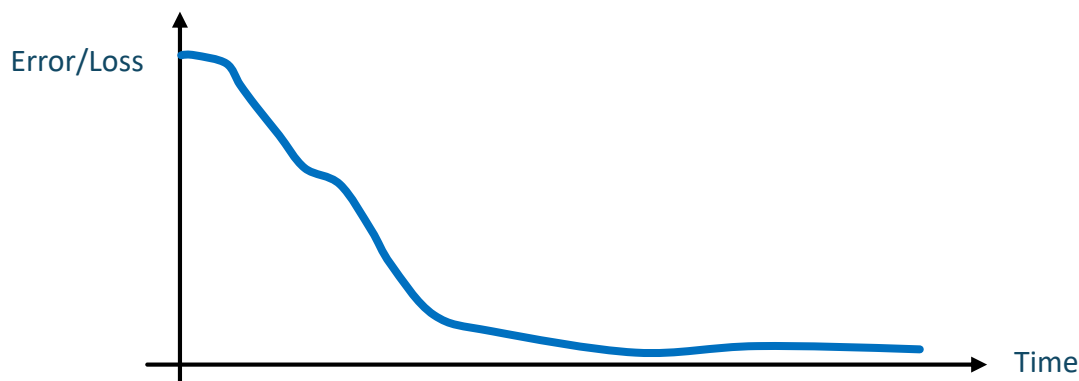
# Training the Network

## How to deal with loss/error functions?

- ▶ They are useful, as they show how far the **actual (expected) output** is from the **current (predicted) outputs** of ANN

### Our goals are

- ▶ To minimize the loss function (as much as possible!)
- ▶ To bring outputs as close as possible to the actual values



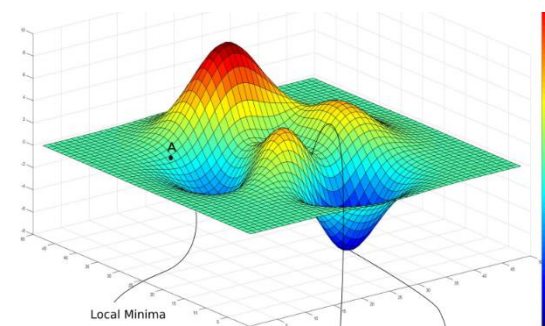
# Training the Network

## Loss Optimization

- ▶ The goal is to **minimize the loss** to achieve the optimal weights
  - ▶ Actually, the lowest possible error
  - ▶ Local Minima and Global Minima
- ▶ If  $W^*$  holds the optimal values of weights vector  $W$ :

$$W^* = \operatorname{argmin} J(W)$$

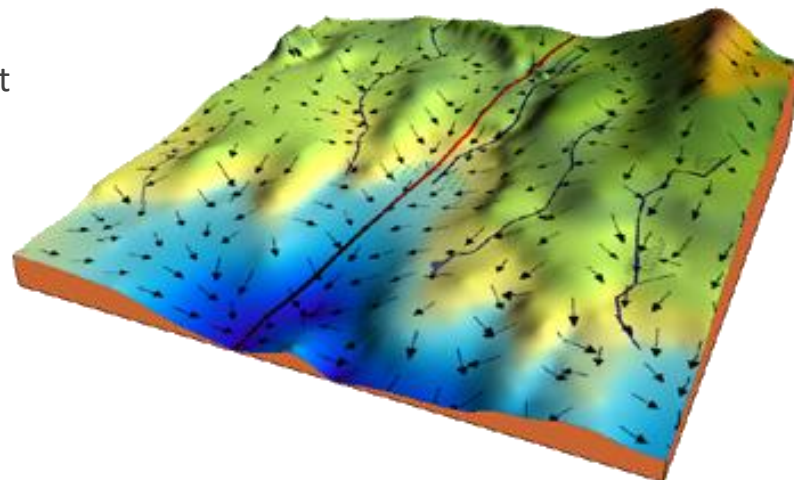
- ▶ It gets more complicated when the number of weights increases
- ▶ How? **Gradient Descent Algorithm**



# Training the Network

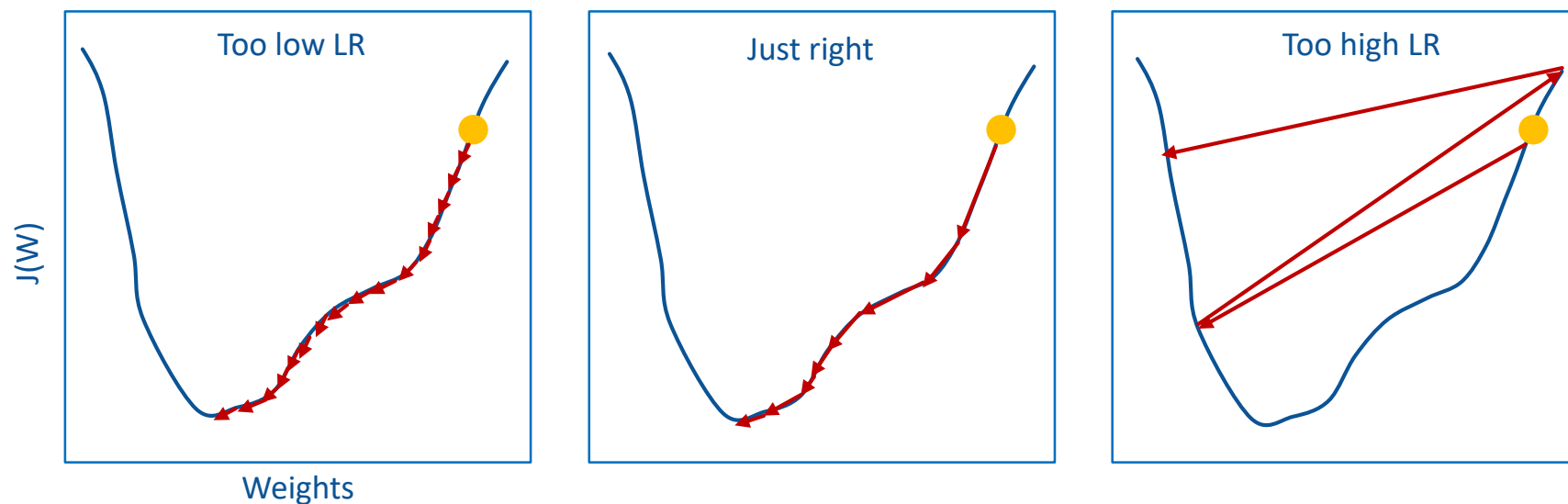
## Gradient Descent Algorithm

- ▶ An iterative optimization algorithm for finding a **local minimum**
  - ▶ By iteratively moving in **the direction of the steepest descent**
- ▶ Just like moving from a mountain towards the sea
  - ▶ A step-to-step downhill in the direction with a **negative gradient**
- ▶ **Learning Rate:** the size of steps
  - ▶ **High LR** → We may miss the optimal point
  - ▶ **Low LR** → We may be too slow



# Training the Network

## Gradient Descent Algorithm



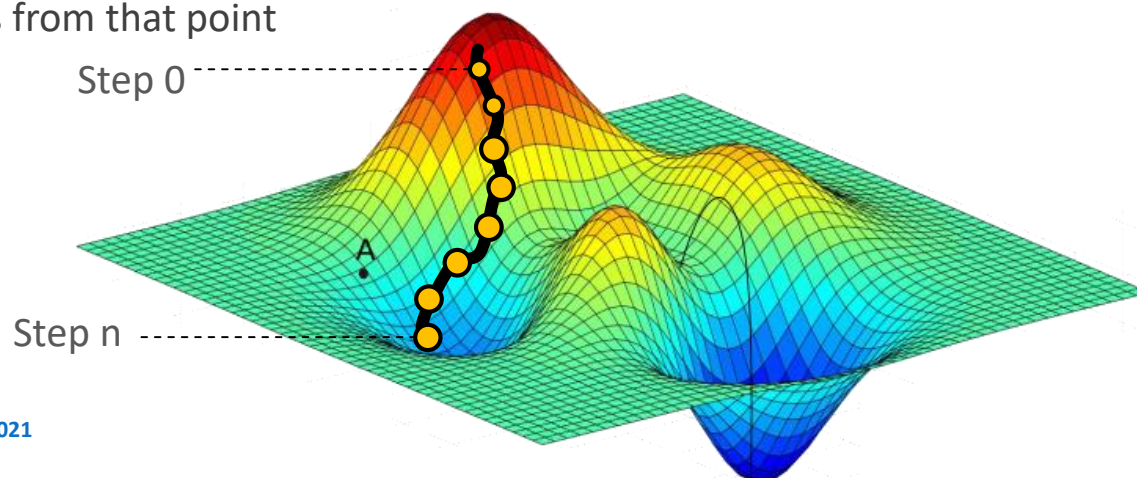
How to find proper LR? **1) trail and error**, **2) Adaptive LR**



# Training the Network

## Gradient Descent Algorithm

- ▶ Pick any random point in the landscape
- ▶ Compute gradient using  $\frac{\partial J(w)}{\partial w}$
- ▶ Take a small step towards the lower stages (steep)
- ▶ Repeat this process until reaching a **Local Minimum**
- ▶ Obtain the weights from that point

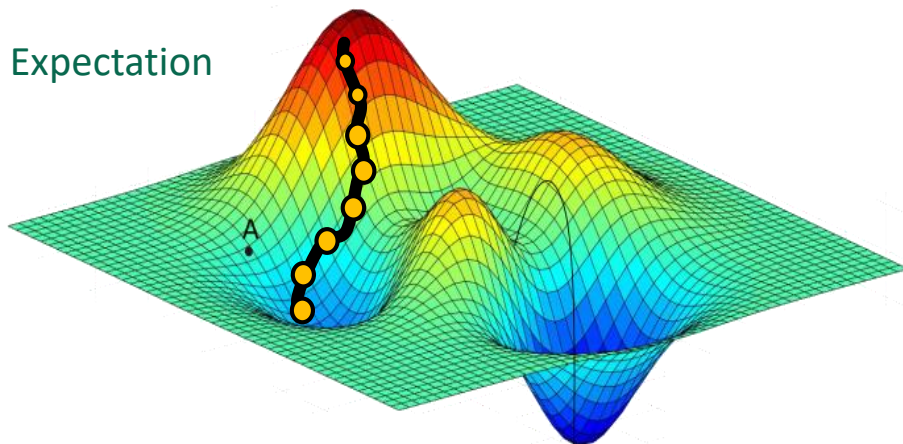


# Training the Network

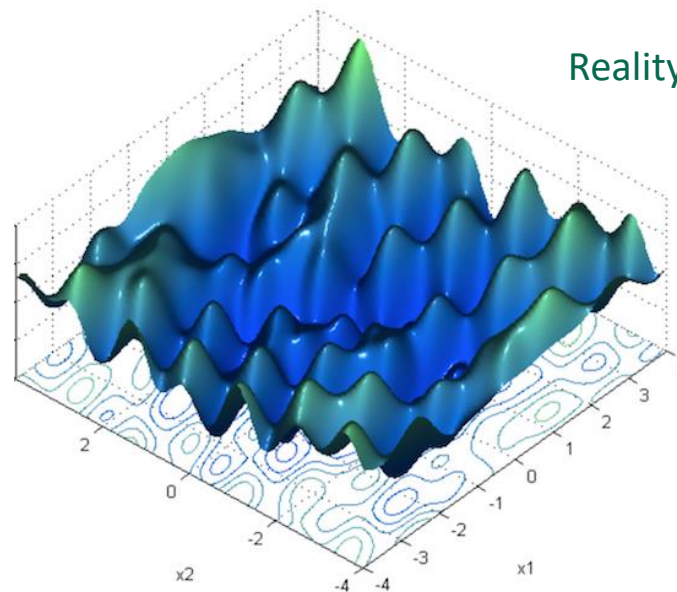
## Gradient Descent Algorithm

- ▶ It is not always that simple!
- ▶ In real-world scenarios, the algorithm may stuck!

Expectation



Reality



# Training the Network

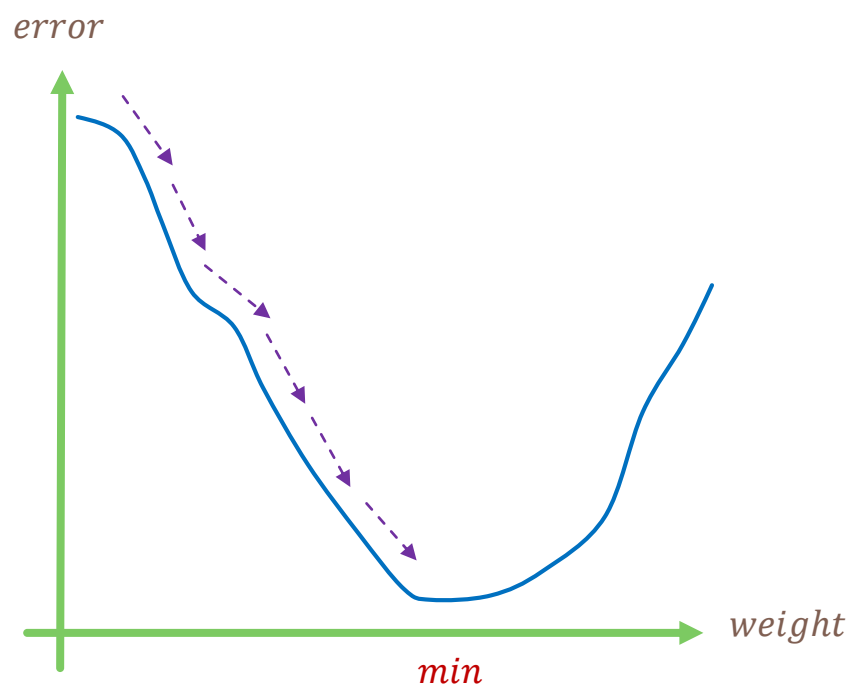
## Gradient Descent Algorithm

- ▶ How to update weights?

$$w_i(t) = w_i(t - 1) + \Delta w_i(t)$$

$$\Delta w_i(t) = \mu \left( -\frac{\delta J(W)}{\delta w} \right)$$

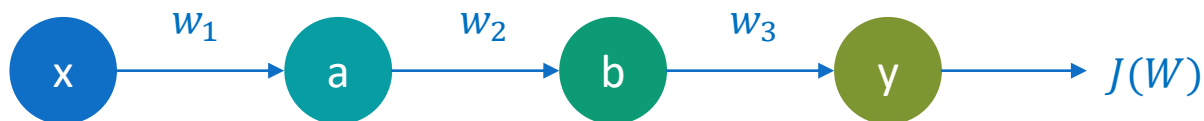
*Learning rate*



# Training the Network

## Gradient Descent Algorithm

- ▶ How do weight changes spread in the network?
  - ▶ **Remember:** the main goal is to **decrease the final loss**
  - ▶ **Recall:** we need to use the **Backpropagation algorithm** (Slide#48)
  - ▶ According to GDA:

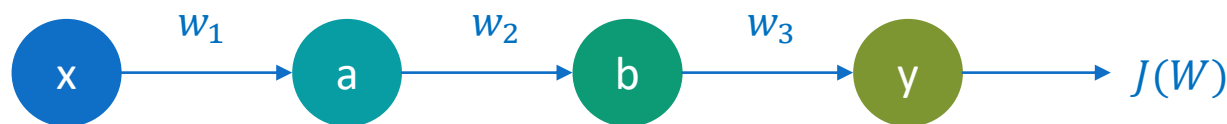


$$\frac{\delta J(W)}{\delta w_1} = \frac{\delta J(W)}{\delta y} * \frac{\delta y}{\delta b} * \frac{\delta b}{\delta a} * \frac{\delta a}{\delta w_1}$$

# Training the Network

## Gradient Descent Algorithm

- ▶ How do weight changes spread in the network?



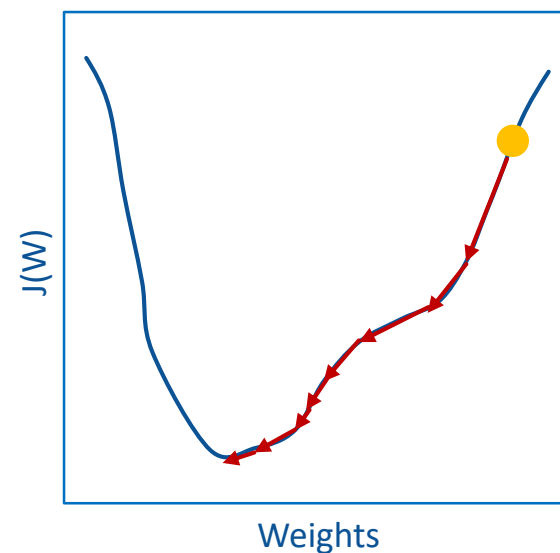
$$\frac{\delta J(W)}{\delta w_1} = \frac{\delta J(W)}{\delta y} * \frac{\delta y}{\delta b} * \frac{\delta b}{\delta a} * \frac{\delta a}{\delta w_1}$$

- ▶ This way, we can calculate all the weights from output to input
- ▶ **Good news:** we do not need to do this practically in our applications!
  - ▶ Backpropagation is implemented in different DL frameworks

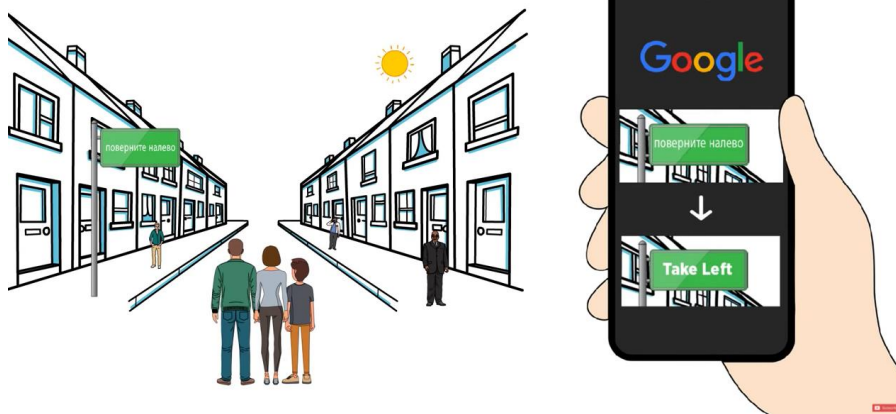
# Training the Network

## Adaptive Learning Rates

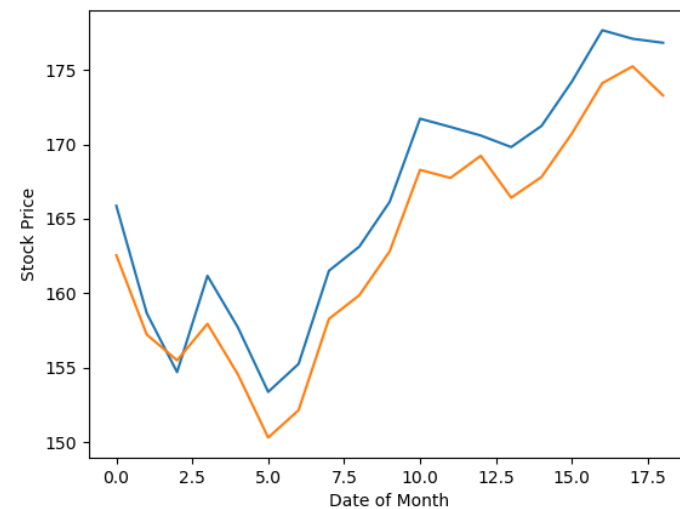
- ▶ In contrast with fixed LR, the value of LR may vary over time
  - ▶ It may become larger or smaller
- ▶ Factors affecting the value of LR
  - ▶ Learning speed (algorithm)
  - ▶ Size of weights
  - ▶ The value of gradient
- ▶ **Algorithms (optimization)**
  - ▶ Adam, AdaDelta, etc.



# Use cases of ANNs



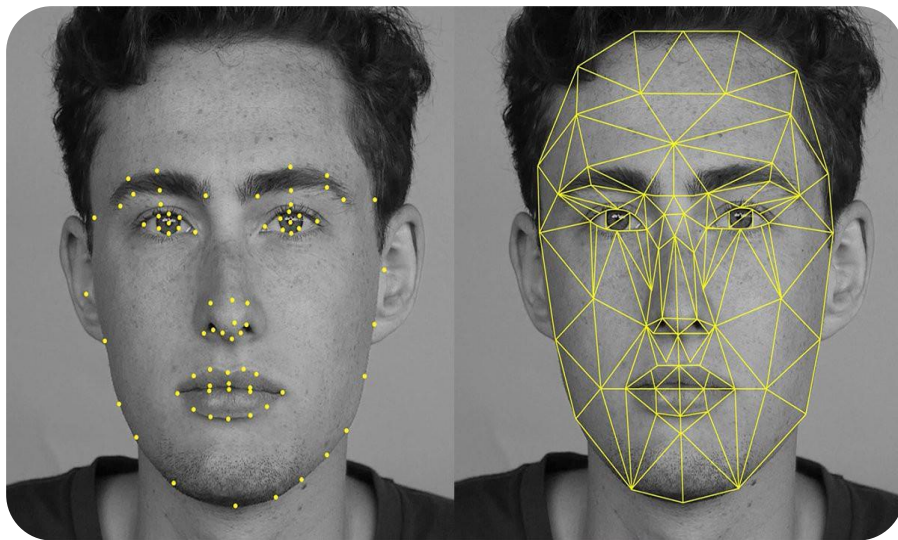
Google online translation tool



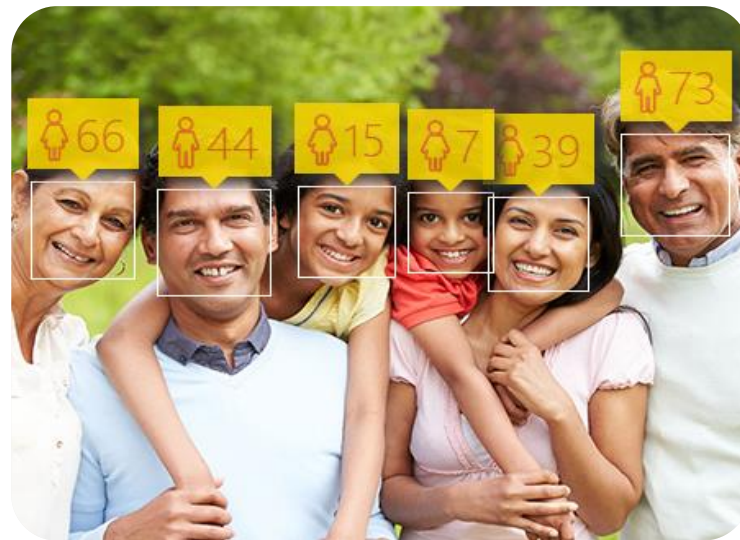
Stock market prediction



# Use cases of ANNs



Face recognition



Age estimation



# Use cases of ANNs



Image completion

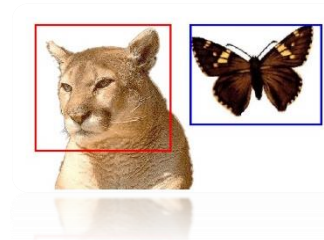
# Use cases of ANNs

## Suggested Projects – General ideas

Beginner

### Image Classification

- ▶ Identify the class of a given image



Beginner

### Visual Tracking

- ▶ Locate moving objects for surveillance



Beginner

### Face Detection

- ▶ Track and visualize human faces



Intermediate

### Image Caption Generation

- ▶ Analyze the context of an image to make captions

# Use cases of ANNs

## Suggested Projects – General ideas

Intermediate

### Text-to-Image

- ▶ Generate images from given texts

Intermediate

### Gender/Age Detection

- ▶ Estimate human age using facial features

Expert

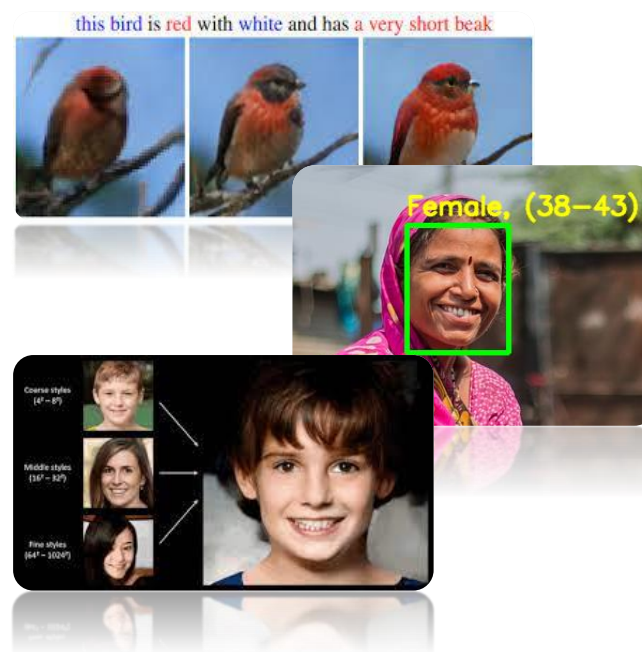
### Human Face Generation

- ▶ Generate unseen human faces

Expert

### Colorize Black & White Images

- ▶ Make B&W images colorful using a deep model



# Use cases of ANNs

## Suggested Projects – Intelligent Transportation Systems Lab

### *Vehicles*

- ▶ Vehicle detection and tracking
- ▶ Vehicle type classification
- ▶ Vehicle pose estimation

### *Drivers*

- ▶ Face pose detection
- ▶ Driver actions detection
- ▶ Driver drowsiness detection



# Assignments

## Find your Field of Interest!

- ▶ Choose a field of interest
- ▶ Find issues that can be resolved using Deep Learning
- ▶ Collect some academic publications about it
- ▶ Share your titles to discuss



# References

- ▶ <http://www.IntroToDeepLearning.com>
- ▶ <https://www.towardsdatascience.com>
- ▶ <https://data-flair.training/blogs/deep-learning-project-ideas/>
- ▶ <https://www.deeplearning.ai/>
- ▶ <https://blog.paperspace.com/a-practical-guide-to-deep-learning-in-6-months/>
- ▶ <https://wiki.pathmind.com/neural-network>

# Questions?

