# Large Language Models (LLMs) Bootcamp

# Deep Learning & Backpropagation
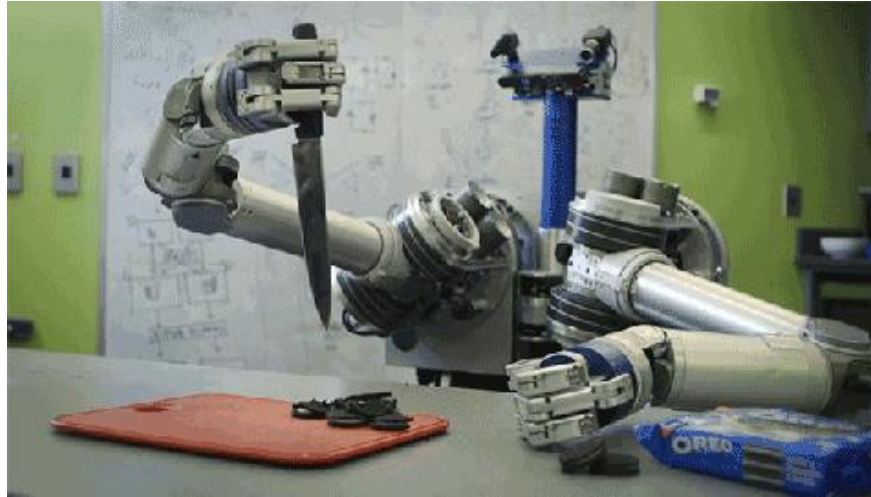
# OUTLINE

- Introduction to Deep Learning
- Applications of Deep Learning
- Neural networks
  - Activation Functions
  - Feed Forward Neural Networks
  - Neural Network Hyperparameters
- How do Neural Networks Learn?

# Introduction to Deep Learning

# Machine Learning

Machine Learning is a type of Artificial Intelligence that provides computers with the ability to learn without being explicitly programmed.
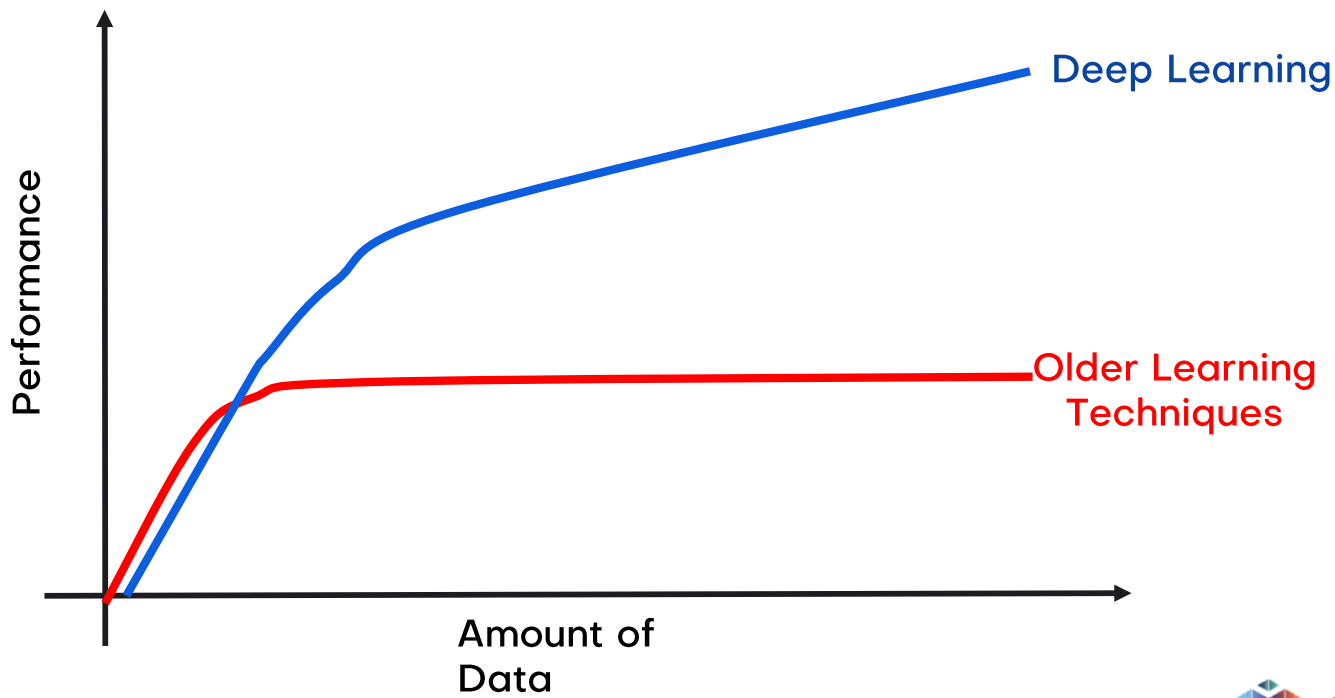
# Deep Learning

- Part of the machine learning field
- Exceptionally effective at identifying and learning patterns.
- Utilizes learning algorithms that derive meaning out of data by using a hierarchy of multiple layers that mimic the neural networks of our brain.
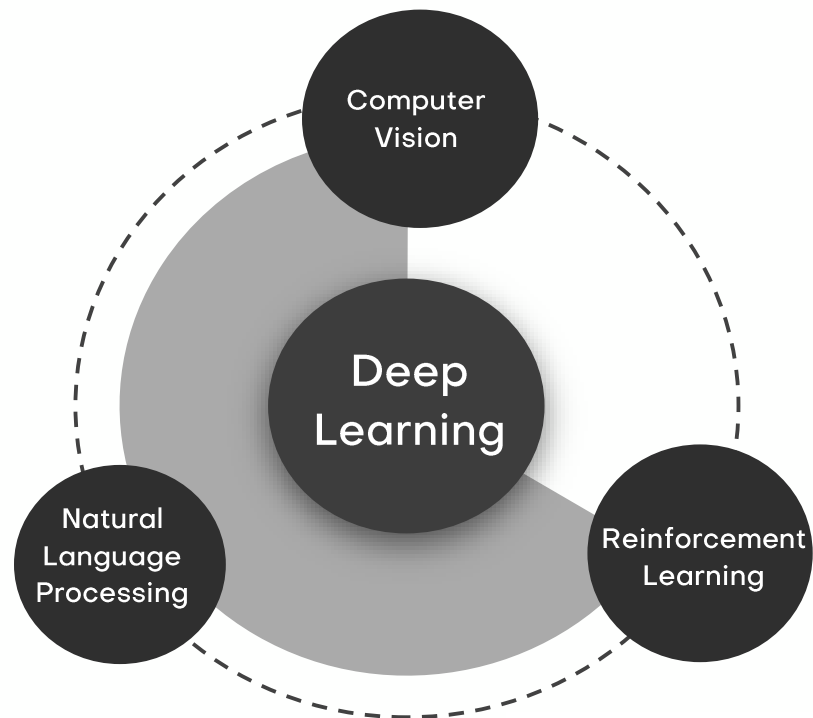
# Comparison
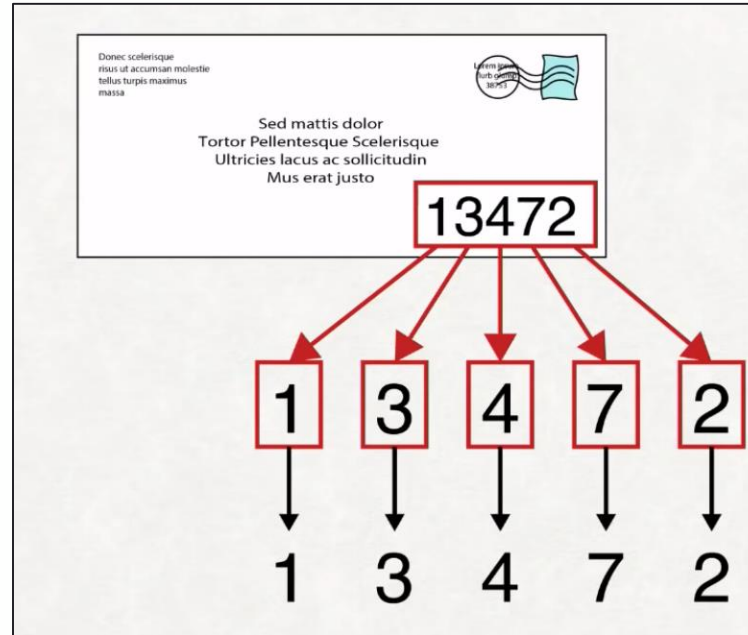
# Applications of Deep Learning

# Applications of Deep Learning
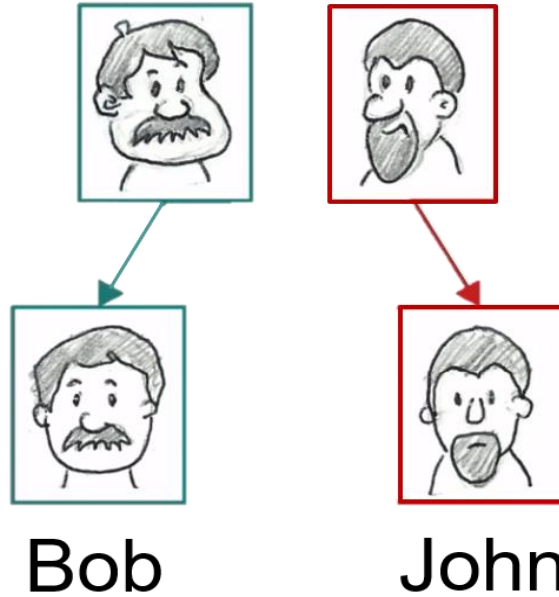
- Computer Vision: image processing, face recognition, object detection, video analytics

- Deep Reinforcement Learning

- Natural Language Processing: sentiment analysis, speech recognition, text-to-speech, conversational AI

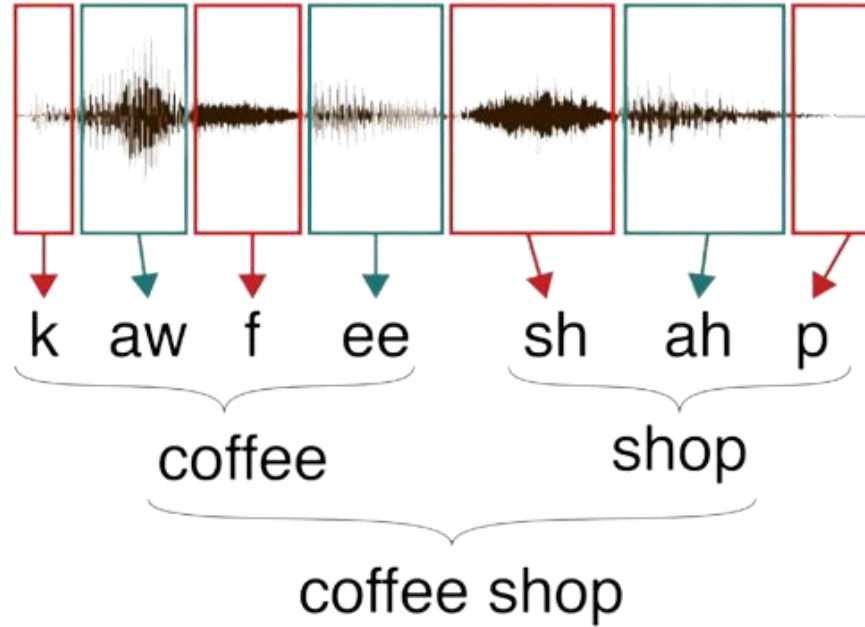# Optical Character Recognition OCR

# Face Recognition



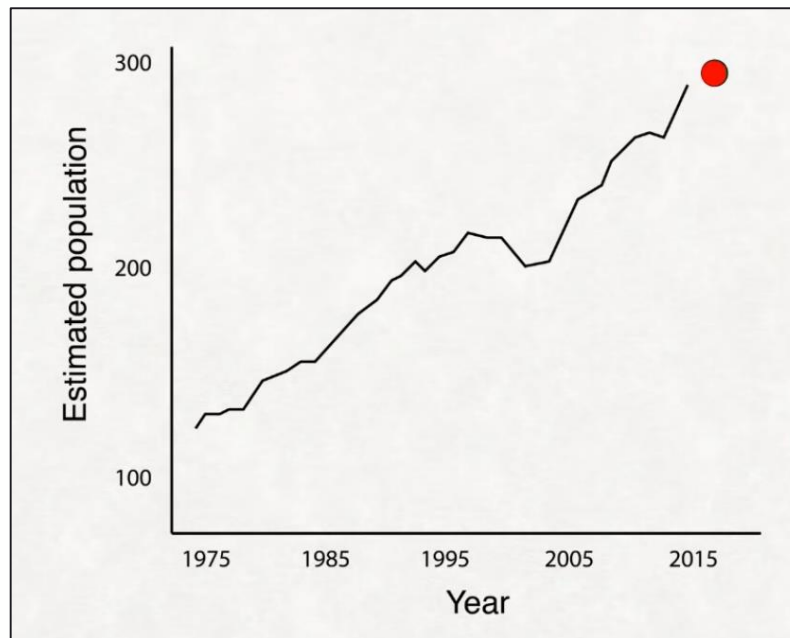Bob          John

# Virtual Assistants

# Time-series Regression

# Image Denoising

# Neural Networks

# Analogy

Our brain has lots of neurons connected together and the strength of the connections between neurons represents long term knowledge.

# Analogy



Inputs

$x_1$
$x_2$
$\vdots$
$x_n$

Dendrite

Axon terminal

Cell body

Myelin sheath

Outputs

$y_1$
$y_2$
$\vdots$
$y_m$

Myelinated axon

Inputs

$x_1$
$x_2$
$x_3$
...
$x_n$

$w_1$
$w_2$
$w_3$
$w_n$

$\Sigma$  $G$

Output

Weighted Sum

Activation Function

# Structure of an Artificial Neuron

The inside of an artificial neuron has 2 fundamental parts:
- The first part computes the weighted sum of the inputs (the net output)
- The second part receives the weighted sum and gives the final output.

# Activation Functions

# Common Activation Functions

The following are the most famous activation functions:

- ❖ Binary step
- ❖ ReLU
- ❖ Sigmoid
- ❖ Tanh



$x_1$ $x_2$ $x_3$

$w_1$ $w_2$ $w_3$

b

$\Sigma$ Weighted Sum

Activation Function f

Output

$$f\left(\sum weight \cdot Input + bias\right)$$

$$\sum weight \cdot Input + bias$$

# Binary Step

It tells if an input value is higher or lower than 0

Y = Sign(x)

Input Value: x



Output Value y

If (x>=0):
        y=1
else:
        y=0

# Rectified Linear Unit (ReLu)

Y = max(0, x)

Input Value: x

Output Value y

$y=x$

$y=0$

3

2

1

-3  -2  -1  0  1  2  3

If (x>=0):
         y=x
else:
         y=0

# Sigmoid Activation Function

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



Input Value: x

Output Value y

If (x>=0):
    0.5 <= y < 1
else:
    0 <= y < 0.5

# Hyperbolic Tangent (tanh)

Y = tanh(x)

Input Value: x

Output Value y

If (x>=0):
0<=y<=1

else:
−1<=y<0

# Feed Forward Neural Networks

# The Perceptron

- Single-output neuron
- Binary step activation function
- Used to classify between linearly separable classes

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

Parameters to learn

# Single Hidden Layer Network

- We have one layer, called 'hidden', between the input and the output layers.



Inputs

- Computes the weighted sum of the outputs of the previous layer.
- Passes sum to an activation function to obtain the final output.

- Computes the weighted sum of the inputs.
- Passes sum to an activation function to obtain output to be passed as input to next layer.

# Deep Neural Networks



| Input Units | Hidden Units | Output Units |

# Neural Network Hyperparameters

# Types of Activation Functions

- There's no Rule of Thumb.
- People tend to use ReLU in their hidden layers.
- The output layer activation should be consistent with the type of problem.
- The step function is not much used because it contains a discontinuity at 0 that leads to problems in derivatives computed during learning phase.



ReLU    ReLU    ReLU    Sigmoid

# Learning Rate

- Choosing the right learning rate is an essential step for reaching convergence.

$$w := w - \alpha \frac{\partial J}{\partial w}$$



$f(w)$ 

$w^*$    $w$

Too small: converge very slowly

$f(w)$

$w^*$    $w$

Too big: overshoot and even diverge

# Batch Size

The amount of data you feed to the network before performing an update in the weights.



All Training Data

Subset

Batch Size

Training Data Size

# Epochs

The number of times we feed ALL the training set examples to our network.

# Batch size vs. Epochs

# Batch size vs. Epochs

# Batch size vs. Epochs

# Batch size vs. Epochs

# Hands-on: Build your first Neural Network

# How do Neural Networks Learn?

# Data Structure

We take a general case where the input data consists of:

- m examples.
- Each of them having n features.
- 1 Output

| | m = 5 Examples | | | | |
|---|---|---|---|---|---|
| x0 | 0.2 | 0.2 | 0.2 | 0.3 | 0.2 |
| x1 | 0.1 | 0.3 | 0.1 | 0.4 | 0.1 |

n = 2 features

Output

| y | 1 | 0 | 1 | 0 | 1 |
|---|---|---|---|---|---|

# Network Structure

In the general case we took, a neural network will consist of:

- Input Layer (of n neurons)
- Hidden Layers
- Output layer (of 1 neuron)

Hidden Layers

Input Layer

Output Layer

n = 2
features

# Weights

A weight $W^l_{(i,j)}$ is defined by 3 characteristics:

- l: Layer index (destination layer)
- i: Index of the neuron in layer l (destination neuron)
- j: Index of the neuron in layer l−1 (source neuron)

Source Layer          Destination Layer

Source neuron  →Weight→  Destination Neuron

# Weights

A weight $W_{(i,j)}^{l}$ is defined by 3 characteristics:

- l: Layer index (destination layer)
- i: Index of the neuron in layer l (destination neuron)
- j: Index of the neuron in layer l–1 (source neuron)



| Layer 0 | Layer 1 | Layer 2 | Layer 3 |
|---------|---------|---------|---------|

$$W_{(2,0)}^{1} \qquad W_{(0,1)}^{2} \qquad W_{(0,0)}^{3}$$

Weights are usually initialized randomly

# Bias

- A bias is a characteristic of one neuron.
- We represent a bias by $b_i^l$ where:
  - i is the index of the given neuron
  - l is the index of the given layer

| Layer 0 | Layer 1 | Layer 2 | Layer 3 |
|---------|---------|---------|---------|



Bias values are usually initialized randomly

# Feed Forward

# Forward Pass

Recall the function that a single neuron having several inputs operates on.



$$Z(x) = b + x_1 * w_1 + x_2 w_2 + x_3 w_3$$

$$A(x) = G(Z(x))$$

# Forward Pass

Now this function will be executed for each neuron in each layer of the network.

- Neurons in each layer operate simultaneously.
- Activation functions are the same in each layer.

# Cost Function

# Loss Function

After the Forward Propagation is complete, we obtain a predicted value y' corresponding to the inputs we provided.

We define the loss function as some error metric between the predicted value (y') and the correct value (y) that we have for a specific input.

# Cost Function

The cost function tells you how much your model is "making mistakes" by averaging the loss functions computed on all individual training samples.

# BackPropagation

This algorithm tries to minimize the cost function J by adjusting the weights.
Given:

- A function J = f(w)

- The function $\frac{dJ}{dw} = \frac{df(w)}{dw}$

- The learning rate $\alpha$

# Gradient Descent

We will use the Gradient Descent algorithm to find the minimum.

- Pick a starting point w0 randomly

- Calculate $\left.\dfrac{dJ}{dw}\right|_{w=w_0}$

- Update w according to the formula

$$w_0 := w_0 - \alpha \left.\frac{dJ}{dw}\right|_{w=w_0}$$

- Get back to step 2

# Variants of Gradient Descent

| Optimization Problem | Samples in each Gradient Step | Number of Updates per epoch |
|---|---|---|
| Batch Gradient Descent | The entire Dataset | 1 |
| Mini-Batch Gradient Descent | Subsets of the dataset | $\frac{N}{\text{size of mini-batch}}$ |
| Stochastic Gradient Descent | Each Sample of the dataset | N |

# Variants of Gradient Descent



— Batch gradient descent
— Mini-batch gradient Descent
— Stochastic gradient descent

# Backward Pass

# Backward Pass

The same GD principle applies to when you're dealing with a function (J) of multiple variables (Weights).

$$\frac{\partial J}{\partial W^3} = \begin{bmatrix} \frac{\partial J}{\partial W_{00}^3} & \frac{\partial J}{\partial W_{01}^3} \end{bmatrix} \quad \frac{\partial J}{\partial b^3} = \begin{bmatrix} \frac{\partial J}{\partial b_0^3} \end{bmatrix}$$

$$\frac{\partial J}{\partial W^2} = \begin{bmatrix} \frac{\partial J}{\partial W_{00}^2} & \frac{\partial J}{\partial W_{01}^2} & \frac{\partial J}{\partial W_{02}^2} \\ \frac{\partial J}{\pa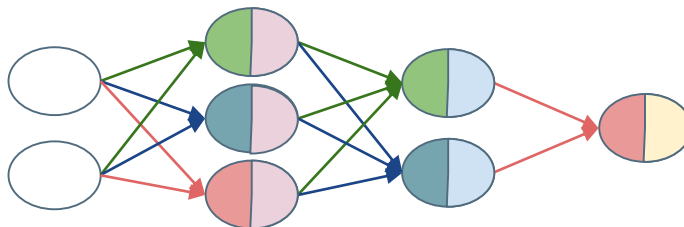rtial W_{10}^2} & \frac{\partial J}{\partial W_{11}^2} & \frac{\partial J}{\partial W_{12}^2} \end{bmatrix} \quad \frac{\partial J}{\partial b^2} = \begin{bmatrix} \frac{\partial J}{\partial b_0^2} \\ \frac{\partial J}{\partial b_1^2} \end{bmatrix}$$

$$\frac{\partial J}{\partial W^1} = \begin{bmatrix} \frac{\partial J}{\partial W_{00}^1} & \frac{\partial J}{\partial W_{01}^1} \\ \frac{\partial J}{\partial W_{10}^1} & \frac{\partial J}{\partial W_{11}^1} \\ \frac{\partial J}{\partial W_{20}^1} & \frac{\partial J}{\partial W_{21}^1} \end{bmatrix} \quad \frac{\partial J}{\partial b^1} = \begin{bmatrix} \frac{\partial J}{\partial b_0^1} \\ \frac{\partial J}{\partial b_1^1} \\ \frac{\partial J}{\partial b_2^1} \end{bmatrix}$$
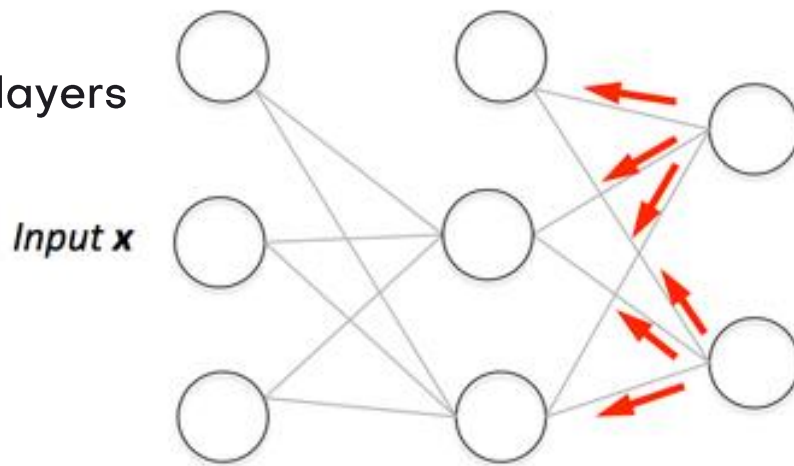
| Layer 3 | Layer 2 | Layer 1 |
|---|---|---|

# Backward Pass Process

- It starts at the output layer
- We compute the gradients of the loss function with respect to the weights
- Weights are updated
- We move backward to the hidden layers

*Input x*

# Weight & Bias Update

# Weight Update

After computing the derivatives, we compute the weights using the Gradient Descent update rule (where: alpha is the learning rate).

$$W^1 := W^1 - \alpha \frac{\partial J}{\partial W^1} \qquad b^1 := b^1 - \alpha \frac{\partial J}{\partial b^1}$$

$$W^2 := W^2 - \alpha \frac{\partial J}{\partial W^2} \qquad b^2 := b^2 - \alpha \frac{\partial J}{\partial b^2}$$

$$W^3 := W^3 - \alpha \frac{\partial J}{\partial W^3} \qquad b^3 := b^3 - \alpha \frac{\partial J}{\partial b^3}$$
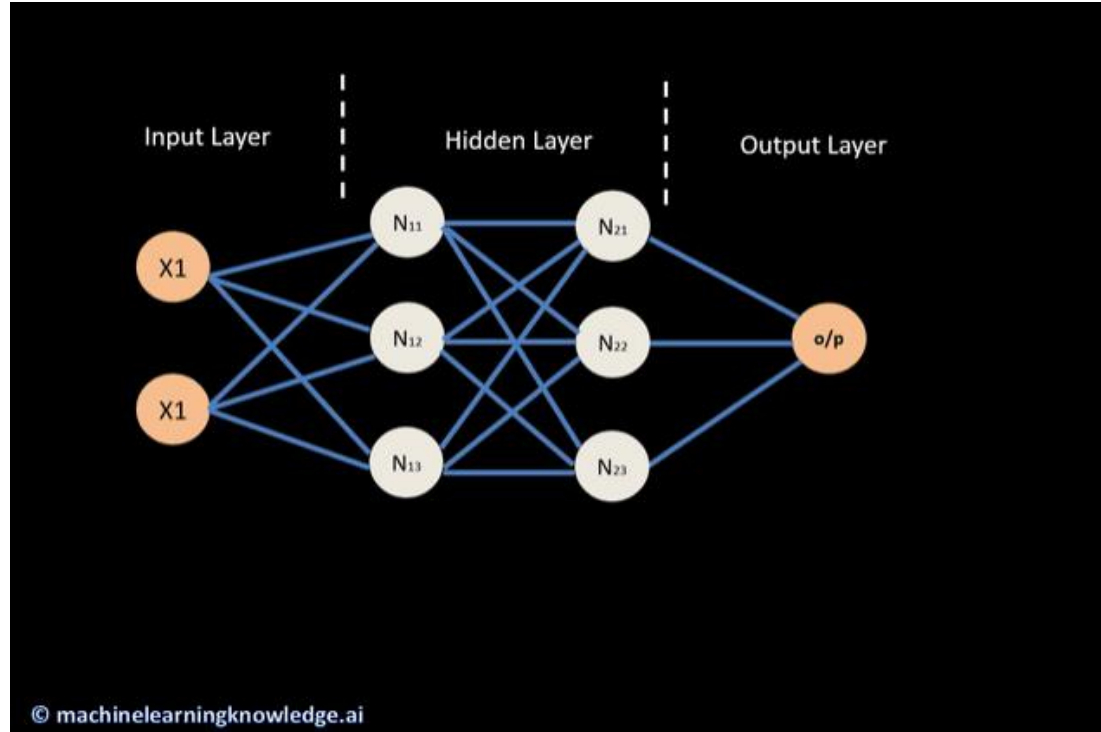
# Repeat the Process again

Now that we've computed the new weights, we repeat the forward-backward process for several iterations.

We stop in 2 cases:

- We reach a maximum number of iterations

- We reach an optimum value for the weights

# Recap



© machinelearningknowledge.ai

شكراً لكم

Thank you

SDAIA
الهيئة السعودية للبيانات
والذكاء الاصطناعي
Saudi Data & AI Authority

# Assignment: Implementing Backpropagation