

# Introduction to Deep Learning

Cristian Pachón García

January 31, 2020

## ① Introduction

## ② Linear regression

## ③ Ingredients for Deep Learning

- Loss function
- Gradient Descent

## ④ Deep Learning

- Introduction
- Perceptron
- Backpropagation

## 1 Introduction

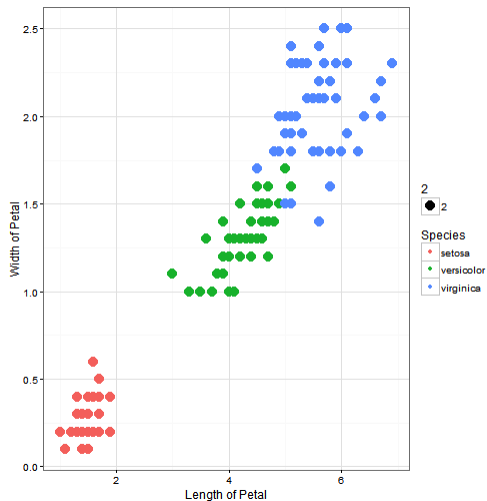
### 3 Ingredients for Deep Learning

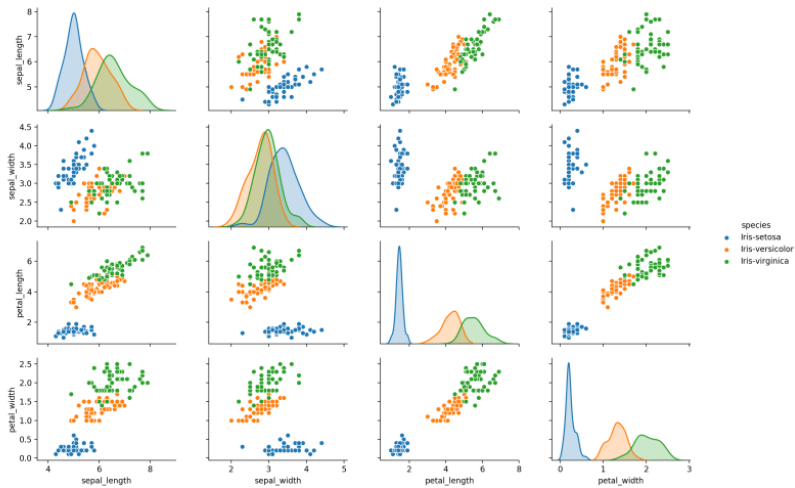
- Loss function
- Gradient Descent

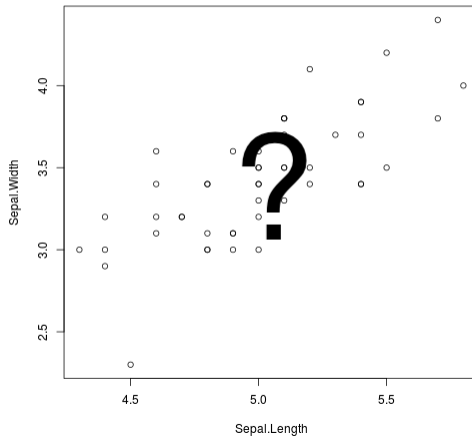
- 4 Deep Learning
  - Introduction
  - Perceptron
  - Backpropagation

# What is Machine Learning?

- Machine learning (ML) is the **scientific study** of **algorithms** and **statistical models** that computer systems use to perform a **specific task without using explicit instructions**, relying on patterns **learn from data**.
- The process of making the machine learn is called **the training process**.





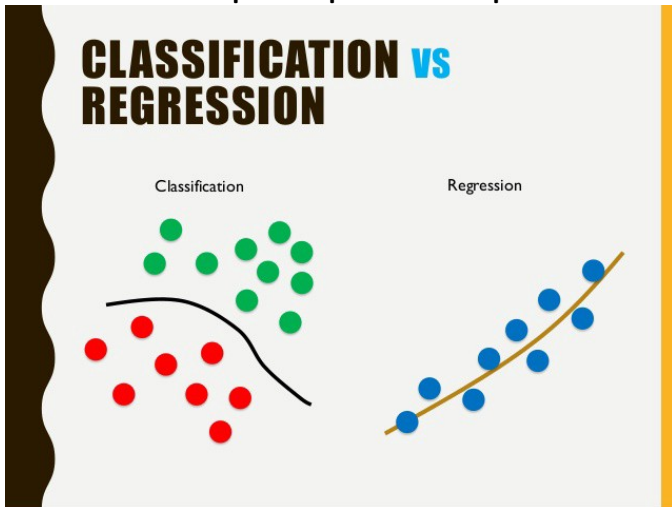


ML algorithms can be classified into two different groups:

- Supervised learning.
- Unsupervised learning.

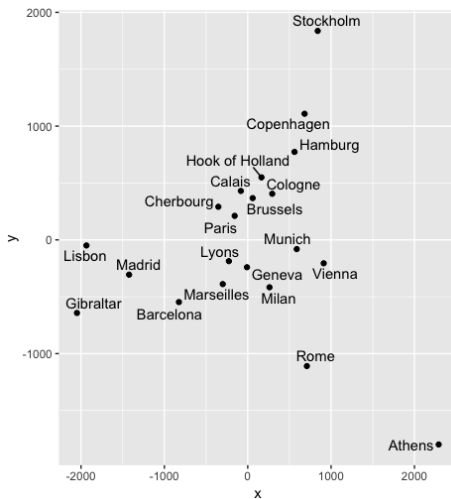


Supervised learning is the machine learning task of learning a function that **maps an input to an output**.



Unsupervised learning is the machine learning task that allows us to discover patterns from the data. Rather than predicting a variable (temperature, flower type, etc.) these algorithms are aimed to discover patterns in the data.

	Athens	Barcelona	Brussels	Calais	Cherbourg	...
Athens	0	3313	2963	3175	3339	...
Barcelona	3313	0	1318	1326	1294	...
Brussels	2963	1318	0	204	583	...
Calais	3175	1326	204	0	460	...
Cherbourg	3339	1294	583	460	0	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮



## 2 Linear regression

### 3 Ingredients for Deep Learning

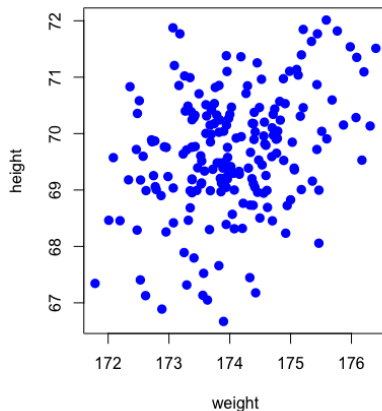
- Loss function
- Gradient Descent

- 4 Deep Learning
  - Introduction
  - Perceptron
  - Backpropagation

# Linear regression

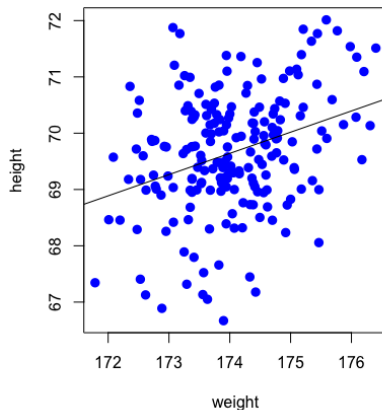
	height	weight
1	173.37	69.76
2	174.18	71.36
3	173.16	70.85
4	175.60	69.91
5	174.33	67.45
6	173.18	71.77
7	174.49	70.46
8	174.74	70.44
9	174.58	69.82
10	173.69	69.99
⋮	⋮	⋮
200	173.61	70.213

Is there any relationship between *height* and *weight*?



- We would like to find two coefficients ( $w$  and  $b$ ) such that  $weight = b + w * height$ .
- In general, given two pairs of variables  $x$  and  $y$ , we would like to find two coefficients ( $w$  and  $b$ ) such that  $y = b + w * x$ .
- $w$  is known as the *weight* and  $b$  is known as the *bias*.





How can we find  $w$  and  $b$  such that  $y = b + w * x$  is a "good approximation" to the data points?

More important question than *how* is: **why** do we know to know these parameters?

- If a system is modeled by an equation, it can help to predict what can happen under certain circumstances.

- Let's suppose  $b = 4.057$  and  $w = 0.3769$ . It means

$$y = 4.057 + 0.3769 * x.$$

What would be the *weight* ( $y$ ) of someone whose height ( $x$ ) is 174 cm?  $y = 4.057 + 0.3769 * 174 = 69.6376$  kg.

### 3 Ingredients for Deep Learning

- Loss function
- Gradient Descent

- 4 Deep Learning
  - Introduction
  - Perceptron
  - Backpropagation

### 3 Ingredients for Deep Learning

- Loss function
- Gradient Descent

- 4 Deep Learning
  - Introduction
  - Perceptron
  - Backpropagation

# Loss function

- Previously, we asked the model to predict the *weight* of a persona whose height is 174 cm.
- The prediction was 69.6376 kg.
- **Notation:** when we use the model to predict, we use a special symbol for the results:  $\hat{y}$ .
- So, in the previous example:  $\hat{y} = 69.6376$

How can we find  $w$  and  $b$  such that  $y = b + w * x$  is a "good approximation" to the data points?

- We need a metric that tells what is "good" and what is "bad".
- We want a metric that is close to 0 when the model is correct.
- We want a metric that increases as long as the model is not correct.

Let's assume we have  $w$  and  $b$ . For instance, at random we choose  $w = 0.2$  and  $b = 3$ :

	height( $x$ )	weight( $y$ )	prediction( $\hat{y}$ )	$error = (y - \hat{y})^2$
1	173.37	69.76	37.67	1029.39
2	174.18	71.36	37.84	1123.97
3	173.16	70.85	37.63	1103.53
⋮	⋮	⋮	⋮	⋮
200	173.6189	70.31279	37.72378	1062.0434

$$loss = MSE = \sqrt{\frac{1}{200}(1029.39 + 1123.97 + \dots + 1062.0434)} = 31.86411$$



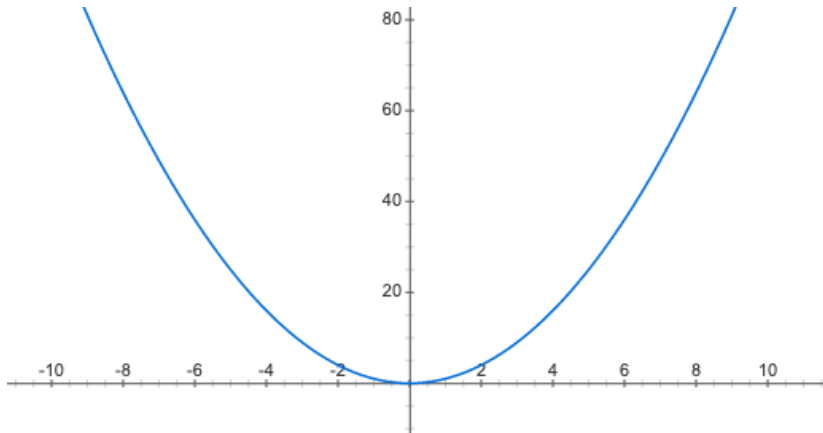
- In general, the formula for loss function (for regression problems) is the following one:

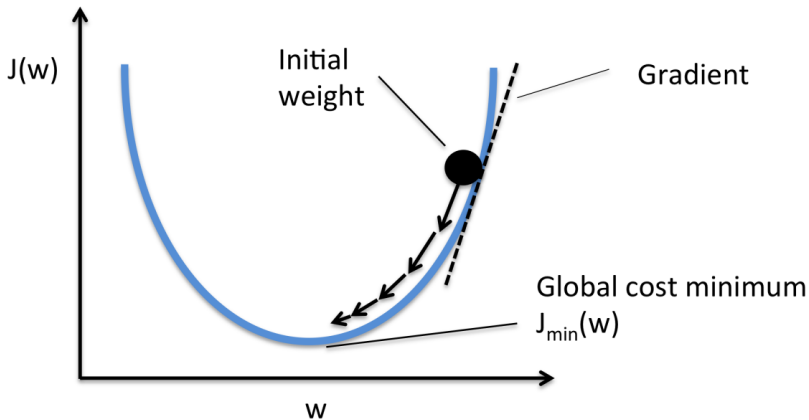
$$loss = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n \left( y_i - (b + wx_i) \right)^2}.$$

- Our goal is to find  $b$  and  $w$  such that the loss is minimum.
- How??? Choosing  $b$  and  $w$  randomly??? We will use **Gradient Descent** algorithm.



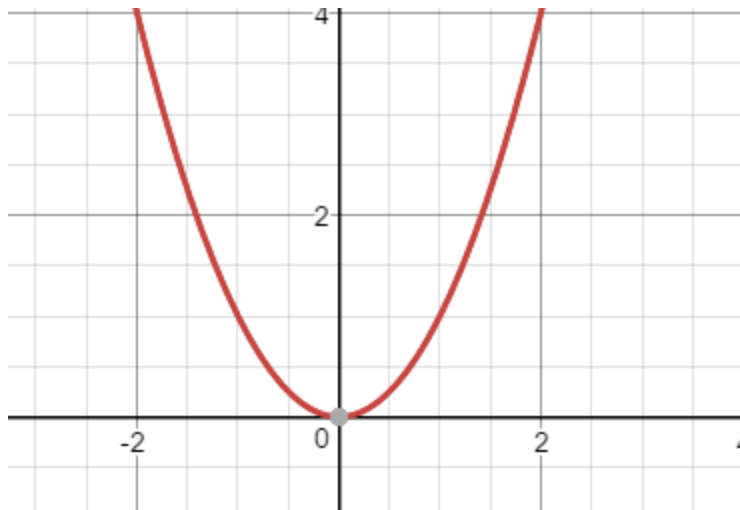








Given  $f(w) = w^2$  and  $w_0 = 1$ , we want to apply GD algorithm to obtain the minimum. We expect to obtain, after some iterations,  $w^* = 0$  (or close to 0).



Remember that:

- $f(w) = w^2$ .
- $derivative = f'(w) = 2w$ .
- $w_1 = w_0 - learning\_rate \cdot derivative$ .

We choose as learning rate a value of 0.01.

iteration	$w_0$	$f(w_0)$	derivative	$w_1$	$f(w_1)$	$ f(w_1) - f(w_0) $
2	1.00	1.00	2.00	0.80	0.64	0.36
3	0.80	0.64	1.60	0.64	0.41	0.23
4	0.64	0.41	1.28	0.51	0.26	0.15
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
26	0.00	0.00	0.01	0.00	0.00	0.00
27	0.00	0.00	0.01	0.00	0.00	0.00





- GD is an algorithm that helps us find the optimal values for the parameters. That is why it is called *optimiser*.
- There are many *optimisers* algorithms:
  - Momentum
  - RMSprop
  - Adam
  - AdaMax
  - Adadelta
  - ...

### 3 Ingredients for Deep Learning

- Loss function
- Gradient Descent

- 4 Deep Learning
  - Introduction
  - Perceptron
  - Backpropagation

## 4 Deep Learning

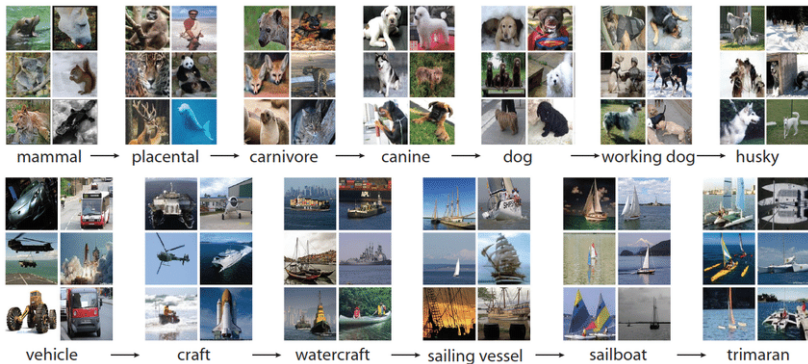
# Introduction

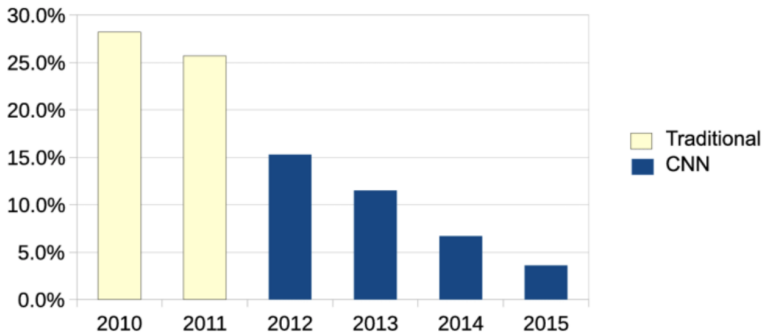
## Why do we need *another* kind of models?

- Unfortunately real life is not linear. We need more *complex/flexible* models.
- Neural networks (Deep Learning models) are very flexible models. They are able to capture very non-linear patterns and model with a high precision.

## ImageNet problem

- The ImageNet project is a large visual database designed for use in visual object recognition software research.
- More than 14 million images have been hand-annotated.
- ImageNet contains more than 20.000 categories with a typical category, such as "balloon" or "strawberry", consisting of several hundred images.

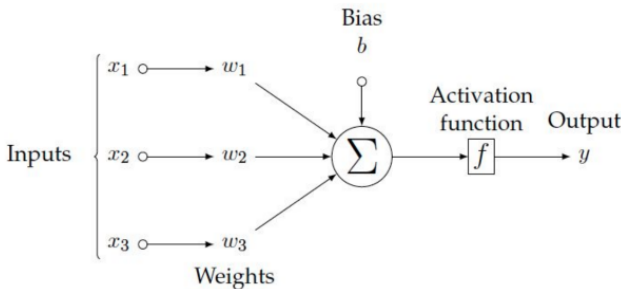






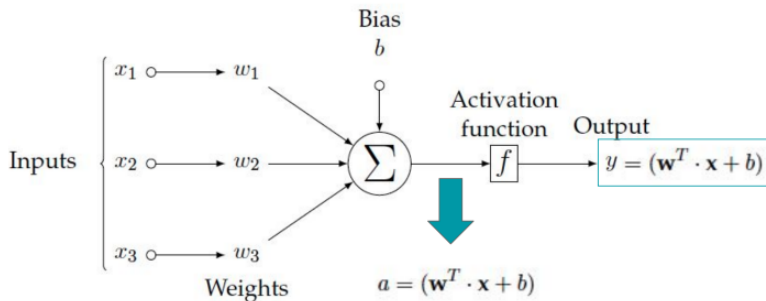


# Single neuron model (perceptron)



- **Weights** and **bias** are the parameters that define the behaviour. They must be estimated during training.
- The **output** ( $y$ ) is derived from a sum of the weighted inputs plus a bias term.
- The **activation** function **introduces non-linearities**.

# Single neuron model: Linear Regression





# Backpropagation

Remember that GD formula is:

$$w = w - \text{learning\_rate} \cdot \text{derivative}.$$

We are going to use the following notation:

$$\text{derivative} = \frac{\partial L}{\partial w}.$$

Our goal is to find a way to calculate the derivative in an efficient way.

# Thank You