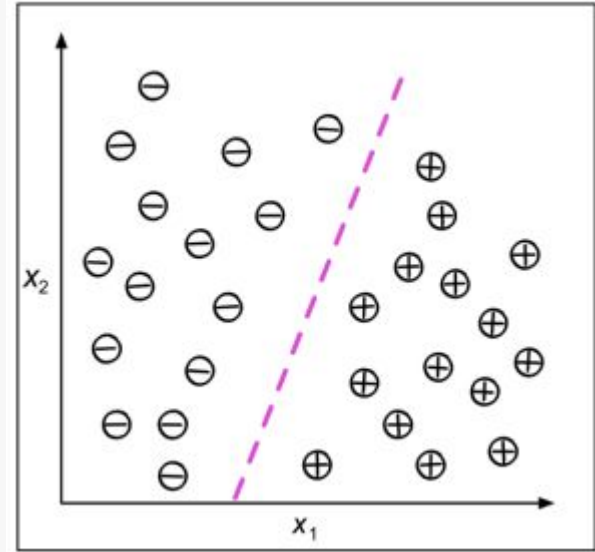# Introduction to Logistic regression
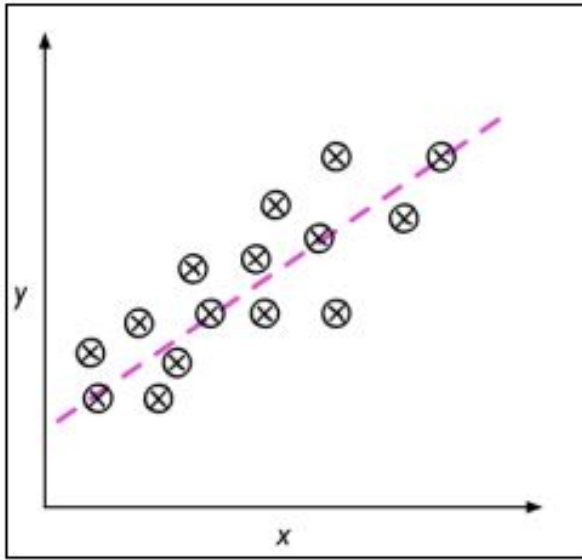
Don't let the name confuse you.  For historical reasons it is called 'regression'.
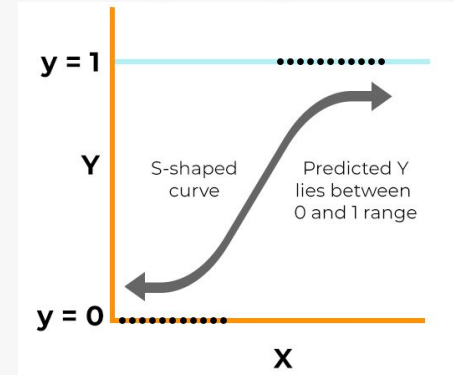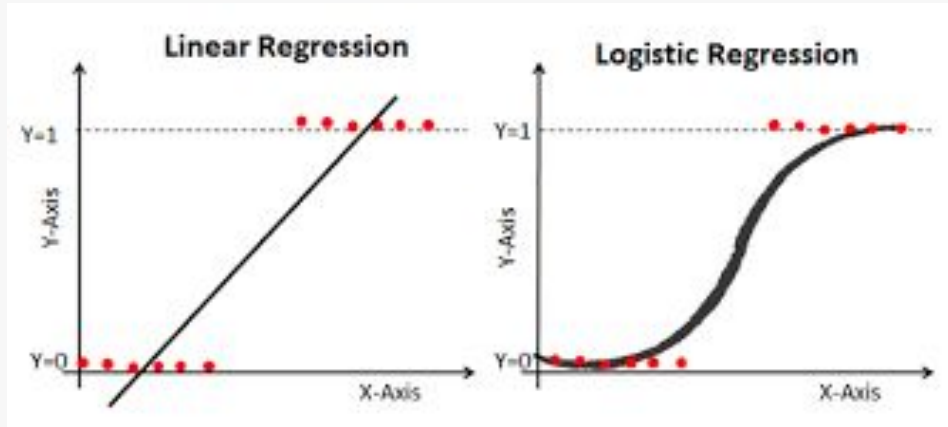
Logistic regression is a _CLASSIFICATION_ algorithm!

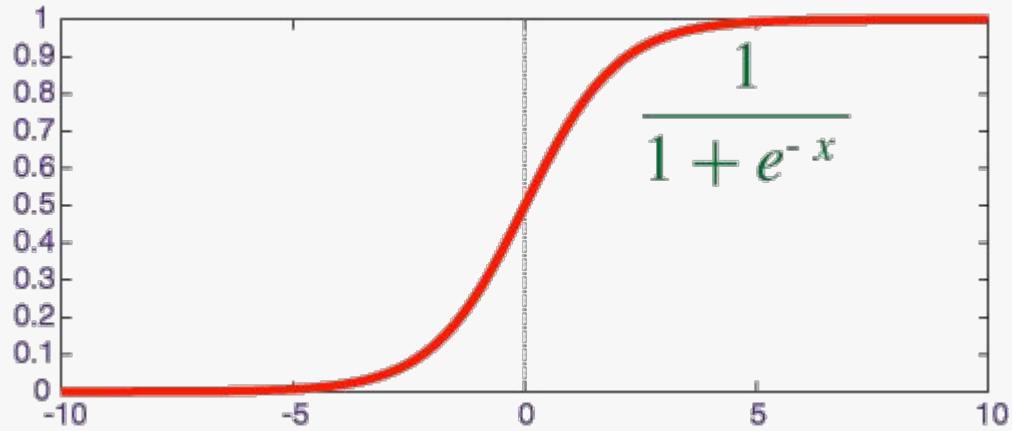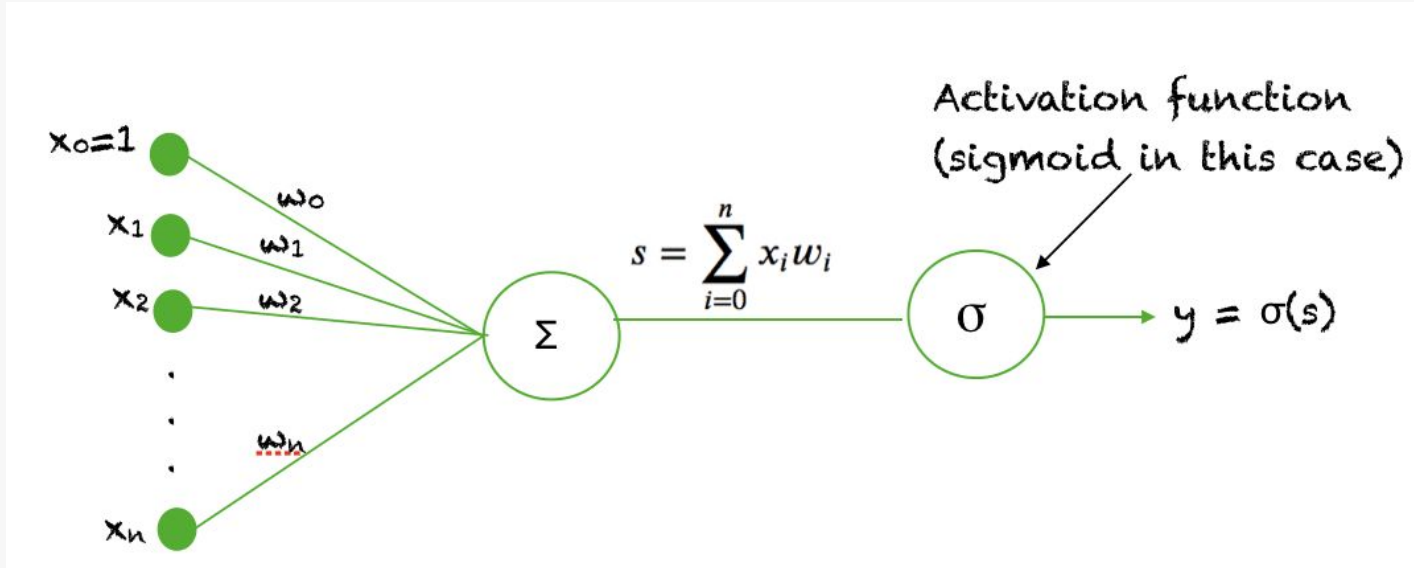# Regression vs Classification

# From straight lines to curves



## Linear Regression

Y=1

Y-Axis

Y=0

X-Axis

## Logistic Regression

Y=1

Y-Axis

Y=0

X-Axis

y = 1

Y

S-shaped curve

Predicted Y lies between 0 and 1 range

y = 0

X

# Sigmoid Function



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

# Turning number predictions into class predictions

# Cost function for logistic regression
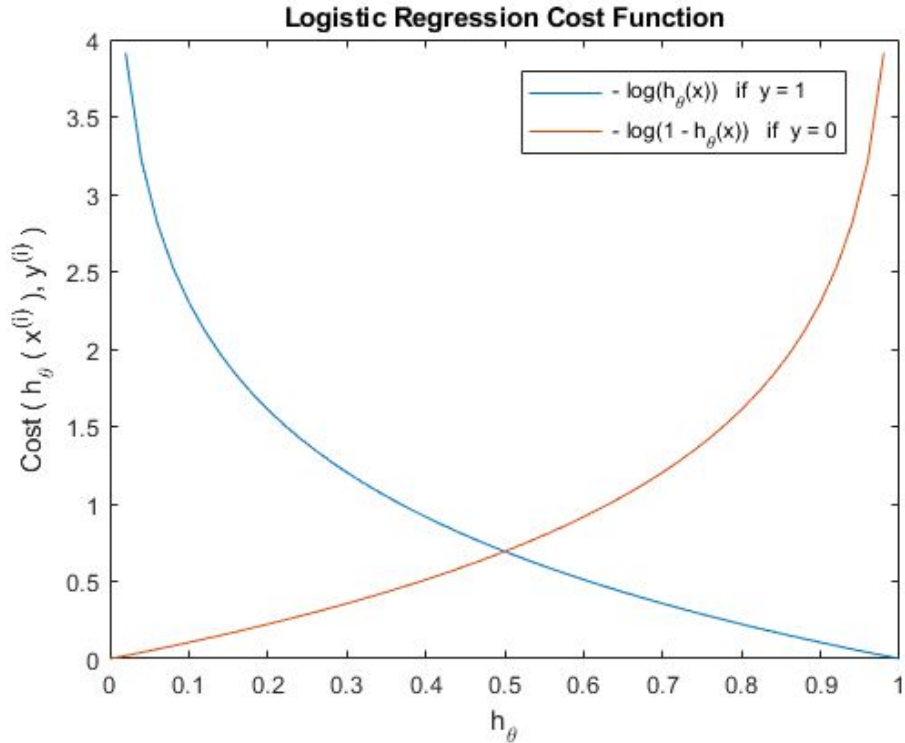
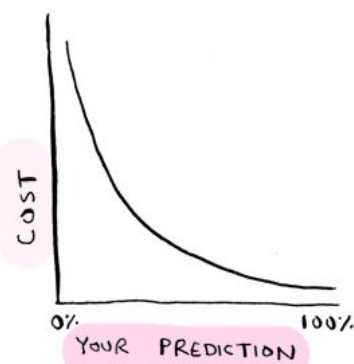$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} Cost(h_\theta(x^{(i)}), y^{(i)})$$

$$J(\theta) = \frac{1}{m} [\sum_{i=1}^{m} -y^{(i)} log(h_\theta(x^{(i)})) + (1 - y^{(i)})log(1 - h_\theta(x^{(i)}))]$$
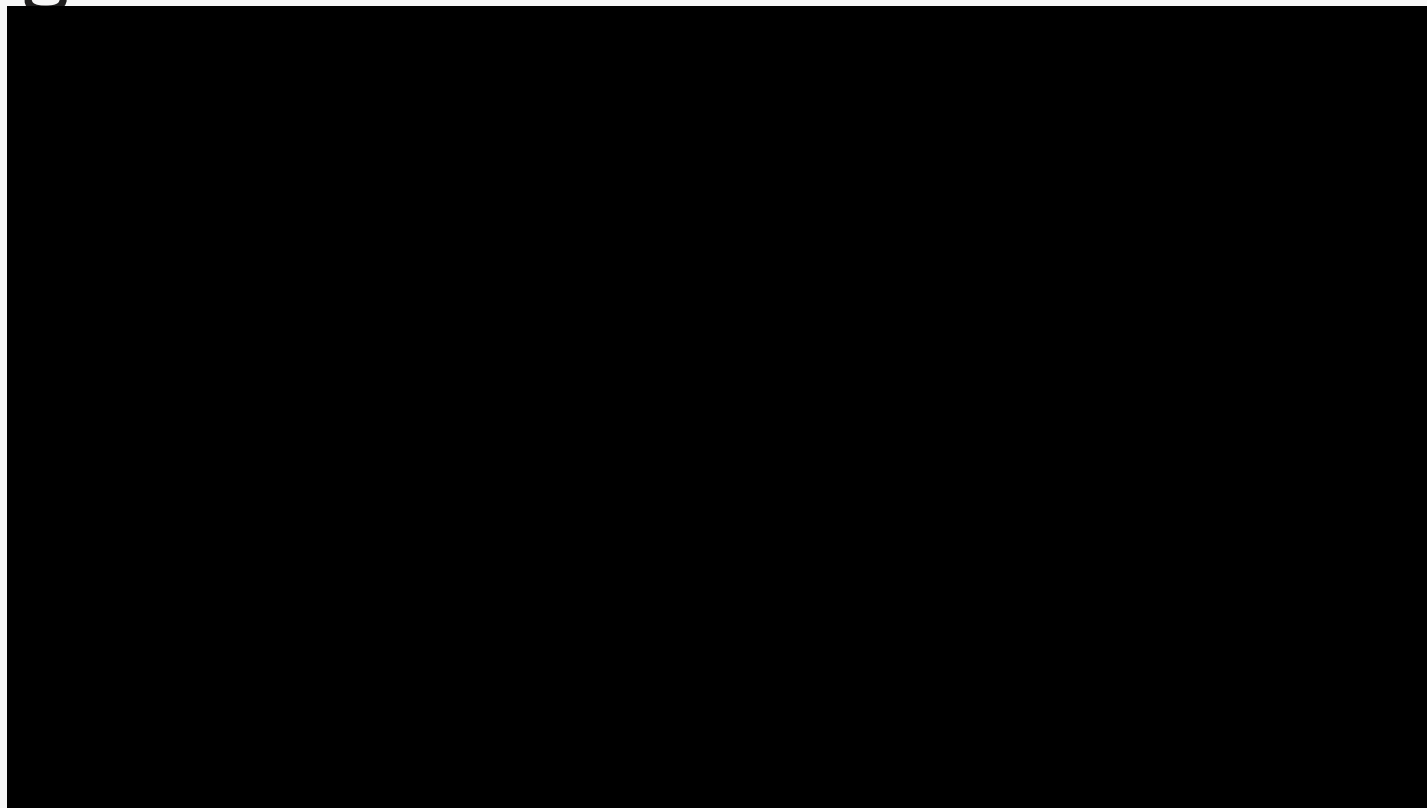
$$m = number\ of\ samples$$

# Why logarithm?



Logistic Regression Cost Function

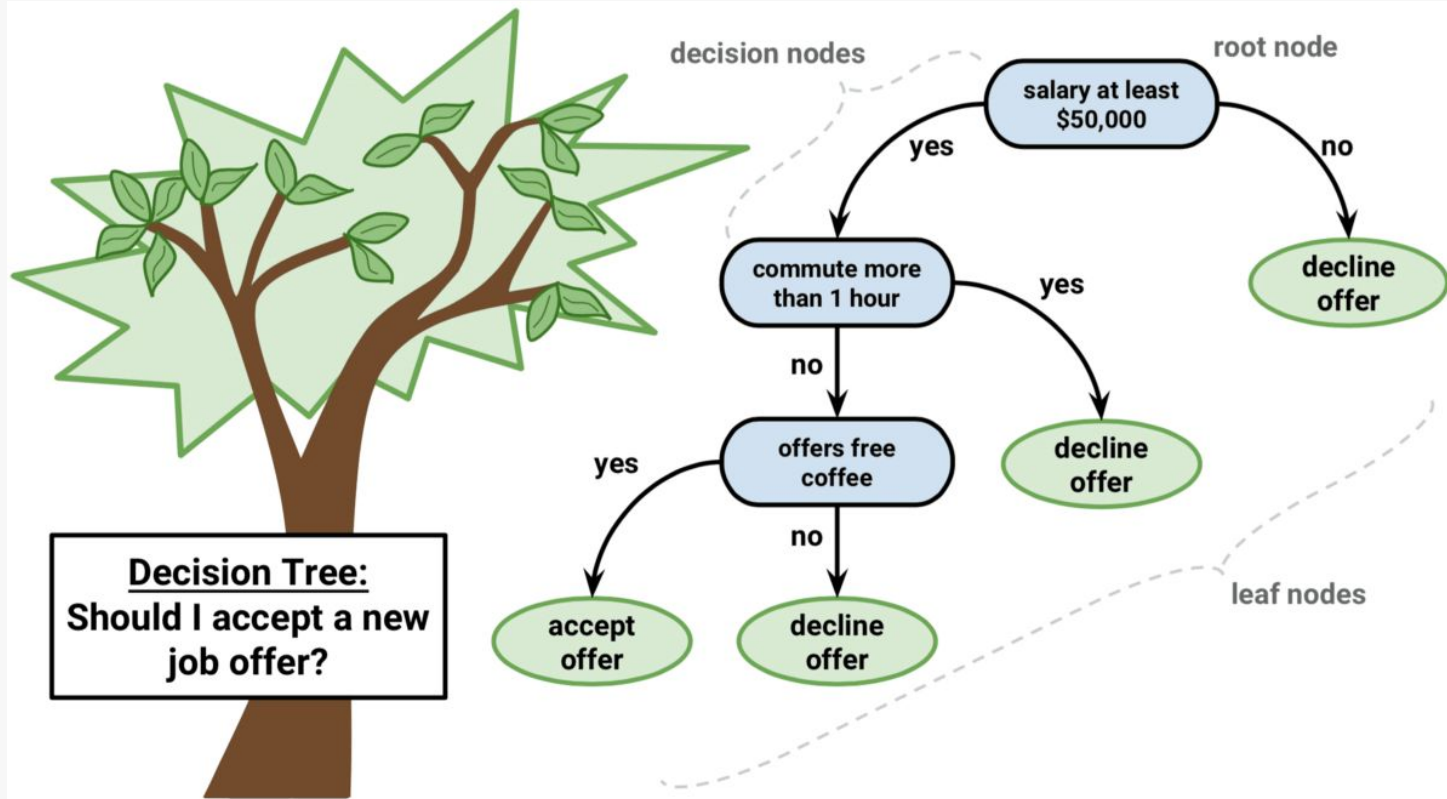$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} Cost\big(h_\theta(x^{(i)}) - y^{(i)}\big)$$

$$Cost(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1 - h_\theta(x)) & \text{if } y=0 \end{cases}$$

# Log loss calculation
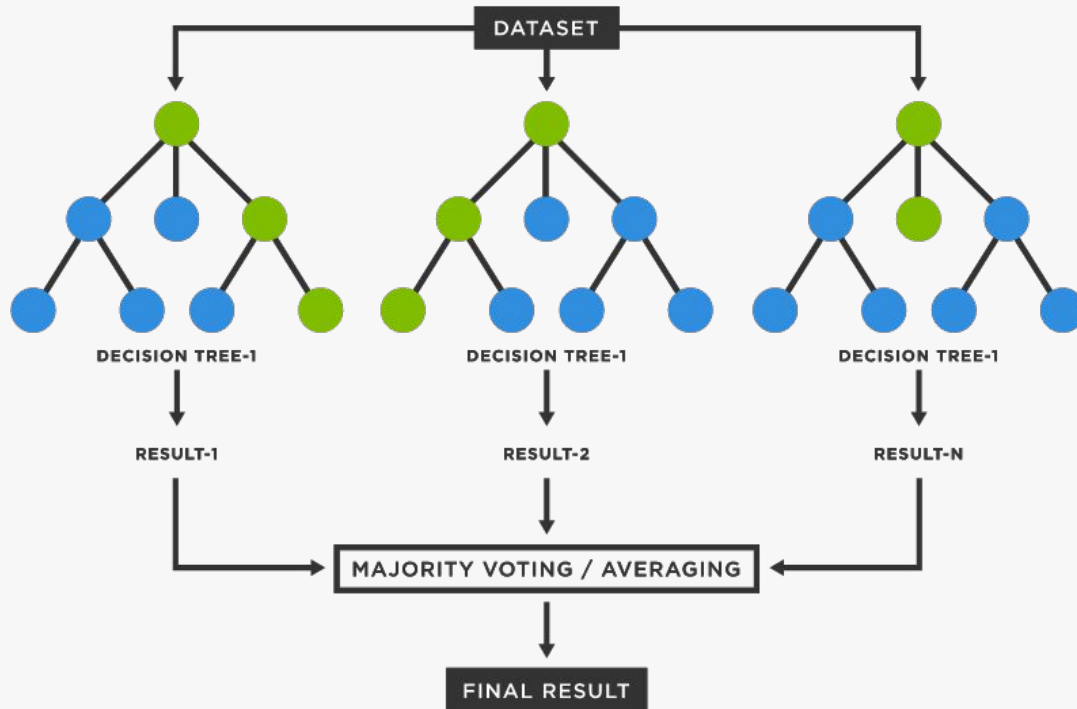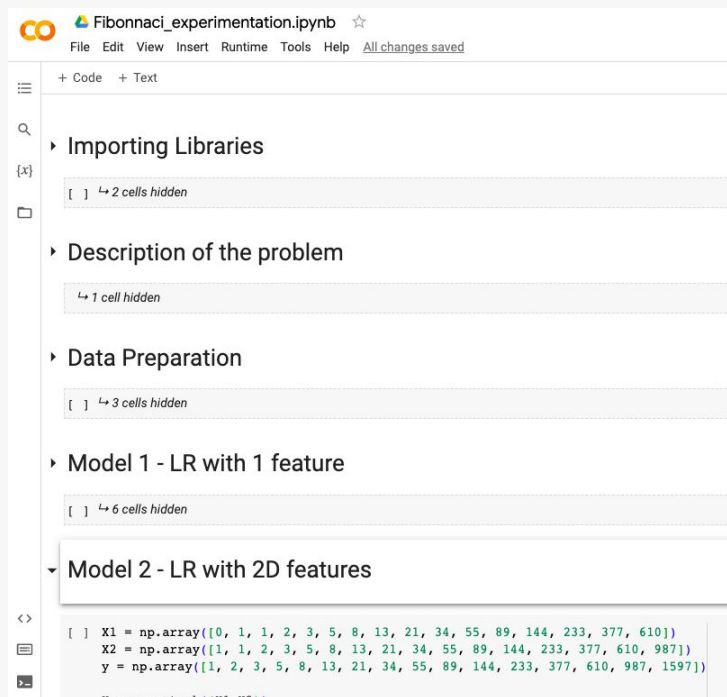
# Non-Parametric models

# Ensemble models

# Module 2

# Data Science process

Data Engineering Vs Data Science : The number and frequency of experiments

Organizing your experiments
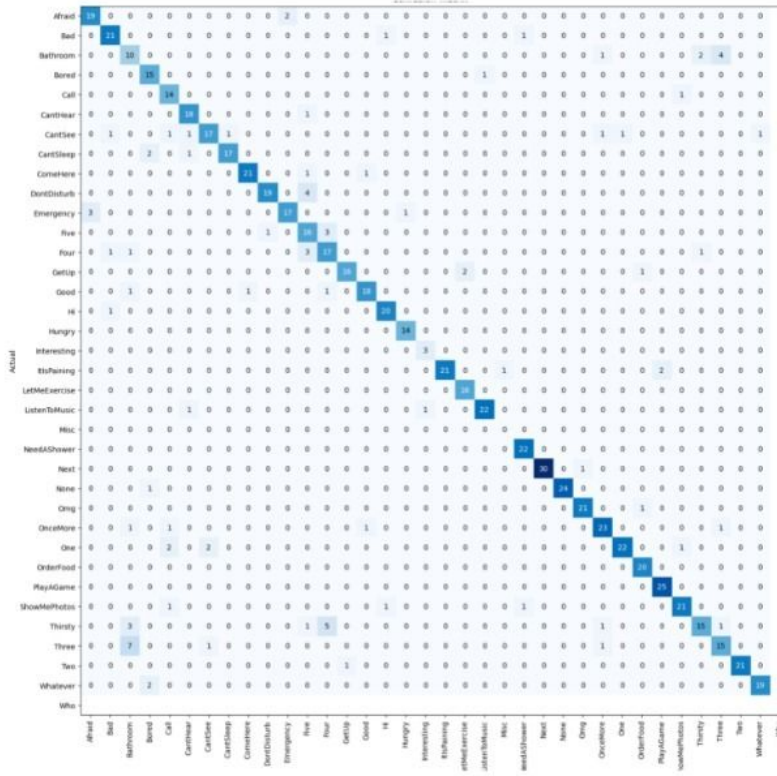
# Introducing ML communities

- Kaggle

- Huggingface  🤗

# How is my model doing?

Metrics notebook demo

# Prediction errors - Confusion Matrix

# Metrics - Measuring model performance

# Going beyond *accuracy*

Accuracy is a limited metric, because models can cheat!!

| Machine learing \ Manual counting | True | False |
|---|---|---|
| True | True Positive (TP) | False Positive (FP) |
| False | False Negative (FN) | True Negative (TN) |

**Equations:**

False positive rate (FPR) = $\dfrac{FP}{FP+TN}$

False negative rate (FNR) = $\dfrac{FN}{FN+TP}$

Sensitivity = $\dfrac{TP}{TP+FN}$

Specificity = $\dfrac{TN}{TN+FP}$

Youden index = Sensitivity + Specificity - 1

Accuracy = $\dfrac{TP+TN}{TP+TN+FP+FN}$
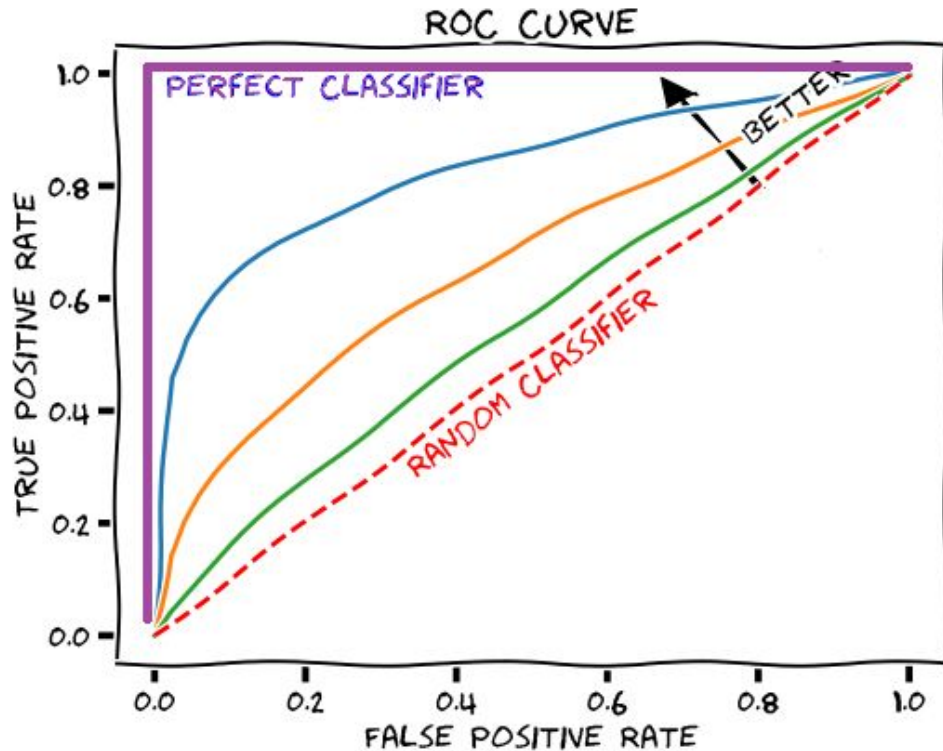
# Going beyond *accuracy*

$$\text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

$$\text{Specificity} = \frac{\text{True Negatives}}{\text{True Negatives} + \text{False Positives}}$$

# AUROC curve



Receiver operating characteristic curve.

# Introduction to Data Science!

# How should we store features?

| ID | First Name | Last Name | Email | Year of Birth |
|---|---|---|---|---|
| 1 | Peter | Lee | plee@university.edu | 1992 |
| 2 | Jonathan | Edwards | jedwards@university.edu | 1994 |
| 3 | Marilyn | Johnson | mjohnson@university.edu | 1993 |
| 6 | Joe | Kim | jkim@university.edu | 1992 |
| 12 | Haley | Martinez | hmartinez@university.edu | 1993 |
| 14 | John | Mfume | jmfume@university.edu | 1991 |
| 15 | David | Letty | dletty@university.edu | 1995 |

**Table: Students**

Feature Vector

*Arrays or vectors!*

# Vectors and Tensors!
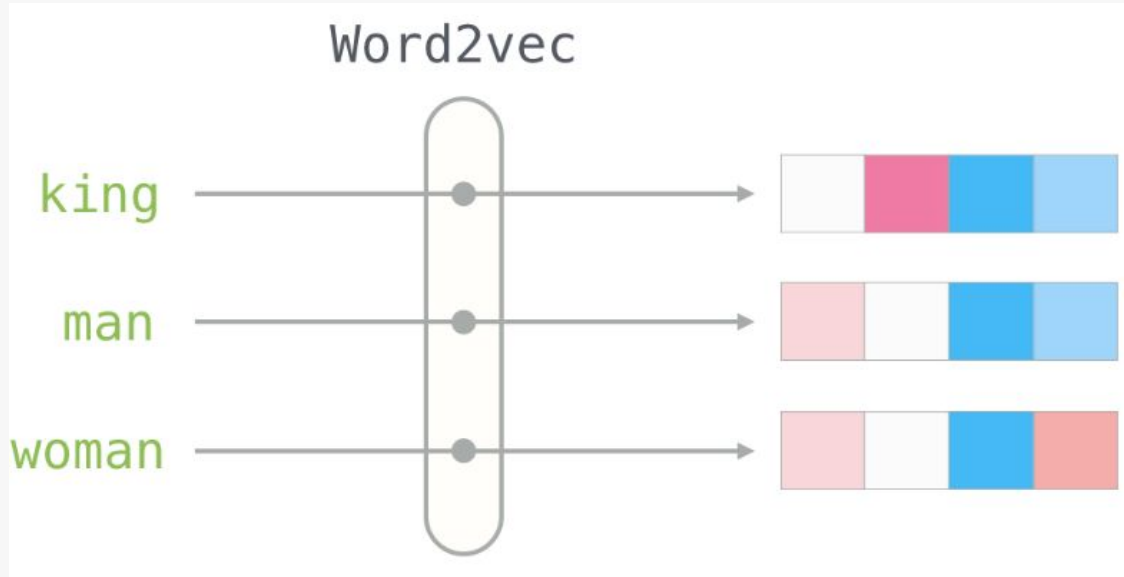


Scalar — Vector — Matrix — Tensor

Scalar: 1

Vector: $\begin{bmatrix} 1 \\ 2 \end{bmatrix}$

Matrix: $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$

Tensor: $\begin{bmatrix} \begin{bmatrix} 1 & 2 \\ 1 & 7 \end{bmatrix} & \begin{bmatrix} 3 & 2 \\ 5 & 4 \end{bmatrix} \end{bmatrix}$

Source:  Hadrien Jean's blog post
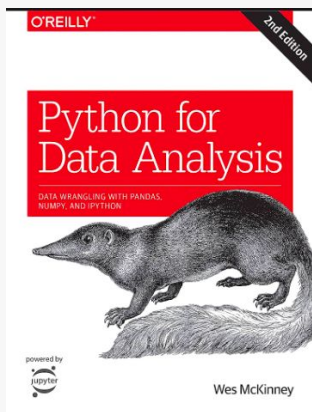
# Vectors Continued…

Source: Jay Alammar's blog post

# Open Source Heroes

Pandas: Wes Mckinney

Python: Guido Van Rossum

Jupyter: Fernando Perez

NumPy: Travis Olyphant