

---

# Práctica 1

Inteligencia de Negocio

Amin Kasrou Aouam



**UNIVERSIDAD  
DE GRANADA**

2020-11-10

## Índice

**Experiments**

**3**

## Experiments

We will first try to gather information about our dataset, by evaluating the statistics of our attributes.

```
1 from pandas import read_csv
2 from sklearn.preprocessing import LabelEncoder
3
4
5 def replace_values(df):
6     columns = ["BI-RADS", "Margin", "Density", "Age"]
7     for column in columns:
8         df[column].fillna(value=df[column].mean(), inplace=True)
9     return df
10
11
12 def process_na(df, action):
13     if action == "drop":
14         return df.dropna()
15     return replace_values(df)
16
17
18 def encode_columns(df):
19     encoder = LabelEncoder()
20     encoder.fit(df["Shape"])
21
22
23 def parse_data(source, action):
24     df = read_csv(filepath_or_buffer=source, na_values="?")
25     processed_df = process_na(df, action)
26     return processed_df
```

```
1 df = parse_data("../data/mamografia.csv", "drop")
2 print(df.describe())
```

	BI-RADS	Age	Margin	Density
count	847.000000	847.000000	847.000000	847.000000
mean	4.322314	55.842975	2.833530	2.909091
std	0.703762	14.603754	1.564049	0.370292
min	0.000000	18.000000	1.000000	1.000000
25%	4.000000	46.000000	1.000000	3.000000
50%	4.000000	57.000000	3.000000	3.000000
75%	5.000000	66.000000	4.000000	3.000000
max	6.000000	96.000000	5.000000	4.000000

We observe that **margin** and **density** are the columns with the most unknown values. The age group of our cohort is middle aged, the BI-RADS score is mostly in the suspicious category, the density is mostly low and the margin belongs to the microlobulated/obscured category.

We'll try to impute values, instead of dropping them, when they're invalid.

```
1 df = parse_data("../data/mamografia.csv", "replace")
2 print(df.describe())
```

	BI-RADS	Age	Margin	Density
1				
2	count	961.000000	961.000000	961.000000
3	mean	4.296142	55.487448	2.796276
4	std	0.705555	14.442373	1.526880
5	min	0.000000	18.000000	1.000000
6	25%	4.000000	45.000000	1.000000
7	50%	4.000000	57.000000	3.000000
8	75%	5.000000	66.000000	3.000000
9	max	6.000000	96.000000	4.000000