# SIMPLE IMPUTER TECHNIQUE

NAME : G.AKHILA

MISSING VALUE TREATMENT

# SIMPLE IMPUTER TECHNIQUES (Numerical Data)

```
            ┌─────────────────────────────┐
            │   SIMPLE IMPUTER            │
            │   TECHNIQUES (Categorical)  │
            └─────────────────────────────┘
```

| LABEL ENCODER | ONEHOT ENCODER | DUMMY VARIABLE |
|---|---|---|

G.AKHILA
SIMPLE IMPUTER
EDA

# SIMPLE IMPUTER

Q. Why do we use simple imputer ?

A. When we have missing values in either categorical data or numerical data we will use simple imputer to impute missing values in the dataset
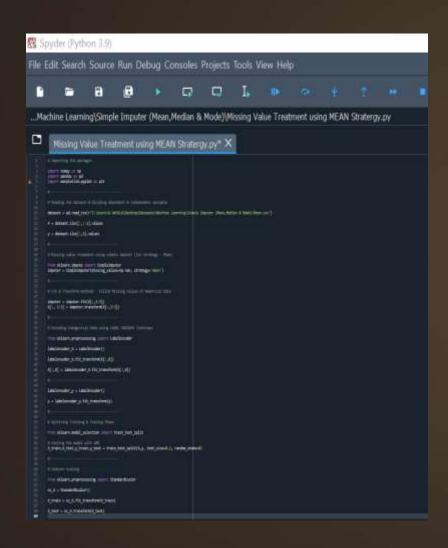
# DATASET FOR PRACTICE

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Country | Age | Salary | Purchased |
| 2 | France | 44 | 72000 | No |
| 3 | Spain | 27 | 48000 | Yes |
| 4 | Germany | 30 | 54000 | No |
| 5 | Spain | 38 | 61000 | No |
| 6 | Germany | 40 | | Yes |
| 7 | France | 45 | 58000 | Yes |
| 8 | Spain | | 52000 | No |
| 9 | France | 48 | 79000 | Yes |
| 10 | Germany | 50 | 83000 | No |
| 11 | France | 37 | 67000 | Yes |

**Missing**

**Missing**

G.AKHILA
SIMPLE IMPUTER
EDA

# PACKAGES WE NEED

➢ numpy

➢ pandas

➢ sklearn (model_selection, preprocessing)

G.AKHILA
SIMPLE IMPUTER
EDA

# CODE FOR IMPUTING MISSING VALUES USING MEAN TECHNIQUE



# MEAN
# TECHNIQUE CODE

G.AKHILA

SIMPLE IMPUTER

EDA

# DATASET (BEFORE & AFTER IMPUTING MEAN TECHNIQUE & LABEL ENCODER)



**BEFORE**

| Index | Country | Age | Salary | Purchased |
|-------|---------|-----|--------|-----------|
| 0 | France | 44 | 72000 | No |
| 1 | Spain | 27 | 48000 | Yes |
| 2 | Germany | 30 | 54000 | No |
| 3 | Spain | 38 | 61000 | No |
| 4 | Germany | 40 | nan | Yes |
| 5 | France | 45 | 58000 | Yes |
| 6 | Spain | nan | 52000 | No |
| 7 | France | 48 | 79000 | Yes |
| 8 | Germany | 50 | 83000 | No |
| 9 | France | 37 | 67000 | Yes |

**AFTER**

X - NumPy object array (read only)

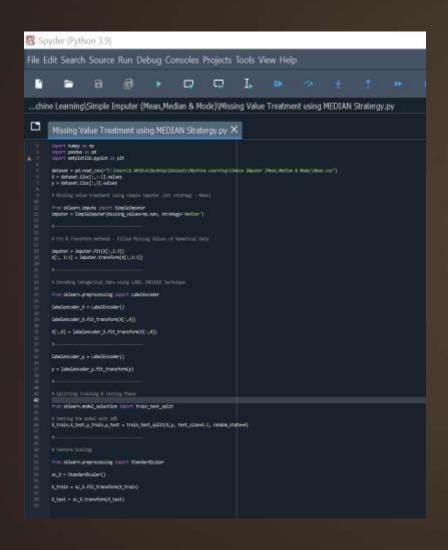| | 0 | 1 | 2 |
|---|---|---|---|
| 0 | 0 | 44.0 | 72000.0 |
| 1 | 2 | 27.0 | 48000.0 |
| 2 | 1 | 30.0 | 54000.0 |
| 3 | 2 | 38.0 | 61000.0 |
| 4 | 1 | 40.0 | 63777.7777777… |
| 5 | 0 | 45.0 | 58000.0 |
| 6 | 2 | 39.8888888888… | 52000.0 |
| 7 | 0 | 48.0 | 79000.0 |
| 8 | 1 | 50.0 | 83000.0 |
| 9 | 0 | 37.0 | 67000.0 |

# MEAN TECHNIQUE

## G.AKHILA
## SIMPLE IMPUTER
## EDA

# CODE FOR IMPUTING MISSING VALUES USING MEDIAN TECHNIQUE



# MEDIAN
# TECHNIQUE CODE

G.AKHILA
SIMPLE IMPUTER
EDA

# DATASET (BEFORE & AFTER IMPUTING MEAN TECHNIQUE & LABEL ENCODER)

**dataset - DataFrame**

| Index | Country | Age | Salary | Purchased |
|-------|---------|-----|--------|-----------|
| 0 | France | 44 | 72000 | No |
| 1 | Spain | 27 | 48000 | Yes |
| 2 | Germany | 30 | 54000 | No |
| 3 | Spain | 38 | 61000 | No |
| 4 | Germany | 40 | nan | Yes |
| 5 | France | 45 | 58000 | Yes |
| 6 | Spain | nan | 52000 | No |
| 7 | France | 48 | 79000 | Yes |
| 8 | Germany | 50 | 83000 | No |
| 9 | France | 37 | 67000 | Yes |

**BEFORE**

**X - NumPy object array (read only)**

| | 0 | 1 | 2 |
|---|---|---|---|
| 0 | 0 | 44.0 | 72000.0 |
| 1 | 2 | 27.0 | 48000.0 |
| 2 | 1 | 30.0 | 54000.0 |
| 3 | 2 | 38.0 | 61000.0 |
| 4 | 1 | 40.0 | 61000.0 |
| 5 | 0 | 45.0 | 58000.0 |
| 6 | 2 | 40.0 | 52000.0 |
| 7 | 0 | 48.0 | 79000.0 |
| 8 | 1 | 50.0 | 83000.0 |
| 9 | 0 | 37.0 | 67000.0 |

**AFTER**

# MEDIAN TECHNIQUE

G.AKHILA
SIMPLE IMPUTER
EDA

# MODE
# TECHNIQUE CODE

G.AKHILA
SIMPLE IMPUTER
EDA

# DATASET (BEFORE & AFTER IMPUTING MEAN TECHNIQUE & LABEL ENCODER)



BEFORE

AFTER

## MODE TECHNIQUE

G.AKHILA
SIMPLE IMPUTER
EDA