

```
[1]: import numpy as np
import pandas as pd
df=pd.read_csv(r"C:\Users\nivet\OneDrive\Documents\Downloads\pre_process_datasample - pre_process_datasample.csv")
df
```

```
[1]:   Country  Age  Salary Purchased
  0 France  44.0  72000.0     No
  1 Spain  27.0  48000.0    Yes
  2 Germany  30.0  54000.0     No
  3 Spain  38.0  61000.0     No
  4 Germany  40.0      NaN    Yes
  5 France  35.0  58000.0    Yes
  6 Spain  NaN  52000.0     No
  7 France  48.0  79000.0    Yes
  8 Germany  50.0  83000.0     No
  9 France  37.0  67000.0    Yes
```

```
[2]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 4 columns):
 #   Column   Non-Null Count  Dtype  
 --- 
  0   Country   10 non-null    object 
  1   Age        9 non-null    float64 
  2   Salary     9 non-null    float64 
  3   Purchased  10 non-null   object 
dtypes: float64(2), object(2)
memory usage: 452.0+ bytes
```

```
[3]: df.Country.mode()
```

```
[3]: 0    France
Name: Country, dtype: object
```

```
[4]: df.Country.mode()[0]
type(df.Country.mode())
```

```
[4]: pandas.core.series.Series
```

```
[5]: df.Country.fillna(df.Country.mode()[0],inplace=True)
df.Age.fillna(df.Age.median(),inplace=True)
df.Salary.fillna(round(df.Salary.mean()),inplace=True)
df
```

```
[5]:   Country  Age  Salary Purchased
  0 France  44.0  72000.0     No
  1 Spain  27.0  48000.0    Yes
  2 Germany  30.0  54000.0     No
  3 Spain  38.0  61000.0     No
  4 Germany  40.0  63778.0    Yes
  5 France  35.0  58000.0    Yes
  6 Spain  38.0  52000.0     No
  7 France  48.0  79000.0    Yes
```

```
[6]: pd.get_dummies(df.Country)
```

```
[6]: France Germany Spain
```

```
0 True False False
1 False False True
2 False True False
3 False False True
4 False True False
5 True False False
6 False False True
7 True False False
8 False True False
9 True False False
```

```
[7]: updated_dataset=pd.concat([pd.get_dummies(df.Country),df.iloc[:,[1,2,3]]],axis=1)
updated_dataset
```

```
[7]: France Germany Spain Age Salary Purchased
```

```
0 True False False 44.0 72000.0 No
1 False False True 27.0 48000.0 Yes
2 False True False 30.0 54000.0 No
3 False False True 38.0 61000.0 No
4 False True False 40.0 63778.0 Yes
5 True False False 35.0 58000.0 Yes
6 False False True 38.0 52000.0 No
7 True False False 48.0 79000.0 Yes
8 False True False 50.0 83000.0 No
9 True False False 37.0 67000.0 Yes
```

```
[8]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype  
 --- 
 0   Country     10 non-null    object 
 1   Age         10 non-null    float64 
 2   Salary       10 non-null    float64 
 3   Purchased    10 non-null    object  
dtypes: float64(2), object(2)
memory usage: 452.0+ bytes
```

```
[9]: updated_dataset.Purchased.replace(['No','Yes'],[0,1],inplace=True)
```

```
C:\Users\nivet\AppData\Local\Temp\ipykernel_18388\725871168.py:1: FutureWarning: Downcasting behavior in `replace` is deprecated and will be removed in a future version. To retain the old behavior, explicitly call `result.infer_objects(copy=False)`. To opt-in to the future behavior, set `pd.set_option('future.no_silent_downcasting', True)`
  updated_dataset.Purchased.replace(['No','Yes'],[0,1],inplace=True)
```

```
[10]: updated_dataset
```

[18]:

	France	Germany	Spain	Age	Salary	Purchased
0	True	False	False	44.0	72000.0	0
1	False	False	True	27.0	48000.0	1
2	False	True	False	30.0	54000.0	0
3	False	False	True	38.0	61000.0	0
4	False	True	False	40.0	63778.0	1
5	True	False	False	35.0	58000.0	1
6	False	False	True	38.0	52000.0	0
7	True	False	False	48.0	79000.0	1
8	False	True	False	50.0	83000.0	0
9	True	False	False	37.0	67000.0	1

[19]: