

Data Preprocessing

Loading dataset.

```
In [1]: import numpy as np
import pandas as pd
```

```
In [2]: df = pd.read_csv('Data.csv')
```

```
df
Out[2]:
```

	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

```
In [3]: df.shape
```

```
Out[3]: (10, 4)
```

```
In [4]: df.dtypes
```

```
Out[4]: Country      object
Age      float64
Salary   float64
Purchased object
dtype: object
```

Observations :

- There are 2 categorical data type and 2 numerical data type.

```
In [5]: df.isnull().sum()
```

```
Out[5]: Country      0
Age          1
Salary        1
Purchased      0
dtype: int64
```

Observations :

- There are some null values in both the numerical features. SO, we can handle them by replacing them with mean.

```
In [6]: df.head()
```

```
Out[6]:
```

	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

Taking care of missing values.

```
In [7]: # handling independent numerical features
from sklearn.impute import SimpleImputer

numerical_features = [feature for feature in df.columns if df[feature].dtype != 'O']

imputer = SimpleImputer(missing_values = np.nan, strategy = 'mean')
df[numerical_features] = imputer.fit_transform(df[numerical_features])

df
```

```
Out[7]:
```

	Country	Age	Salary	Purchased
0	France	44.000000	72000.000000	No
1	Spain	27.000000	48000.000000	Yes
2	Germany	30.000000	54000.000000	No
3	Spain	38.000000	61000.000000	No
4	Germany	40.000000	63777.777778	Yes
5	France	35.000000	58000.000000	Yes
6	Spain	38.777778	52000.000000	No
7	France	48.000000	79000.000000	Yes
8	Germany	50.000000	83000.000000	No
9	France	37.000000	67000.000000	Yes

```
In [18]: df.isnull().sum()
```

```
Out[18]: Age          0
Salary        0
Purchased      0
dtype: int64
```

Observations :

- There is no missing values.

Taking care of categorical features.

```
In [8]: # handling independent categorical features
from sklearn.preprocessing import OneHotEncoder

encoder = OneHotEncoder()
encoded = encoder.fit_transform(df[['Country']]).toarray()
encoded_df = pd.DataFrame(encoded, columns = encoder.get_feature_names_out())

df.drop(columns = ['Country'], inplace = True)
new_df = pd.concat([df, encoded_df], axis = 1)

new_df
```

```
Out[8]:
```

	Age	Salary	Purchased	Country_France	Country_Germany	Country_Spain
0	44.000000	72000.000000	No	1.0	0.0	0.0
1	27.000000	48000.000000	Yes	0.0	0.0	1.0
2	30.000000	54000.000000	No	0.0	1.0	0.0
3	38.000000	61000.000000	No	0.0	0.0	1.0
4	40.000000	63777.777778	Yes	0.0	1.0	0.0
5	35.000000	58000.000000	Yes	1.0	0.0	0.0
6	38.777778	52000.000000	No	0.0	0.0	1.0
7	48.000000	79000.000000	Yes	1.0	0.0	0.0
8	50.000000	83000.000000	No	0.0	1.0	0.0
9	37.000000	67000.000000	Yes	1.0	0.0	0.0

```
In [9]: # handling dependent categorical features
from sklearn.preprocessing import LabelEncoder
```

```
encoder = LabelEncoder()
new_df['Purchased'] = encoder.fit_transform(new_df['Purchased'])

new_df
```

```
Out[9]:
```

	Age	Salary	Purchased	Country_France	Country_Germany	Country_Spain
0	44.000000	72000.000000	0	1.0	0.0	0.0
1	27.000000	48000.000000	1	0.0	0.0	1.0
2	30.000000	54000.000000	0	0.0	1.0	0.0
3	38.000000	61000.000000	0	0.0	0.0	1.0
4	40.000000	63777.777778	1	0.0	1.0	0.0
5	35.000000	58000.000000	1	1.0	0.0	0.0
6	38.777778	52000.000000	0	0.0	0.0	1.0
7	48.000000	79000.000000	1	1.0	0.0	0.0
8	50.000000	83000.000000	0	0.0	1.0	0.0
9	37.000000	67000.000000	1	1.0	0.0	0.0

Splitting dataset into independent and dependent features

```
In [10]: from sklearn.model_selection import train_test_split
```

```
x = new_df.drop(columns = ['Purchased'])
y = new_df['Purchased']

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state = 42)
```

```
In [11]: x_train
```

```
Out[11]:
```

	Age	Salary	Country_France	Country_Germany	Country_Spain
5	35.000000	58000.000000	1.0	0.0	0.0
0	44.000000	72000.000000	1.0	0.0	0.0
7	48.000000	79000.000000	1.0	0.0	0.0
2	30.000000	54000.000000	0.0	1.0	0.0
9	37.000000	67000.000000	1.0	0.0	0.0
4	40.000000	63777.777778	0.0	1.0	0.0
3	38.000000	61000.000000	0.0	0.0	1.0
6	38.777778	52000.000000	0.0	0.0	1.0

```
In [12]: y_train
```

```
Out[12]: 5      1
0      0
7      1
2      0
9      1
4      1
3      0
6      0
Name: Purchased, dtype: int32
```

```
In [13]: x_test
```

```
Out[13]:
```

	Age	Salary	Country_France	Country_Germany	Country_Spain
8	50.0	83000.0	0.0	1.0	0.0
1	27.0	48000.0	0.0	0.0	1.0

```
In [14]: y_test
```

```
Out[14]: 8      0
1      1
Name: Purchased, dtype: int32
```

Standardization

```
In [15]: from sklearn.preprocessing import StandardScaler
```

```
scaler = StandardScaler()
x_train = scaler.fit_transform(x_train)
x_test = scaler.transform(x_test) # to avoid data leakage we just do transform
```

```
In [16]: x_train
```

```
Out[16]: array([[ -0.7529426, -0.62603778,  1.          , -0.57735027, -0.57735027],
 [  1.00845381,  1.01304295,  1.          , -0.57735027, -0.57735027],
 [  1.79129666,  1.83258331,  1.          , -0.57735027, -0.57735027],
 [-1.73149616, -1.09434656, -1.          ,  1.73205081, -0.57735027],
 [-0.36152118,  0.42765698,  1.          , -0.57735027, -0.57735027],
 [  0.22561096,  0.05940924, -1.          ,  1.73205081, -0.57735027],
 [  0.16591046, -0.27480619, -1.          , -0.57735027,  1.73205081],
 [-0.01359102, -1.32850095, -1.          , -0.57735027,  1.73205081]])
```

```
In [17]: x_test
```

```
Out[17]: array([[  2.18271898,  2.30089209, -1.          ,  1.73205081, -0.57735027],
 [-2.3186283 , -1.79680973, -1.          , -0.57735027,  1.73205081]])
```


