```python
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
```

```python
df = pd.read_csv("Lab1.csv")
df.head()
```

|   | country | age | salary | purchased |
|---|---------|-----|--------|-----------|
| 0 | France  | NaN | 7200   | no        |
| 1 | Spain   | 27.0| 4800   | yes       |
| 2 | Germany | 30.0| 5400   | yes       |
| 3 | UK      | 49.0| 98000  | no        |

```python
df.tail()
```

|   | country | age | salary | purchased |
|---|---------|-----|--------|-----------|
| 0 | France  | NaN | 7200   | no        |
| 1 | Spain   | 27.0| 4800   | yes       |
| 2 | Germany | 30.0| 5400   | yes       |
| 3 | UK      | 49.0| 98000  | no        |

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4 entries, 0 to 3
Data columns (total 4 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   country    4 non-null      object
 1   age        3 non-null      float64
 2   salary     4 non-null      int64
 3   purchased  4 non-null      object
dtypes: float64(1), int64(1), object(2)
memory usage: 256.0+ bytes
```

```python
df.describe()
```

|       | age       | salary       |
|-------|-----------|--------------|
| count | 3.000000  | 4.000000     |
| mean  | 35.333333 | 28850.000000 |
| std   | 11.930353 | 46111.278447 |
| min   | 27.000000 | 4800.000000  |
| 25%   | 28.500000 | 5250.000000  |
| 50%   | 30.000000 | 6300.000000  |
| 75%   | 39.500000 | 29900.000000 |
| max   | 49.000000 | 98000.000000 |

```python
df.isnull().sum()
```

```
country      0
age          1
salary       0
purchased    0
dtype: int64
```

```python
df['age'].fillna(df['age'].mean(), inplace = True)
df['salary'].fillna(df['salary'].mean(), inplace=True)
```

```
df.isnull().sum()
```

```
country    0
age        0
salary     0
purchased  0
dtype: int64
```

```
from sklearn.impute import SimpleImputer
x = df.iloc[:,:-1].values
x
```

```
array([['France', 35.333333333333336, 7200],
       ['Spain', 27.0, 4800],
       ['Germany', 30.0, 5400],
       ['UK', 49.0, 98000]], dtype=object)
```

```
y = df.iloc[:,3: ].values
y
```

```
array([['no'],
       ['yes'],
       ['yes'],
       ['no']], dtype=object)
```

```
imp = SimpleImputer(missing_values =np.nan, strategy = "mean")
x[:, 1:3] = imp.fit_transform(x[:, 1:3])
x
```

```
array([['France', 35.333333333333336, 7200.0],
       ['Spain', 27.0, 4800.0],
       ['Germany', 30.0, 5400.0],
       ['UK', 49.0, 98000.0]], dtype=object)
```

```
from sklearn.preprocessing import LabelEncoder
```

```
le = LabelEncoder()
h = le.fit_transform(x[:,0])
h
```

```
array([0, 2, 1, 3])
```

```
y = le.fit_transform(y)
y
```

```
C:\Users\25LAB-2BCA\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A column-vector y was passed when a 1d array
was expected. Please change the shape of y to (n_samples, ), for example using ravel().
  return f(*args, **kwargs)
```

```
array([0, 1, 1, 0])
```

```
from sklearn.utils import column_or_1d
y = column_or_1d(y, warn = True)
y
```

```
array([0, 1, 1, 0])
```

```
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
```

```
transform = ColumnTransformer([('norm1', OneHotEncoder(), [0])], remainder ="passthrough")
x = transform.fit_transform(x)
x
```

```
array([[1.0, 0.0, 0.0, 0.0, 35.333333333333336, 7200.0],
       [0.0, 0.0, 1.0, 0.0, 27.0, 4800.0],
       [0.0, 1.0, 0.0, 0.0, 30.0, 5400.0],
       [0.0, 0.0, 0.0, 1.0, 49.0, 98000.0]], dtype=object)
```

```python
from sklearn.model_selection import train_test_split
x_train,x_test,y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state = 0)
```

```python
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
```
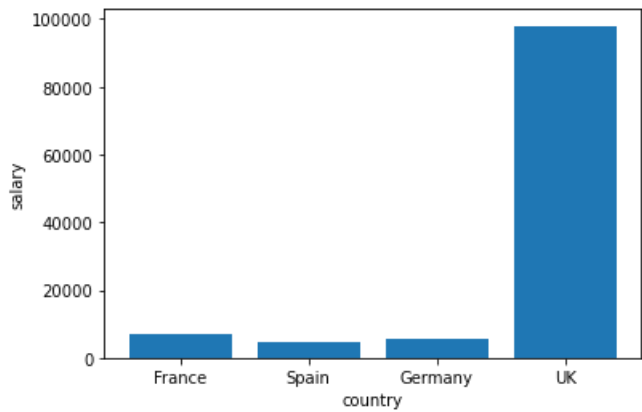
```python
x_train[:,4:6] = sc.fit_transform(x_train[:, 4:6])
x_train
```

```
array([[0.0, 0.0, 0.0, 1.0, 1.3109359202840398, 1.4138527953056175],
       [0.0, 0.0, 1.0, 0.0, -1.1149081191200718, -0.7345887349522664],
       [1.0, 0.0, 0.0, 0.0, -0.1960278011639687, -0.679264060353351]],
      dtype=object)
```

```python
plt.bar(df['country '],df['salary'])
plt.xlabel('country')
plt.ylabel('salary')
plt.show()
```

```python
import seaborn as sns
sns.pairplot(df)
```

```
<seaborn.axisgrid.PairGrid at 0x18b7dac21f0>
```



Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js