

## Data Preprocessing Tools

### Importing the libraries

```
In [ ]: 1 import numpy as np
        2 import matplotlib.pyplot as plt
        3 import pandas as pd
```

### Importing the dataset

```
In [ ]: 1 dataset = pd.read_csv('Data.csv')
        2 X = dataset.iloc[:, :-1].values
        3 y = dataset.iloc[:, -1].values
```

```
In [ ]: 1 print(X)

[['France' 44.0 72000.0]
 ['Spain' 27.0 48000.0]
 ['Germany' 30.0 54000.0]
 ['Spain' 38.0 61000.0]
 ['Germany' 40.0 nan]
 ['France' 35.0 58000.0]
 ['Spain' nan 52000.0]
 ['France' 48.0 79000.0]
 ['Germany' 50.0 83000.0]
 ['France' 37.0 67000.0]]
```

```
In [ ]: 1 print(y)

['No' 'Yes' 'No' 'No' 'Yes' 'Yes' 'No' 'Yes' 'No' 'Yes']
```

### Taking care of missing data

```
In [ ]: 1 from sklearn.impute import SimpleImputer
        2 imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
        3 imputer.fit(X[:, 1:3])
        4 X[:, 1:3] = imputer.transform(X[:, 1:3])
```

```
In [ ]: 1 print(X)

[['France' 44.0 72000.0]
 ['Spain' 27.0 48000.0]
 ['Germany' 30.0 54000.0]
 ['Spain' 38.0 61000.0]
 ['Germany' 40.0 63777.77777777778]
 ['France' 35.0 58000.0]
 ['Spain' 38.77777777777778 52000.0]
 ['France' 48.0 79000.0]
 ['Germany' 50.0 83000.0]
 ['France' 37.0 67000.0]]
```

### Encoding categorical data

#### Encoding the Independent Variable

```
In [ ]: 1 from sklearn.compose import ColumnTransformer
        2 from sklearn.preprocessing import OneHotEncoder
        3 ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(), [0])], remainder='passthrough')
        4 X = np.array(ct.fit_transform(X))
```

```
In [ ]: 1 print(X)

[[1.0 0.0 0.0 44.0 72000.0]
 [0.0 0.0 1.0 27.0 48000.0]
 [0.0 1.0 0.0 30.0 54000.0]
 [0.0 0.0 1.0 38.0 61000.0]
 [0.0 1.0 0.0 40.0 63777.77777777778]
 [1.0 0.0 0.0 35.0 58000.0]
 [0.0 0.0 1.0 38.77777777777778 52000.0]
 [1.0 0.0 0.0 48.0 79000.0]
 [0.0 1.0 0.0 50.0 83000.0]
 [1.0 0.0 0.0 37.0 67000.0]]
```

#### Encoding the Dependent Variable

```
In [ ]: 1 from sklearn.preprocessing import LabelEncoder
        2 le = LabelEncoder()
        3 y = le.fit_transform(y)
```

```
In [ ]: 1 print(y)

[0 1 0 0 1 1 0 1 0 1]
```

### Splitting the dataset into the Training set and Test set

```
In [ ]: 1 from sklearn.model_selection import train_test_split
        2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 1)
```

```
In [ ]: 1 print(X_train)

[[0.0 0.0 1.0 38.77777777777778 52000.0]
 [0.0 1.0 0.0 40.0 63777.77777777778]
 [1.0 0.0 0.0 44.0 72000.0]
 [0.0 0.0 1.0 38.0 61000.0]
 [0.0 0.0 1.0 27.0 48000.0]
 [1.0 0.0 0.0 48.0 79000.0]
 [0.0 1.0 0.0 50.0 83000.0]
 [1.0 0.0 0.0 35.0 58000.0]]
```

```
In [ ]: 1 print(X_test)

[[0.0 1.0 0.0 30.0 54000.0]
 [1.0 0.0 0.0 37.0 67000.0]]
```

```
In [ ]: 1 print(y_train)

[0 1 0 0 1 1 0 1]
```

```
In [ ]: 1 print(y_test)

[0 1]
```

### Feature Scaling

```
In [ ]: 1 from sklearn.preprocessing import StandardScaler
        2 sc = StandardScaler()
        3 X_train[:, 3:] = sc.fit_transform(X_train[:, 3:])
        4 X_test[:, 3:] = sc.transform(X_test[:, 3:])
```

```
In [ ]: 1 print(X_train)

[[0.0 0.0 1.0 -0.19159184384578545 -1.0781259408412425]
```

```
`\n[0.0 1.0 0.0 -0.014117293757057777 -0.07013167641635372]\n[1.0 0.0 0.0 0.56670850653324 0.633562432710455]\n[0.0 0.0 1.0 -0.30453019390224867 -0.30786617274297867]\n[0.0 0.0 1.0 -1.9018011447007988 -1.420463615551582]\n[1.0 0.0 0.0 1.1475343068237058 1.232653363453549]\n[0.0 1.0 0.0 1.4379472009688968 1.5749910381638885]\n[1.0 0.0 0.0 -0.7401495441200351 -0.5646194287757332]]
```

```
In [ ]: 1 print(X_test)
```

```
[[0.0 1.0 0.0 -1.4661817944830124 -0.9069571034860727]\n [1.0 0.0 0.0 -0.44973664397484414 0.2056403393225306]]
```