

1. The Dataset was Successfully downloaded.

```
In [1]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
from sklearn.preprocessing import LabelEncoder, MinMaxScaler
from sklearn.model_selection import train_test_split
```

1. Load the Dataset

```
In [4]: data = pd.read_csv("Downloads/Churn_Modelling.csv")
```

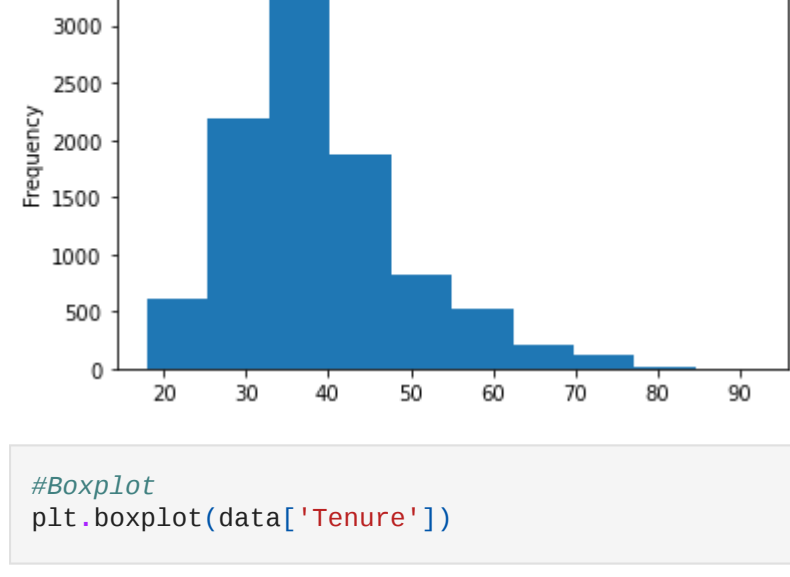
```
In [5]: data.head()
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	1	15634602	Hargrave	619	France	Female	42	2	0.00	1	1	1	101348.88	1
1	2	15647311	Hill	608	Spain	Female	41	1	83807.86	1	0	1	112542.58	0
2	3	15619304	Onio	502	France	Female	42	8	159660.80	3	1	0	113931.57	1
3	4	15701354	Boni	699	France	Female	39	1	0.00	2	0	0	93826.63	0
4	5	15737688	Mitchell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

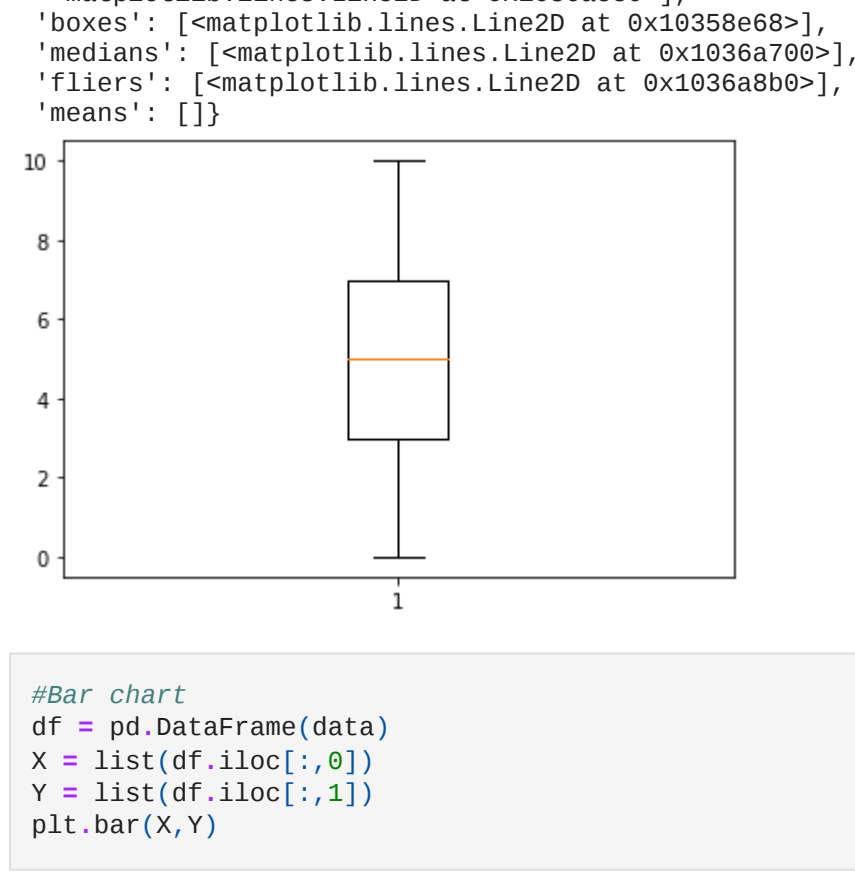
1. Perform Below Visualizations 3.1 Univariate Analysis

```
In [6]: #Histogram
data['Age'].plot(kind='hist')
```

```
Out[6]: <AxesSubplot:ylabel='Frequency'>
```

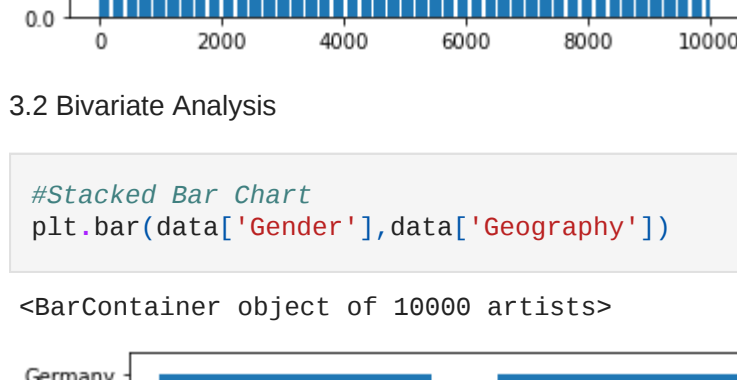


```
In [7]: #Boxplot
plt.boxplot(data['Tenure'])
```



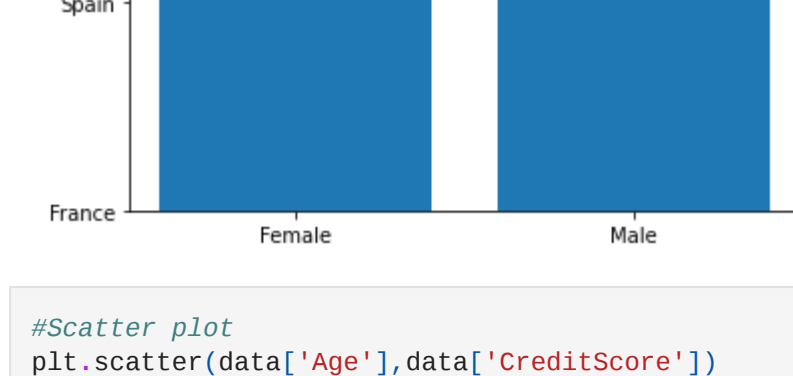
```
In [9]: #Bar chart
df = pd.DataFrame(data)
X = list(df.iloc[:,0])
Y = list(df.iloc[:,1])
plt.bar(X,Y)
```

```
Out[9]: <BarContainer object of 10000 artists>
```



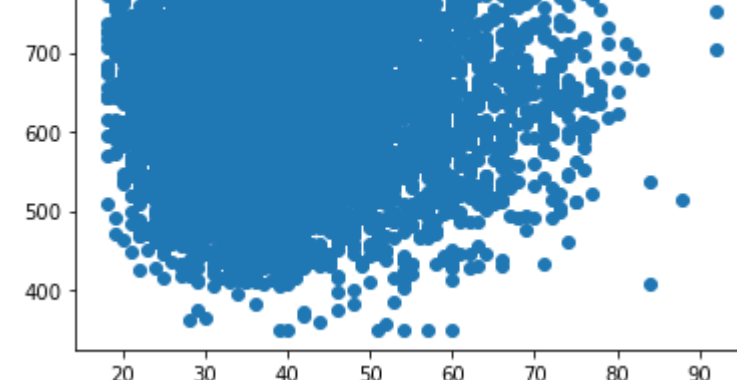
```
In [13]: #Stacked Bar Chart
plt.bar(data['Gender'],data['Geography'])
```

```
Out[13]: <BarContainer object of 10000 artists>
```



```
In [12]: #Scatter plot
plt.scatter(data['Age'],data['CreditScore'])
```

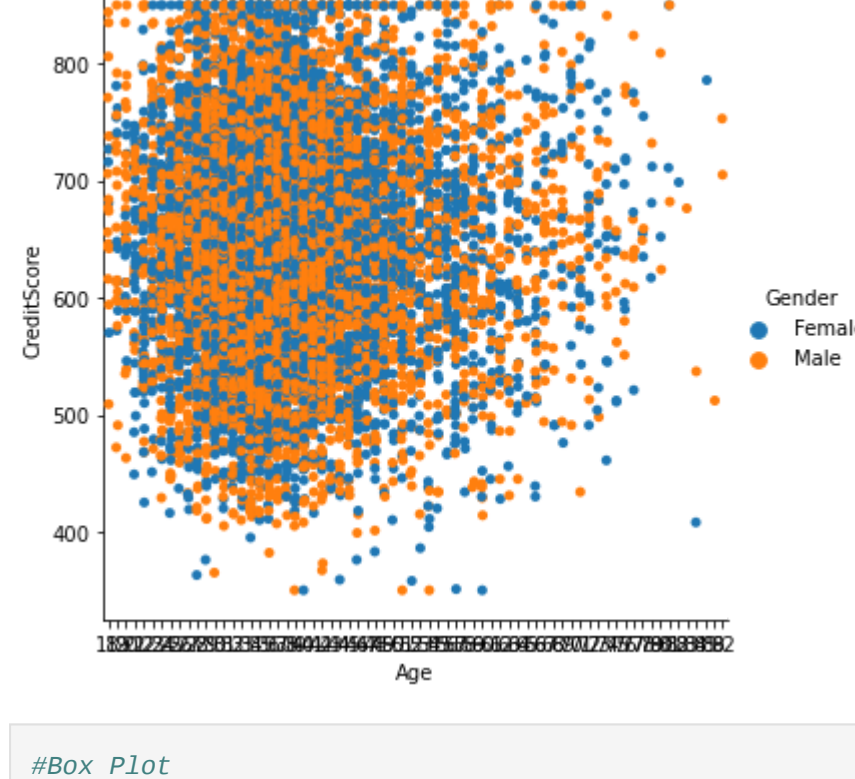
```
Out[12]: <matplotlib.collections.PathCollection at 0x21d13a98>
```



3.3 Multivariate Analysis

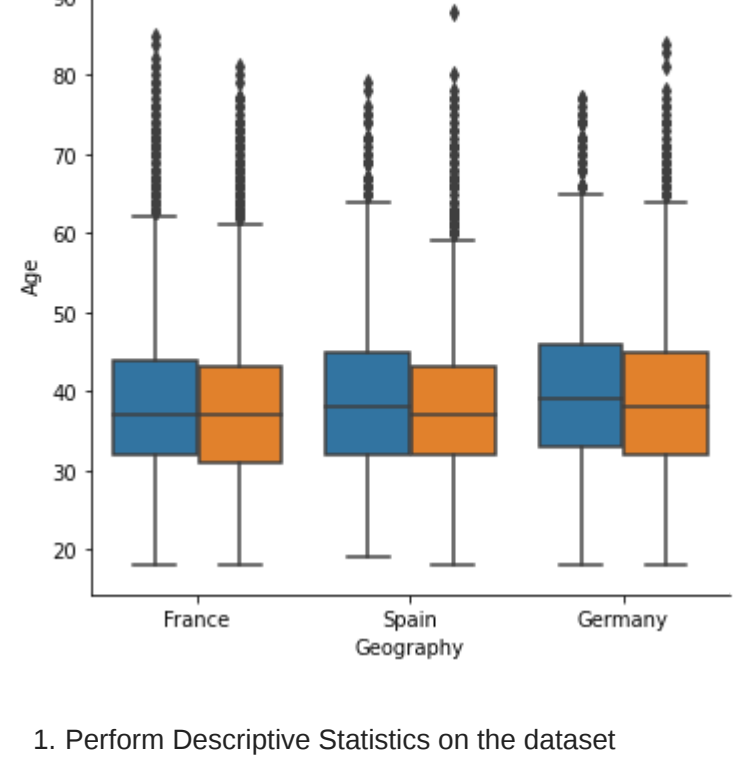
```
In [14]: sns.catplot(data=data,x='Age',y='CreditScore',hue='Gender')
```

```
Out[14]: <seaborn.axisgrid.FacetGrid at 0xa21ca7e8>
```



```
In [15]: #Box Plot
sns.catplot(data=data,x='Geography',y='Age',hue='Gender',kind='box')
```

```
Out[15]: <seaborn.axisgrid.FacetGrid at 0xe1e32c678>
```



1. Perform Descriptive Statistics on the dataset

```
In [16]: data.mean()
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
count	10000.00000	1.000000e+04	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	5000.50000	1.56094e+07	650.528800	38.921800	5.012800	76485.889288	1.530200	0.70550	0.515100	100090.239881	0.203700
std	2886.89568	7.19361e+04	96.653299	10.487906	2.892174	62397.405202	0.581654	0.45584	0.499797	57510.492818	0.402769
min	1.00000	1.556570e+07	350.000000	18.000000	0.000000	0.000000	1.000000	0.000000	0.000000	51002.11000	0.000000
25%	2500.75000	1.562635e+07	584.000000	32.000000	3.000000	0.000000	1.000000	0.000000	0.000000	51002.110000	0.000000
50%	5000.50000	1.56074e+07	652.000000	37.000000	5.000000	97198.540000	1.000000	1.000000	1.000000	100193.915000	0.000000
75%	7500.25000	1.575323e+07	716.000000	44.000000	7.000000	127644.240000	2.000000	1.000000	1.000000	149388.247500	0.000000
max	10000.00000	1.581568e+07	850.000000	92.000000	10.000000	250898.090000	4.000000	1.000000	1.000000	199992.480000	1.000000

```
In [17]: data.median()
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
count	10000.00000	1.000000e+04	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	5000.50000	1.56094e+07	650.528800	38.921800	5.012800	76485.889288	1.530200	0.70550	0.515100	100090.239881	0.203700
std	2886.89568	7.19361e+04	96.653299	10.487906	2.892174	62397.405202	0.581654	0.45584	0.499797	57510.492818	0.402769
min	1.00000	1.556570e+07	350.000000	18.000000	0.000000	0.000000	1.000000	0.000000	0.000000	51002.11000	0.000000
25%	2500.75000	1.562635e+07	584.000000	32.000000	3.000000	0.000000	1.000000	0.000000	0.000000	51002.110000	0.000000
50%	5000.50000	1.56074e+07	652.000000	37.000000	5.000000	97198.540000	1.000000	1.000000	1.000000	100193.915000	0.000000
75%	7500.25000	1.575323e+07	716.000000	44.000000	7.000000	127644.240000	2.000000	1.000000	1.000000	149388.247500	0.000000
max	10000.00000	1.581568e+07	850.000000	92.000000	10.000000	250898.090000	4.000000	1.000000	1.000000	199992.480000	1.000000

```
In [18]: data.describe()
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
count	10000.00000	1.000000e+04	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	5000.50000	1.56094e+07	650.528800	38.921800	5.012800	76485.889288	1.530200	0.70550	0.515100	100090.239881	0.203700
std	2886.89568	7.19361e+04	96.653299	10.487906	2.892174	62397.405202	0.581654	0.45584	0.499797	57510.492818	0.402769
min	1.00000	1.556570e+07	350.000000	18.000000	0.000000	0.000000	1.000000	0.000000	0.000000	51002.11000	0.000000
25%	2500.75000	1.562635e+07	584.000000	32.000000	3.000000	0.000000	1.000000	0.000000	0.000000	51002.110000	0.000000
50%	5000.50000	1.56074e+07	652.000000	37.000000	5.000000	97198.540000	1.000000	1.000000	1.000000	100193.915000	0.000000
75%	7500.25000	1.575323e+07	716.000000	44.000000	7.000000	127644.240000	2.000000	1.000000	1.000000	149388.247500	0.000000
max	10000.00000	1.581568e+07	850.000000	92.000000	10.000000	250898.090000	4.000000	1.000000	1.000000	199992.480000	1.000000

```
In [19]: data.shape
```

```
Out[19]: (10000, 14)
```

1. Handle the missing values

```
In [20]: data.isnull().sum()
```

```
Out[20]: RowNumber      0
CustomerId      0
Surname         0
CreditScore     0
Geography       0
Gender          0
Age            0
Tenure         0
Balance        0
NumOfProducts  0
HasCrCard      0
IsActiveMember 0
EstimatedSalary 0
Exited         0
dtype: int64
```

1. Find the Outliers and replace the Outliers

```
In [21]: sns.boxplot(data['Age'])
```

C:\Users\sssl\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be data, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
Out[21]: <AxesSubplot:ylabel='Age'>
```



```
In [22]: qnt = data.quantile(q=[0.25,0.75])
qnt
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0.25	2500.75	1562628.25	584.0	32.0	3.0	0.00	1.0	0.0	0.0	51002.1100	0.0
0.75	7500.25	15753233.75	716.0	44.0	7.0	127644.24	2.0	1.0	1.0	149388.2475	0.0

```
In [23]: IQR = qnt.loc[0.75] - qnt.loc[0.25]
IQR
```

```
Out[23]: RowNumber      4999.5000
CustomerId      124795.5000
CreditScore      134.0000
Age              12.0000
Tenure           4.0000
Balance         127644.2400
NumOfProducts     1.0000
HasCrCard        1.0000
IsActiveMember    1.0000
EstimatedSalary  98386.1375
Exited           0.0000
dtype: float64
```

```
In [24]: upper_extreme = qnt.loc[0.75]+1.5*IQR
upper_extreme
```

```
Out[24]: RowNumber      1.499950e+04
CustomerId      1.594029e+07
CreditScore      3.838000e+02
Age              1.400000e+01
Tenure          -3.000000e+00
Balance         -1.914650e+05
NumOfProducts    -5.000000e-01
HasCrCard        -1.500000e+00
IsActiveMember    -1.500000e+00
EstimatedSalary  -9.657710e+04
Exited           0.000000e+00
dtype: float64
```

```
In [25]: lower_extreme = qnt.loc[0.25]-1.5*IQR
lower_extreme
```

```
Out[25]: RowNumber      -4.998500e+03
CustomerId      1.544147e+07
CreditScore      3.838000e+02
Age              1.400000e+01
Tenure          -3.000000e+00
Balance         -1.914650e+05
NumOfProducts    -5.000000e-01
HasCrCard        -1.500000e+00
IsActiveMember    -1.500000e+00
EstimatedSalary  -9.657710e+04
Exited           0.000000e+00
dtype: float64
```

```
In [26]: df2 = data[(data['Age']<upper_extreme['Age']) & (data['Age']>lower_extreme['Age'])]
df2.shape
```

```
Out[26]: (18000, 14)
```

```
In [27]: df2.shape
```

```
Out[27]: (9589, 14)
```

```
In [28]: sns.boxplot(df2['Age'])
```

C:\Users\sssl\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be data, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
Out[28]: <AxesSubplot:ylabel='Age'>
```



1. Check the Categorical columns and perform Encoding

```
In [29]: le = LabelEncoder()
df2['Geography'] = le.fit_transform(df2['Geography'])
df2['Gender'] = le.fit_transform(df2['Gender'])
df2.head()
```

```
<ipython-input-29-7fbcaf02a1ee>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df2['Geography'] = le.fit_transform(df2['Geography'])
<ipython-input-29-7fbcaf02a1ee>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df2['Gender'] = le.fit_transform(df2['Gender'])
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	1	15634602	Hargrave	619	0	0	42	2	0.00	1	1	1	101348.88	1
1	2	15647311	Hill	608	2	0	41	1	83807.86	1	0	1	112542.58	0
2	3	15619304	Onio	502	0	0	42	8	159660.80	3	1	0	113931.57	1
3	4	15701354	Boni	699	0	0	39	1	0.00	2	0	0	93826.63	0
4	5	15737688	Mitchell	850	2	0	43	2	125510.82	1	1	1	79084.10	0

1. Split the data into dependent and independent variables

```
In [30]: x = df.iloc[:, :-1].values
x
```

```
Out[30]: array([[1, 15634602, 'Hargrave', ..., 1, 1, 101348.88],
[2, 15647311, 'Hill', ..., 0, 1, 112542.58],
[3, 15619304, 'Onio', ..., 1, 0, 113931.57],
...,
[9999, 15684532, 'Liu', ..., 0, 1, 42895.58],
[9999, 15682355, 'Sabbatini', ..., 1, 0, 92888.52],
[10000, 15628319, 'Walker', ..., 1, 0, 38190.78]], dtype=object)
```

```
In [31]: y = df.iloc[:, -1].values
y
```

```
Out[31]: array([1, 0, 1, ..., 1, 1, 0], dtype=int64)
```

1. Scale the independent variables

```
In [32]: scaler = MinMaxScaler()
df[['CustomerId']] = scaler.fit_transform(df[['CustomerId']])
df
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	1	0.265462	Hargrave	619	France	Female	42	2	0.00	1	1	1	101348.88	1
1	2	0.325454	Hill	608	Spain	Female	41	1	83807.86	1	0	1	112542.58	0
2	3	0.214421	Onio	502	France	Female	42	8	159660.80	3	1	0	113931.57	1
3	4	0.542636	Boni	699	France	Female	39	1	0.00	2	0	0	93826.63	0
4	5	0.608778	Mitchell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

9998	9999	0.466637	Sabbatini	772	Germany	Male	42	3	75075.31	2	1	0	92888.52
9999	10000	0.250483	Walker	792	France	Female	28	4	130142.79	1	1	0	38190.78

10000 rows x 14 columns

1. Split the data into training and testing

```
train_size=0.8
X = df.drop(columns=['Tenure']).copy()
```