Assignment -2 Data visualization and Preprocessing

Assignment Date	21 September 2022
Student Name	Dhanishkumar.R
Student Roll Number	610519104014
Maximum Marks	2 Marks

1.DOWNLOAD THE DATASET 2.LOAD THE DATASET

							4 3 k		_					
impo	r t panda	ıs as pd											I	n [1]:
	r t nump	-												
impo	r t matpl	otlib.py _l	plot as	plt										
impo	rt se	aborn	as s	ns									I	n [2]:
df = p	d.read	d_csv('/con	tent/C	hurn_N	Model	lin	g.cs	v')					
a e	In [6]:													
df													Ο	ut[6]:
													O	սւլоյ.
	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A g e	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMem ber	Estima tedSal ary	Ex ite d
0	1	1563 4602	Har gra ve	619	Fran ce	Fe ma le	4 2	2	0.00	1	1	1	101348 .88	1
1	2	1564 7311	Hill	608	Spai n	Fe ma le	4 1	1	838 07.8 6	1	0	1	112542 .58	0
2	3	1561 9304	Oni o	502	Fran ce	Fe ma le	4 2	8	159 660. 80	3	1	0	113931 .57	1
3	4	1570 1354	Bon i	699	Fran ce	Fe ma le	3 9	1	0.00	2	0	0	93826. 63	0

125

2 510.

79084.

Fe

ma 3

Spai

850

Mit

chel

1

1573

7888

9 9996 1560 Obi jiak 771 Fran Ma 3 5 0.00 2 1 0 96270. 64 0

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A g e	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMem ber	Estima tedSal ary	Ex ite d
9 9 9 6	9997	1556 9892	Joh nsto ne	516	Fran ce	Ma le	3 5	10	573 69.6 1	1	1	1	101699 .77	0
9 9 9 7	9998	1558 4532	Liu	709	Fran ce	Fe ma le	3 6	7	0.00	1	0	1	42085. 58	1
9 9 9 8	9999	1568 2355	Sab bati ni	772	Ger man y	Ma le	4 2	3	750 75.3 1	2	1	0	92888. 52	1
9 9 9	1000	1562 8319	Wal ker	792	Fran ce	Fe ma le	2 8	4	130 142. 79	1	1	0	38190. 78	0

 $10000 \ rows \times 14 \ columns$

In [3]:

df.head()

Out[3]:

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A g e	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMemb er	Estima tedSala ry	Ex ite d
0	1	1563 4602	Har gra ve	619	Fran ce	Fe ma le	4 2	2	0.00	1	1	1	101348 .88	1
1	2	1564 7311	Hill	608	Spai n	Fe ma le	4	1	838 07.8 6	1	0	1	112542 .58	0
2	3	1561 9304	Oni o	502	Fran ce	Fe ma le	4 2	8	159 660. 80	3	1	0	113931 .57	1
3	4	1570 1354	Bon i	699	Fran ce	Fe ma le	3 9	1	0.00	2	0	0	93826. 63	0

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A g e	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMemb er	Estima tedSala ry	Ex ite d	
4	5	1573 7888	Mit chel l	850	Spai n	Fe ma le	4 3	2	125 510. 82	1	1	1	79084. 10	0	
df.shape												I	n [4]:		
														ut[4]:	

3.Univariate,Bivariate & MultiVariate Analysis

Univariate Analysis

```
df_france=df.loc[df['Geography']=='France']
df_spain=df.loc[df['Geography']=='Spain']
df_germany=df.loc[df['Geography']=='Germany']

In [9]:

plt.plot(df_france['Balance'],np.zeros_like(df_france['Balance']),'o')
plt.plot(df_spain['Balance'],np.zeros_like(df_spain['Balance']),'o')
plt.plot(df_germany['Balance'],np.zeros_like(df_germany['Balance']),'o')
plt.xlabel('Age')
plt.show()
```

Bivariate Analysis

```
In [18]:
sns.FacetGrid(df,hue="Geography",size=5).map(plt.scatter,"Age","Balance").a
dd_legend();
plt.show()
/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:337: UserWarning
: The `size` parameter has been renamed to `height`; please update your cod
e.
    warnings.warn(msg, UserWarning)
```

Multivariate Analysis

```
In [24]:
```

/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:2076: UserWarnin g: The `size` parameter has been renamed to `height`; please update your co de.

warnings.warn(msg, UserWarning)

Out[24]:

<seaborn.axisgrid.PairGrid at 0x7f9a9f3029d0>

4.Descriptive Statistics

df.head()

Out[29]:

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A g e	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMemb er	Estima tedSala ry	Ex ite d
0	1	1563 4602	Har gra ve	619	Fran ce	Fe ma le	4 2	2	0.00	1	1	1	101348 .88	1
1	2	1564 7311	Hill	608	Spai n	Fe ma le	4	1	838 07.8 6	1	0	1	112542 .58	0
2	3	1561 9304	Oni o	502	Fran ce	Fe ma le	4 2	8	159 660. 80	3	1	0	113931 .57	1
3	4	1570 1354	Bon i	699	Fran ce	Fe ma le	3 9	1	0.00	2	0	0	93826. 63	0
4	5	1573 7888	Mit chel l	850	Spai n	Fe ma le	4 3	2	125 510. 82	1	1	1	79084. 10	0

In [30]:

df.mean() # Get the mean of each column

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

"""Entry point for launching an IPython kernel.

Out[30]:

RowNumber	5.000500e+03
CustomerId	1.569094e+07
CreditScore	6.505288e+02
Age	3.892180e+01
Tenure	5.012800e+00

```
Balance7.648589e+04NumOfProducts1.530200e+00HasCrCard7.055000e-01IsActiveMember5.151000e-01EstimatedSalary1.000902e+05Exited2.037000e-01
```

dtype: float64

In [31]:

```
df.mean(axis=1) # Get the mean of each row
```

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

"""Entry point for launching an IPython kernel.

```
Out[31]:
0
       1.430602e+06
1
       1.440392e+06
       1.444860e+06
       1.435993e+06
3
      1.449399e+06
9995 1.428483e+06
9996
      1.430866e+06
9997
      1.421579e+06
9998
      1.441922e+06
9999
      1.437044e+06
Length: 10000, dtype: float64
```

df.median()

Get the median of each column

In [32]:

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

"""Entry point for launching an IPython kernel.

Out[32]:

```
5.000500e+03
RowNumber
CustomerId
                 1.569074e+07
CreditScore
                6.520000e+02
                3.700000e+01
                5.000000e+00
Tenure
Balance
                 9.719854e+04
NumOfProducts
               1.000000e+00
HasCrCard
                 1.000000e+00
IsActiveMember
               1.000000e+00
EstimatedSalary 1.001939e+05
                 0.000000e+00
Exited
```

dtype: float64

```
norm data = pd.DataFrame(np.random.normal(size=100000))
```

In [39]:

```
norm data.plot(kind="density",
```

figsize=(10,10)); plt.vlines(norm_data.mean(), #Plot black line at

mean

```
ymin=0,
ymax=0.4,
```

```
linewidth=5.0);
plt.vlines(norm data.median(),  # Plot red line at median
           ymin=0,
            ymax=0.4,
            linewidth=2.0,
            color="red");
                                                                            In [36]:
skewed data = pd.DataFrame(np.random.exponential(size=100000))
skewed data.plot(kind="density",
               figsize=(10,10),
               xlim=(-1,5));
plt.vlines(skewed_data.mean(),
                                    # Plot black line at mean
            ymin=0,
            ymax=0.8,
            linewidth=5.0);
plt.vlines(skewed data.median(),  # Plot red line at median
           ymin=0,
           ymax=0.8,
           linewidth=2.0,
           color="red");
                                                                            In [40]:
norm data = np.random.normal(size=50)
outliers = np.random.normal(15, size=3)
combined data = pd.DataFrame(np.concatenate((norm data, outliers), axis=0))
combined data.plot(kind="density",
               figsize=(10,10),
               xlim=(-5,20);
plt.vlines(combined_data.mean(),
                                       # Plot black line at mean
            ymin=0,
            ymax=0.2,
            linewidth=5.0);
plt.vlines(combined_data.median(),
                                      # Plot red line at median
           ymin=0,
           ymax=0.2,
           linewidth=2.0,
           color="red");
                                                                            In [42]:
df.mode()
```

Out[42]:

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A g e	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMem ber	Estima tedSala ry	Ex ite d
0	1	1556 5701	Smi th	850.0	Fran ce	Ma le	3 7. 0	2.0	0.0	1.0	1.0	1.0	24924. 92	0. 0
1	2	1556 5706	Na N	NaN	NaN	Na N	N a N	Na N	Na N	NaN	NaN	NaN	NaN	Na N
2	3	1556 5714	Na N	NaN	NaN	Na N	N a N	Na N	Na N	NaN	NaN	NaN	NaN	Na N
3	4	1556 5779	Na N	NaN	NaN	Na N	N a N	Na N	Na N	NaN	NaN	NaN	NaN	Na N
4	5	1556 5796	Na N	NaN	NaN	Na N	N a N	Na N	Na N	NaN	NaN	NaN	NaN	Na N
9 9 9 5	9996	1581 5628	Na N	NaN	NaN	Na N	N a N	Na N	Na N	NaN	NaN	NaN	NaN	Na N
9 9 9 6	9997	1581 5645	Na N	NaN	NaN	Na N	N a N	Na N	Na N	NaN	NaN	NaN	NaN	Na N
9 9 9 7	9998	1581 5656	Na N	NaN	NaN	Na N	N a N	Na N	Na N	NaN	NaN	NaN	NaN	Na N
9 9 9 8	9999	1581 5660	Na N	NaN	NaN	Na N	N a N	Na N	Na N	NaN	NaN	NaN	NaN	Na N
9 9 9	1000	1581 5690	Na N	NaN	NaN	Na N	N a N	Na N	Na N	NaN	NaN	NaN	NaN	Na N

Measures of Spread

```
In [43]:
max(df["Age"]) - min(df["Age"])
                                                                         Out[43]:
74
                                                                         In [45]:
five num = [df["Age"].quantile(0),
            df["Age"].quantile(0.25),
            df["Age"].quantile(0.50),
            df["Age"].quantile(0.75),
            df["Age"].quantile(1)]
five num
[18.0, 32.0, 37.0, 44.0, 92.0]
                                                                         Out[45]:
df["Age"].describe()
                                                                          In [46]:
count 10000.000000
mean 38.921800
std 10.487806
                                                                         Out[46]:
          18.000000
min
           32.000000
25%
          37.000000
50%
75%
max
           44.000000
           92.000000
Name: Age, dtype: float64
df["Age"].quantile(0.75) - df["Age"].quantile(0.25)
                                                                          In [47]:
12.0
df.boxplot(column="Age",
                                                                         Out[47]:
               return type='axes',
               figsize=(8,8))
                                                                         In [49]:
plt.text(x=0.74, y=22.25, s="3rd Quartile")
plt.text(x=0.8, y=18.75, s="Median")
plt.text(x=0.75, y=15.5, s="1st Quartile")
plt.text(x=0.9, y=10, s="Min")
plt.text(x=0.9, y=33.5, s="Max")
plt.text(x=0.7, y=19.5, s="IQR", rotation=90, size=25);
                                                                          In [50]:
df["Age"].var()
                                                                         Out[50]:
109.99408416841683
                                                                          In [51]:
df["Age"].std()
```

```
Out[51]:
10.487806451704609
                                                                            In [52]:
abs median devs = abs(df["Age"] - df["Age"].median())
abs median devs.median() * 1.4826
                                                                           Out[52]:
8.8956
Skewness and Kurtosis
df["Age"].skew() # Check skewness
                                                                            In [53]:
1.0113202630234552
                                                                           Out[53]:
df["Age"].kurt() # Check kurtosis
                                                                            In [54]:
1.3953470615086956
                                                                           Out[54]:
norm data = np.random.normal(size=100000)
                                                                            In [55]:
skewed_data = np.concatenate((np.random.normal(size=35000)+2,
                               np.random.exponential(size=65000)),
                               axis=0)
uniform data = np.random.uniform(0,2, size=100000)
peaked data = np.concatenate((np.random.exponential(size=50000),
                               np.random.exponential(size=50000)*(-1)),
                               axis=0)
data_df = pd.DataFrame({"norm":norm_data,
                         "skewed":skewed data,
                         "uniform":uniform data,
                         "peaked":peaked data})
                                                                            In [56]:
data df.plot(kind="density",
             figsize=(10,10),
             xlim=(-5,5));
                                                                            In [57]:
data df.skew()
                                                                           Out[57]:
norm -0.007037
skewed 1.002549
          1.002549
uniform -0.004434
peaked 0.018058
dtype: float64
data df.kurt()
                                                                            In [58]:
         -0.009914
skewed
          1.314497
                                                                           Out[58]:
```

uniform -1.201740 peaked 2.971592

dtype: float64

5.Handle the Missing values

Values

In [83]:

df=pd.read_csv('/content/Churn_Modelling.csv')

In [84]:

df.head()

Out[84]:

Row Cust Sur Cred Geo Ge A Te Bal NumO Has IsActiv Estima Ex Num omer na itSco grap nd g nu anc fProdu CrC eMemb tedSala ite bar Id me re by or a re of cts ord are recorded.

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A g e	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMemb er	Estima tedSala ry	Ex ite d
0	1	1563 4602	Har gra ve	619	Fran ce	Fe ma le	4 2	2	0.00	1	1	1	101348 .88	1
1	2	1564 7311	Hill	608	Spai n	Fe ma le	4	1	838 07.8 6	1	0	1	112542 .58	0
2	3	1561 9304	Oni o	502	Fran ce	Fe ma le	4 2	8	159 660. 80	3	1	0	113931 .57	1
3	4	1570 1354	Bon i	699	Fran ce	Fe ma le	3 9	1	0.00	2	0	0	93826. 63	0
4	5	1573 7888	Mit chel l	850	Spai n	Fe ma le	4 3	2	125 510. 82	1	1	1	79084. 10	0

df.isnull()

Out[86]:

In [86]:

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A ge	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMem ber	Estima tedSala ry	Ex ite d	
0	False	False	Fals e	False	False	Fal se	F al	Fal se	Fal se	False	False	False	False	Fa lse	

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A ge	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMem ber	Estima tedSala ry	Ex ite d
1	False	False	Fals e	False	False	Fal se	F al se	Fal se	Fal se	False	False	False	False	Fa lse
2	False	False	Fals e	False	False	Fal se	F al se	Fal se	Fal se	False	False	False	False	Fa Ise
3	False	False	Fals e	False	False	Fal se	F al se	Fal se	Fal se	False	False	False	False	Fa lse
4	False	False	Fals e	False	False	Fal se	F al se	Fal se	Fal se	False	False	False	False	Fa lse
•••														
9 9 9 5	False	False	Fals e	False	False	Fal se	F al se	Fal se	Fal se	False	False	False	False	Fa lse
9 9 9 6	False	False	Fals e	False	False	Fal se	F al se	Fal se	Fal se	False	False	False	False	Fa lse
9 9 9 7	False	False	Fals e	False	False	Fal se	F al se	Fal se	Fal se	False	False	False	False	Fa lse
9 9 9 8	False	False	Fals e	False	False	Fal se	F al se	Fal se	Fal se	False	False	False	False	Fa lse
9 9 9	False	False	Fals e	False	False	Fal se	F al se	Fal se	Fal se	False	False	False	False	Fa lse

 $10000 \ rows \times 14 \ columns$

```
Out[89]:
<matplotlib.axes. subplots.AxesSubplot at 0x7f9a987d8290>
                                                                         In [93]:
sns.set style('whitegrid')
sns.countplot(x='Geography',data=df)
<matplotlib.axes._subplots.AxesSubplot at 0x7f9a92a88850>
                                                                        Out[93]:
sns.set style('whitegrid')
                                                                        In [94]:
sns.countplot(x='Geography', hue='Gender', data=df, palette='RdBu r')
<matplotlib.axes. subplots.AxesSubplot at 0x7f9a92ec10d0>
                                                                        Out[94]:
sns.set style('whitegrid')
sns.countplot(x='Geography', hue='Gender', data=df, palette='rainbow')
                                                                         In [96]:
<matplotlib.axes. subplots.AxesSubplot at 0x7f9a92afac50>
                                                                        Out[96]:
sns.distplot(df['Age'].dropna(),kde=False,color='darkred',bins=40)
                                                                         In [97]:
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: Futur
eWarning: `distplot` is a deprecated function and will be removed in a futu
re version. Please adapt your code to use either `displot` (a figure-level
function with similar flexibility) or `histplot` (an axes-level function fo
r histograms).
  warnings.warn(msg, FutureWarning)
                                                                        Out[97]:
<matplotlib.axes. subplots.AxesSubplot at 0x7f9a98787590>
                                                                         In [98]:
df['Age'].hist(bins=30,color='darkred',alpha=0.3)
<matplotlib.axes._subplots.AxesSubplot at 0x7f9a92d64c10>
                                                                        Out[98]:
sns.countplot(x='NumOfProducts',data=df)
                                                                        In [100]:
<matplotlib.axes. subplots.AxesSubplot at 0x7f9a9306f790>
                                                                       Out[100]:
df['Age'].hist(color='green',bins=40,figsize=(8,4))
<matplotlib.axes. subplots.AxesSubplot at 0x7f9a90f52d90>
                                                                        In [101]:
```

Cufflinks for plots

Out[101]:

```
In [102]:
import cufflinks as cf
cf.go_offline()
                                                                                       In []:
df['Age'].iplot(kind='hist',bins=30,color='green')
Data Cleaning
                                                                                    In [107]:
plt.figure(figsize=(12, 7))
sns.boxplot(x='Gender',y='Age',data=df,palette='winter')
                                                                                   Out[107]:
<matplotlib.axes. subplots.AxesSubplot at 0x7f9a90f59450>
def impute age(cols):
                                                                                    In [307]:
    Age = cols[0]
    Pclass = cols[1]
    if pd.isnull(Age):
         if Pclass == 1:
                return 37
              elif Pclass == 2:
                return 29
         else:
                return 24
     else:
                                                                                    In [122]:
         return Age sns.heatmap(df.isnull(),yticklabels=False,cbar=False,cmap='viridis')
                                                                                   Out[122]:
<matplotlib.axes. subplots.AxesSubplot at 0x7f9a8aa699d0>
                                                                                    In [112]:
df.drop('Gender',axis=1,inplace=True)
df.head()
                                                                                    In [114]:
                                                                                   Out[114]:
                                                                      IsActive
    RowN
                   Sur
                         Credi
                                                      NumOf
                                                               HasC
                                                                               Estimat
                                                                                        Ex
           Custo
                                Geog
                                           Te
                                                Bala
                                                                               edSalar
                                           nu
                                                               rCar
                                                                      Membe
    umbe
            merI
                  nam
                         tScor
                                raph
                                                      Product
                                                                                        ite
                                       \mathbf{g}
                                                nce
```

re

 \mathbf{S}

d

d

e

0 1 15634 Har grav 619 Franc 4 2 0.00 1 1 1 1 101348. 1

	RowN umbe r	Custo merI d	Sur nam e	Credi tScor e	Geog raph y	A g e	Te nu re	Bala nce	NumOf Product s	HasC rCar d	IsActive Membe r	Estimat edSalar y	Ex ite d
1	2	15647 311	Hill	608	Spain	4 1	1	8380 7.86	1	0	1	112542. 58	0
2	3	15619 304	Oni o	502	Franc e	4 2	8	1596 60.8 0	3	1	0	113931. 57	1
3	4	15701 354	Bon i	699	Franc e	3 9	1	0.00	2	0	0	93826.6	0
4	5	15737 888	Mitc hell	850	Spain	4 3	2	1255 10.8 2	1	1	1	79084.1 0	0

Converting Categorical Features

In [116]:

```
df.info()
```

<class 'pandas.core.frame.DataFrame'> RangeIndex: 10000 entries, 0 to 9999 Data columns (total 13 columns):

#	Column	Non-Null Count	Dtype	
0	RowNumber	10000 non-null	int64	
1	CustomerId	10000 non-null	int64	
2	Surname	10000 non-null	object	
3	CreditScore	10000 non-null	int64	
4	Geography	10000 non-null	object	
5	Age	10000 non-null	int64	
6	Tenure	10000 non-null	int64	
7	Balance	10000 non-null	float64	
8	NumOfProducts	10000 non-null	int64	
9	HasCrCard	10000 non-null	int64	
10	IsActiveMember	10000 non-null	int64	
11	EstimatedSalary	10000 non-null	float64	
12	Exited	10000 non-null	int64	
dtype	es: float64(2), in	nt64(9), object(2)	
memo	ry usage: 1015.8+	KB		

pd.get_dummies(df['Geography'],drop_first=True).head()

In [118]:

Out[118]: Germany Spain

0 0

1	0 1								
2	0 0								
3	0 0								
4	0 1								
								I	n [124]:
df.in	fo								
								0	ut[124]:
		taFrame.info		RowNumber	Cust	comerId	S	urname	e Cre
ditSc 0	ore Geograp 1	hy Age Ten 15634602	ure \		610	Eman	~~	42	2
1	2	15647311	Hargrave Hill		619 608	Fran Spa		41	1
2	3	15619304	Onio		502	Fran		42	8
3	4	15701354	Boni		699	Fran		39	1
4	5	15737888	Mitchell		850	Spa		43	2
	• • •	• • •							
9995	9996	15606229	Obijiaku		771	Fran		39	5
9996	9997	15569892	Johnstone		516	Fran		35	10
9997	9998	15584532	Liu		709	Fran		36	7
9998	9999	15682355	Sabbatini		772	Germa		42	3
9999	10000	15628319	Walker		792	Fran	_	28	4
\	Balance	NumOfProduc	ts HasCrC	ard IsAc	ctiveMe	ember :	Estim	atedSa	alary
0	0.00		1	1		1		10134	18.88
1	83807.86		1	0		1			12.58
2	159660.80		3	1		0			31.57
3	0.00		2	0		0			26.63
4	125510.82		1	1		1			34.10
9995	0.00		2	1		0		962	70.64
9996	57369.61		1	1		1		10169	99.77
9997	0.00		1	0		1		4208	35.58
9998	75075.31		2	1		0			38.52
9999	130142.79		1	1		0		3819	90.78
	Exited								
0	1								
1	0								
2	1								
3	0								
4	0								
9995	0								
9996	0								
9997	1								
9998	1								

Germany Spain

```
9999 0
```

[10000 rows x 13 columns]>

sex = pd.get_dummies(df['Age'],drop_first=True)
embark = pd.get_dummies(df['Balance'],drop_first=True)

In [127]:

df.drop(['Age','HasCrCard','Surname','CustomerId'],axis=1,inplace=True)

In [129]:

df.head()

Out[129]:

	RowNum ber	CreditSc ore	Geogra phy	Tenu re	Balanc e	NumOfProd ucts	IsActiveMe mber	EstimatedSa lary	Exit ed
0	1	619	France	2	0.00	1	1	101348.88	1
1	2	608	Spain	1	83807.8 6	1	1	112542.58	0
2	3	502	France	8	159660. 80	3	0	113931.57	1
3	4	699	France	1	0.00	2	0	93826.63	0
4	5	850	Spain	2	125510. 82	1	1	79084.10	0

In [130]:

train = pd.concat([df,sex,embark],axis=1)

train.head()

Out[131]:

																					- 1
	Ro											2	2	2	2	2	2	2	2	2	2
		C	•	T _N	В	Nu	IsA	Est	E			1	1	1	1	1	1	2	2 2	3 8	5
	w ed	Cr eo	G e u	itS	al	mO fPr	ctiv	ima	x it			2	2	2	3	4	6	1	2	3	0 8
	gr	n	m	co	a n	odu	eM em	ted Sal	e	1		6	6	7	1	3	1	•	6 7.	8 7.	9 8.
	ap be	re	hy	r	ce	cts	ber	ary	d	9	•	9	9	7	4	4	0	5	6	5	0
		10	щ	e							•	2.	6.	8.	6.	6.	9.	3			
	r			_										0.	υ.			2.			
												9	3	2	2	9 6	8 8			_	
												7	2			6	8	8	3	6	9
			Fr		0.			101													
			an	2	0.	1	1	348	1	0	•	0	0	0	0	0	0	0	0	0	0
		61	ce	_	0		_	.88	_	-											
0	1	9																			
			Sp		8			112													
			ai	1	3	1	1	542	0	0		0	0	0	0	0	0	0	0	0	0

1 69 8

6. Find the outliers and replace the outliers

In [147]:

dataset= [11,10,12,14,12,15,14,13,15,102,12,14,17,19,107, 10,13,12,14,12,108,12,11,14,13,15,10,15,12,10,14,13,15,10]

Detecting outlier using Z score

Using Z score

In [148]:

outliers=[]
def detect_outliers(data):
 threshold=3
 mean = np.mean(data)

```
std =np.std(data)
   for i in data:
       z_score= (i - mean)/std
        if np.abs(z_score) > threshold:
            outliers.append(y)
    return outliers
                                                                      In [151]:
outlier pt=detect outliers(dataset)
                                                                      In [152]:
outlier pt
                                                                     Out[152]:
[0
       101348.88
 1
        112542.58
        113931.57
 3
        93826.63
         79084.10
        96270.64
 9995
      101699.77
 9996
 9997
        42085.58
 9998
        92888.52
      38190.78
 Name: EstimatedSalary, Length: 10000, dtype: float64, 0 101348.88
       112542.58
        113931.57
        93826.63
 3
         79084.10
        96270.64
 9995
      101699.77
 9996
 9997
         42085.58
 9998
        92888.52
       38190.78
 Name: EstimatedSalary, Length: 10000, dtype: float64, 0 101348.88
        112542.58
        113931.57
 3
         93826.63
         79084.10
          . . .
 9995
        96270.64
 9996
       101699.77
 9997
         42085.58
         92888.52
 9998
 9999
         38190.78
 Name: EstimatedSalary, Length: 10000, dtype: float64]
                                                                      In [153]:
## Perform all the steps of IQR
sorted(dataset)
                                                                     Out[153]:
[10,
10,
 10,
```

```
10,
 10,
 11,
 11,
 12,
 12,
 12,
 12,
 12,
 12,
 12,
 13,
 13,
 13,
 13,
 14,
 14,
 14,
 14,
 14,
 14,
 15,
 15,
 15,
 15,
 15,
 17,
 19,
 102,
 107,
 108]
                                                                               In [155]:
quantile1, quantile3= np.percentile(dataset,[25,75])
print(quantile1,quantile3)
                                                                               In [156]:
12.0 15.0
## Find the IQR
                                                                               In [157]:
iqr value=quantile3-quantile1
print(iqr_value)
3.0
## Find the lower bound value and the higher bound value
                                                                               In [159]:
lower_bound_val = quantile1 -(1.5 * iqr_value)
upper_bound_val = quantile3 +(1.5 * iqr_value)
print(lower_bound_val,upper_bound_val)
                                                                               In [160]:
7.5 19.5
```

7. Check for Categorical columns andperform encoding

df=1	pd.rea	.d_csv	('/co	ntent/	Churn	_Mode	ellin	g.cs	sv')				ln	[161]:
df.	head())											In	[162]:
													Out	[162]:
	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	nd	A g e	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMemb er	Estima tedSala ry	Ex ite d
0	1	1563 4602	Har gra ve	619	Fran ce	ma	4 2	2	0.00	1	1	1	101348 .88	1
1	2	1564 7311	Hill	608	Spai n		4 1	1	838 07.8 6	1	0	1	112542 .58	0
2	3	1561 9304	Oni o	502	Fran ce	ma	4 2	8	159 660. 80	3	1	0	113931 .57	1
3	4	1570 1354	Bon i	699	Fran ce	111121	3 9	1	0.00	2	0	0	93826. 63	0
4	5	1573 7888	Mit chel l	850	Spai n	m	4 3	2	125 510. 82	1	1	1	79084. 10	0
df	numer '	ic = c	₹[['	R∩wNiim	ber!	. 'Cı	ıstom	erTo	1'. '(CreditSc	core!.	'Age'.		[163]:
'Ba 'Nu	lance' mOfPro	', oducts	s','H	asCrCa	ırd',	'IsAc	ctive	Memk	per',	' Estimat	tedSala			
_	catego numer:			I[['Su	ırnam	e', '	Geog	raph	ту',	'Gender'	.]]		In	[164]:
α ₁ _	. 1 WING L -	- O • 1166	()										Out	[164]:
	RowNu mber			Credit Score	A ge	Ten ure	Balan ce		umOfPr oducts		IsActiv en	veM Es aber	stimated Salary	Exi ted
0	1	1563	346 02	619	42	2	0.00		1	1		1 10	01348.88	1

	RowNu mber	Custo merId	Credit Score	A ge	Ten ure	Balan ce	NumOfPr oducts	HasCr Card	IsActiveM ember	Estimated Salary	Exi ted
1	2	156473 11	608	41	1	83807 .86	1	0	1	112542.58	0
2	3	156193 04	502	42	8	15966 0.80	3	1	0	113931.57	1
3	4	157013 54	699	39	1	0.00	2	0	0	93826.63	0
4	5	157378 88	850	43	2	12551 0.82	1	1	1	79084.10	0
df_	categori	ical.head	d()							In	[165]:
	Surname	Geography	v Gen	der						Out	:[165]:
0	Hargrave	France	e Fen	nale							
1	Hill	Spair	n Fen	nale							
2	Onio	France	e Fen	nale							
3	Boni	France	e Fen	nale							
4	Mitchell	Spair	n Fen	nale							
pri	nt(df[' G	Surname'] eography' Gender'].].uni	que (ln	[166]:
[' F		'Spain' '				ashiwa	gi' 'Ald	ridge' '	'Burbidge']	
fro	m sklear	n.prepro	cess	ing i	mpor	t Labe	lEncoder				
mar	ry_encod	der = Lak	pelEnd	coder	()					In	[167]:
mar	ry_encod	er.fit(df	_cate	egori	cal['	Gender	'])				
Lab	elEncode	er()								In	[168]:
										Out	:[168]:

In [169]:

```
marry values = marry encoder.transform(df categorical['Gender'])
                                                                                                                                                             In [170]:
print("Before Encoding:", list(df categorical['Gender'][-10:]))
print("After Encoding:", marry_values[-10:])
print("The inverse from the encoding result:",
marry encoder.inverse transform(marry values[-10:]))
Before Encoding: ['Male', 'Female', 'Male', 'Male', 'Female', 'Male', 
', 'Female', 'Male', 'Female']
After Encoding: [1 0 1 1 0 1 1 0 1 0]
The inverse from the encoding result: ['Male' 'Female' 'Male' 'Femal
e' 'Male' 'Male' 'Female' 'Male'
  'Female']
                                                                                                                                                              In [171]:
residence encoder = LabelEncoder()
residence values =
residence_encoder.fit_transform(df_categorical['Geography'])
print("Before Encoding:", list(df_categorical['Geography'][:5]))
print("After Encoding:", residence_values[:5])
print("The inverse from the encoding result:",
residence encoder.inverse transform(residence values[:5]))
Before Encoding: ['France', 'Spain', 'France', 'France', 'Spain']
After Encoding: [0 2 0 0 2]
The inverse from the encoding result: ['France' 'Spain' 'France' 'France' '
Spain']
                                                                                                                                                              In [172]:
from sklearn.preprocessing import OneHotEncoder
gender encoder = OneHotEncoder()
                                                                                                                                                             In [174]:
from sklearn.preprocessing import OneHotEncoder
import numpy as np
gender encoder = OneHotEncoder()
gender reshaped = np.array(df categorical['Gender']).reshape(-1, 1)
gender values = gender encoder.fit transform(gender reshaped)
print(df_categorical['Gender'][:5])
print()
print(gender values.toarray()[:5])
print(gender encoder.inverse transform(gender values)[:5])
\cap
          Female
1
          Female
2
          Female
         Female
         Female
Name: Gender, dtype: object
[[1. 0.]
  [1. 0.]
  [1. 0.]
   [1. 0.]
   [1. 0.]]
```

```
[['Female']
 ['Female']
 ['Female']
 ['Female']
 ['Female']]
                                                                        In [175]:
smoke encoder = OneHotEncoder()
smoke reshaped = np.array(df categorical['Surname']).reshape(-1, 1)
smoke values = smoke encoder.fit transform(smoke reshaped)
print(df_categorical['Surname'][:5])
print()
print(smoke values.toarray()[:5])
print()
print(smoke encoder.inverse transform(smoke values)[:5])
     Hargrave
0
1
         Hill
2
         Onio
3
         Boni
    Mitchell
Name: Surname, dtype: object
[[0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]]
[['Hargrave']
 ['Hill']
 ['Onio']
 ['Boni']
 ['Mitchell']]
work encoder = OneHotEncoder()
                                                                        In [176]:
work_reshaped = np.array(df_categorical['Geography']).reshape(-1, 1)
work values = work encoder.fit transform(work reshaped)
print(df categorical['Geography'][:5])
print()
print(work_values.toarray()[:5])
print()
print(work_encoder.inverse_transform(work values)[:5])
0
     France
1
      Spain
2
     France
3
    France
      Spain
Name: Geography, dtype: object
[[1. 0. 0.]
 [0. 0. 1.]
 [1. 0. 0.]
 [1. 0. 0.]
 [0. 0. 1.]]
```

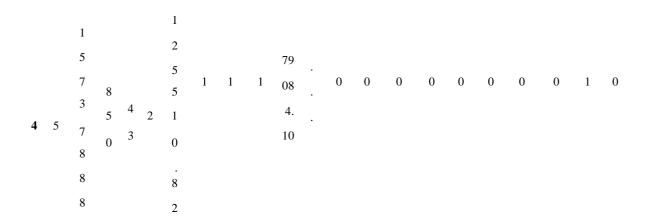
```
['Spain']
  ['France']
  ['France']
  ['Spain']]
                                                                                                                           In [178]:
df_categorical_encoded = pd.get_dummies(df_categorical, drop_first=True)
df_categorical_encoded.head()
                                                                                                                         Out[178]:
                                                                                                         \mathbf{S}
                                                                                                                S
      \mathbf{S}
            \mathbf{S}
                                                                           \mathbf{S}
                                                                                                    \mathbf{S}
                                                                                       \mathbf{S}
                                                                                            Su
                                                                                                    u
                                                                                      ur
                                                                                            rn
                                                                                                    r
                  S
                               \mathbf{S}
                                    Su
                                            \mathbf{S}
                                                 Su
                                                        Su
                                                               Su
                                                                                                                             G
            u
      u
                         \mathbf{S}
                                                                                                                     Ge
                                                                                                                                  G
                                                                                      na
                                                                                             a
                                                                                                    n
                                                               rn
            r
                 ur
                                                  rn
                                                                                      m
                                                                                             m
      r
                                     rn
                                           ur
                                                                                                                      og
                        ur
                                                               a
                                                                           n
                                                                                                   m
                                                               me
                                                                                       \mathbf{Z}
                                                                                             \mathbf{Z}
      n
            n
                 na
                                                                                                                            gr
                                                                                                    e
                                                                                                                n
                                                                                                                      ra
                                                                                       u
                                                                                            ub
                                                               _A
                                                                                                   ar{\mathbf{z}}
                                                  m
                                                        me
                                           m
      a
            a
                  m
                               n
                                     m
                                                                                                                             a
                                                                                                                     ph
                                                                                      ba
                                                                                             ar
                                                               br
                                                                          m
                                                                                      re
                                                                                             ev
                                                                                                    u
                                                                                                         m
                                     \mathbf{e}_{-}
                                                  \mathbf{e}_{-}
                                                        _A
     m
           m
                  e_
                                           e_
                                                               a
                               a
                                                                                                               m
                                                                                                                     \mathbf{y}_{-}
                                                                                                                             p
                                                                                m
                                                               m
                                                                          e_
                                                                                                          e
                                                  A
                                                         br
                                            A
           \mathbf{e}_{-}
                  \mathbf{A}
                              m
                                                                                                               e_
                                                                                                                     Ge
                                                                                                                             h
                                                               ow
                                                                                 e
                                                                           \mathbf{Z}
                                                               itz
            A
                                     be
                                           br
                                                  br
                                                          a
                                                                                                               \mathbf{Z}
                  b
                               e
                                                                                                                       r
                                                                          ot
      \mathbf{A}
            b
                                                   a
                                                         m
                                                                                                                u
                                                                                                                      m
                                                                                                                             \mathbf{S}
                  d
                                                                                                                                  M
                                                                                \mathbf{Z}
                                                                                                          u
                                                         ov
      b
            b
                                     at
                                           m
                                                                                                               ye
                                                                                                                     an
                  ul
                              \mathbf{A}
                                                                                                                             p
                                                         ic
                                                  \mathbf{ov}
     bi
                                    hy
                                                                                                                       y
                  la
                               b
                                                                                                                             ai
                                                                                                                                  le
                        ov
                                                                                                          e
                                                          h
      e
            t
                                                                                                                a
                  h
                              el
                                                                                                                             n
                                                                           0
                                                                 0
                                                   0
                                                          0
                                                                           0
 1 0
                               0
                                                   0
                                                          0
                                                                           0
 2 0
            0
                  0
                         0
                               0
                                      0
                                            0
                                                   0
                                                          0
                                                                 0
                         0
                               0
                                            0
                                                   0
                                                         0
                                                                 0
                                      0
```

[['France']

.

Out[179]:

	R o w N u m b er	C u st o m e rI d	C r e di t S c o r e	A g e	T e n u r e	B a l a n c	N u m Of Pr od uc ts	H a s C r C a r d	Is Ac tiv e M e m be r	Es ti m at ed Sa la ry	 Su rn a m e_ Zo to va	S u r n a m e_ Z o x	Su rn a me _Z ub ar ev	Su rn am e_ Zu ba rev a	S ur na m e_Z ue v	Su rn a m e_ Z uy ev	Su rn a m e_ Zu ye va	Ge ogr ap hy _G er ma ny	Ge og ra ph y_ Sp ai n	G e n d er - M al e
0				4 2	2	0 0 0	1	1	1	10 13 48 .8 8	0	0	0	0	0	0	0	0	0	0
1	2	1 5 6 4 7 3 1	6 0 8	4	1	8 3 8 0 7 8 6	1	0	1	11 25 42 .5	 0	0	0	0	0	0	0	0	1	0
2	3	1 5 6 1 9 3 0 4	5 0 2	4 2	8	5 9 6 6 0 8 0	3	1	0	11 39 31 .5 7	 0	0	0	0	0	0	0	0	0	0
3	4	1 5 7 0 1 3 5 4	6 9 9	3 9	1	0 0 0	2	0	0	93 82 6. 63	 0	0	0	0	0	0	0	0	0	0



 $5 \text{ rows} \times 2945 \text{ columns}$

8. Split the data into dependent and independent variables.

```
In [180]:
 df=pd.read csv('/content/Churn Modelling.csv')
                                                                                                                                           In [182]:
 print(df["Balance"].min())
 print(df["Balance"].max())
 print(df["Balance"].mean())
 0.0
 250898.09
 76485.889288
 print(df.count(0))

      Print (df.count(0))

      RowNumber
      10000

      CustomerId
      10000

      Surname
      10000

      CreditScore
      10000

      Geography
      10000

      Age
      10000

      Tenure
      10000

      Balance
      10000

      NumOfProducts
      10000

      HasCrCard
      10000

      IsActiveMember
      10000

      EstimatedSalary
      10000

      Exited
      10000

                                                                                                                                           In [183]:
                                    10000
 Exited
 dtype: int64
 print(df.shape)
 (10000, 14)
                                                                                                                                           In [184]:
 print(df.size)
 140000
                                                                                                                                           In [185]:
 X = df.iloc[:, :-1].values
 print(X)
 [[1 15634602 'Hargrave' ... 1 1 101348.88]
                                                                                                                                           In [187]:
   [2 15647311 'Hill' ... 0 1 112542.58]
[3 15619304 'Onio' ... 1 0 113931.57]
   [9998 15584532 'Liu' ... 0 1 42085.58]
   [9999 15682355 'Sabbatini' ... 1 0 92888.52]
   [10000 15628319 'Walker' ... 1 0 38190.78]]
 Y = df.iloc[:, -1].values
 print(Y)
 [1 0 1 ... 1 1 0]
                                                                                                                                           In [271]:
```

9. Scale the independent variables

In [215]: df = pd.read csv('/content/Churn Modelling.csv') x = df[['Age', 'Tenure']].values y = df['Gender'].values fig, ax = plt.subplots(ncols=2, figsize=(12, 4)) ax[0].scatter(x[:,0], y)ax[1].scatter(x[:,1], y)plt.show() In [216]: fig, ax = plt.subplots(figsize=(12, 4))ax.scatter(x[:,0], y)ax.scatter(x[:,1], y)<matplotlib.collections.PathCollection at 0x7f9a8a854ad0> Out[216]: fig, ax = plt.subplots(figsize=(12, 4)) In [217]: ax.hist(x[:,0])ax.hist(x[:,1])Out[217]: (array([413., 1035., 1048., 1009., 989., 1012., 967., 1028., 1025., array([0., 1., 2., 3., 4., 5., 6., 7., 8., 9., 10.]), <a list of 10 Patch objects>) In [220]: from sklearn.preprocessing import StandardScalerfrom sklearn.preprocessing import MinMaxScaler fig, ax = plt.subplots(figsize=(12, 4)) scaler = StandardScaler() x_std = scaler.fit_transform(x) ax.hist(x std[:,0])ax.hist(x std[:,1])Out[220]: (array([413., 1035., 1048., 1009., 2001., 0., 1995., 0., 1025., 1474.]), array([-1.73331549, -1.38753759, -1.04175968, -0.69598177, -0.35020386,

-0.00442596, 0.34135195, 0.68712986, 1.03290776, 1.37868567,

1.724463581),

<a list of 10 Patch objects>)

```
In [219]:
fig, ax = plt.subplots(figsize=(12, 4))
scaler = StandardScaler()
x std = scaler.fit transform(x)
ax.scatter(x_std[:,0], y)
ax.scatter(x std[:,1], y)
                                                                      Out[219]:
<matplotlib.collections.PathCollection at 0x7f9a8a2fde50>
fig, ax = plt.subplots(figsize=(12, 4))
                                                                       In [221]:
scaler = MinMaxScaler()
x_minmax = scaler.fit_transform(x)
ax.hist(x minmax [:,0])
ax.hist(x minmax [:,1])
                                                                       Out[221]:
(array([ 413., 1035., 1048., 1009., 989., 1012., 967., 1028., 1025.,
        1474.]),
 array([0., 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.]),
<a list of 10 Patch objects>)
                                                                       In [222]:
fig, ax = plt.subplots(figsize=(12, 4))
scaler = MinMaxScaler()
x minmax = scaler.fit transform(x)
ax.scatter(x minmax [:,0], y)
ax.scatter(x_minmax [:,1], y)
<matplotlib.collections.PathCollection at 0x7f9a8a0cae10>
                                                                       Out[222]:
fig, ax = plt.subplots(figsize=(12, 4))
                                                                       In [223]:
scaler = MinMaxScaler()
x minmax = scaler.fit transform(x)
ax.scatter(x minmax [:,0], y)
<matplotlib.collections.PathCollection at 0x7f9a8a0caf10>
                                                                       Out[223]:
fig, ax = plt.subplots(figsize=(12, 4))
scaler = MinMaxScaler()
                                                                        In [224]:
x minmax = scaler.fit transform(x)
ax.hist(x minmax [:,0])
```

```
Out[224]:
```

from sklearn.model_selection import train_test_split

from sklearn.pipeline import Pipeline

from sklearn.linear_model import SGDRegressor from sklearn.preprocessing import StandardScalerfrom sklearn.preprocessing import MinMaxScaler from sklearn.metrics import mean_absolute_error import sklearn.metrics as metrics

Mean Absolute Error: 57120.533393590835

Score 0.0004207814312172653

import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

Import Data

```
df = pd.read csv('/content/Churn Modelling.csv')
x = df[['Age', 'Tenure']].values
y = df['Balance'].values
# Split into a training and testing set
X train, X test, Y train, Y test = train test split(x, y)
# Define the pipeline for scaling and model fitting
pipeline = Pipeline([
    ("MinMax Scaling", MinMaxScaler()),
    ("SGD Regression", SGDRegressor())
])
# Scale the data and fit the model
pipeline.fit(X train, Y train)
# Evaluate the model
Y pred = pipeline.predict(X test)
print('Mean Absolute Error: ', mean absolute error(Y pred, Y test))
print('Score', pipeline.score(X_test, Y_test))
```

10.Split the data into training and testing

dataset = pd.read_csv('/content/Churn_Modelling.csv')
print(dataset)

In [267]:

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age
\							
0	1	15634602	Hargrave	619	France	Female	42
1	2	15647311	Hill	608	Spain	Female	41
2	3	15619304	Onio	502	France	Female	42
3	4	15701354	Boni	699	France	Female	39

4		5 1573	37888 M	itchell	850	-	Femal	e 43
9995 9996 9997 9998 9999	9 <u>9</u> 9 <u>9</u> 9 <u>9</u>	997 1556 998 1558 999 1568	06229 Ol 59892 Jol 84532 82355 Sal	nnstone Liu	772 516 709 772 792	France France France Germany	Male Female Male	e 39 e 35 e 36 e 42
0 1 2 3 4	2 1 8 1 2	0.00 83807.86 159660.80 0.00 125510.82		1 1 3 2 1	1 0 1 0 1	IsActiveMemk	1 1 0 0 1	
9995 9996 9997 9998 9999	7	0.00 57369.61 0.00 75075.31 130142.79		2 1 1 2 1	1 1 0 1		0 1 1 0 0	
datase	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	edSalary .01348.88 .12542.58 .13931.57 .93826.63 .79084.10 .96270.64 .01699.77 .42085.58 .92888.52 .38190.78 .14 column ["HasCrCard	1 0 1 0 0 0 0 1 1 0 0 ss]	,inplace=	=True)			In [287]: In [288]:
0 1 2 3 4 5 6 7 8	stomerId 15634602 15647311 15619304 15701354 15737888 15574012 15592531 15656148 15792365 15592389 timatedSa		core Age 619 42 608 41 502 42 699 39 850 43 645 44 822 50 376 29 501 44 684 27		Balance 0.00 83807.86 159660.80 0.00 125510.82 113755.78 0.00 115046.74 142051.07 134603.88	IsActiveMen	nber 1 1 0 0 0 1 0 1 1 1	

```
2
        113931.57
3
          93826.63
          79084.10
4
5
        149756.71
6
         10062.80
7
        119346.88
          74940.50
8
9
          71725.73
                                                                         In [289]:
X=dataset.iloc[:,:-1].values
Χ
                                                                        Out[289]:
array([[1.5634602e+07, 6.1900000e+02, 4.2000000e+01, 2.0000000e+00,
        0.0000000e+00, 1.0000000e+00],
       [1.5647311e+07, 6.0800000e+02, 4.1000000e+01, 1.0000000e+00,
        8.3807860e+04, 1.0000000e+00],
       [1.5619304e+07, 5.0200000e+02, 4.2000000e+01, 8.0000000e+00,
        1.5966080e+05, 0.0000000e+00],
       [1.5584532e+07, 7.0900000e+02, 3.6000000e+01, 7.0000000e+00,
        0.0000000e+00, 1.0000000e+00],
       [1.5682355e+07, 7.7200000e+02, 4.2000000e+01, 3.0000000e+00,
        7.5075310e+04, 0.0000000e+00],
       [1.5628319e+07, 7.9200000e+02, 2.8000000e+01, 4.0000000e+00,
        1.3014279e+05, 0.0000000e+00]])
                                                                         In [290]:
Y=dataset.iloc[:,-1].values
array([101348.88, 112542.58, 113931.57, ..., 42085.58, 92888.52,
                                                                        Out[290]:
        38190.781)
from sklearn.model_selection import train_test_split
                                                                         In [291]:
X train, X test, Y train, Y test = train test split( X, Y, test size = 0.25,
random state = 0 )
                                                                         In [306]:
from sklearn.preprocessing import StandardScaler
sc=StandardScaler()
X train = sc.fit transform(X train)
X test = sc.transform(X test)
print(X train)
[[-1.34333028 - 0.73550706 \ 0.01526571 \ 0.00886037 \ 0.67316003 - 1.03446007]
 [-0.65515619 \quad 0.80829492 \quad -0.46178778 \quad 1.39329338 \quad -0.35693706 \quad 0.96668786]
 . . .
  [-1.63542994 \quad 0.90092304 \quad -0.36637708 \quad 0.00886037 \quad 1.36657199 \quad -1.03446007] 
 [-0.38540456 -0.62229491 -0.08014499 1.39329338 -1.20772417 0.96668786]
 [-1.37829524 -0.28265848 0.87396199 -1.37557264 0.51741687 -1.03446007]]
                                                                         In [305]:
print(X_test)
 \begin{bmatrix} [-1.05852196 & -0.55025082 & -0.36637708 & 1.04718513 & 0.88494297 & 0.96668786] \end{bmatrix} 
 [-0.51554728 -1.31185979 \ 0.11067641 -1.02946438 \ 0.43586703 -1.03446007]
 [-0.8058485 \quad 0.57157862 \quad 0.3014978 \quad 1.04718513 \quad 0.31486378 \quad 0.96668786]
```