Assignment - 2 Data Visualization and pre-processing

Assignment submission	26 September 2022
Student Name	Sebastian prabu.M
Student Roll Number	951919CS087
Maximum Marks	2 Marks

1. Download the Dataset

2. Import required library

import pandas as pd import seaborn as sns import matplotlib.pyplot as plt import numpy as np sns.set_style('darkgrid') sns.set(font_scale=1.3)

3. Read dataset and do pre-processing

df=pd.read_csv("/content/drive/MyDrive/IBM/Assignment - 2 /Churn_Modelling.csv") df.head()

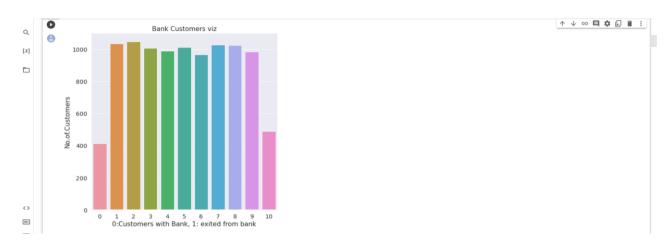
	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	1	15634602	Hargrave	619	France	Female	42	2	0.00	1	1	1	101348.88	1
1	2	15647311	Hill	608	Spain	Female	41	1	83807.86	1	0	1	112542.58	0
2	! 3	15619304	Onio	502	France	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	15737888	Mitchell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

Drop the columns that are not required for the neural network. df.drop(["RowNumber","CustomerId","Surname"],axis=1,inplace=True) df.info()

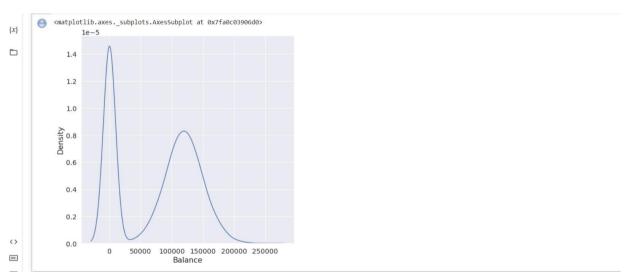
```
Data columns (total 11 columns):
            # Column
                                Non-Null Count Dtype
                                10000 non-null int64
            0 CreditScore
                                 10000 non-null object
             1 Geography
                                 10000 non-null object
                                 10000 non-null int64
             3 Age
                                 10000 non-null int64
            5 Balance 10000 non-null float64
6 NumOfProducts 10000 non-null int64
            7 HasCrCard 10000 non-null int64
8 IsActiveMember 10000 non-null int64
               EstimatedSalary 10000 non-null float64
           10 Exited 10000 non-null ir dtypes: float64(2), int64(7), object(2)
                                 10000 non-null int64
1
            memory usage: 859.5+ KB
```

4.A. Perform Univariate Analysis

```
plt.figure(figsize=(8,8))
sns.countplot(x='Tenure',data=df)
plt.xlabel('0:Customers with Bank, 1: exited from bank')
plt.ylabel('No.of.Customers')
plt.title("Bank Customers viz")
plt.show()
```



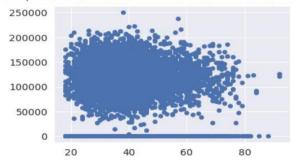
plt.figure(figsize=(8,8))
sns.kdeplot(x=df['Balance'])



4.B. Perform Bi-variate Analysis

plt.scatter(df.Age,df.Balance)

<matplotlib.collections.PathCollection at 0x7fa0d35a7dd0>



df.corr()

	CreditScore	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
CreditScore	1.000000	0.007888	-0.003965	0.000842	0.006268	0.012238	-0.005458	0.025651	-0.001384	-0.027094
Gender	0.007888	1.000000	0.022812	0.003739	0.069408	0.003972	-0.008523	0.006724	-0.001369	0.035943
Age	-0.003965	0.022812	1.000000	-0.009997	0.028308	-0.030680	-0.011721	0.085472	-0.007201	0.285323
Tenure	0.000842	0.003739	-0.009997	1.000000	-0.012254	0.013444	0.022583	-0.028362	0.007784	-0.014001
Balance	0.006268	0.069408	0.028308	-0.012254	1.000000	-0.304180	-0.014858	-0.010084	0.012797	0.118533
NumOfProducts	0.012238	0.003972	-0.030680	0.013444	-0.304180	1.000000	0.003183	0.009612	0.014204	-0.047820
HasCrCard	-0.005458	-0.008523	-0.011721	0.022583	-0.014858	0.003183	1.000000	-0.011866	-0.009933	-0.007138
IsActiveMember	0.025651	0.006724	0.085472	-0.028362	-0.010084	0.009612	-0.011866	1.000000	-0.011421	-0.156128
EstimatedSalary	-0.001384	-0.001369	-0.007201	0.007784	0.012797	0.014204	-0.009933	-0.011421	1.000000	0.012097
Exited	-0.027094	0.035943	0.285323	-0.014001	0.118533	-0.047820	-0.007138	-0.156128	0.012097	1.000000

#Perform Bivariate Analysis import statsmodels.api as sm

#define response variable y = df['CreditScore']

#define explanatory variable x = df[['EstimatedSalary']]

#add constant to predictor variables x = sm.add_constant(x)

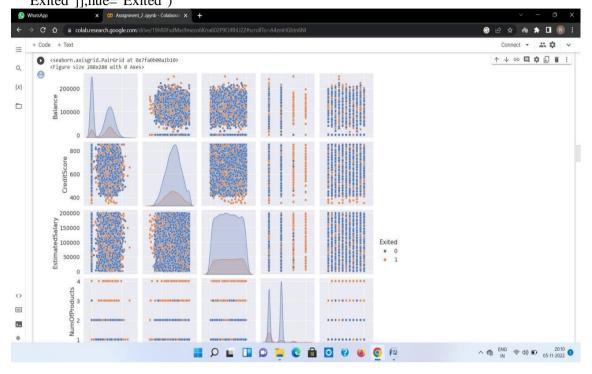
#fit linear regression model model = sm.OLS(y, x).fit()

#view model summary
print(model.summary())

4.C. Perform Multi-variate Analysis

#Perform Multivariate Analysis plt.figure(figsize=(4,4))

sns.pairplot(data=df[["Balance","CreditScore","EstimatedSalary","NumOfProducts","Tenure", "Exited"]],hue="Exited")



5. Perform descriptive statistics on the datasets:

df=pd.DataFrame(df)
print(df.sum())



```
#Perform Descriptive Statistics
print("----Sum Value-----")
print(df.sum(1))
print("-----Product Value-----")
print(df.prod())
print("______")
```

```
{x}
                   ----Sum Value----
0 102015.88
1 197002.44
2 274149.37
            0
94567.63
205492.92
                   " 203492,92

9995 97088.64

9996 159633.38

9997 42840.58

9998 168784.83

9999 169159.57

Length: 10000, dtype: float64
                    -----Product Value-
CreditScore
                    Creditscore
Age
Tenure
Balance
NumOfProducts
HasCrCard
IsActiveMember
EstimatedSalary
Exited
dtype: float64
                   /usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:3: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a futur This is separate from the ipykernel package so we can avoid doing imports until /usr/local/lib/python3.7/dist-packages/numpy/core/_methods.py:52: RuntimeWarning: overflow encountered in reduce return umr-prod(a, axis, dtype, out, keepdims, initial, where) /usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:6: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a futur
#Perform Descriptive Statistics
    print("-----")
    print(df.mean())
    print("____")
    print("----- Median Value ----- ")
    print(df.median())
    print("_____")
print("-----Mode Value -----")
    print(df.mode())
    ↑ ↓ ∞ □ ‡ ☐ i :
            Creditscore 650.528800

Age 38.921800

Fenure 5.012800

Balance 76485.889288

NumOFProducts 1.530200

HasCrCard 0.765500

IsActiveNember 0.5151000

EstimatedSalary 100090.239881
{x}
Exited
dtype: float64
                                                           0.203700
                  Creditscore
Age
Tenure
Balance
NumofProducts
HasCrCard
IsActiveMember
Estimatedsalary
Exited
dtype: float64
                        CreditScore Geography Gender Age Tenure Balance NumOfProducts \
850 France Male 37 2 0.0 1
                    HasCrCard IsActiveMember EstimatedSalary Exited 0 1 1 24924.92 0
                   /usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:3: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a futur This is separate from the ipykernel package so we can avoid doing imports until /usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:6: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a futur
<>
```

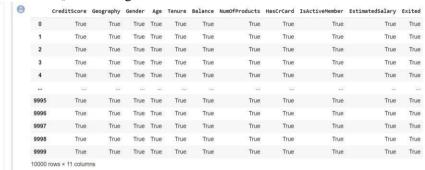
6. Handle the missing values:

df.isnull()#Checking values are null

0	False											
4	False	False			False							
3	raise											
2	False											
3	False											
4	False											
	(888)	****	(89)	300	***	111	(***)	***	3399	(888)	15.500	
9995	False											
9996	False											
9997	False											
9998	False											
9999	False											

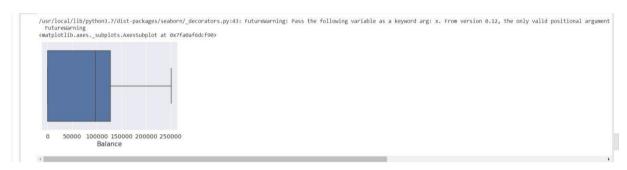
#Handling with missing Values

df.notnull()#Checking values are not null



7. Find outlier and replace the outlier:

sns.boxplot(df['Balance'])



print(np.where(df['Balance']>100000))

```
Q (array([ 2, 4, 5, ..., 9987, 9993, 9999]),)
```

#Find outliers & replace the outliers from scipy import stats import numpy as np

z = np.abs(stats.zscore(df["EstimatedSalary"]))
print(z)

```
0 0.021886

1 0.216534

2 0.240687

3 0.108918

4 0.365276

....

9995 0.066419

9996 0.027988

9997 1.080643

9998 0.12531

9999 1.076370

Name: EstimatedSalary, Length: 10000, dtype: float64
```

8. Check for categorical columns & performs encoding:

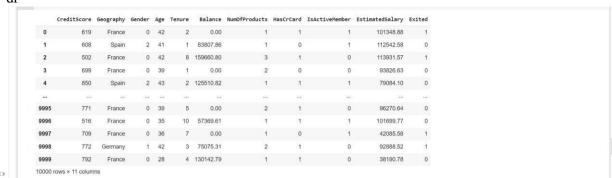
from sklearn.preprocessing import LabelEncoder df['Gender'].unique()

```
array(['Female', 'Male'], dtype=object)
```

df['Gender'].value_counts()

```
Male 5457
Female 4543
Name: Gender, dtype: int64
```

```
#Check for categorical columns & performs encoding encoding=LabelEncoder()
df["Gender"]=encoding.fit_transform(df.iloc[:,1].values)
df
```



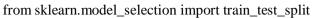
9. Split the data into Dependent & Independent Variables:

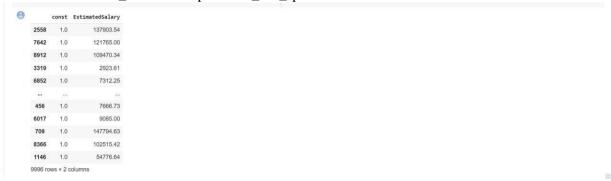
10. Scale the independent variables:

```
from sklearn.preprocessing import StandardScaler
object= StandardScaler()
# standardization
scale = object.fit_transform(df)
print(scale)

© [[-0.3262142 0.29351742 -1.04175968 ... 0.97024255 0.02188649
1.97716468]
[-0.44003595 0.19816383 -1.38753759 ... 0.97024255 0.021653375
-0.96577476]
[-1.53679418 0.29351742 1.03290776 ... -1.03067011 0.2406869
1.97716468]
...
[ 0.66498839 -0.27860412 0.68712986 ... 0.97024255 -1.08064308
1.97716468]
[ 1.25683526 0.29351742 -0.69598177 ... -1.03067011 -1.07636976
-0.56577476]]
```

11. Split the data into training & testing:





y_train ② 2558 727 7642 811 8912 623 3319 439 6852 680 ... 456 733 6017 487 709 686 8366 637 1146 614 Name: CreditScore, Length: 9996, dtype: int64

```
y_test

1603 576

8713 786

4561 562

6600 505

Name: CreditScore, dtype: int64
```