

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

In [2]:

```
#load the dataset
santhos=pd.read_csv("/content/Churn_Modelling.csv")
santhos
```

Out[2]:

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geog raph y	Ge nd er	A ge	Te nu re	Bala nce	NumOf Produc ts	Has CrC ard	IsActiv eMemb er	Estimat edSalar y	Ex ite d
0	1	15634602	Har grav e	619	Fran ce	Fe mal e	42	2	0.00	1	1	1	101348.88	1
1	2	15647311	Hill	608	Spai n	Fe mal e	41	1	83807.86	1	0	1	112542.58	0
2	3	15619304	Oni o	502	Fran ce	Fe mal e	42	8	159660.80	3	1	0	113931.57	1
3	4	15701354	Bon i	699	Fran ce	Fe mal e	39	1	0.00	2	0	0	93826.63	0
4	5	15737888	Mit chel l	850	Spai n	Fe mal e	43	2	125510.82	1	1	1	79084.10	0
...
9995	9996	15606229	Obij iaku	771	Fran ce	Ma le	39	5	0.00	2	1	0	96270.64	0
9996	9997	15569892	Joh nstone	516	Fran ce	Ma le	35	10	57369.61	1	1	1	101699.77	0

	Row Number	Customer Id	Sur name	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOf Products	Has CrCard	IsActiveMember	EstimatedSalary	Exited
9997	9998	15584532	Liu	709	France	Female	36	7	0.00	1	0	1	42085.58	1
9998	9999	15682355	Sab bati ni	772	Germany	Male	42	3	75075.31	2	1	0	92888.52	1
9999	10000	15628319	Wal ker	792	France	Female	28	4	130142.79	1	1	0	38190.78	0

10000 rows × 14 columns

In [3]:
santhose.head

Out[3]:
In [4]:

santhose.shape

(10000, 14)

Out[4]:

In [5]:
santhose.std()
/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[5]:

```

RowNumber      2886.895680
CustomerId     71936.186123
CreditScore     96.653299
Age            10.487806
Tenure         2.892174
Balance       62397.405202
NumOfProducts   0.581654
HasCrCard      0.455840
IsActiveMember  0.499797
EstimatedSalary 57510.492818
Exited         0.402769
dtype: float64

```

In [6]:

```
santhose.median()
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1: FutureWarning
: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=No
ne') is deprecated; in a future version this will raise TypeError.  Select on
ly valid columns before calling the reduction.
```

```
"""Entry point for launching an IPython kernel.
```

Out[6]:

```
RowNumber      5.000500e+03
CustomerId     1.569074e+07
CreditScore    6.520000e+02
Age            3.700000e+01
Tenure         5.000000e+00
Balance        9.719854e+04
NumOfProducts  1.000000e+00
HasCrCard      1.000000e+00
IsActiveMember 1.000000e+00
EstimatedSalary 1.001939e+05
Exited         0.000000e+00
dtype: float64
```

In [7]:

```
sns.scatterplot(x=santhose.index,y=santhose['EstimatedSalary'],hue=santhose['
Age'])
```

Out[7]:

In [8]:

```
santhose.columns
```

Out[8]:

```
Index(['RowNumber', 'CustomerId', 'Surname', 'CreditScore', 'Geography',
      'Gender', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard',
      'IsActiveMember', 'EstimatedSalary', 'Exited'],
      dtype='object')
```

In [9]:

```
santhose.info()
```

```
RangeIndex: 10000 entries, 0 to 9999
```

```
Data columns (total 14 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	Gender	10000 non-null	object
6	Age	10000 non-null	int64
7	Tenure	10000 non-null	int64
8	Balance	10000 non-null	float64
9	NumOfProducts	10000 non-null	int64
10	HasCrCard	10000 non-null	int64
11	IsActiveMember	10000 non-null	int64
12	EstimatedSalary	10000 non-null	float64

```
13 Exited          10000 non-null int64
dtypes: float64(2), int64(9), object(3)
memory usage: 1.1+ MB
```

In [10]:

```
#statistical description of the dataset
santhose.describe()
```

Out[10]:

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCreditCard	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	1.569094e+07	650.528800	38.921800	5.012800	76485.889288	1.530200	0.70550	0.515100	100090.239881	0.203700
std	2886.89568	7.193619e+04	96.653299	10.487806	2.892174	62397.405202	0.581654	0.45584	0.499797	57510.492818	0.402769
min	1.000000	1.556570e+07	350.000000	18.000000	0.000000	0.000000	1.000000	0.000000	0.000000	11.580000	0.000000
25%	2500.75000	1.562853e+07	584.000000	32.000000	3.000000	0.000000	1.000000	0.000000	0.000000	51002.110000	0.000000
50%	5000.50000	1.569074e+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	718.000000	44.000000	7.000000	127644.240000	2.000000	1.000000	1.000000	149388.247500	0.000000
max	10000.00000	1.581569e+07	850.000000	92.000000	10.000000	250898.090000	4.000000	1.000000	1.000000	199992.480000	1.000000

In [11]:

```
sns.histplot(x='EstimatedSalary',data=santhose,color='skyblue') #univariate analysis
```

Out[11]:

In [12]:

```
santhose.mean()

/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[12]:

```
RowNumber      5.000500e+03
CustomerId     1.569094e+07
CreditScore    6.505288e+02
Age            3.892180e+01
Tenure         5.012800e+00
Balance        7.648589e+04
NumOfProducts  1.530200e+00
HasCrCard      7.055000e-01
IsActiveMember 5.151000e-01
EstimatedSalary 1.000902e+05
Exited         2.037000e-01
dtype: float64
```

In [13]:

```
santhose.dtypes
```

Out[13]:

```
RowNumber      int64
CustomerId     int64
Surname        object
CreditScore    int64
Geography      object
Gender         object
Age            int64
Tenure         int64
Balance        float64
NumOfProducts  int64
HasCrCard      int64
IsActiveMember int64
EstimatedSalary float64
Exited         int64
dtype: object
```

In [14]:

```
santhose.corr() #bivariate analysis
```

Out[14]:

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
RowNumber	1.000000	0.004202	0.005840	0.000783	-0.006495	-0.009067	0.007246	0.000599	0.012044	-0.005988	-0.016571

	RowN umber	Custo merId	Credit Score	Age	Ten ure	Bala nce	NumOfP roducts	HasC rCard	IsActive Member	Estimate dSalary	Exit ed
CustomerId	0.004202	1.000000	0.005308	0.009497	-0.014883	-0.012419	0.016972	-0.014025	0.001665	0.015271	-0.006248
CreditScore	0.005840	0.005308	1.000000	-0.003965	0.000842	0.006268	0.012238	-0.005458	0.025651	-0.001384	-0.027094
Age	0.000783	0.009497	-0.003965	1.000000	-0.009997	0.028308	-0.030680	-0.011721	0.085472	-0.007201	0.285323
Tenure	-0.006495	-0.014883	0.000842	-0.009997	1.000000	-0.012254	0.013444	0.022583	-0.028362	0.007784	-0.014001
Balance	-0.009067	-0.012419	0.006268	0.028308	-0.012254	1.000000	-0.304180	-0.014858	-0.010084	0.012797	0.118533
NumOfProducts	0.007246	0.016972	0.012238	-0.030680	0.013444	-0.304180	1.000000	0.003183	0.009612	0.014204	-0.047820
HasCrCard	0.000599	-0.014025	-0.005458	-0.011721	0.022583	-0.014858	0.003183	1.000000	-0.011866	-0.009933	-0.007138
IsActiveMember	0.012044	0.001665	0.025651	0.085472	-0.028362	0.010084	0.009612	-0.011866	1.000000	-0.011421	-0.156128
EstimatedSalary	-0.005988	0.015271	-0.001384	-0.007201	0.000842	0.012797	0.014204	-0.009933	-0.011421	1.000000	0.012097
Exited	-0.016571	-0.006248	-0.027094	0.285323	-0.014001	0.118533	-0.047820	-0.007138	-0.156128	0.012097	1.000000

In [15]:

```
santhose.duplicated().sum()
```

Out[15]:

0

In [17]:

```
#statistical descripton of the dataset

df = pd.DataFrame({'customer id': [1, 1, 1, 2, 2, 2, 3, 3, 3, 3,
                                   3, 4, 4, 5, 5, 6, 6, 6, 7, 8],
                   'Estimated salary (in thousands)': [75, 66, 68, 74, 78,
72, 85, 82, 90, 82,
                                   80, 88, 85, 90, 92, 94, 94, 88, 91, 96]})
```

```
print(df)
```

	customer id	Estimated salary (in thousands)
0	1	75
1	1	66
2	1	68
3	2	74
4	2	78
5	2	72
6	3	85
7	3	82
8	3	90
9	3	82
10	3	80
11	4	88
12	4	85
13	5	90
14	5	92
15	6	94
16	6	94
17	6	88
18	7	91
19	8	96

In [18]:

```
sns.regplot(x='HasCrCard',y='EstimatedSalary',data=santhose)
plt.ylim(0,10000)
```

Out[18]:

```
(0.0, 10000.0)
```

In [19]:

```
#multivariant
sns.heatmap(santhose.corr(),annot=True)
```

Out[19]:

In [20]:

```
santhose.isnull()
```

Out[20]:

	Row Number	Customer Id	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
9999	False	False	False	False	False	False	False	False	False	False	False	False	False	False

10000 rows × 14 columns

In [21]:

`santhose.dropna()`

Out[21]:

	Row Number	Customer Id	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	Michell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0
...
99995	9996	15606229	Obijaku	771	France	Male	39	5	0.00	2	1	0	96270.64	0

	Row Number	Customer Id	Sur name	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOf Products	Has CrCard	IsActiveMember	EstimatedSalary	Exited
9996	9997	15569892	Johnstone	516	France	Male	35	10	57369.61	1	1	1	101699.77	0
9997	9998	15584532	Liu	709	France	Female	36	7	0.00	1	0	1	42085.58	1
9998	9999	15682355	Sabatinini	772	Germany	Male	42	3	75075.31	2	1	0	92888.52	1
9999	10000	15628319	Walker	792	France	Female	28	4	130142.79	1	1	0	38190.78	0

10000 rows × 14 columns

In [22]:
santhose.fillna(0)

	Row Number	Customer Id	Sur name	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOf Products	Has CrCard	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

	Row Number	Cust omer Id	Sur na me	Cred itSco re	Geog raph y	Ge nd er	A ge	Te nu re	Bala nce	NumOf Produc ts	Has CrC ard	IsActiv eMemb er	Estimat edSalar y	Ex ite d
3	4	1570 1354	Bon i	699	Fran ce	Fe mal e	3 9	1	0.00	2	0	0	93826.6 3	0
4	5	1573 7888	Mit chel l	850	Spai n	Fe mal e	4 3	2	1255 10.8 2	1	1	1	79084.1 0	0
...
9 9 9 5	9996	1560 6229	Obij iaku	771	Fran ce	Ma le	3 9	5	0.00	2	1	0	96270.6 4	0
9 9 9 6	9997	1556 9892	Joh nsto ne	516	Fran ce	Ma le	3 5	10	5736 9.61	1	1	1	101699. 77	0
9 9 9 7	9998	1558 4532	Liu	709	Fran ce	Fe mal e	3 6	7	0.00	1	0	1	42085.5 8	1
9 9 9 8	9999	1568 2355	Sab bati ni	772	Ger many	Ma le	4 2	3	7507 5.31	2	1	0	92888.5 2	1
9 9 9 9	10000	1562 8319	Wal ker	792	Fran ce	Fe mal e	2 8	4	1301 42.7 9	1	1	0	38190.7 8	0

10000 rows × 14 columns

In [23]:

```
#finding the outerlier
sns.boxplot(y='EstimatedSalary',data=santhose)
```

Out[23]:

In [34]:

```
outliers_low = (df<santhose)
```

```
print(santhose)
```

	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	15737888	Mitchell	850	Spain	Female	43	
...	
9995	9996	15606229	Obijiaku	771	France	Male	39	
9996	9997	15569892	Johnstone	516	France	Male	35	
9997	9998	15584532	Liu	709	France	Female	36	
9998	9999	15682355	Sabbatini	772	Germany	Male	42	
9999	10000	15628319	Walker	792	France	Female	28	

	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	\
0	2	0.00	1	1	1	
1	1	83807.86	1	0	1	
2	8	159660.80	3	1	0	
3	1	0.00	2	0	0	
4	2	125510.82	1	1	1	
...	
9995	5	0.00	2	1	0	
9996	10	57369.61	1	1	1	
9997	7	0.00	1	0	1	
9998	3	75075.31	2	1	0	
9999	4	130142.79	1	1	0	

	EstimatedSalary	Exited
0	101348.88	1
1	112542.58	0
2	113931.57	1
3	93826.63	0
4	79084.10	0
...
9995	96270.64	0
9996	101699.77	0
9997	42085.58	1
9998	92888.52	1
9999	38190.78	0

```
[10000 rows x 14 columns]
```

In [28]:

```
#categorical columns
import numpy as np

import pandas as pd

# Categorical using dtype

c = pd.Series(["a", "b", "d", "a", "d"], dtype="category")

print ("\nCategorical without pandas.Categorical() : \n", c)
```

```

c1 = pd.Categorical([1, 2, 3, 1, 2, 3])

print ("\n\nc1 : ", c1)

c2 = pd.Categorical(['e', 'm', 'f', 'i',
                    'f', 'e', 'h', 'm' ])

print ("\nc2 : ", c2)

Categorical without pandas.Categorical() :
0      a
1      b
2      d
3      a
4      d
dtype: category
Categories (3, object): ['a', 'b', 'd']

c1 :  [1, 2, 3, 1, 2, 3]
Categories (3, int64): [1, 2, 3]

c2 :  ['e', 'm', 'f', 'i', 'f', 'e', 'h', 'm']
Categories (5, object): ['e', 'f', 'h', 'i', 'm']

#split the data
X = df.iloc[:, :-1].values
print(X)

[[1]
 [1]
 [1]
 [2]
 [2]
 [2]
 [3]
 [3]
 [3]
 [3]
 [3]
 [4]
 [4]
 [5]
 [5]
 [6]
 [6]
 [6]
 [7]
 [8]]

```

In [29]:

In [30]:

```
#scale the independent
import pandas as pd

from sklearn.preprocessing import StandardScaler

# Read Data from CSV

santhose =pd.read_csv('Churn_Modelling.csv')
santhose.head()

# Initialise the Scaler

scaler = StandardScaler()

print(santhose)
```

	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	15737888	Mitchell	850	Spain	Female	43	
...	
9995	9996	15606229	Obijiaku	771	France	Male	39	
9996	9997	15569892	Johnstone	516	France	Male	35	
9997	9998	15584532	Liu	709	France	Female	36	
9998	9999	15682355	Sabbatini	772	Germany	Male	42	
9999	10000	15628319	Walker	792	France	Female	28	

	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	\
0	2	0.00	1	1	1	
1	1	83807.86	1	0	1	
2	8	159660.80	3	1	0	
3	1	0.00	2	0	0	
4	2	125510.82	1	1	1	
...	
9995	5	0.00	2	1	0	
9996	10	57369.61	1	1	1	
9997	7	0.00	1	0	1	
9998	3	75075.31	2	1	0	
9999	4	130142.79	1	1	0	

	EstimatedSalary	Exited
0	101348.88	1
1	112542.58	0
2	113931.57	1
3	93826.63	0
4	79084.10	0
...
9995	96270.64	0
9996	101699.77	0
9997	42085.58	1

```
9998          92888.52          1
9999          38190.78          0
```

```
[10000 rows x 14 columns]
```

In [31]:

```
#split the data into training
import pandas as pd

from sklearn.linear_model import LinearRegression

from sklearn.model_selection import train_test_split

# read the dataset

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

# get the locations

X = df.iloc[:, :-1]

y = df.iloc[:, -1]

# split the dataset

X_train, X_test, y_train, y_test = train_test_split(

    X, y, test_size=0.05, random_state=0)
print(df)
```

	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	15737888	Mitchell	850	Spain	Female	43	
...	
9995	9996	15606229	Obijiaku	771	France	Male	39	
9996	9997	15569892	Johnstone	516	France	Male	35	
9997	9998	15584532	Liu	709	France	Female	36	
9998	9999	15682355	Sabbatini	772	Germany	Male	42	
9999	10000	15628319	Walker	792	France	Female	28	

	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	\
0	2	0.00	1	1	1	
1	1	83807.86	1	0	1	
2	8	159660.80	3	1	0	
3	1	0.00	2	0	0	
4	2	125510.82	1	1	1	
...	
9995	5	0.00	2	1	0	
9996	10	57369.61	1	1	1	
9997	7	0.00	1	0	1	
9998	3	75075.31	2	1	0	

9999	4	130142.79	1	1	0
------	---	-----------	---	---	---

	EstimatedSalary	Exited
0	101348.88	1
1	112542.58	0
2	113931.57	1
3	93826.63	0
4	79084.10	0
...
9995	96270.64	0
9996	101699.77	0
9997	42085.58	1
9998	92888.52	1
9999	38190.78	0

[10000 rows x 14 columns]