

Steps for preprocessing

1.Import the necessary libraries 2.Import the dataset 3.Handling the null values(churn-leaving the company) 4.Separate dependent and independent variables(dep,output-input) 5.Encoding 6.Splitting into training and testing set 7.Feature scaling

D.Harshita VIT-AP

```
In [1]: pip install seaborn

Requirement already satisfied: seaborn in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (0.12.2)
Requirement already satisfied: numpy<24.0, >=1.17 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from seaborn) (1.24.3)
Requirement already satisfied: pandas<0.25 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from seaborn) (2.0.3)
Requirement already satisfied: matplotlib<3.6.1, >=3.1 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from seaborn) (3.7.2)
Requirement already satisfied: contourpy<1.0.1 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from matplotlib<3.6.1, >=3.1->seaborn) (1.1.0)
Requirement already satisfied: cycler<0.10.0, >=0.1 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from matplotlib<3.6.1, >=3.1->seaborn) (0.11.0)
Requirement already satisfied: fonttools<=4.22.0 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from matplotlib<3.6.1, >=3.1->seaborn) (4.41.0)
Requirement already satisfied: pyparsing<=3.0.2023.1 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from matplotlib<3.6.1, >=3.1->seaborn) (3.1.4)
Requirement already satisfied: pillow<=6.2.0 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from matplotlib<3.6.1, >=3.1->seaborn) (10.0.0)
Requirement already satisfied: pytz<=2023.1 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from pandas<0.25->seaborn) (2023.3)
Requirement already satisfied: tzdata<=2022.1 in c:\users\harsh\appdata\local\programs\python\python311\lib\site-packages (from pandas<0.25->seaborn) (2023.3)
Note: you may need to restart the kernel to use updated packages.
```

1.Import the necessary libraries

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

2.Import the dataset

```
In [4]: # (csv, .csv, .json, .excel)
dataset=pd.read_csv("Titanic-Dataset.csv")
```

```
In [5]: dataset

Out[5]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	3	Johnson, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W/C 6607	23.4500	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows x 12 columns

```
In [6]: dataset.head(3)

Out[6]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	NaN	S

```
In [7]: dataset.tail() #default last 5 rows

Out[7]:
```

886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00	B42	S
888	889	0	3	Johnson, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W/C 6607	23.45	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75	NaN	Q

```
In [8]: dataset.shape

Out[8]: (891, 12)
```

```
In [9]: dataset.info()

Out[9]:
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 0   Column      Non-Null Count  Dtype
---  --
 0   PassengerId   891 non-null    int64
 1   Survived      891 non-null    int64
 2   Pclass        891 non-null    int64
 3   Name          891 non-null    object
 4   Sex           891 non-null    object
 5   Age           714 non-null    float64
 6   SibSp         891 non-null    int64
 7   Parch         891 non-null    int64
 8   Ticket        891 non-null    object
 9   Fare          891 non-null    float64
10   Cabin         284 non-null    object
11   Embarked      889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

```
In [10]: dataset.describe()

Out[10]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.698118	0.523008	0.381594	32.204208
std	257.353842	0.486982	0.836071	14.526487	1.102743	0.800587	49.493429
min	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000	0.000000	7.914000
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	668.500000	1.000000	3.000000	30.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.339200

colted dependent variable -output

inputs-independent variables

```
In [11]: numeric_dataset = dataset.select_dtypes(include=[np.number])
numeric_dataset corr()
#we can write dataset.corr() or dataset.select_dtypes(include=[np.number]).corr()
```

```
Out[11]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
PassengerId	1.000000	-0.050007	-0.038144	0.036647	-0.057527	-0.001852	0.012658
Survived	-0.050007	1.000000	-0.334881	-0.077221	-0.035322	0.081629	0.257307
Pclass	-0.038144	-0.334881	1.000000	-0.309226	0.093981	0.013443	-0.549500
Age	0.036647	-0.077221	-0.309226	1.000000	-0.300347	0.389119	0.096907
SibSp	-0.057527	-0.035322	0.093981	-0.300347	1.000000	0.414838	0.159951
Parch	-0.001852	0.081629	0.013443	-0.189118	0.414838	1.000000	0.216225
Fare	0.012658	0.257307	-0.549500	0.096907	0.159951	0.216225	1.000000

```
In [16]: plt.subplots(figsize=(20,15))
sns.heatmap(numeric_dataset.corr(),annot = True)
```



```
In [17]: dataset.Parch.value_counts()

Out[17]:
```

Parch	count
0	678
1	118
2	88
3	5
4	5
5	4
6	1

Name: count, dtype: int64

3.Handling null values

```
In [18]: dataset.isnull().any()

Out[18]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	False	False	False	True	True	True	True	True	True	True	False	True
1	False	True	False	True	True	True	True	True	True	True	True	True
2	False	True	False	True	True	True	True	True	True	True	True	True
3	False	True	False	True	True	True	True	True	True	True	True	True
4	False	False	False	True	True	True	True	True	True	True	False	True
...
886	False	False	False	True	True	True	True	True	True	True	False	True
887	False	True	False	True	True	True	True	True	True	True	True	True
888	False	False	False	True	True	True	True	True	True	True	False	True
889	False	True	False	True	True	True	True	True	True	True	True	True
890	False	True	False	True	True	True	True	True	True	True	True	True
891	False	False	False	True	True	True	True	True	True	True	False	True

```
In [19]: dataset.isnull().sum()

Out[19]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0
...
886	0	0	0	0	0	0	0
887	0	0	0	0	0	0	0
888	0	0	0	0	0	0	0
889	0	0	0	0	0	0	0
890	0	0	0	0	0	0	0
891	0	0	0	0	0	0	0

```
In [20]: dataset.head()

Out[20]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

4.outliers

```
In [21]: sns.boxplot(dataset.Pclass)

Out[21]: <Axes: >
```



5.Separate dependent and independent variables

```
In [26]: y=dataset.iloc[:,3:13]
X=dataset.iloc[:,0:3:14]
print(x.head())
print(y.tail())
```

```
Out[26]:
```

	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	NaN	S
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
Out[26]:
```

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S

```
Out[26]:
```

	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
886	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	Johnson, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W/C 6607	23.4500	NaN	S
889	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C
890	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

```
Out[26]:
```

	Ticket	Fare	Cabin	Embarked
886	211536	13.00	NaN	S
887	112053	30.00	B42	S
888	W/C 6607	23.45	NaN	S
889	111369	30.00	C148	C
890	370376	7.75	NaN	Q

```
In [27]: print(y.head())
print(y.tail())
```

```
Out[27]:
```

	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282	7.9250	NaN	S
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
Out[27]:
```

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S

```
Out[27]:
```

	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
886	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	Johnson, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W/C 6607	23.4500	NaN	S
889	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C
890	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

```
Out[27]:
```

	Ticket	Fare	Cabin	Embarked
886	211536	13.00	NaN	S
887	112053			

