

21BCE7139 VIT-AP

Import the Libraries

```
In [3]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

Import the Dataset

In [4]: df=pd.read_csv("Titanic-Dataset.csv")

In [5]: df.head()

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  Cumings, Mrs. John Bradley Florence Briggs Th...  female  38.0    1    0  PC17599   71.2833   C85    C
2            3         1       3  Heikinen, Mrs. Lena   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   373650   8.0500   NaN    S

In [6]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   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           204 non-null    object
11  Embarked         891 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB

In [7]: df.describe()

Out[7]:
   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   29.699118   0.523008   0.381594   32.204208
std     257.353842   0.486592   0.838671   14.526497   1.102743   0.806067   49.693429
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   7.910400
50%     446.000000   0.000000   3.000000   28.000000   0.000000   0.000000   14.454200
75%     698.500000   1.000000   3.000000   38.000000   1.000000   0.000000   31.000000
max     891.000000   1.000000   3.000000   80.000000   8.000000   6.000000   512.329200

In [8]: df.corr()

C:\Users\AKARSHA\AppData\Local\Temp\ipykernel_9232\1134722445.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
df.corr()

Out[8]:
   PassengerId  Survived  Pclass    Age    SibSp    Parch    Fare
PassengerId   1.000000  -0.05007  -0.03544  0.03847  -0.05737  -0.00162  0.012658
Survived       -0.05007  1.000000  -0.338481  -0.07721  -0.03322  0.081629  -0.257397
Pclass         -0.03544  -0.338481  1.000000  -0.369226  0.083081  0.018443  -0.549500
Age            0.03847  -0.07721  -0.369226  1.000000  -0.308247  0.189119  -0.096067
SibSp          -0.05737  -0.03322  0.083081  -0.308247  1.000000  0.414838  -0.159651
Parch         -0.00162  0.081629  0.018443  -0.189119  0.414838  1.000000  0.216225
Fare           0.012658  -0.257397  -0.549500  -0.096067  -0.159651  0.216225  1.000000

In [9]: df.corr().Survived.sort_values(ascending = False)

C:\Users\AKARSHA\AppData\Local\Temp\ipykernel_9232\1287823212.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
df.corr().Survived.sort_values(ascending = False)

Out[9]:
Survived
PassengerId    -0.050069
Fare             0.012658
Parch           -0.001622
SibSp           -0.057327
Age             -0.035322
Sex             -0.077221
Pclass          -0.338481
Name: Survived, dtype: float64

Handling Missing/Null Values

In [10]: df.isnull().any()

PassengerId    False
Survived        False
Pclass          False
Name            False
Sex             False
Age             True
SibSp           False
Parch           False
Ticket          False
Fare            False
Cabin           True
Embarked        True
dtype: bool

In [11]: sum(df.Cabin.isnull())

687

Out[11]:

In [12]: sum(df.Age.isnull())

In [13]: df["Age"].fillna(df["Age"].mean(),inplace=True)

In [14]: sum(df.Embarked.isnull())

Out[14]:

In [15]: df["Embarked"].fillna(df["Embarked"].mode()[0],inplace=True)

In [16]: df.describe()

Out[16]:
   PassengerId  Survived  Pclass    Age    SibSp    Parch    Fare
count  891.000000   891.000000   891.000000   891.000000   891.000000   891.000000   891.000000
mean    446.000000   0.383838   2.308642   29.699118   0.523008   0.381594   32.204208
std     257.353842   0.486592   0.838671   14.526497   1.102743   0.806067   49.693429
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   7.910400
50%     446.000000   0.000000   3.000000   29.699118   0.000000   0.000000   14.454200
75%     698.500000   1.000000   3.000000   35.000000   1.000000   0.000000   31.000000
max     891.000000   1.000000   3.000000   80.000000   8.000000   6.000000   512.329200

Data Visualization

In [17]: plt.scatter(df["Fare"],df["Survived"])

Out[17]:
<matplotlib.collections.PathCollection at 0x22164e91298>

Out[17]:
<matplotlib.collections.PathCollection at 0x22164e91298>

In [18]: sns.heatmap(df.corr(),annot=True)

C:\User\AKARSHA\AppData\Local\Temp\ipykernel_9232\4277794465.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
sns.heatmap(df.corr(),annot=True)

Out[18]:
<Axes: >

PassengerId    1
Survived       -0.005
Pclass         -0.035
Age            0.033
SibSp          -0.058
Parch          -0.0017
Fare           0.013

Survived        1
Pclass         -0.34
Age            -0.33
SibSp          -0.035
Parch          0.082
Fare           0.26

Pclass          1
Age            -0.33
SibSp          0.083
Parch          0.018
Fare          -0.55

Age             1
SibSp          -0.23
Parch          -0.18
Fare           0.092

SibSp           1
Parch          0.41
Fare           0.16

Parch           1
Fare           0.22

Fare            1

In [19]: sns.pairplot(df)

Out[19]:
<seaborn.axisgrid.PairGrid at 0x22164ea8850>

Out[19]:
<seaborn.axisgrid.PairGrid at 0x22164ea8850>

In [20]: sns.barplot(x=df["Sex"],y=df["Survived"],ci=0)

C:\Users\AKARSHA\AppData\Local\Temp\ipykernel_9232\2482959670.py:1: FutureWarning:
The 'ci' parameter is deprecated. Use 'errorbar=(\'ci\', 0)' for the same effect.
sns.barplot(x=df["Sex"],y=df["Survived"],ci=0)
<Axes: xlabel='Sex', ylabel='Survived'>

Out[20]:
<Axes: xlabel='Sex', ylabel='Survived'>

In [21]: sns.barplot(x=df["Embarked"],y=df["Survived"],ci=0)

C:\Users\AKARSHA\AppData\Local\Temp\ipykernel_9232\2962821836.py:1: FutureWarning:
The 'ci' parameter is deprecated. Use 'errorbar=(\'ci\', 0)' for the same effect.
sns.barplot(x=df["Embarked"],y=df["Survived"],ci=0)
<Axes: xlabel='Embarked', ylabel='Survived'>

Out[21]:
<Axes: xlabel='Embarked', ylabel='Survived'>

In [22]: sns.barplot(x=df["Parch"],y=df["Survived"],ci=0)

C:\Users\AKARSHA\AppData\Local\Temp\ipykernel_9232\2317442107.py:1: FutureWarning:
The 'ci' parameter is deprecated. Use 'errorbar=(\'ci\', 0)' for the same effect.
sns.barplot(x=df["Parch"],y=df["Survived"],ci=0)
<Axes: xlabel='Parch', ylabel='Survived'>

Out[22]:
<Axes: xlabel='Parch', ylabel='Survived'>

Outlier Detection

In [23]: sns.boxplot(df)

<Axes: >

Out[23]:
<Axes: >

In [24]: sns.boxplot(df.Age)

<Axes: >

Out[24]:
<Axes: >

In [52]: sns.boxplot(df.Age)

<Axes: >

Out[52]:
<Axes: >

In [53]: sns.boxplot(df.SibSp)

<Axes: >

Out[53]:
<Axes: >

In [54]: p99 = df.SibSp.quantile(0.99)

In [55]: df = df[df.SibSp <= p99]

In [56]: sns.boxplot(df["SibSp"])

<Axes: >

Out[56]:
<Axes: >

In [51]: sns.boxplot(df["Parch"])

<Axes: >

Out[51]:
<Axes: >

In [57]: p99 = df.Parch.quantile(0.99)

In [58]: df = df[df.Parch <= p99]

In [59]: sns.boxplot(df["Parch"])

<Axes: >

Out[59]:
<Axes: >

In [60]: sns.boxplot(df["Fare"])

<Axes: >

Out[60]:
<Axes: >

In [62]: sns.boxplot(df.Fare)

<Axes: >

Out[62]:
<Axes: >

Splitting Dependent and Independent Variables

In [63]: x = df.drop(columns=["Survived","PassengerId","Name","Ticket","Cabin"],axis=1)
#Independent variables should be in df or 2d array

Out[63]:
<Axes: >

In [64]: x.head()

Out[64]:
   Pclass  Sex    Age  SibSp  Parch    Fare  Embarked
0      3  male  22.000000   1    0    7.2500    S
1      3  female  26.000000   0    0    7.9250    S
2      3  male  35.000000   0    0    8.0500    S
3      3  male  29.699118   0    0    8.4583    Q
4      3  female  14.000000   1    0    30.0708    C

In [65]: y = pd.Series(df["Survived"])

In [66]: y.head()

Out[66]:
0    0
1    1
2    1
3    0
4    1
dtype: object
Name: Survived, dtype: int64

Encoding

In [67]: from sklearn.preprocessing import LabelEncoder

In [68]: le = LabelEncoder()

In [69]: x["Sex"] = le.fit_transform(x["Sex"])

In [70]: x.head()

Out[70]:
   Pclass  Sex    Age  SibSp  Parch    Fare  Embarked
0      3    1  22.000000   1    0    7.2500    S
1      3    0  26.000000   0    0    7.9250    S
2      3    1  35.000000   0    0    8.0500    S
3      3    1  29.699118   0    0    8.4583    Q
4      3    0  14.000000   1    0    30.0708    C

In [72]: print(le.classes_)

['female' 'male']

In [73]: mapping=dict(zip(le.classes_,range(len(le.classes_))))

In [74]: mapping

Out[74]:
{'female': 0, 'male': 1}

In [75]: le1 = LabelEncoder()

In [76]: x["Embarked"] = le1.fit_transform(x["Embarked"])

In [77]: x.head()

Out[77]:
   Pclass  Sex    Age  SibSp  Parch    Fare  Embarked
0      3    1  22.000000   1    0    7.2500    2
1      3    0  26.000000   0    0    7.9250    1
2      3    1  35.000000   0    0    8.0500    2
3      3    1  29.699118   0    0    8.4583    1
4      3    0  14.000000   1    0    30.0708    0

In [78]: print(le1.classes_)

['C' 'Q' 'S']

In [79]: mapping=dict(zip(le1.classes_,range(len(le1.classes_))))

In [80]: mapping

Out[80]:
{'C': 0, 'Q': 1, 'S': 2}

Feature Scaling

In [81]: from sklearn.preprocessing import MinMaxScaler
ms = MinMaxScaler()

In [82]: x_Scaled = pd.DataFrame(ms.fit_transform(x),columns = x.columns)

In [83]: x_Scaled.head()

Out[83]:
   Pclass  Sex    Age  SibSp  Parch    Fare  Embarked
0      1.0  1.0  0.358974  0.333333  0.0  0.195024    1.0
1      1.0  0.0  0.461538  0.000000   0.0  0.214165    1.0
2      1.0  1.0  0.692308  0.000000   0.0  0.217543    1.0
3      1.0  1.0  0.556388  0.000000   0.0  0.228577    0.5
4      0.5  0.0  0.153846  0.333333  0.0  0.812632    0.0

Splitting Training and Testing Data

In [84]: from sklearn.model_selection import train_test_split

In [85]: x_train,x_test,y_train,y_test = train_test_split(x_Scaled,y,test_size = 0.2,random_state = 9)

In [86]: print(x_train.shape,x_test.shape,y_train.shape,y_test.shape)

(453, 7) (114, 7) (453,) (114, )
```