

In []: Text

TITANIC Data Preprocessing

```
In [1]: import pandas as pd  
import numpy as np
```

```
In [3]: titanic = pd.read_csv(r"E:\Data Science & AI\Dataset files\titanic dataset.csv")
```

```
In [5]: titanic
```

Out[5]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599 7
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803 5
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536 1
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053 3
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607 2
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369 3
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376

891 rows × 12 columns



In [7]: titanic.tail()

Out[7]:


	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7



In [9]: `titanic.head()`

Out[9]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.25
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.2833
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.92
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.05



Performing Data Cleaning and Analysis Understanding meaning of each column: Data Dictionary: Variable Description
Survived - Survived (1) or died (0) Pclass - Passenger's class (1 = 1st, 2 = 2nd, 3 = 3rd) Name - Passenger's name Sex - Passenger's sex Age - Passenger's age SibSp - Number of siblings/spouses aboard Parch - Number of parents/children aboard

(Some children travelled only with a nanny, therefore parch=0 for them.) Ticket - Ticket number Fare - Fare Cabin - Cabin Embarked - Port of embarkation (C = Cherbourg, Q = Queenstown, S = Southampton)

```
In [12]: titanic.describe()
```

Out[12]:

	PassengerId	Survived	Pclass	Age	SibSp	Parch	
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	32.204167
std	257.353842	0.486592	0.836071	14.526497	1.102743	0.806057	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	7.910000
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454167
75%	668.500000	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	512.320833

```
In [14]: #Name column can never decide survival of a person, hence we can safely delete it  
del titanic["Name"]
```

```
In [16]: del titanic["Ticket"]
```

```
In [18]: del titanic["Fare"]
```

```
In [20]: del titanic['Cabin']
```

```
In [22]: titanic.head()
```

Out[22]:

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Embarked
0	1	0	3	male	22.0	1	0	S
1	2	1	1	female	38.0	1	0	C
2	3	1	3	female	26.0	0	0	S
3	4	1	1	female	35.0	1	0	S
4	5	0	3	male	35.0	0	0	S

```
In [24]: # Changing Value for "Male, Female" string values to numeric values , male=1 and female=2  
def getNumber(str):  
    if str=="male":  
        return 1  
    else:  
        return 2  
titanic["Gender"]=titanic["Sex"].apply(getNumber)
```

```
In [26]: titanic.head()
```

```
Out[26]:
```

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Embarked	Gender
0	1	0	3	male	22.0	1	0	S	1
1	2	1	1	female	38.0	1	0	C	2
2	3	1	3	female	26.0	0	0	S	2
3	4	1	1	female	35.0	1	0	S	2
4	5	0	3	male	35.0	0	0	S	1

```
In [28]: #Deleting Sex column, since no use of it now
del titanic["Sex"]
```

```
In [30]: titanic.head()
```

```
Out[30]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Embarked	Gender
0	1	0	3	22.0	1	0	S	1
1	2	1	1	38.0	1	0	C	2
2	3	1	3	26.0	0	0	S	2
3	4	1	1	35.0	1	0	S	2
4	5	0	3	35.0	0	0	S	1

```
In [32]: titanic.isnull().sum()
```

```
Out[32]: PassengerId      0
Survived      0
Pclass        0
Age          177
SibSp         0
Parch         0
Embarked       2
Gender        0
dtype: int64
```

```
In [34]: means= titanic[titanic.Survived==1].Age.mean()
means
```

```
Out[34]: 28.343689655172415
```

```
In [36]: titanic["age"]=np.where(pd.isnull(titanic.Age) & titanic["Survived"]==1 ,means,
titanic.head())
```

```
Out[36]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Embarked	Gender	age
0	1	0	3	22.0	1	0	S	1	22.0
1	2	1	1	38.0	1	0	C	2	38.0
2	3	1	3	26.0	0	0	S	2	26.0
3	4	1	1	35.0	1	0	S	2	35.0
4	5	0	3	35.0	0	0	S	1	35.0

```
In [38]: titanic.isnull().sum()
```

```
Out[38]: PassengerId      0
Survived      0
Pclass        0
Age          177
SibSp         0
Parch         0
Embarked      2
Gender        0
age          125
dtype: int64
```

```
In [40]: # Finding the mean age of "Not Survived" people

means1=titanic[titanic.Survived==0].Age.mean()
means1
```

```
Out[40]: 30.62617924528302
```

```
In [44]: titanic.age.fillna(means1,inplace=True)
```

C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\2890245544.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
titanic.age.fillna(means1,inplace=True)
```

```
In [46]: titanic.head()
```

```
Out[46]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Embarked	Gender	age
0	1	0	3	22.0	1	0	S	1	22.0
1	2	1	1	38.0	1	0	C	2	38.0
2	3	1	3	26.0	0	0	S	2	26.0
3	4	1	1	35.0	1	0	S	2	35.0
4	5	0	3	35.0	0	0	S	1	35.0

```
In [48]: titanic.isnull().sum()
```

```
Out[48]: PassengerId      0
Survived      0
Pclass        0
Age          177
SibSp         0
Parch         0
Embarked      2
Gender        0
age           0
dtype: int64
```

```
In [50]: del titanic['Age']
titanic.head()
```

```
Out[50]:
```

	PassengerId	Survived	Pclass	SibSp	Parch	Embarked	Gender	age
0	1	0	3	1	0	S	1	22.0
1	2	1	1	1	0	C	2	38.0
2	3	1	3	0	0	S	2	26.0
3	4	1	1	1	0	S	2	35.0
4	5	0	3	0	0	S	1	35.0

```
In [52]: # Finding the number of people who have survived
# given that they have embarked or boarded from a particular port

survivedQ = titanic[titanic.Embarked == 'Q'][titanic.Survived == 1].shape[0]
survivedC = titanic[titanic.Embarked == 'C'][titanic.Survived == 1].shape[0]
survivedS = titanic[titanic.Embarked == 'S'][titanic.Survived == 1].shape[0]
print(survivedQ)
print(survivedC)
print(survivedS)
```

```
30
93
217
```

```
C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\3300902897.py:4: UserWarning: Boolean Series key will be reindexed to match DataFrame index.
  survivedQ = titanic[titanic.Embarked == 'Q'][titanic.Survived == 1].shape[0]
C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\3300902897.py:5: UserWarning: Boolean Series key will be reindexed to match DataFrame index.
  survivedC = titanic[titanic.Embarked == 'C'][titanic.Survived == 1].shape[0]
C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\3300902897.py:6: UserWarning: Boolean Series key will be reindexed to match DataFrame index.
  survivedS = titanic[titanic.Embarked == 'S'][titanic.Survived == 1].shape[0]
```

```
In [54]: survivedQ = titanic[titanic.Embarked == 'Q'][titanic.Survived == 0].shape[0]
survivedC = titanic[titanic.Embarked == 'C'][titanic.Survived == 0].shape[0]
survivedS = titanic[titanic.Embarked == 'S'][titanic.Survived == 0].shape[0]
print(survivedQ)
print(survivedC)
print(survivedS)
```

```
47
75
427
```

```
C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\3240960939.py:1: UserWarning: Boolean Series key will be reindexed to match DataFrame index.
  survivedQ = titanic[titanic.Embarked == 'Q'][titanic.Survived == 0].shape[0]
C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\3240960939.py:2: UserWarning: Boolean Series key will be reindexed to match DataFrame index.
  survivedC = titanic[titanic.Embarked == 'C'][titanic.Survived == 0].shape[0]
C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\3240960939.py:3: UserWarning: Boolean Series key will be reindexed to match DataFrame index.
  survivedS = titanic[titanic.Embarked == 'S'][titanic.Survived == 0].shape[0]
```

```
In [58]: titanic.dropna(inplace=True)
titanic.head()
```

```
Out[58]:
```

	PassengerId	Survived	Pclass	SibSp	Parch	Embarked	Gender	age
0	1	0	3	1	0	S	1	22.0
1	2	1	1	1	0	C	2	38.0
2	3	1	3	0	0	S	2	26.0
3	4	1	1	1	0	S	2	35.0
4	5	0	3	0	0	S	1	35.0

```
In [60]: titanic.isnull().sum()
```

```
Out[60]: PassengerId    0
Survived              0
Pclass               0
SibSp                0
Parch                0
Embarked             0
Gender               0
age                  0
dtype: int64
```

```
In [62]: #Renaming "age" and "gender" columns
```



```
titanic.rename(columns={'age':'Age'}, inplace=True)
titanic.head()
```

```
Out[62]:
```

	PassengerId	Survived	Pclass	SibSp	Parch	Embarked	Gender	Age
0	1	0	3	1	0	S	1	22.0
1	2	1	1	1	0	C	2	38.0
2	3	1	3	0	0	S	2	26.0
3	4	1	1	1	0	S	2	35.0
4	5	0	3	0	0	S	1	35.0

```
In [64]: titanic.rename(columns={'Gender':'Sex'}, inplace=True)
```

```
In [66]: titanic.head()
```

```
Out[66]:
```

	PassengerId	Survived	Pclass	SibSp	Parch	Embarked	Sex	Age
0	1	0	3	1	0	S	1	22.0
1	2	1	1	1	0	C	2	38.0
2	3	1	3	0	0	S	2	26.0
3	4	1	1	1	0	S	2	35.0
4	5	0	3	0	0	S	1	35.0

```
In [68]: def getEmb(str):
          if str=="S":
              return 1
          elif str=='Q':
              return 2
          else:
              return 3
titanic["Embark"]=titanic["Embarked"].apply(getEmb)
titanic.head()
```

```
Out[68]:
```

	PassengerId	Survived	Pclass	SibSp	Parch	Embarked	Sex	Age	Embark
0	1	0	3	1	0	S	1	22.0	1
1	2	1	1	1	0	C	2	38.0	3
2	3	1	3	0	0	S	2	26.0	1
3	4	1	1	1	0	S	2	35.0	1
4	5	0	3	0	0	S	1	35.0	1

```
In [70]: del titanic['Embarked']
titanic.rename(columns={'Embark':'Embarked'}, inplace=True)
titanic.head()
```

Out[70]:

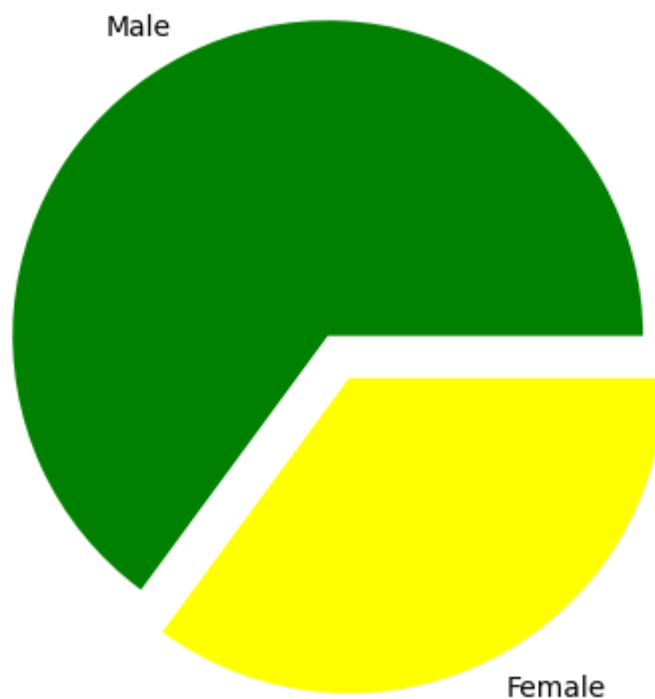
	PassengerId	Survived	Pclass	SibSp	Parch	Sex	Age	Embarked
0	1	0	3	1	0	1	22.0	1
1	2	1	1	1	0	2	38.0	3
2	3	1	3	0	0	2	26.0	1
3	4	1	1	1	0	2	35.0	1
4	5	0	3	0	0	1	35.0	1

```
In [72]: #Drawing a pie chart for number of males and females aboard

import matplotlib.pyplot as plt
from matplotlib import style

males = (titanic['Sex'] == 1).sum()
#Summing up all the values of column gender with a
#condition for male and similary for females
females = (titanic['Sex'] == 2).sum()
print(males)
print(females)
p = [males, females]
plt.pie(p,      #giving array
        labels = ['Male', 'Female'], #Correspndingly giving labels
        colors = ['green', 'yellow'], # Corresponding colors
        explode = (0.15, 0),         #How much the gap should be there between the pie
        startangle = 0) #what start angle should be given
plt.axis('equal')
plt.show()
```

577
312



```
In [74]: # More Precise Pie Chart
```

```

MaleS=titanic[titanic.Sex==1][titanic.Survived==1].shape[0]
print(MaleS)
MaleN=titanic[titanic.Sex==1][titanic.Survived==0].shape[0]
print(MaleN)
FemaleS=titanic[titanic.Sex==2][titanic.Survived==1].shape[0]
print(FemaleS)
FemaleN=titanic[titanic.Sex==2][titanic.Survived==0].shape[0]
print(FemaleN)

```

109

468

231

81

C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\2810620594.py:3: UserWarning: Boolean Series key will be reindexed to match DataFrame index.

```
MaleS=titanic[titanic.Sex==1][titanic.Survived==1].shape[0]
```

C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\2810620594.py:5: UserWarning: Boolean Series key will be reindexed to match DataFrame index.

```
MaleN=titanic[titanic.Sex==1][titanic.Survived==0].shape[0]
```

C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\2810620594.py:7: UserWarning: Boolean Series key will be reindexed to match DataFrame index.

```
FemaleS=titanic[titanic.Sex==2][titanic.Survived==1].shape[0]
```

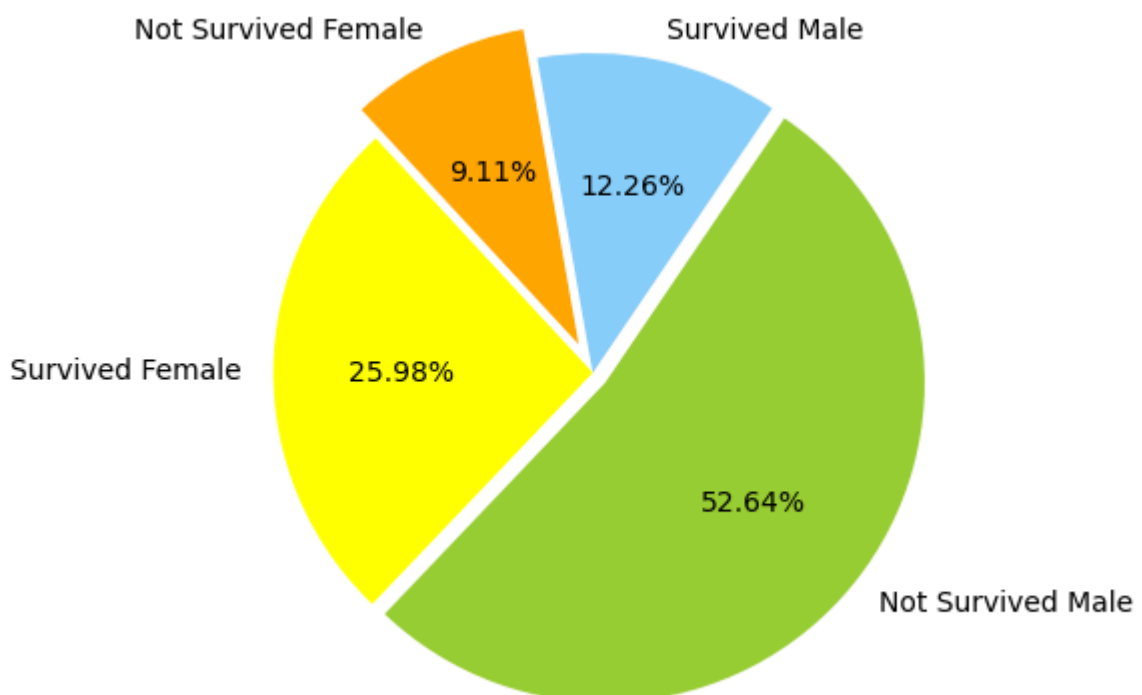
C:\Users\roy62\AppData\Local\Temp\ipykernel_11016\2810620594.py:9: UserWarning: Boolean Series key will be reindexed to match DataFrame index.

```
FemaleN=titanic[titanic.Sex==2][titanic.Survived==0].shape[0]
```

```

In [76]: chart=[MaleS, MaleN, FemaleS, FemaleN]
         colors=['lightskyblue', 'yellowgreen', 'Yellow', 'Orange']
         labels=["Survived Male", "Not Survived Male", "Survived Female", "Not Survived Female"]
         explode=[0, 0.05, 0, 0.1]
         plt.pie(chart, labels=labels, colors=colors, explode=explode, startangle=100, counter
         plt.axis("equal")
         plt.show()

```



In []: