

Data Preprocessing for Beginners

```
In [3]: import numpy as np
import pandas as pd
```

```
In [89]: titanic = pd.read_csv(r"C:\Users\gadel\OneDrive\Desktop\Nareshit DataScience by Pr
titanic.tail()
```

Out[89]:

	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

Performing Data Cleaning and Analysis

1. 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)
 2. Analysing which columns are completely useless in predicting the survival and deleting them
 Note - Don't just delete the columns because you are not finding it useful. Or focus is not on deleting the columns. Our focus is on analysing how each column is affecting the result or the prediction and in accordance with that deciding whether to keep the column or to delete the column or fill the null values of the column by some values and if yes, then what values.

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

Out[10]:

	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	32.200000
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	7.912500
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.450000
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	512.320000

In [12]: *#Name column can never decide survival of a person, hence we can safely delete it*
`del titanic["Name"]`
`titanic.head()`

Out[12]:

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	male	22.0	1	0	A/5 21171	7.2500	NaN
1	2	1	1	female	38.0	1	0	PC 17599	71.2833	C85
2	3	1	3	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN
3	4	1	1	female	35.0	1	0	113803	53.1000	C123
4	5	0	3	male	35.0	0	0	373450	8.0500	NaN

In [14]: `del titanic["Ticket"]`
`titanic.head()`

Out[14]:

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Cabin	Embarked
0	1	0	3	male	22.0	1	0	7.2500	NaN	S
1	2	1	1	female	38.0	1	0	71.2833	C85	C
2	3	1	3	female	26.0	0	0	7.9250	NaN	S
3	4	1	1	female	35.0	1	0	53.1000	C123	S
4	5	0	3	male	35.0	0	0	8.0500	NaN	S

In [16]: `del titanic["Fare"]`
`titanic.head()`

Out[16]:

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Cabin	Embarked
0	1	0	3	male	22.0	1	0	NaN	S
1	2	1	1	female	38.0	1	0	C85	C
2	3	1	3	female	26.0	0	0	NaN	S
3	4	1	1	female	35.0	1	0	C123	S
4	5	0	3	male	35.0	0	0	NaN	S

In [18]: `del titanic["Cabin"]`
`titanic.head()`

Out[18]:

	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 [20]: *# 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)

#We have created a new column called "Gender" and
#filling it with values 1,2 based on the values of sex column

titanic.head()
```

Out[20]:

	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 [22]: *#Deleting Sex column, since no use of it now*

```
del titanic["Sex"]
titanic.head()
```

Out[22]:

	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 [24]: `titanic.isnull().sum()`

Out[24]:

PassengerId	0
Survived	0
Pclass	0
Age	177
SibSp	0
Parch	0
Embarked	2
Gender	0

dtype: int64

Fill the null values of the Age column. Fill mean Survived age(mean age of the survived people) in the column where the person has survived and mean not Survived age (mean age of the people who have not survived) in the column where person has not survived

In [28]: `meanS = titanic[titanic.Survived==1].Age.mean()
meanS`

Out[28]: 28.343689655172415

Creating a new "Age" column , filling values in it with a condition if goes True then given values (here meanS) is put in place of last values else nothing happens, simply the values are copied from the "Age" column of the dataset

In [31]: `titanic["age"]=np.where(pd.isnull(titanic.Age) & titanic["Survived"]==1 ,meanS, titanic.Age)
titanic.head()`

Out[31]:

	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 [33]: `titanic.isnull().sum()`

```
Out[33]: PassengerId    0
         Survived      0
         Pclass       0
         Age         177
         SibSp       0
         Parch       0
         Embarked     2
         Gender      0
         age         125
         dtype: int64
```

```
In [35]: # Finding the mean age of "Not Survived" people

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

```
Out[35]: 30.62617924528302
```

```
In [49]: import warnings
         warnings.filterwarnings('ignore')
```

```
In [51]: titanic.age.fillna(meanNS,inplace=True)
         titanic.head()
```

```
Out[51]:
```

	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 [53]: titanic.isnull().sum()
```

```
Out[53]: PassengerId    0
         Survived      0
         Pclass       0
         Age         177
         SibSp       0
         Parch       0
         Embarked     2
         Gender      0
         age         0
         dtype: int64
```

```
In [55]: del titanic["Age"]
         titanic.head()
```

```
Out[55]:
```

	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

We want to check if "Embarked" column is important for analysis or not, that is whether survival of the person depends on the Embarked column value or not

```
In [62]: # 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
```

```
In [64]: 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
```

As there are significant changes in the survival rate based on which port the passengers aboard the ship. We cannot delete the whole embarked column(It is useful). Now the Embarked column has some null values in it and hence we can safely say that deleting some rows from total rows will not affect the result. So rather than trying to fill those null values with some vales. We can simply remove them.

```
In [67]: titanic.dropna(inplace=True)
titanic.head()
```

Out[67]:

	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 [69]: `titanic.isnull().sum()`

Out[69]:

```

PassengerId    0
Survived        0
Pclass          0
SibSp           0
Parch           0
Embarked        0
Gender          0
age             0
dtype: int64

```

In [71]: `#Renaming "age" and "gender" columns`
`titanic.rename(columns={'age':'Age'}, inplace=True)`
`titanic.head()`

Out[71]:

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

Out[73]:

	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 [77]:

```

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[77]:

	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 [79]:

```

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

```

Out[79]:

	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 [83]: *#Drawing a pie chart for number of males and females aboard*

```

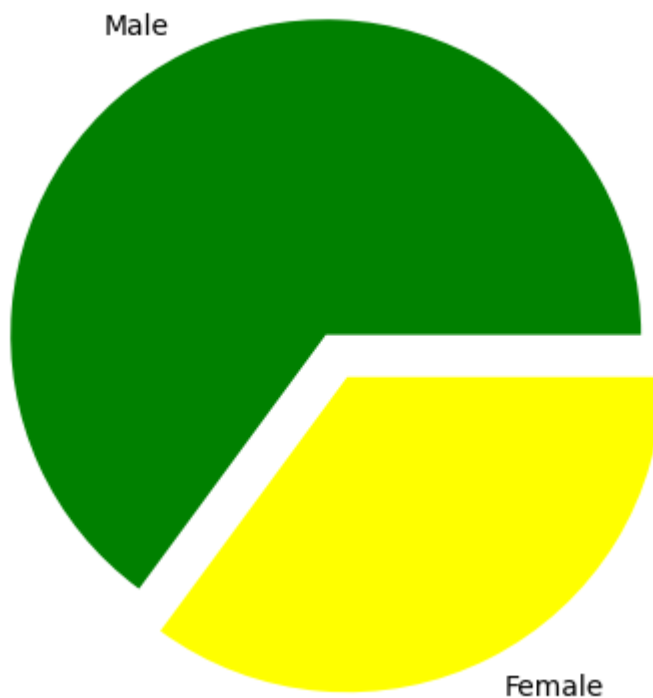
import matplotlib.pyplot as plt
from matplotlib import style

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 pies
        startangle = 0) # what start angle should be given
plt.axis('equal')
plt.show()

```

577

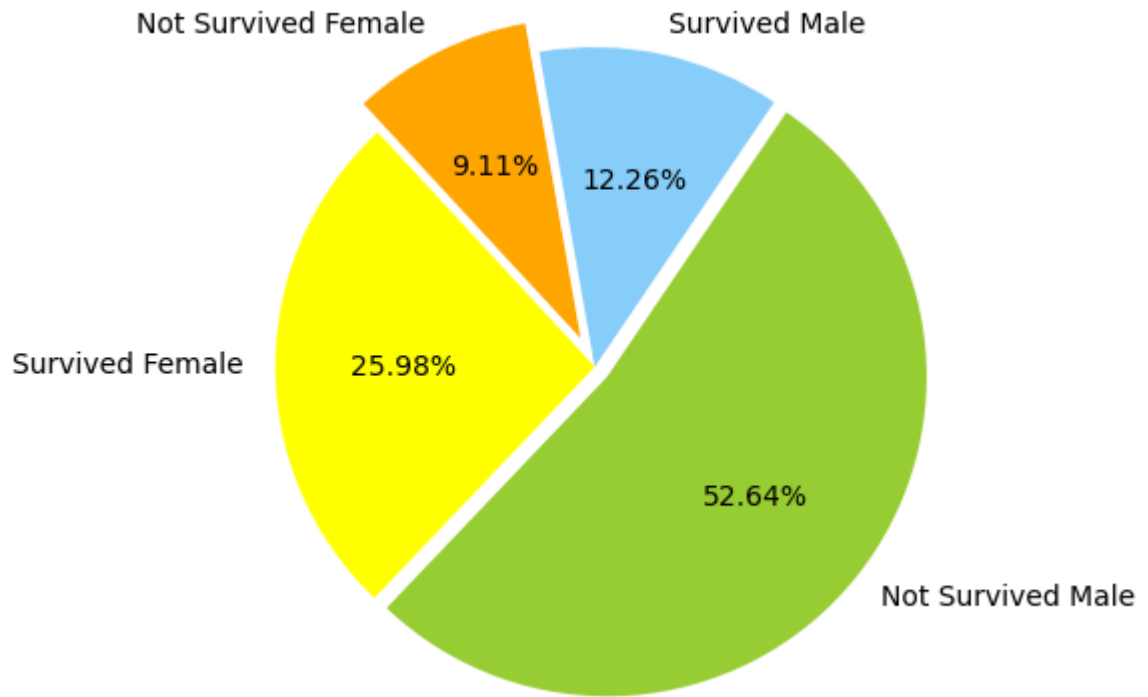
312



```
In [85]: # 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
```

```
In [87]: 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, counterclockwise=True)
plt.axis("equal")
plt.show()
```



In []: