

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

from warnings import filterwarnings
filterwarnings(action='ignore')
```

```
In [2]: train = pd.read_csv(r"C:\Users\Admin\Desktop\DataScience(ProdigyInfotech)\Task2_Titanic\train1.csv")
trial = pd.read_csv(r"C:\Users\Admin\Desktop\DataScience(ProdigyInfotech)\Trial.csv")
```

```
In [3]: print(train.head())
print(trial.head())
```

	PassengerId	Survived	Pclass	Name	Gender	Age	SibSp	Parch	\
0	1	0	3	shivani	female	22.0	1	0	
1	2	1	1	isha	female	38.0	1	0	
2	3	1	3	bablu	male	26.0	0	0	
3	4	1	1	pinki	male	35.0	1	0	
4	5	0	3	surender singh	male	35.0	0	0	

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

	student	name	std	div	roll	no
0	Aria	2	A		21	
1	Veda	3	A		29	
2	Beena	7	C		12	
3	Seema	1	C		34	
4	Arohi	3	B		19	

```
In [4]: train.head()
```

Out[4]:

	PassengerId	Survived	Pclass	Name	Gender	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	shivani	female	22.0	1	0	A/5 21171	7.2500	C85	S
1	2	1	1	isha	female	38.0	1	0	PC 17599	71.2833	NaN	C
2	3	1	3	bablu	male	26.0	0	0	STON/O2. 3101282	7.9250	C123	S
3	4	1	1	pinki	male	35.0	1	0	113803	53.1000	NaN	S
4	5	0	3	surender singh	male	35.0	0	0	373450	8.0500	NaN	S

In [5]: `train.tail()`

Out[5]:

	PassengerId	Survived	Pclass	Name	Gender	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
59	60	0	3	mohd. aslam	male	11.0	5	2	113509	46.9000	NaN	S
60	61	0	3	bharat damor	male	22.0	0	0	113509	7.2292	B28	C
61	62	1	1	alok kumar thakur	male	38.0	0	0	113509	80.0000	C83	NaN
62	63	0	1	najim urf babbu	male	45.0	1	0	113509	83.4750	NaN	S
63	64	0	3	amar singh	male	4.0	3	2	113509	27.9000	NaN	S

In [6]: `train`

Out[6]:

	PassengerId	Survived	Pclass	Name	Gender	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	shivani	female	22.0	1	0	A/5 21171	7.2500	C85	S
1	2	1	1	isha	female	38.0	1	0	PC 17599	71.2833	NaN	C
2	3	1	3	bablu	male	26.0	0	0	STON/O2. 3101282	7.9250	C123	S
3	4	1	1	pinki	male	35.0	1	0	113803	53.1000	NaN	S
4	5	0	3	surender singh	male	35.0	0	0	373450	8.0500	NaN	S
...	...	...	...	...	...	...	...	...	...	...	...	...
59	60	0	3	mohd. aslam	male	11.0	5	2	113509	46.9000	NaN	S
60	61	0	3	bharat damor	male	22.0	0	0	113509	7.2292	B28	C
61	62	1	1	alok kumar thakur	male	38.0	0	0	113509	80.0000	C83	NaN
62	63	0	1	najim urf babbu	male	45.0	1	0	113509	83.4750	NaN	S
63	64	0	3	amar singh	male	4.0	3	2	113509	27.9000	NaN	S

64 rows × 12 columns

In [7]: `train.shape`Out[7]: `(64, 12)`In [8]: `#Checking for Null values  
train.isnull().sum()`

```
Out[8]: PassengerId      0
Survived          0
Pclass            0
Name              0
Gender            0
Age               15
SibSp             0
Parch             0
Ticket            1
Fare              0
Cabin            50
Embarked          1
dtype: int64
```

```
In [9]: #Description of dataset
train.describe(include="all")
```

	PassengerId	Survived	Pclass	Name	Gender	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
<b>count</b>	64.000000	64.0000	64.000000	64	64	49.000000	64.000000	64.000000	63	64.000000	14	63
<b>unique</b>	Nan	Nan	Nan	62	2	Nan	Nan	Nan	54	Nan	14	3
<b>top</b>	Nan	Nan	Nan	anita	male	Nan	Nan	Nan	113509	Nan	C85	S
<b>freq</b>	Nan	Nan	Nan	2	40	Nan	Nan	Nan	9	Nan	1	40
<b>mean</b>	32.500000	0.4375	2.359375	Nan	Nan	27.234694	0.781250	0.437500	Nan	30.210158	Nan	Nan
<b>std</b>	18.618987	0.5000	0.842656	Nan	Nan	16.528732	1.147375	1.037013	Nan	39.276369	Nan	Nan
<b>min</b>	1.000000	0.0000	1.000000	Nan	Nan	2.000000	0.000000	0.000000	Nan	7.225000	Nan	Nan
<b>25%</b>	16.750000	0.0000	2.000000	Nan	Nan	15.000000	0.000000	0.000000	Nan	8.044800	Nan	Nan
<b>50%</b>	32.500000	0.0000	3.000000	Nan	Nan	27.000000	0.000000	0.000000	Nan	17.900000	Nan	Nan
<b>75%</b>	48.250000	1.0000	3.000000	Nan	Nan	38.000000	1.000000	0.000000	Nan	32.415625	Nan	Nan
<b>max</b>	64.000000	1.0000	3.000000	Nan	Nan	66.000000	5.000000	5.000000	Nan	263.000000	Nan	Nan

```
In [10]: #Display shape
train.shape
```

Out[10]: (64, 12)

```
In [11]: train.corr
```

				PassengerId	Survived	Pclass	Name	Gender	Age	SibSp	\
0	1	0	3	shivani	female	22.0					1
1	2	1	1	isha	female	38.0					1
2	3	1	3	bablu	male	26.0					0
3	4	1	1	pinki	male	35.0					1
4	5	0	3	surrender singh	male	35.0					0
..	...	...	...	...	...	...	...	...	...	...	...
59	60	0	3	mohd. aslam	male	11.0					5
60	61	0	3	bharat damor	male	22.0					0
61	62	1	1	alok kumar thakur	male	38.0					0
62	63	0	1	najim urf babbu	male	45.0					1
63	64	0	3	amar singh	male	4.0					3
	Parch		Ticket	Fare	Cabin	Embarked					
0	0	A/5 21171	7.2500	C85		S					
1	0	PC 17599	71.2833	Nan		C					
2	0	STON/O2. 3101282	7.9250	C123		S					
3	0	113803	53.1000	Nan		S					
4	0	373450	8.0500	Nan		S					
..	...	...	...	...	...	...	...	...	...	...	...
59	2	113509	46.9000	Nan		S					
60	0	113509	7.2292	B28		C					
61	0	113509	80.0000	C83		Nan					
62	0	113509	83.4750	Nan		S					
63	2	113509	27.9000	Nan		S					

[64 rows x 12 columns]>

```
In [12]: train.groupby('Survived').mean
```

Out[12]: <bound method GroupBy.mean of <pandas.core.groupby.generic.DataFrameGroupBy object at 0x000002A652890800>>

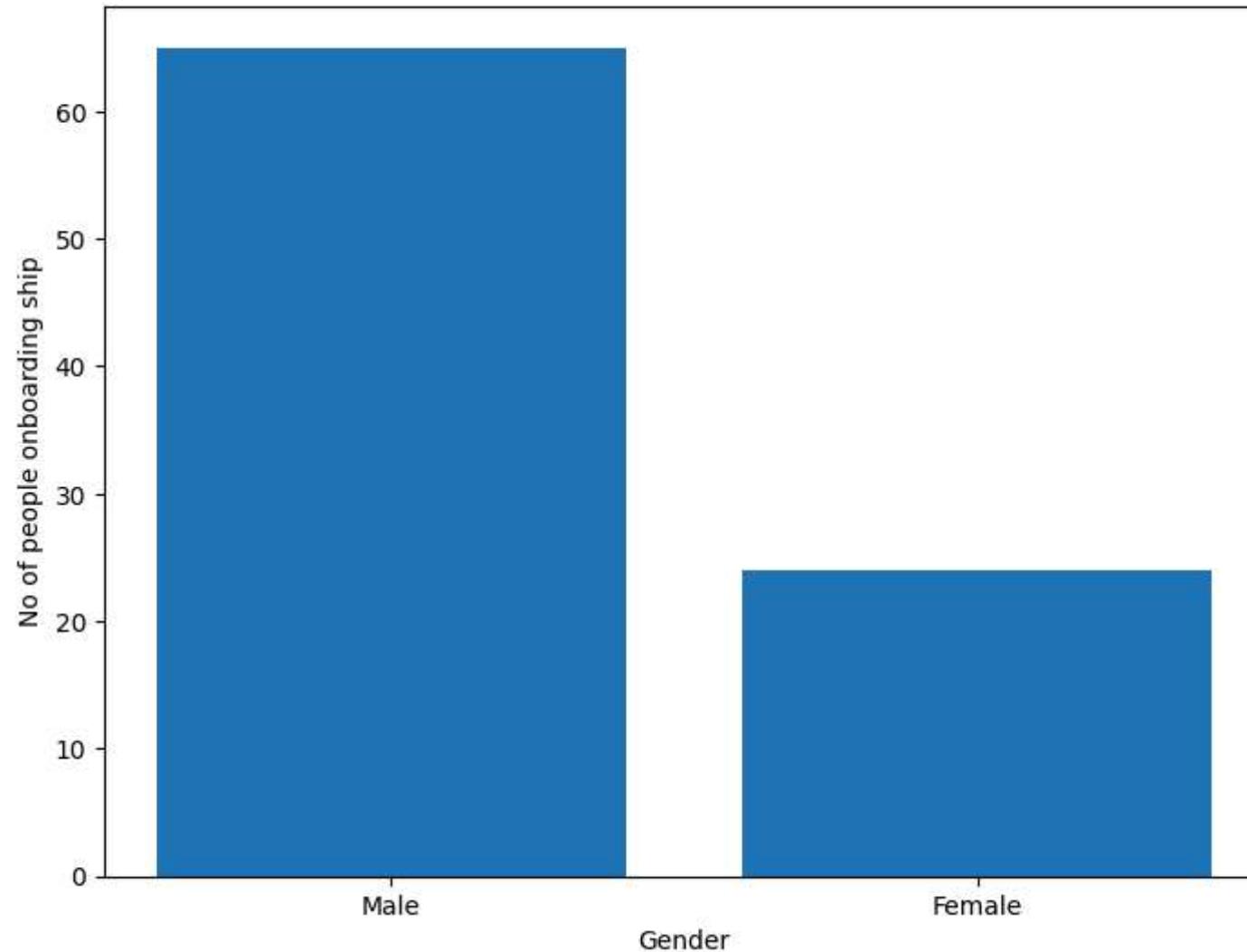
```
In [13]: male_ind = len(train[train['Gender'] == 'male'])
print("No of Males in Titanic:", male_ind)
```

No of Males in Titanic: 40

```
In [14]: female_ind = len(train[train['Gender'] == 'female'])
print("No of Females in Titanic:", female_ind)
```

No of Females in Titanic: 24

```
In [15]: #Plotting
fig = plt.figure()
ax = fig.add_axes([0,0,1,1])
gender = ['Male', 'Female']
index = [65,24]
ax.bar(gender, index)
plt.xlabel("Gender")
plt.ylabel("No of people onboarding ship")
plt.show()
```



```
In [16]: alive = len(train[train['Survived'] == 1])
dead = len(train[train['Survived'] == 0])
```

```
In [17]: train.groupby('Gender')[['Survived']].mean()
```

Out[17]:

**Survived**

Gender	
female	0.5
male	0.4

In [18]:

```
train.groupby('Age')[['Survived']].mean()
```

Out[18]:

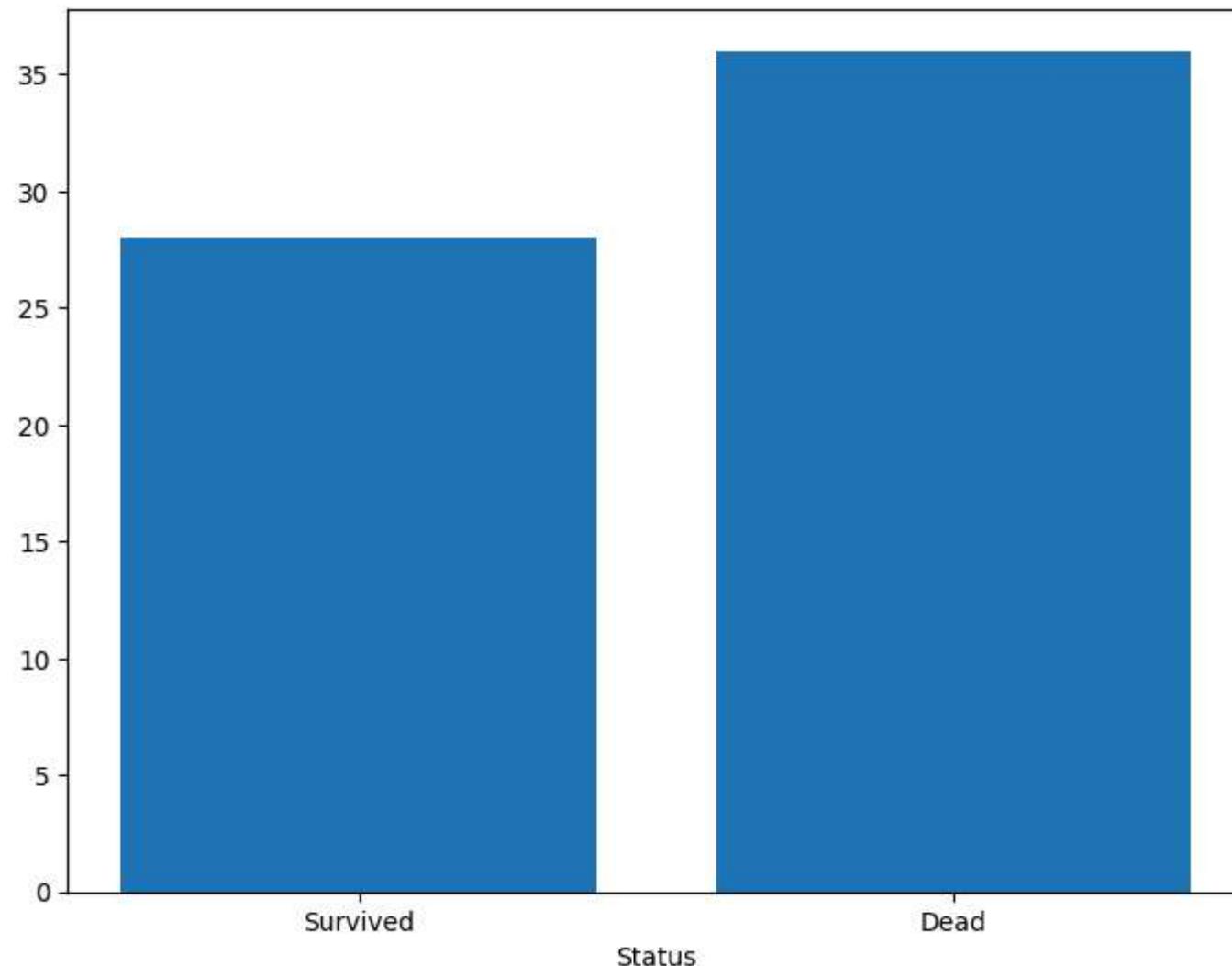
**Survived**

<b>Age</b>	
<b>2.0</b>	0.000000
<b>3.0</b>	1.000000
<b>4.0</b>	0.500000
<b>5.0</b>	1.000000
<b>7.0</b>	0.000000
<b>8.0</b>	0.000000
<b>11.0</b>	0.000000
<b>14.0</b>	0.666667
<b>15.0</b>	1.000000
<b>18.0</b>	0.000000
<b>19.0</b>	0.500000
<b>20.0</b>	0.000000
<b>21.0</b>	0.333333
<b>22.0</b>	0.000000
<b>26.0</b>	1.000000
<b>27.0</b>	0.500000
<b>28.0</b>	0.500000
<b>28.5</b>	0.000000
<b>29.0</b>	1.000000
<b>31.0</b>	0.000000
<b>34.0</b>	1.000000

**Survived**

Age	
<b>35.0</b>	0.333333
<b>38.0</b>	1.000000
<b>39.0</b>	0.000000
<b>40.0</b>	0.000000
<b>42.0</b>	0.000000
<b>45.0</b>	0.000000
<b>49.0</b>	1.000000
<b>54.0</b>	0.000000
<b>55.0</b>	1.000000
<b>58.0</b>	1.000000
<b>65.0</b>	0.000000
<b>66.0</b>	0.000000

```
In [19]: fig = plt.figure()
ax = fig.add_axes([0,0,1,1])
status = ['Survived', 'Dead']
ind = [alive,dead]
ax.bar(status,ind)
plt.xlabel("Status")
plt.show()
```



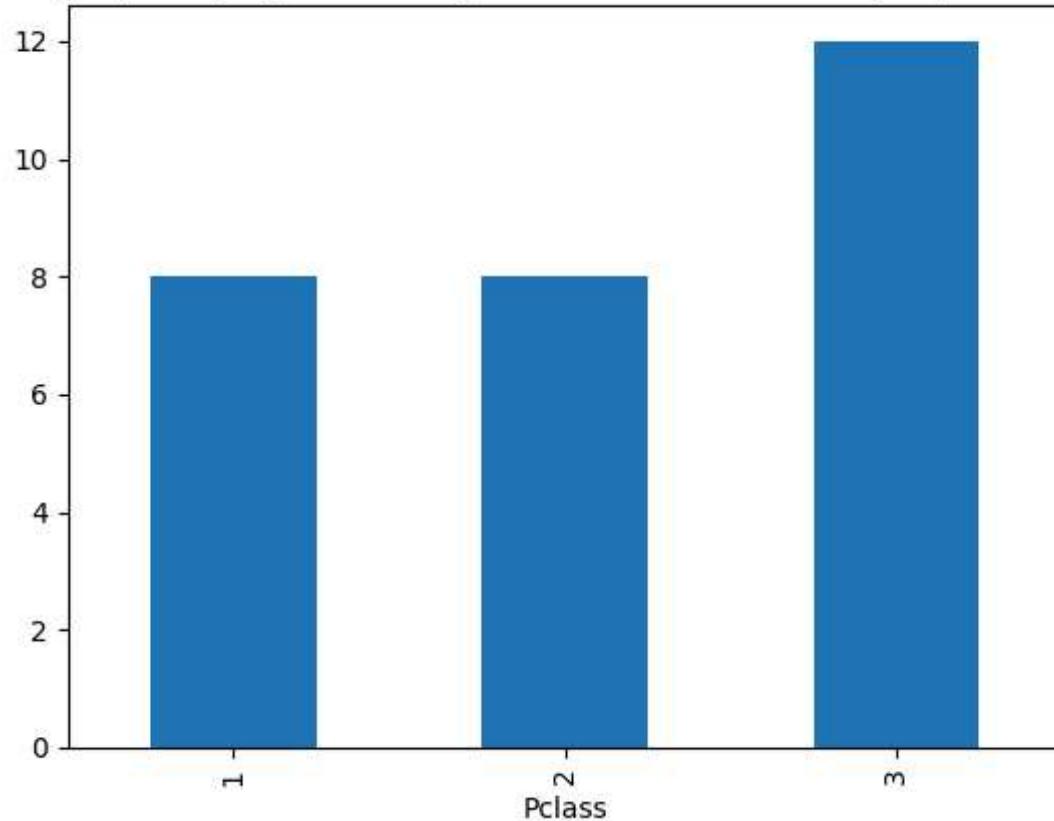
```
In [20]: plt.figure(1)
train.loc[train['Survived'] == 1, 'Pclass'].value_counts().sort_index().plot.bar()
plt.title('Bar graph of people according to ticket class in which people survived')

plt.figure(2)
```

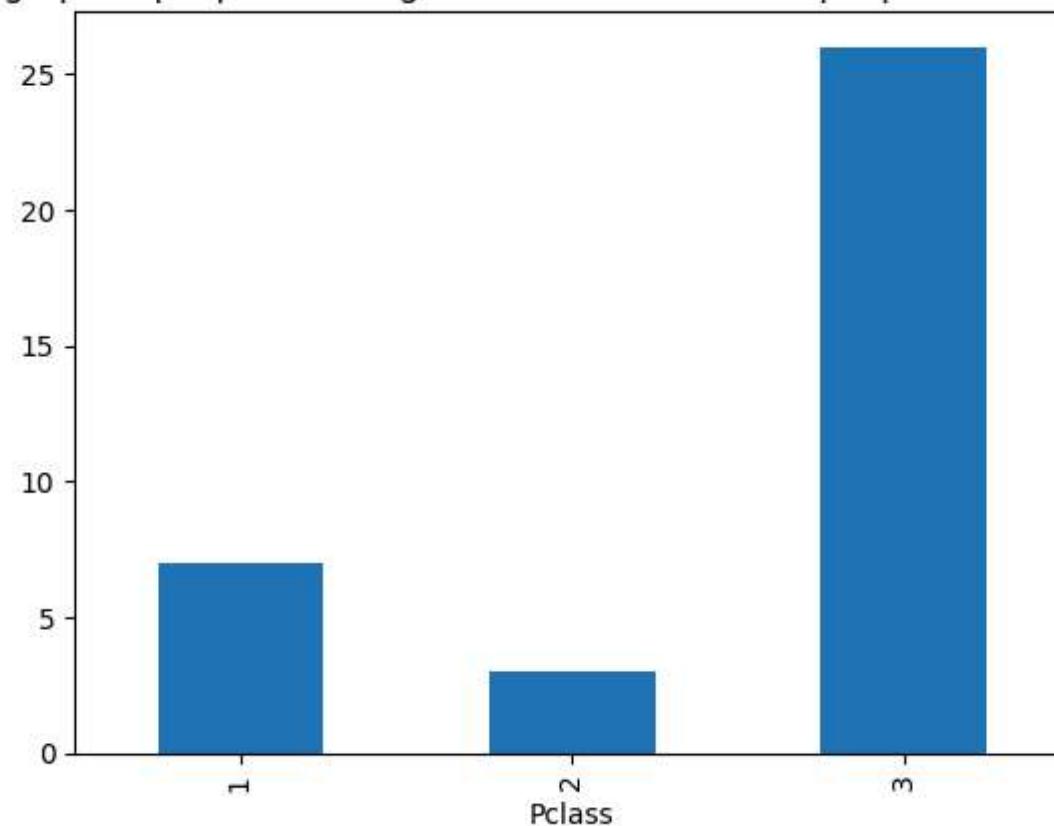
```
train.loc[train['Survived'] == 0, 'Pclass'].value_counts().sort_index().plot.bar()  
plt.title('Bar graph of people according to ticket class in which people couldn\'t survive')
```

Out[20]: Text(0.5, 1.0, "Bar graph of people according to ticket class in which people couldn't survive")

Bar graph of people according to ticket class in which people survived



Bar graph of people according to ticket class in which people couldn't survive

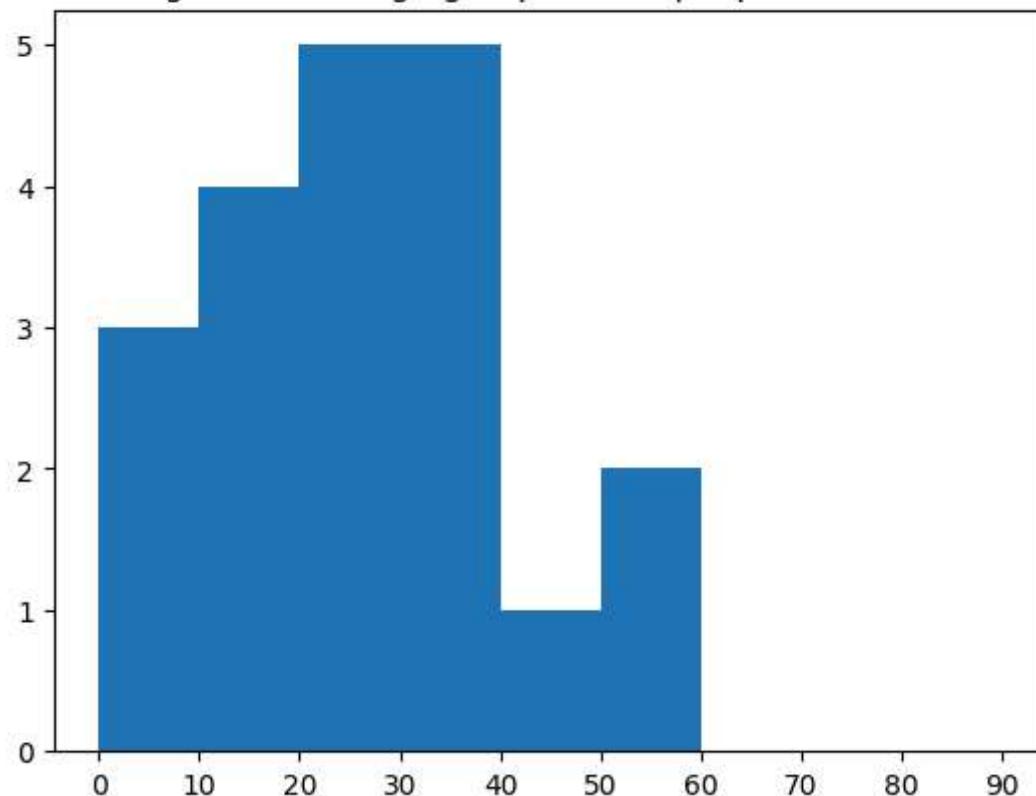


```
In [21]: plt.figure(1)
age = train.loc[train.Survived == 1, 'Age']
plt.title('The histogram of the age groups of the people that had survived')
plt.hist(age, np.arange(0,100,10))
plt.xticks(np.arange(0,100,10))

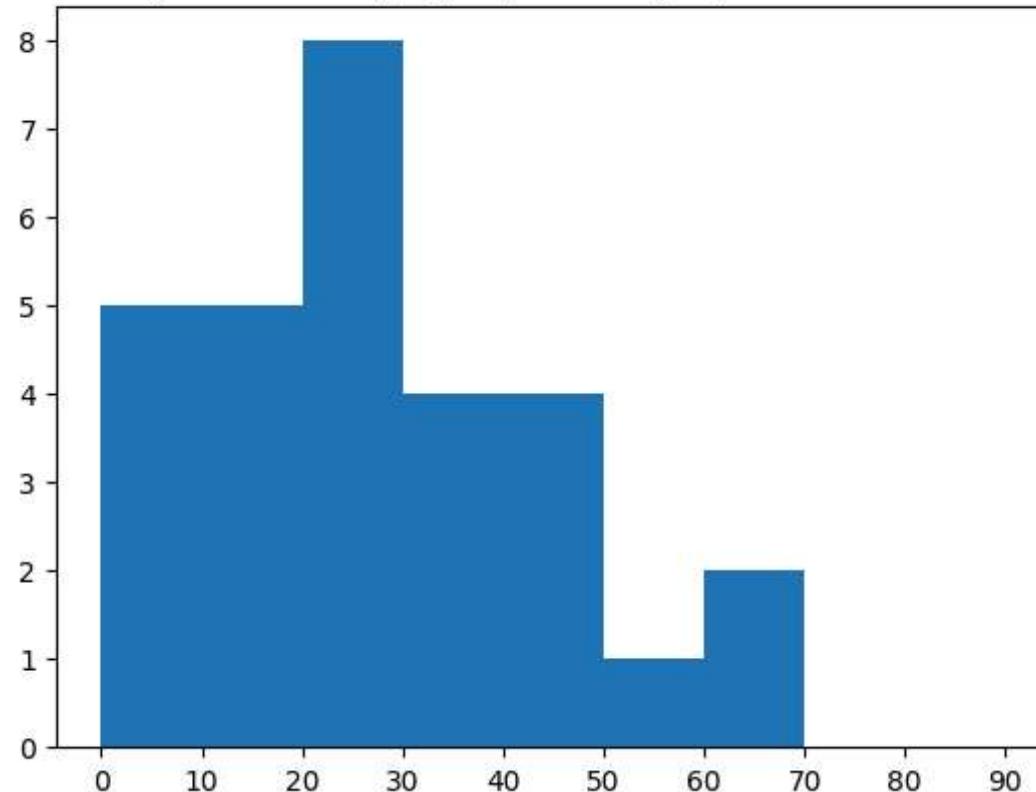
plt.figure(2)
age = train.loc[train.Survived == 0, 'Age']
plt.title('The histogram of the age groups of the people that couldn\'t survive')
plt.hist(age, np.arange(0,100,10))
plt.xticks(np.arange(0,100,10))
```

```
Out[21]: ([<matplotlib.axis.XTick at 0x2a654cbede0>,
    <matplotlib.axis.XTick at 0x2a654d4e7e0>,
    <matplotlib.axis.XTick at 0x2a654d4c950>,
    <matplotlib.axis.XTick at 0x2a6550bf7d0>,
    <matplotlib.axis.XTick at 0x2a6550bfda0>,
    <matplotlib.axis.XTick at 0x2a6550f0a70>,
    <matplotlib.axis.XTick at 0x2a654d4c410>,
    <matplotlib.axis.XTick at 0x2a6550f0650>,
    <matplotlib.axis.XTick at 0x2a6550f1d60>,
    <matplotlib.axis.XTick at 0x2a6550f2660>],
[Text(0, 0, '0'),
 Text(10, 0, '10'),
 Text(20, 0, '20'),
 Text(30, 0, '30'),
 Text(40, 0, '40'),
 Text(50, 0, '50'),
 Text(60, 0, '60'),
 Text(70, 0, '70'),
 Text(80, 0, '80'),
 Text(90, 0, '90')])
```

The histogram of the age groups of the people that had survived



The histogram of the age groups of the people that coudn't survive



```
In [22]: train[["SibSp", "Survived"]].groupby(['SibSp'], as_index=False).mean().sort_values(by='Survived', ascending=False)
```

Out[22]:

	SibSp	Survived
1	1	0.52381
0	0	0.50000
2	2	0.00000
3	3	0.00000
4	4	0.00000
5	5	0.00000

In [23]: `train[['Pclass", "Survived']].groupby(['Pclass'], as_index=False).mean().sort_values(by='Survived', ascending=False)`

Out[23]:

	Pclass	Survived
1	2	0.727273
0	1	0.533333
2	3	0.315789

In [24]: `train[['Age", "Survived']].groupby(['Age'], as_index=False).mean().sort_values(by='Age', ascending=True)`

Out[24]:

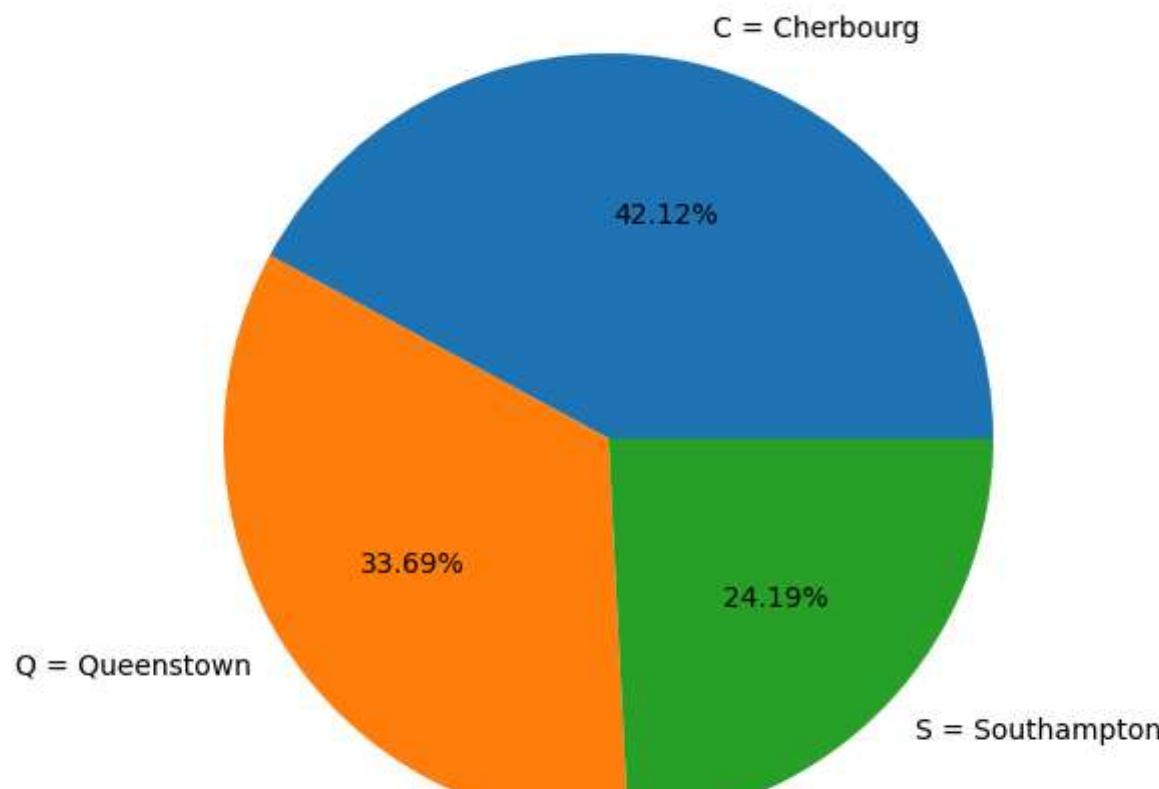
	Age	Survived
<b>0</b>	2.0	0.000000
<b>1</b>	3.0	1.000000
<b>2</b>	4.0	0.500000
<b>3</b>	5.0	1.000000
<b>4</b>	7.0	0.000000
<b>5</b>	8.0	0.000000
<b>6</b>	11.0	0.000000
<b>7</b>	14.0	0.666667
<b>8</b>	15.0	1.000000
<b>9</b>	18.0	0.000000
<b>10</b>	19.0	0.500000
<b>11</b>	20.0	0.000000
<b>12</b>	21.0	0.333333
<b>13</b>	22.0	0.000000
<b>14</b>	26.0	1.000000
<b>15</b>	27.0	0.500000
<b>16</b>	28.0	0.500000
<b>17</b>	28.5	0.000000
<b>18</b>	29.0	1.000000
<b>19</b>	31.0	0.000000
<b>20</b>	34.0	1.000000
<b>21</b>	35.0	0.333333

	Age	Survived
22	38.0	1.000000
23	39.0	0.000000
24	40.0	0.000000
25	42.0	0.000000
26	45.0	0.000000
27	49.0	1.000000
28	54.0	0.000000
29	55.0	1.000000
30	58.0	1.000000
31	65.0	0.000000
32	66.0	0.000000

```
In [25]: train[["Embarked", "Survived"]].groupby(['Embarked'], as_index=False).mean().sort_values(by='Survived', ascending=False)
```

	Embarked	Survived
1	Q	0.571429
0	C	0.500000
2	S	0.375000

```
In [26]: fig = plt.figure()
ax = fig.add_axes([0,0,1,1])
ax.axis('equal')
l = ['C = Cherbourg', 'Q = Queenstown', 'S = Southampton']
s = [0.625000,0.500000,0.358974]
ax.pie(s, labels = l, autopct='%1.2f%')
plt.show()
```



```
In [27]: train.describe(include="all")
```

Out[27]:

	PassengerId	Survived	Pclass	Name	Gender	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
<b>count</b>	64.000000	64.0000	64.000000	64	64	49.000000	64.000000	64.000000	63	64.000000	14	63
<b>unique</b>	Nan	Nan	Nan	62	2	Nan	Nan	Nan	54	Nan	14	3
<b>top</b>	Nan	Nan	Nan	anita	male	Nan	Nan	Nan	113509	Nan	C85	S
<b>freq</b>	Nan	Nan	Nan	2	40	Nan	Nan	Nan	9	Nan	1	40
<b>mean</b>	32.500000	0.4375	2.359375	Nan	Nan	27.234694	0.781250	0.437500	Nan	30.210158	Nan	Nan
<b>std</b>	18.618987	0.5000	0.842656	Nan	Nan	16.528732	1.147375	1.037013	Nan	39.276369	Nan	Nan
<b>min</b>	1.000000	0.0000	1.000000	Nan	Nan	2.000000	0.000000	0.000000	Nan	7.225000	Nan	Nan
<b>25%</b>	16.750000	0.0000	2.000000	Nan	Nan	15.000000	0.000000	0.000000	Nan	8.044800	Nan	Nan
<b>50%</b>	32.500000	0.0000	3.000000	Nan	Nan	27.000000	0.000000	0.000000	Nan	17.900000	Nan	Nan
<b>75%</b>	48.250000	1.0000	3.000000	Nan	Nan	38.000000	1.000000	0.000000	Nan	32.415625	Nan	Nan
<b>max</b>	64.000000	1.0000	3.000000	Nan	Nan	66.000000	5.000000	5.000000	Nan	263.000000	Nan	Nan

In [28]:

```
#Feature Selection
column_train=['Age','Pclass','SibSp','Parch','Fare','Gender','Embarked']
#training values
X=train[column_train]
#target value
Y=train['Survived']
```

In [29]:

```
X['Age'].isnull().sum()
X['Pclass'].isnull().sum()
X['SibSp'].isnull().sum()
X['Parch'].isnull().sum()
X['Fare'].isnull().sum()
X['Gender'].isnull().sum()
X['Embarked'].isnull().sum()
```

Out[29]: np.int64(1)

```
In [30]: X['Age']=X['Age'].fillna(X['Age'].median())
X['Age'].isnull().sum()
```

```
Out[30]: np.int64(0)
```

```
In [31]: X['Embarked'] = train['Embarked'].fillna(method ='pad')
X['Embarked'].isnull().sum()
```

```
Out[31]: np.int64(0)
```

```
In [32]: d={'male':0, 'female':1}
X['Gender']=X['Gender'].apply(lambda x:d[x])
X['Gender'].head()
```

```
Out[32]: 0      1
         1      1
         2      0
         3      0
         4      0
Name: Gender, dtype: int64
```

```
In [33]: e={'C':0, 'Q':1 , 'S':2}
X['Embarked']=X['Embarked'].apply(lambda x:e[x])
X['Embarked'].head()
```

```
Out[33]: 0      2
         1      0
         2      2
         3      2
         4      2
Name: Embarked, dtype: int64
```

```
In [37]: results = pd.DataFrame({
    'Age': ['Logistic Regression', 'Survived', 'Dead', 'No injuries' , 'Deep injuries'],
    'Ticket': [0.75,0.66,0.76,0.66,0.74]})

result_df = results.sort_values(by='Ticket', ascending=False)
result_df = result_df.set_index('Ticket')
result_df.head(9)
```

Out[37]:

Ticket	Age
<b>0.76</b>	Dead
<b>0.75</b>	Logistic Regression
<b>0.74</b>	Deep injuries
<b>0.66</b>	Survived
<b>0.66</b>	No injuries

In [ ]: