Practical 1

Problem Statement - To apply Preprocessing techniques
            on raw dataset
            To perform E.DA.

      Libraries used - Numpy
                       Pandas
                       Matplot
                       scalearn
                       SKlearn

      data set - titanic

      Plots used - histograms & Boxplots

      train to test split ratio :> 80 : 20

# Practical No.1

## Data Science and Visualization (Honors Course)

**Name:Manjunath GB**

**PRN: 72018269H**

**Class: TE ENTC 'B'**

## In this practical we will access an open source dataset 'titanic.csv' and apply pre-processing techniques on the raw dataset.

In [1]:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

*We will now check the current version of all the packages which we imported.*

In [2]:

```python
pd.__version__
```

Out[2]:

```
'1.2.4'
```

In [3]:

```python
np.__version__
```

Out[3]:

```
'1.20.1'
```

In [4]:

```python
sns.__version__
```

Out[4]:

```
'0.11.1'
```

*We will now get the datasets which are already inbuilt in the packages.*

In [5]:

```python
sns.get_dataset_names()
```

Out[5]:

```
['anagrams',
 'anscombe',
 'attention',
 'brain_networks',
 'car_crashes',
 'diamonds',
 'dots',
 'exercise',
 'flights',
 'fmri'
```

```
        'gammas',
        'geyser',
        'iris',
        'mpg',
        'penguins',
        'planets',
        'taxis',
        'tips',
        'titanic']
```

**_There are various ways to import a dataset which are as follows:_**

In [6]:

```
dataset = sns.load_dataset('titanic')
```

In [7]:

```
dataset
```

Out[7]:

| | survived | pclass | sex | age | sibsp | parch | fare | embarked | class | who | adult_male | deck | embark_town | alive |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 3 | male | 22.0 | 1 | 0 | 7.2500 | S | Third | man | True | NaN | Southampton | no |
| 1 | 1 | 1 | female | 38.0 | 1 | 0 | 71.2833 | C | First | woman | False | C | Cherbourg | yes |
| 2 | 1 | 3 | female | 26.0 | 0 | 0 | 7.9250 | S | Third | woman | False | NaN | Southampton | yes |
| 3 | 1 | 1 | female | 35.0 | 1 | 0 | 53.1000 | S | First | woman | False | C | Southampton | yes |
| 4 | 0 | 3 | male | 35.0 | 0 | 0 | 8.0500 | S | Third | man | True | NaN | Southampton | no |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | .. |
| 886 | 0 | 2 | male | 27.0 | 0 | 0 | 13.0000 | S | Second | man | True | NaN | Southampton | no |
| 887 | 1 | 1 | female | 19.0 | 0 | 0 | 30.0000 | S | First | woman | False | B | Southampton | yes |
| 888 | 0 | 3 | female | NaN | 1 | 2 | 23.4500 | S | Third | woman | False | NaN | Southampton | no |
| 889 | 1 | 1 | male | 26.0 | 0 | 0 | 30.0000 | C | First | man | True | C | Cherbourg | yes |
| 890 | 0 | 3 | male | 32.0 | 0 | 0 | 7.7500 | Q | Third | man | True | NaN | Queenstown | no |

**891 rows × 15 columns**

In [8]:

```
df = pd.read_csv('https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/stuff/titani
c.csv')
```

In [9]:

```
df
```

Out[9]:

| | Survived | Pclass | Name | Sex | Age | Siblings/Spouses Aboard | Parents/Children Aboard | Fare |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 3 | Mr. Owen Harris Braund | male | 22.0 | 1 | 0 | 7.2500 |
| 1 | 1 | 1 | Mrs. John Bradley (Florence Briggs Thayer) Cum... | female | 38.0 | 1 | 0 | 71.2833 |
| 2 | 1 | 3 | Miss. Laina Heikkinen | female | 26.0 | 0 | 0 | 7.9250 |
| 3 | 1 | 1 | Mrs. Jacques Heath (Lily May Peel) Futrelle | female | 35.0 | 1 | 0 | 53.1000 |
| 4 | 0 | 3 | Mr. William Henry Allen | male | 35.0 | 0 | 0 | 8.0500 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |

| | Survived | Pclass | Name | Sex | Age | Siblings/Spouses Aboard | Parents/Children Aboard | Fare |
|---|---|---|---|---|---|---|---|---|
| 882 | 0 | 2 | Rev. Juozas Montvila | male | 27.0 | 0 | 0 | 13.0000 |
| 883 | 1 | 1 | Miss. Margaret Edith Graham | female | 19.0 | 0 | 0 | 30.0000 |
| 884 | 0 | 3 | Miss. Catherine Helen Johnston | female | 7.0 | 1 | 2 | 23.4500 |
| 885 | 1 | 1 | Mr. Karl Howell Behr | male | 26.0 | 0 | 0 | 30.0000 |
| 886 | 0 | 3 | Mr. Patrick Dooley | male | 32.0 | 0 | 0 | 7.7500 |

**887 rows × 8 columns**

*We will now perform certain pre processing operations on our dataset.*

In [11]:

```
df.columns   #The title of all the columns in the dataset.
```

Out[11]:

```
Index(['Survived', 'Pclass', 'Name', 'Sex', 'Age', 'Siblings/Spouses Aboard',
       'Parents/Children Aboard', 'Fare'],
      dtype='object')
```

In [12]:

```
df.shape
```

Out[12]:

```
(887, 8)
```

In [13]:

```
dataset.shape
```

Out[13]:

```
(891, 15)
```

In [14]:

```
dataset.columns
```

Out[14]:

```
Index(['survived', 'pclass', 'sex', 'age', 'sibsp', 'parch', 'fare',
       'embarked', 'class', 'who', 'adult_male', 'deck', 'embark_town',
       'alive', 'alone'],
      dtype='object')
```

In [16]:

```
df.head() #the .head() function returns the first five rows of dataset by default.
```

Out[16]:

| | Survived | Pclass | Name | Sex | Age | Siblings/Spouses Aboard | Parents/Children Aboard | Fare |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 3 | Mr. Owen Harris Braund | male | 22.0 | 1 | 0 | 7.2500 |
| 1 | 1 | 1 | Mrs. John Bradley (Florence Briggs Thayer) Cum... | female | 38.0 | 1 | 0 | 71.2833 |
| 2 | 1 | 3 | Miss. Laina Heikkinen | female | 26.0 | 0 | 0 | 7.9250 |
| 3 | 1 | 1 | Mrs. Jacques Heath (Lily May Peel) Futrelle | female | 35.0 | 1 | 0 | 53.1000 |
| 4 | 0 | 3 | Mr. William Henry Allen | male | 35.0 | 0 | 0 | 8.0500 |

In [17]:

```
dataset.head()
```

Out[17]:

| | survived | pclass | sex | age | sibsp | parch | fare | embarked | class | who | adult_male | deck | embark_town | alive | alo |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 3 | male | 22.0 | 1 | 0 | 7.2500 | S | Third | man | True | NaN | Southampton | no | Fa |
| 1 | 1 | 1 | female | 38.0 | 1 | 0 | 71.2833 | C | First | woman | False | C | Cherbourg | yes | Fa |
| 2 | 1 | 3 | female | 26.0 | 0 | 0 | 7.9250 | S | Third | woman | False | NaN | Southampton | yes | Tı |
| 3 | 1 | 1 | female | 35.0 | 1 | 0 | 53.1000 | S | First | woman | False | C | Southampton | yes | Fa |
| 4 | 0 | 3 | male | 35.0 | 0 | 0 | 8.0500 | S | Third | man | True | NaN | Southampton | no | Tı |

In [19]:

```
df.tail() #the .tail() function returns the last five rows of dataset by default.
```

Out[19]:

| | Survived | Pclass | Name | Sex | Age | Siblings/Spouses Aboard | Parents/Children Aboard | Fare |
|---|---|---|---|---|---|---|---|---|
| 882 | 0 | 2 | Rev. Juozas Montvila | male | 27.0 | 0 | 0 | 13.00 |
| 883 | 1 | 1 | Miss. Margaret Edith Graham | female | 19.0 | 0 | 0 | 30.00 |
| 884 | 0 | 3 | Miss. Catherine Helen Johnston | female | 7.0 | 1 | 2 | 23.45 |
| 885 | 1 | 1 | Mr. Karl Howell Behr | male | 26.0 | 0 | 0 | 30.00 |
| 886 | 0 | 3 | Mr. Patrick Dooley | male | 32.0 | 0 | 0 | 7.75 |

In [20]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 887 entries, 0 to 886
Data columns (total 8 columns):
 #   Column                   Non-Null Count  Dtype
---  ------                   --------------  -----
 0   Survived                 887 non-null    int64
 1   Pclass                   887 non-null    int64
 2   Name                     887 non-null    object
 3   Sex                      887 non-null    object
 4   Age                      887 non-null    float64
 5   Siblings/Spouses Aboard  887 non-null    int64
 6   Parents/Children Aboard  887 non-null    int64
 7   Fare                     887 non-null    float64
dtypes: float64(2), int64(4), object(2)
memory usage: 55.6+ KB
```

In [21]:

```
dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 15 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   survived     891 non-null    int64
 1   pclass       891 non-null    int64
 2   sex          891 non-null    object
 3   age          714 non-null    float64
 4   sibsp        891 non-null    int64
 5   parch        891 non-null    int64
 6   fare         891 non-null    float64
 7   embarked     889 non-null    object
 8   class        891 non-null    category
```

```
 9   who          891 non-null    object
 10  adult_male   891 non-null    bool
 11  deck         203 non-null    category
 12  embark_town  889 non-null    object
 13  alive        891 non-null    object
 14  alone        891 non-null    bool
dtypes: bool(2), category(2), float64(2), int64(4), object(5)
memory usage: 80.7+ KB
```

In [22]:

```
dataset.describe()
```

Out[22]:

|       | survived   | pclass     | age        | sibsp      | parch      | fare       |
|-------|------------|------------|------------|------------|------------|------------|
| count | 891.000000 | 891.000000 | 714.000000 | 891.000000 | 891.000000 | 891.000000 |
| mean  | 0.383838   | 2.308642   | 29.699118  | 0.523008   | 0.381594   | 32.204208  |
| std   | 0.486592   | 0.836071   | 14.526497  | 1.102743   | 0.806057   | 49.693429  |
| min   | 0.000000   | 1.000000   | 0.420000   | 0.000000   | 0.000000   | 0.000000   |
| 25%   | 0.000000   | 2.000000   | 20.125000  | 0.000000   | 0.000000   | 7.910400   |
| 50%   | 0.000000   | 3.000000   | 28.000000  | 0.000000   | 0.000000   | 14.454200  |
| 75%   | 1.000000   | 3.000000   | 38.000000  | 1.000000   | 0.000000   | 31.000000  |
| max   | 1.000000   | 3.000000   | 80.000000  | 8.000000   | 6.000000   | 512.329200 |

In [23]:

```
df.describe()
```

Out[23]:

|       | Survived   | Pclass     | Age        | Siblings/Spouses Aboard | Parents/Children Aboard | Fare      |
|-------|------------|------------|------------|-------------------------|-------------------------|-----------|
| count | 887.000000 | 887.000000 | 887.000000 | 887.000000              | 887.000000              | 887.00000 |
| mean  | 0.385569   | 2.305524   | 29.471443  | 0.525366                | 0.383315                | 32.30542  |
| std   | 0.487004   | 0.836662   | 14.121908  | 1.104669                | 0.807466                | 49.78204  |
| min   | 0.000000   | 1.000000   | 0.420000   | 0.000000                | 0.000000                | 0.00000   |
| 25%   | 0.000000   | 2.000000   | 20.250000  | 0.000000                | 0.000000                | 7.92500   |
| 50%   | 0.000000   | 3.000000   | 28.000000  | 0.000000                | 0.000000                | 14.45420  |
| 75%   | 1.000000   | 3.000000   | 38.000000  | 1.000000                | 0.000000                | 31.13750  |
| max   | 1.000000   | 3.000000   | 80.000000  | 8.000000                | 6.000000                | 512.32920 |

In [24]:

```
df.count()
```

Out[24]:

```
Survived                   887
Pclass                     887
Name                       887
Sex                        887
Age                        887
Siblings/Spouses Aboard    887
Parents/Children Aboard    887
Fare                       887
dtype: int64
```

In [25]:

```
dataset.count()
```

```
Out[25]:
survived        891
pclass          891
sex             891
age             714
sibsp           891
parch           891
fare            891
embarked        889
class           891
who             891
adult_male      891
deck            203
embark_town     889
alive           891
alone           891
dtype: int64
```

In [26]:

```python
dataset.isnull().sum()
```

Out[26]:

```
survived          0
pclass            0
sex               0
age             177
sibsp             0
parch             0
fare              0
embarked          2
class             0
who               0
adult_male        0
deck            688
embark_town       2
alive             0
alone             0
dtype: int64
```

In [27]:

```python
dataset = dataset.drop('deck', axis = 1)
```

In [28]:

```python
dataset.isnull().sum()
```

Out[28]:

```
survived          0
pclass            0
sex               0
age             177
sibsp             0
parch             0
fare              0
embarked          2
class             0
who               0
adult_male        0
embark_town       2
alive             0
alone             0
dtype: int64
```
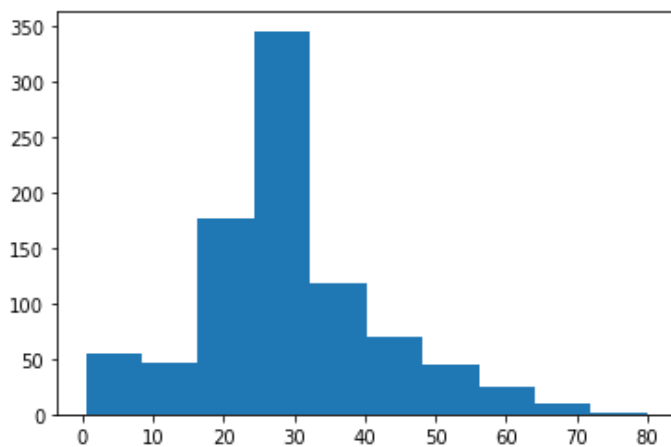
In [29]:

```python
dataset['age'] = dataset['age'].fillna(dataset['age'].median())
```

```
In [30]:
```

```
dataset.isnull().sum()
```

```
Out[30]:
```

```
survived        0
pclass          0
sex             0
age             0
sibsp           0
parch           0
fare            0
embarked        2
class           0
who             0
adult_male      0
embark_town     2
alive           0
alone           0
dtype: int64
```

```
In [31]:
```

```
dataset['embarked'].mode()[0]
```

```
Out[31]:
```

```
'S'
```

```
In [32]:
```

```
dataset['embark_town'].mode()[0]
```

```
Out[32]:
```

```
'Southampton'
```

```
In [33]:
```

```
dataset['embarked'] = dataset['embarked'].fillna(
    dataset['embarked'].mode()[0])
```

```
In [34]:
```

```
dataset['embark_town'] = dataset['embark_town'].fillna(
    dataset['embark_town'].mode()[0])
```

```
In [35]:
```

```
dataset.isnull().sum()
```

```
Out[35]:
```

```
survived        0
pclass          0
sex             0
age             0
sibsp           0
parch           0
fare            0
embarked        0
class           0
who             0
adult_male      0
embark_town     0
alive           0
alone           0
dtype: int64
```

```
In [36]:
```

```
dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 14 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   survived     891 non-null    int64
 1   pclass       891 non-null    int64
 2   sex          891 non-null    object
 3   age          891 non-null    float64
 4   sibsp        891 non-null    int64
 5   parch        891 non-null    int64
 6   fare         891 non-null    float64
 7   embarked     891 non-null    object
 8   class        891 non-null    category
 9   who          891 non-null    object
 10  adult_male   891 non-null    bool
 11  embark_town  891 non-null    object
 12  alive        891 non-null    object
 13  alone        891 non-null    bool
dtypes: bool(2), category(1), float64(2), int64(4), object(5)
memory usage: 79.4+ KB
```

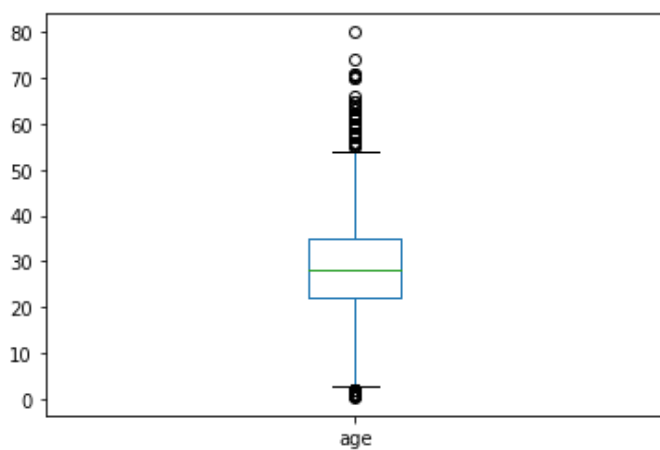## Visualization of dataset

In [37]:

```
plt.hist(dataset['age']);
```



In [38]:

```
dataset['age'].plot(kind='box')
```
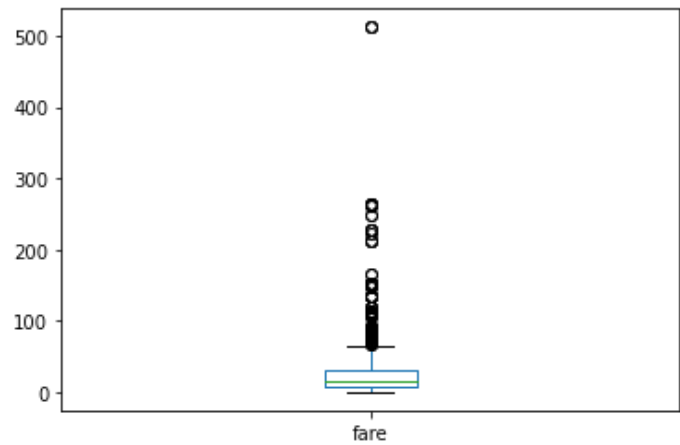
Out[38]:

```
<AxesSubplot:>
```



In [39]:

```
dataset['fare'].plot(kind='box')
```

```
<AxesSubplot:>
```



In [40]:

```
dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 14 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   survived     891 non-null    int64
 1   pclass       891 non-null    int64
 2   sex          891 non-null    object
 3   age          891 non-null    float64
 4   sibsp        891 non-null    int64
 5   parch        891 non-null    int64
 6   fare         891 non-null    float64
 7   embarked     891 non-null    object
 8   class        891 non-null    category
 9   who          891 non-null    object
 10  adult_male   891 non-null    bool
 11  embark_town  891 non-null    object
 12  alive        891 non-null    object
 13  alone        891 non-null    bool
dtypes: bool(2), category(1), float64(2), int64(4), object(5)
memory usage: 79.4+ KB
```

In [41]:

```
pd.get_dummies(dataset).head()
```

Out[41]:

| | survived | pclass | age | sibsp | parch | fare | adult_male | alone | sex_female | sex_male | ... | class_Second | class_Third | who |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 3 | 22.0 | 1 | 0 | 7.2500 | True | False | 0 | 1 | ... | 0 | 1 | |
| 1 | 1 | 1 | 38.0 | 1 | 0 | 71.2833 | False | False | 1 | 0 | ... | 0 | 0 | |
| 2 | 1 | 3 | 26.0 | 0 | 0 | 7.9250 | False | True | 1 | 0 | ... | 0 | 1 | |
| 3 | 1 | 1 | 35.0 | 1 | 0 | 53.1000 | False | False | 1 | 0 | ... | 0 | 0 | |
| 4 | 0 | 3 | 35.0 | 0 | 0 | 8.0500 | True | True | 0 | 1 | ... | 0 | 1 | |

**5 rows × 24 columns**

# Training our Dataset

*Importing the required library.*

In [42]:

```python
from sklearn.model_selection import train_test_split
```

In [43]:

```python
train, test = train_test_split(dataset,test_size=0.20)
```

In [44]:

```python
len(dataset)
```

Out[44]:

891

In [45]:

```python
len(train)
```

Out[45]:

712

In [46]:

```python
len(test)
```

Out[46]:

179

In [ ]:

In [ ]:

In [ ]: