# Assignment 1

```python
#import required library
import pandas as pd
import numpy as np
```

```python
#read the dataset
dataset = pd.read_csv('Titanic.csv')
dataset.head()
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fai |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 892 | 0 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.82ξ |
| 1 | 893 | 1 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.00( |

```python
#check the dataset information
dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 12 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   PassengerId  418 non-null    int64
 1   Survived     418 non-null    int64
 2   Pclass       418 non-null    int64
 3   Name         418 non-null    object
 4   Sex          418 non-null    object
 5   Age          332 non-null    float64
 6   SibSp        418 non-null    int64
 7   Parch        418 non-null    int64
 8   Ticket       418 non-null    object
 9   Fare         417 non-null    float64
 10  Cabin        91 non-null     object
 11  Embarked     418 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 39.3+ KB
```

```python
#describe the dataset
dataset.describe()
```

| | PassengerId | Survived | Pclass | Age | SibSp | Parch | F |
|---|---|---|---|---|---|---|---|
| count | 418.000000 | 418.000000 | 418.000000 | 332.000000 | 418.000000 | 418.000000 | 417.000 |
| mean | 1100.500000 | 0.363636 | 2.265550 | 30.272590 | 0.447368 | 0.392344 | 35.627 |
| std | 120.810458 | 0.481622 | 0.841838 | 14.181209 | 0.896760 | 0.981429 | 55.907 |
| min | 892.000000 | 0.000000 | 1.000000 | 0.170000 | 0.000000 | 0.000000 | 0.000 |
| 25% | 996.250000 | 0.000000 | 1.000000 | 21.000000 | 0.000000 | 0.000000 | 7.895 |
| 50% | 1100.500000 | 0.000000 | 3.000000 | 27.000000 | 0.000000 | 0.000000 | 14.454 |
| 75% | 1204.750000 | 1.000000 | 3.000000 | 39.000000 | 1.000000 | 0.000000 | 31.500 |
| max | 1309.000000 | 1.000000 | 3.000000 | 76.000000 | 8.000000 | 9.000000 | 512.329 |

We will check the null values in dataset

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

```
PassengerId      0
Survived         0
Pclass           0
Name             0
Sex              0
Age             86
SibSp            0
Parch            0
Ticket           0
Fare             1
Cabin          327
```

```
Embarked          0
dtype: int64
```

## Imputing missing values of 'Age' column

```
#lets find the mean of the 'Age' column
d1 = dataset['Age'].mean()
d1
```

```
    30.272590361445783
```

```
#find the round value of d1
d1 = round(d1)
d1
```

```
    30
```

```
#fill the rounded value with missing value using fillna() method
dataset['Age'] = dataset['Age'].fillna(d1)
```

```
#check the result is affected or not
dataset.isnull().sum()
```

```
    PassengerId       0
    Survived          0
    Pclass            0
    Name              0
    Sex               0
    Age               0
    SibSp             0
    Parch             0
    Ticket            0
    Fare              1
    Cabin           327
    Embarked          0
    dtype: int64
```

```
#do the same process for anothe columns that are missing value
d2 = round(dataset['Fare'].mean())
dataset['Fare'] = dataset['Fare'].fillna(d2)
dataset.isnull().sum()
```

```
    PassengerId       0
    Survived          0
    Pclass            0
    Name              0
    Sex               0
    Age               0
    SibSp             0
    Parch             0
    Ticket            0
    Fare              0
    Cabin           327
    Embarked          0
    dtype: int64
```

```
#same for 'Cabin'
dataset['Cabin'].unique()
```

```
    array([nan, 'B45', 'E31', 'B57 B59 B63 B66', 'B36', 'A21', 'C78', 'D34',
           'D19', 'A9', 'D15', 'C31', 'C23 C25 C27', 'F G63', 'B61', 'C53',
           'D43', 'C130', 'C132', 'C101', 'C55 C57', 'B71', 'C46', 'C116',
           'F', 'A29', 'G6', 'C6', 'C28', 'C51', 'E46', 'C54', 'C97', 'D22',
           'B10', 'F4', 'E45', 'E52', 'D30', 'B58 B60', 'E34', 'C62 C64',
           'A11', 'B11', 'C80', 'F33', 'C85', 'D37', 'C86', 'D21', 'C89',
           'F E46', 'A34', 'D', 'B26', 'C22 C26', 'B69', 'C32', 'B78',
           'F E57', 'F2', 'A18', 'C106', 'B51 B53 B55', 'D10 D12', 'E60',
           'E50', 'E39 E41', 'B52 B54 B56', 'C39', 'B24', 'D28', 'B41', 'C7',
           'D40', 'D38', 'C105'], dtype=object)
```

```
#lets count the value
dataset['Cabin'].value_counts()
```

```
    B57 B59 B63 B66    3
    B45                2
    C89                2
    C55 C57            2
    A34                2
                      ..
    E52                1
    D30                1
```

```
              1
              1
              1

          E31
          C62 C64
          C105
          Name: Cabin, Length: 76, dtype: int64
```

#'ffill' stands for 'forward fill' and will propagate last valid observation forward.
dataset['Cabin'] = dataset['Cabin'].ffill()
dataset.isnull().sum()

```
      PassengerId    0
      Survived       0
      Pclass         0
      Name           0
      Sex            0
      Age            0
      SibSp          0
      Parch          0
      Ticket         0
      Fare           0
      Cabin          12
      Embarked       0
      dtype: int64
```

#bfill() will backward fill the NaN values that are present in the pandas dataframe.
dataset['Cabin'] = dataset['Cabin'].bfill()
dataset.isnull().sum()

```
      PassengerId    0
      Survived       0
      Pclass         0
      Name           0
      Sex            0
      Age            0
      SibSp          0
      Parch          0
      Ticket         0
      Fare           0
      Cabin          0
      Embarked       0
      dtype: int64
```

#lets check our preprocessed dataset without any missing values
dataset.head()

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fai |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 892 | 0 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.82ξ |
| **1** | 893 | 1 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.00( |