

Group A

Assignment 1

Data Wrangling I

Unsupported Cell Type. Double-Click to inspect/edit the content.

Import all the required Python Libraries.

```
import numpy as np
import pandas as pd
```

Load the Dataset into pandas dataframe.

```
df = pd.read_csv('titanic.csv')
```

```
df.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	F
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs)	female	38.0	1	0	PC 17599	71.2

```
df.tail()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00

```
df.sample()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fa
				Brown,						

Data Preprocessing

check for missing values in the data using pandas isnull()

```
df.isnull().sum()
```

PassengerId	0
Survived	0
Pclass	0
Name	0
Sex	0
Age	177
SibSp	0
Parch	0
Ticket	0
Fare	0
Cabin	687

```
Embarked      2
dtype: int64
```

```
df['Age'].fillna(df['Age'].mean(), inplace = True)
```

```
df['Age'].isna().sum()
```

```
0
```

```
df['Embarked'].value_counts()
```

```
S    644
C    168
Q     77
Name: Embarked, dtype: int64
```

```
df['Embarked'].fillna('S',inplace = True)
```

```
df['Embarked'].isna().sum()
```

```
0
```

```
df.drop(columns = ['Cabin'],axis=1,inplace=True)
```

```
df.isnull().sum()
```

```
PassengerId    0
Survived        0
Pclass          0
Name            0
Sex             0
Age            77
SibSp           0
Parch           0
Ticket          0
Fare            0
Embarked        0
dtype: int64
```

✓ describe() function to get some initial statistics. Provide variable descriptions.

```
df.describe()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
std	257.353842	0.486592	0.836071	13.002015	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	22.000000	0.000000	0.000000	7.910461
50%	446.000000	0.000000	3.000000	29.699118	0.000000	0.000000	14.454295
75%	668.500000	1.000000	3.000000	35.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.3291

✓ Types of variables

```
df.dtypes
```

```
PassengerId    int64
Survived        int64
Pclass          int64
Name            object
Sex             object
Age            float64
SibSp           int64
Parch           int64
Ticket          object
Fare            float64
Embarked        object
dtype: object
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 11 columns):
 #   Column        Non-Null Count  Dtype
---  --
 0   PassengerId   891 non-null    int64
 1   Survived      891 non-null    int64
 2   Pclass        891 non-null    int64
 3   Name          891 non-null    object
 4   Sex           891 non-null    object
 5   Age           891 non-null    float64
 6   SibSp         891 non-null    int64
 7   Parch         891 non-null    int64
 8   Ticket        891 non-null    object
 9   Fare          891 non-null    float64
10   Embarked      891 non-null    object
dtypes: float64(2), int64(5), object(4)
memory usage: 76.7+ KB
```

✓ Check the dimensions of the data frame

```
df.shape
```

```
(891, 11)
```

```
df.shape[0]
```

```
891
```

✓ Data Formatting and Data Normalization

- ✓ Summarize the types of variables by checking the data types (i.e., character, numeric, integer, factor, and logical) of the variables in the data set.

```
df.nunique()
```

```
PassengerId    891
Survived         2
Pclass          3
Name            891
Sex              2
Age             89
SibSp           7
Parch           7
Ticket         681
Fare           248
Embarked        3
dtype: int64
```

```
df['Survived'].value_counts()
```

```
0    549
1    342
Name: Survived, dtype: int64
```

```
df['Pclass'].value_counts()
```

```
3    491
1    216
2    184
Name: Pclass, dtype: int64
```

```
df['Sex'].value_counts()
```

```
male      577
female    314
Name: Sex, dtype: int64
```

```
df['SibSp'].value_counts()
```

```
0    608
1    209
2     28
4     18
3     16
```

```
8      7
5      5
Name: SibSp, dtype: int64
```

```
df['Parch'].value_counts()
```

```
0      678
1      118
2       80
5         5
3         5
4         4
6         1
Name: Parch, dtype: int64
```

```
df['Embarked'].value_counts()
```

```
S      646
C      168
Q       77
Name: Embarked, dtype: int64
```

✓ If variables are not in the correct data type, apply proper type conversions.

```
df.dtypes
```

```
PassengerId    int64
Survived        int64
Pclass         int64
Name           object
Sex            object
Age           float64
SibSp          int64
Parch          int64
Ticket         object
Fare           float64
Embarked       object
dtype: object
```

```
df['Age'] = df['Age'].astype('int64')
```

✓ Turn categorical variables into quantitative variables in Python.

```
df["Sex"].replace(['female','male'],[0,1],inplace = True)
```

```
df['Sex'].value_counts()
```

```
1      577
0      314
Name: Sex, dtype: int64
```

```
df['Embarked'].replace(['C','Q','S'],[1,2,3],inplace= True)
```

```
df['Embarked'].value_counts()
```

```
3      646
1      168
2       77
Name: Embarked, dtype: int64
```

```
df.dtypes
```

```
PassengerId    int64
Survived        int64
Pclass         int64
Name           object
Sex            int64
Age            int64
SibSp          int64
Parch          int64
Ticket         object
Fare           float64
Embarked       int64
dtype: object
```

```
df.drop(columns=['Name','PassengerId','Ticket'],axis = 1,inplace = True)
```

```
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 8 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Survived    891 non-null    int64
1   Pclass      891 non-null    int64
2   Sex         891 non-null    int64
3   Age         891 non-null    int64
4   SibSp       891 non-null    int64
5   Parch       891 non-null    int64
6   Fare        891 non-null    float64
7   Embarked    891 non-null    int64
dtypes: float64(1), int64(7)
memory usage: 55.8 KB
```