

Import Modules

```
In [1]: import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

In [2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
matplotlib inline

In [3]: #loading DATA SET
train = pd.read_csv("TITANIC.csv")
test = pd.read_csv("TITANIC.csv")
train.head()
```

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked		
0	892	1	0	3	Kelly, Mr. James	male	34.5	0	0	330911	7.0292	NaN	Q
1	893	1	3	Wiles, Mrs. James (Ellen Needs)	female	47.0	1	0	363272	7.0200	NaN	S	
2	894	0	2	Wiles, Mr. Thomas Francis	male	62.0	0	0	240276	9.6875	NaN	Q	
3	895	0	3	Witz, Mr. Albert	male	27.0	0	0	315154	8.6625	NaN	S	
4	896	1	3	Hivonen, Mrs. Alexander (Helga E Lindqvist)	female	22.0	1	1	3101298	12.2875	NaN	S	

```
In [4]: # statistical info
train.describe()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	418.000000	418.000000	418.000000	332.000000	418.000000	418.000000	417.000000
mean	1105.500000	0.363636	2.265550	30.277590	0.473661	0.382344	55.627388
std	120.819458	0.481622	0.841838	14.181209	0.896760	0.981429	55.627388
min	892.000000	0.000000	1.000000	0.170000	0.000000	0.000000	0.000000
25%	996.250000	0.000000	1.000000	21.000000	0.000000	0.000000	7.895680
50%	1105.500000	0.000000	3.000000	27.000000	0.000000	0.000000	14.454200
75%	1204.750000	1.000000	3.000000	39.000000	1.000000	0.000000	31.500000
max	1309.000000	1.000000	3.000000	76.000000	0.000000	9.000000	512.329200

```
In [5]: # datatype info
train.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                  417 non-null    object
 5   Age                  332 non-null    float64
 6   SibSp               418 non-null    int64
 7   Parch              418 non-null    int64
 8   Ticket              417 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: 59.3+ KB
```

```
In [6]: #EXPLORATORY DATA ANALYSIS
# categorical attributes
sns.countplot(train['Survived'])
```

Survived	count
0	140
1	278

```
In [7]: sns.countplot(train['Pclass'])
```

Pclass	count
1	103
2	97
3	218

```
In [8]: sns.countplot(train['Sex'])
```

Sex	count
male	260
female	158

```
In [9]: sns.countplot(train['SibSp'])
```

SibSp	count
0	280
1	103
2	10
3	1
4	0
5	0
6	0
8	0

```
In [10]: sns.countplot(train['Parch'])
```

Parch	count
0	300
1	50
2	10
3	0
4	0
5	0
6	0
9	0

```
In [11]: sns.countplot(train['Embarked'])
```

Embarked	count
Q	40
S	270
C	90

```
In [12]: # numerical attributes
sns.distplot(train['Age'])
```

Age	Density
0	0.000
10	0.005
20	0.035
30	0.025
40	0.015
50	0.005
60	0.002
70	0.001
80	0.000

```
In [13]: sns.distplot(train['Fare'])
```

Fare	Density
0	0.000
50	0.035
100	0.005
200	0.002
300	0.001
400	0.000
500	0.000

```
In [14]: class_fare = train.pivot_table(index='Pclass', values='Fare', values='Fare')
class_fare.plot(kind='bar')
plt.xlabel('Pclass')
plt.ylabel('Avg. Fare')
plt.xticks(rotation=90)
plt.show()
```

Pclass	Avg. fare
1	51.2
2	20.9
3	14.5

```
In [15]: class_fare = train.pivot_table(index='Pclass', values='Fare', aggfunc=np.sum)
class_fare.plot(kind='bar')
plt.xlabel('Pclass')
plt.ylabel('Total Fare')
plt.xticks(rotation=90)
plt.show()
```

Pclass	Total fare
1	5170.64
2	1980.75
3	2369.75

```
In [16]: sns.barplot(data=train, x='Pclass', y='Fare', hue='Survived')
```

Pclass	Survived	Fare
1	0	7.0292
1	1	51.2108
2	0	9.6875
2	1	20.6619
3	0	8.6625
3	1	23.9854

```
In [17]: sns.barplot(data=train, x='Survived', y='Fare', hue='Pclass')
```

Survived	Pclass	Fare
0	1	7.0292
0	2	9.6875
0	3	8.6625
1	1	51.2108
1	2	20.6619
1	3	23.9854

Preprocessing of Data

```
In [18]: train_len = len(train)
# combine two dataframes
df = pd.concat([train, test], axis=0)
df = df.reset_index(drop=True)
df.head()
```

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	
0	892	1	3	Kelly, Mr. James	male	34.5	0	0	330911	7.0292	NaN	Q
1	893	1	3	Wiles, Mrs. James (Ellen Needs)	female	47.0	1	0	363272	7.0200	NaN	S
2	894	0	2	Wiles, Mr. Thomas Francis	male	62.0	0	0	240276	9.6875	NaN	Q
3	895	0	3	Witz, Mr. Albert	male	27.0	0	0	315154	8.6625	NaN	S
4	896	1	3	Hivonen, Mrs. Alexander (Helga E Lindqvist)	female	22.0	1	1	3101298	12.2875	NaN	S

```
In [19]: # find the null values
df.isnull().sum()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
PassengerId	0	0	0	0	0	0	0	0	0	0	0	0
Survived	0	0	0	0	0	0	0	0	0	0	0	0
Pclass	0	0	0	0	0	0	0	0	0	0	0	0
Name	0	0	0	0	0	0	0	0	0	0	0	0
Sex	0	0	0	0	0	0	0	0	0	0	0	0
Age	0	0	0	0	0	0	0	0	0	0	0	0
SibSp	0	0	0	0	0	0	0	0	0	0	0	0
Parch	0	0	0	0	0	0	0	0	0	0	0	0
Ticket	0	0	0	0	0	0	0	0	0	0	0	0
Fare	0	0	0	0	0	0	0	0	0	0	0	0
Cabin	0	0	0	0	0	0	0	0	0	0	0	0
Embarked	0	0	0	0	0	0	0	0	0	0	0	0
dtype:	int64	int64	int64	object	object	float64	int64	int64	object	float64	object	object

```
In [20]: # drop or delete the column
df = df.drop(columns='Cabin', axis=1)
```

```
Out[20]: 30.2759381445783
```

```
In [22]: # fill missing values using mean of the numerical column
df['Age'] = df['Age'].fillna(df['Age'].mean())
df['Fare'] = df['Fare'].fillna(df['Fare'].mean())
```

```
In [23]: df['Embarked'].mode()[0]
```

```
Out[23]: 'S'
```

```
In [24]: # fill missing values using mode of the categorical column
df['Embarked'] = df['Embarked'].fillna(df['Embarked'].mode()[0])
```

Log transformation for uniform data distribution

```
In [25]: sns.distplot(df['Fare'])
```

Fare	Density
0	0.000
50	0.035
100	0.005
200	0.002
300	0.001
400	0.000
500	0.000

```
In [26]: df['Fare'] = np.log(df['Fare']+1)
sns.distplot(df['Fare'])
```

Fare	Density
-1	0.000
0	0.035
1	0.005
2	0.002
3	0.001
4	0.000
5	0.000
6	0.000

Correlation Matrix

```
In [27]: corr = df.corr()
plt.figure(figsize=(15, 9))
sns.heatmap(corr, annot=True, cmap='coolwarm')
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
PassengerId	1	-0.023	-0.027	-0.031	0.0038	0.043	0.021
Survived	-0.023	1	-0.11	-1.1e-05	0.1	0.16	0.22
Pclass	-0.027	-0.11	1	-0.44	0.0011	0.019	-0.74
Age	-0.031	-1.1e-05	-0.44	1	-0.08	-0.045	0.33
SibSp	0.0038	0.1	0.0011	-0.08	1	0.31	0.3
Parch	0.043	0.16	0.019	-0.045	0.31	1	0.31
Fare	0.021	0.22	-0.74	0.33	0.3	0.31	1

```
In [28]: df.head()
```

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked	
0	892	0	3	Kelly, Mr. James	male	34.5	0	0	330911	2.178064	Q
1	893	1	3	Wiles, Mrs. James (Ellen Needs)	female	47.0	1	0	363272	2.078442	S
2	894	0	2	Wiles, Mr. Thomas Francis	male	62.0	0	0	240276	2.369075	Q
3	895	0	3	Witz, Mr. Albert	male	27.0	0	0	315154	2.268252	S
4	896	1	3	Hivonen, Mrs. Alexander (Helga E Lindqvist)	female	22.0	1	1	3101298	2.586824	S

```
In [29]: # drop unnecessary columns
df = df.drop(columns='Name', 'Ticket', axis=1)
df.head()
```

PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	
0	892	0	3	male	34.5	0	0	2.178064	Q
1	893	1	3	female	47.0	1	0	2.078442	S
2	894	0	2	male	62.0	0	0	2.369075	Q
3	895	0	3	male	27.0	0	0	2.268252	S
4	896	1	3	female	22.0	1	1	2.586824	S

Label Encoding

```
In [30]: from sklearn.preprocessing import LabelEncoder
cols = ['Sex', 'Embarked']
le = LabelEncoder()

for col in cols:
    df[col] = le.fit_transform(df[col])

df.head()
```

PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	
0	892	0	3	1	34.5	0	0	2.178064	1
1	893	1	3	0	47.0	1	0	2.078442	2
2	894	0	2	1	62.0	0	0	2.369075	1
3	895	0	3	1	27.0	0	0	2.268252	2
4	896	1	3	0	22.0	1	1	2.586824	2

Train-Test Split

```
In [31]: train = df.iloc[:train_len, :]
test = df.iloc[train_len, :]
train.head()
```

PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	
0	892	0	3	1	34.5	0	0	2.178064	1
1	893	1	3	0	47.0	1	0	2.078442	2
2	894	0	2	1	62.0	0	0	2.369075	1
3	895	0	3	1	27.0	0	0	2.268252	2
4	896	1	3	0	22.0	1	1	2.586824	2

```
In [32]: test.head()
```

PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	
418	892	0	3	1	34.5	0	0	2.178064	1
419	893	1	3	0	47.0	1	0	2.078442	2
420	894	0	2	1	62.0	0	0	2.369075	1
421	895	0	3	1	27.0	0	0	2.268252	2
422	896	1	3	0	22.0	1	1	2.586824	2

Model Training

```
In [34]: from sklearn.model_selection import train_test_split, cross_val_score
# classify column
def classify(y_train, y_test):
    x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=42)
    model = LogisticRegression()
    print('Accuracy:', model.score(x_test, y_test))

score = cross_val_score(model, X, y, cv=5)
print('CV Score:', np.mean(score))
```

```
In [35]: from sklearn.linear_model import LogisticRegression
model = LogisticRegression()
classify(model)
```

```
Accuracy: 1.8
CV Score: 1.8
```

```
In [36]: from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()
classify(model)
```

```
Accuracy: 1.8
CV Score: 1.8
```

```
In [37]: from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
classify(model)
```

```
Accuracy: 1.8
CV Score: 1.8
```

```
In [38]: from sklearn.ensemble import ExtraTreesClassifier
model = ExtraTreesClassifier()
classify(model)
```

```
Accuracy: 1.8
CV Score: 1.8
```

Complete Model Training with Full Data

```
In [49]: model = RandomForestClassifier()
model.fit(X, y)
```

PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	
418	892	0	3	1	34.5	0	0	2.178064	1
419	893	1	3	0	47.0	1	0	2.078442	2
420	894	0	2	1	62.0	0	0	2.369075	1
421	895	0	3	1	27.0	0	0	2.268252	2
422	896	1	3	0	22.0	1	1	2.586824	2

```
In [50]: # input split for test data
X_test = test.drop(columns='PassengerId', 'Survived', axis=1)
X_test.head()
```

Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	
418	3	1	34.5	0	0	2.178064	1
419	3	0	47.0	1	0	2.078442	2
420	2	1	62.0	0	0	2.369075	1
421	3	1	27.0	0	0	2.268252	2
422	3	0	22.0	1	1	2.586824	2

```
In [51]: pred = model.predict(X_test)
pred
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
--	---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----