

ds-task-1

January 4, 2024

1 Import necessary libraries

```
[1]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report
```

2 loading data

```
[2]: df = pd.read_csv('E:\\codsoft Data science\\titanic survey task_1\\Titanic-Dataset.csv')
df.head()
```

```
[2]:
```

	PassengerId	Survived	Pclass	\
0	1	0	3	
1	2	1	1	
2	3	1	3	
3	4	1	1	
4	5	0	3	

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.0	1	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	
2	Heikkinen, Miss. Laina	female	26.0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	
4	Allen, Mr. William Henry	male	35.0	0	

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

3 preprocessing

```
[3]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 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            714 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   Cabin          204 non-null    object
11   Embarked        889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

```
[4]: df.describe()
```

```
[4]:
```

	PassengerId	Survived	Pclass	Age	SibSp	\
count	891.000000	891.000000	891.000000	714.000000	891.000000	
mean	446.000000	0.383838	2.308642	29.699118	0.523008	
std	257.353842	0.486592	0.836071	14.526497	1.102743	
min	1.000000	0.000000	1.000000	0.420000	0.000000	
25%	223.500000	0.000000	2.000000	20.125000	0.000000	
50%	446.000000	0.000000	3.000000	28.000000	0.000000	
75%	668.500000	1.000000	3.000000	38.000000	1.000000	
max	891.000000	1.000000	3.000000	80.000000	8.000000	

	Parch	Fare
count	891.000000	891.000000
mean	0.381594	32.204208
std	0.806057	49.693429
min	0.000000	0.000000
25%	0.000000	7.910400
50%	0.000000	14.454200
75%	0.000000	31.000000
max	6.000000	512.329200

```
[5]: df.shape
```

```
[5]: (891, 12)
```

```
[6]: df.isnull().sum()
```

```
[6]: 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
```

```
[7]: df.isnull().sum().sum()
```

```
[7]: 866
```

```
[8]: # Handle missing values
     df['Age'].fillna(df['Age'].median(), inplace=True)
     df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)
```

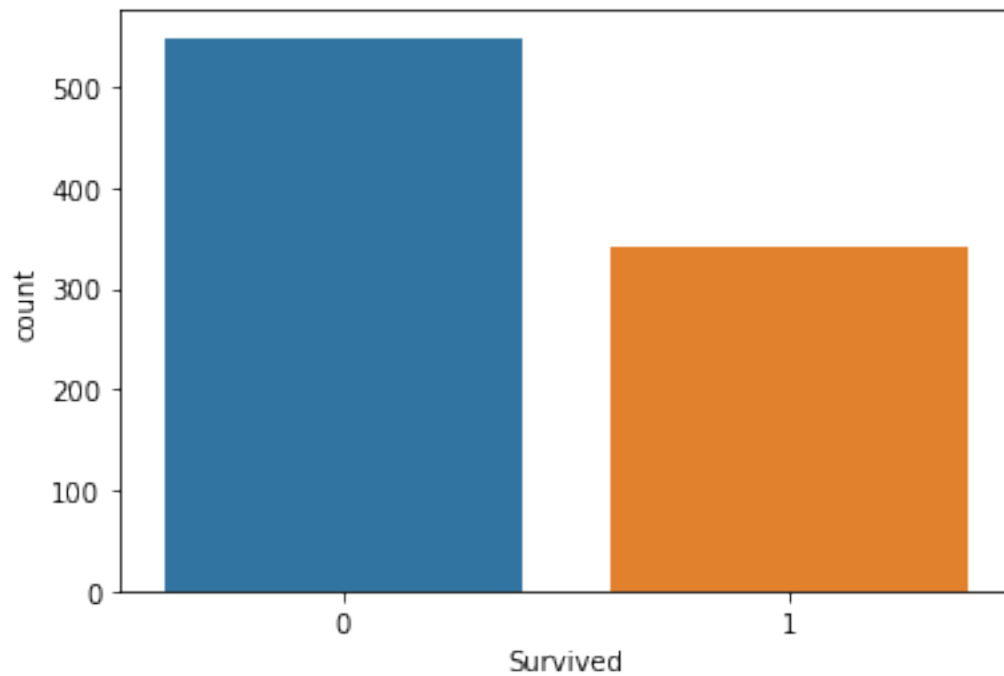
```
[9]: # Convert categorical features to numerical
     df = pd.get_dummies(df, columns=['Sex', 'Embarked'], drop_first=True)
```

4 visualization

```
[13]: import matplotlib.pyplot as plt
     import seaborn as sns
```

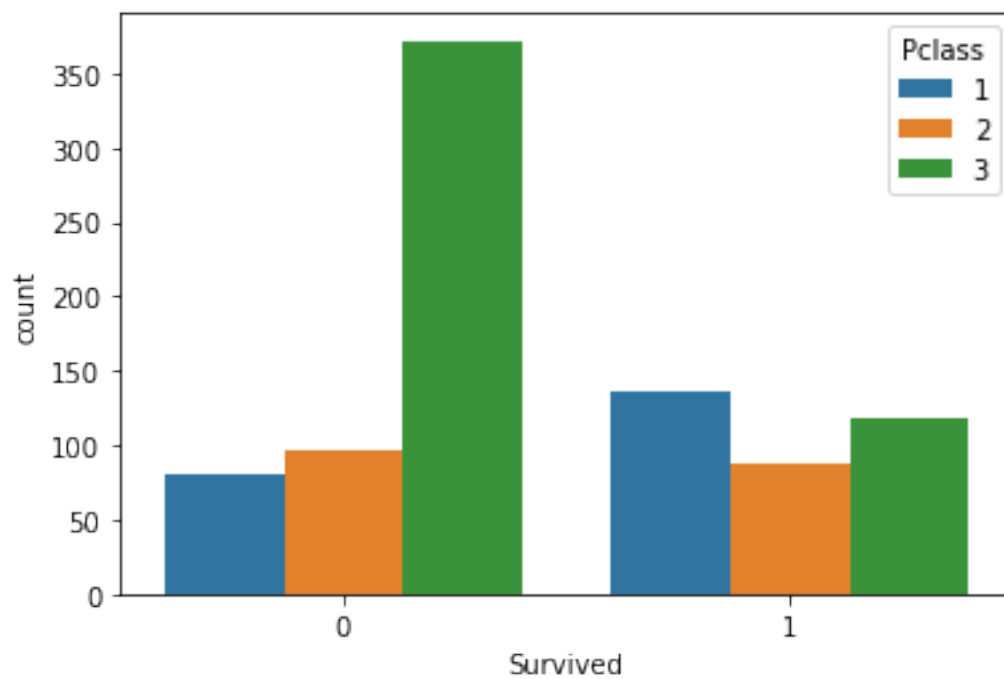
```
[14]: sns.countplot(x="Survived", data=df)
```

```
[14]: <AxesSubplot:xlabel='Survived', ylabel='count'>
```



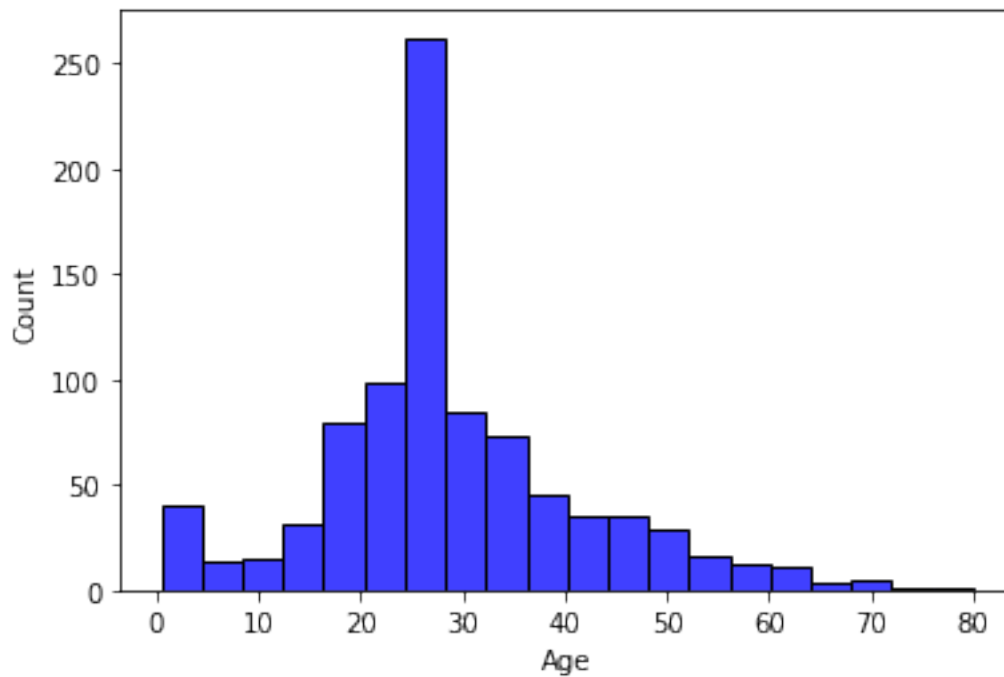
```
[15]: sns.countplot(x="Survived",data=df,hue= 'Pclass')
```

```
[15]: <AxesSubplot:xlabel='Survived', ylabel='count'>
```



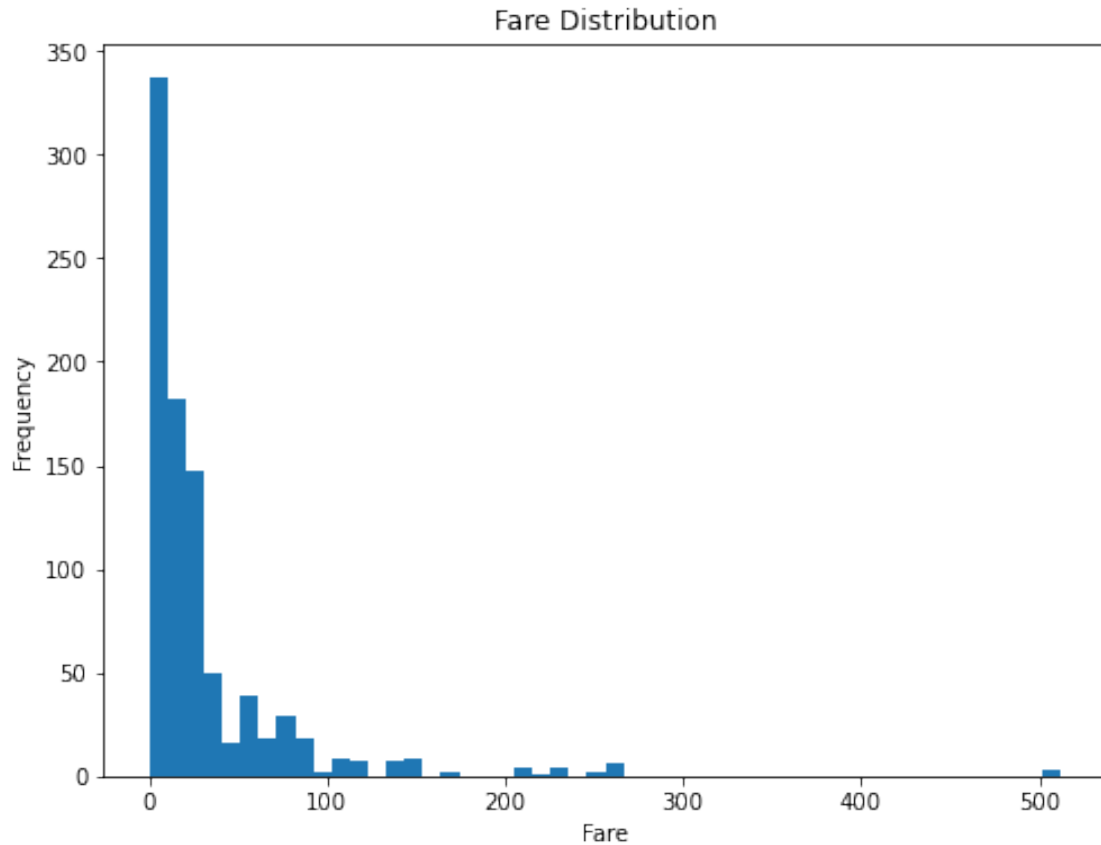
```
[16]: sns.histplot(x="Age",data=df,bins=20,color="blue")
```

```
[16]: <AxesSubplot:xlabel='Age', ylabel='Count'>
```



```
[17]: plt.figure(figsize=(8,6))  
plt.hist(df['Fare'], bins = 50)  
plt.title("Fare Distribution")  
plt.xlabel('Fare')  
plt.ylabel('Frequency')
```

```
[17]: Text(0, 0.5, 'Frequency')
```



```
[18]: # Define features (X) and target variable (y)
X = df.drop('Survived', axis=1)
y = df['Survived']

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)

# Initialize and train the model
model = RandomForestClassifier(random_state=42)
model.fit(X_train, y_train)

# Make predictions on the test set
predictions = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, predictions)
report = classification_report(y_test, predictions)

print(f"Accuracy: {accuracy:.2f}")
```

```
print("Classification Report:\n", report)
```

Accuracy: 0.82

Classification Report:

	precision	recall	f1-score	support
0	0.83	0.87	0.85	105
1	0.80	0.76	0.78	74
accuracy			0.82	179
macro avg	0.82	0.81	0.81	179
weighted avg	0.82	0.82	0.82	179

```
[ ]:
```