

```
[2]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.metrics import accuracy_score
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.svm import SVC
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.naive_bayes import GaussianNB
8 from sklearn.svm import LinearSVC
9 from sklearn.linear_model import SGDClassifier
10 from sklearn.tree import DecisionTreeClassifier
11 from sklearn.svm import SVC
12 from sklearn.ensemble import RandomForestClassifier
13 from matplotlib import pylab
14 import xgboost as Xgb
15 import seaborn as sns
16 import matplotlib.pyplot as plt
17 import numpy as np
18 import pandas as pd
19 import missingno as msno
20 import warnings
21 warnings.filterwarnings('ignore')
22 %matplotlib inline
23
24 import os
25 for dirname, _, filenames in os.walk('/kaggle/input'):
26     for filename in filenames:
27         print(os.path.join(dirname, filename))
28
```

/kaggle/input/titanic-dataset/Titanic-Dataset.csv

```
[3]: 1 df = pd.read_csv("/kaggle/input/titanic-dataset/Titanic-Dataset.csv")
2 df.head().style.set_properties(
3     **{
4         'background-color': 'LightBlue',
5         'color': 'Black',
6         'border-color': 'darkblack'
7     })
```

```
[3]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.000000	1	0	A/5 21171	7.250000	nan	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	38.000000	1	0	PC 17599	71.283300	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.000000	0	0	STON/O2. 3101282	7.925000	nan	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.000000	1	0	113803	53.100000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.000000	0	0	373450	8.050000	nan	S

Shape of the dataset

```
[4]: 1 print("Shape of the Dataset is:", df.shape)
```

Shape of the Dataset is: (891, 12)

Getting information about dataset

```
[5]: 1 df.describe()
```

[5]:

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	714.000000	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	14.526497	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	20.125000	0.000000	0.000000	7.910400
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	668.500000	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

[6]:

```
1 df.info()
```

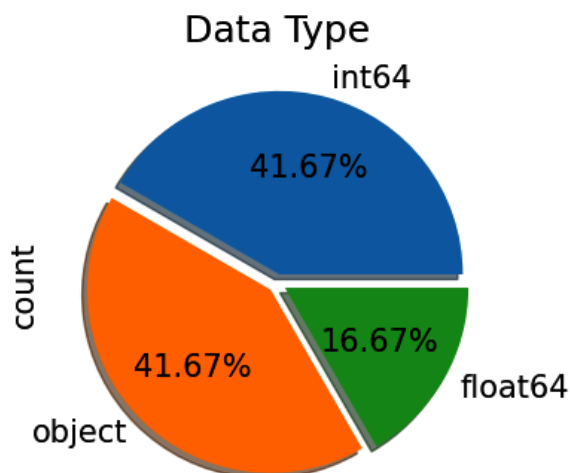
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   PassengerId  891 non-null    int64
1   Survived     891 non-null    int64
2   Pclass       891 non-null    int64
3   Name         891 non-null    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
```

Visualisation on Data Type

[14]:

```
1 plt.figure(figsize=(6,4))
2 plt.rcParams.update({'font.size':15})
3
4 df.dtypes.value_counts().plot.pie(explode=[0.05,0.05,0.05],
5                                   autopct='%1.2f%%',
6                                   shadow=True)
7
8 plt.title('Data Type',
9           color='Black',
10          loc='center')
```

[14... Text(0.5, 1.0, 'Data Type')

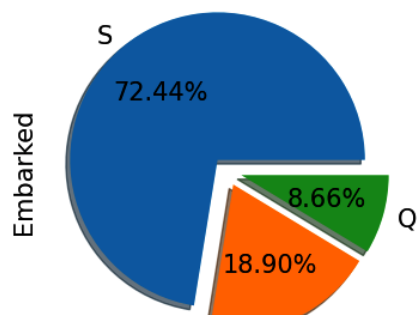


Visualisation on "Embarked" column

```
[13]: 1 plt.figure(figsize=(6,4))
2 plt.rcParams.update({'font.size': 15})
3 df['Embarked'].value_counts().plot.pie(explode=[0.1, 0.1, 0.1],
4                                         autopct='%1.2f%%',
5                                         shadow=True)
6 plt.title('Embarked Distribution',color='Black',loc='center')
7 plt.ylabel('Embarked')
```

```
[13... Text(0, 0.5, 'Embarked')
```

Embarked Distribution

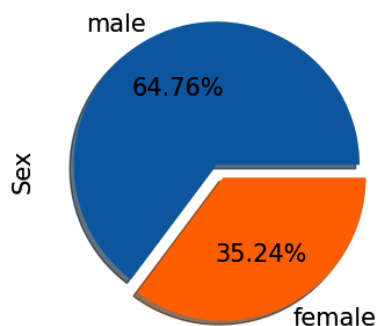


Visualisation on Gender Distribution

```
[15]: 1 plt.figure(figsize=(6,4))
2 plt.rcParams.update({'font.size': 15})
3 df['Sex'].value_counts().plot.pie(explode=[.05,.05],
4                                    autopct='%1.2f%%',
5                                    shadow=True)
6 plt.title('Gender Distribution')
7 plt.ylabel('Sex')
```

```
[15... Text(0, 0.5, 'Sex')
```

Gender Distribution

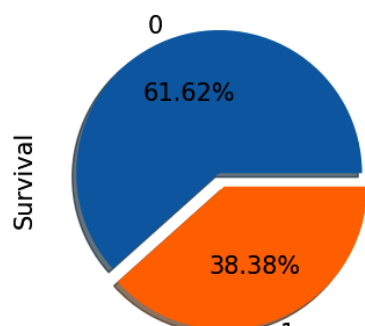


Visualisation on Survival Status

```
[16]: 1 plt.figure(figsize=(6,4))
2 plt.rcParams.update({'font.size': 15})
3 df['Survived'].value_counts().plot.pie(explode=[.05,.05],
4                                         autopct='%1.2f%%',
5                                         shadow=True)
6 plt.title('Survival Status', loc='center')
7 plt.ylabel('Survival')
8 print('0 = Died\n1 = Survived')
```

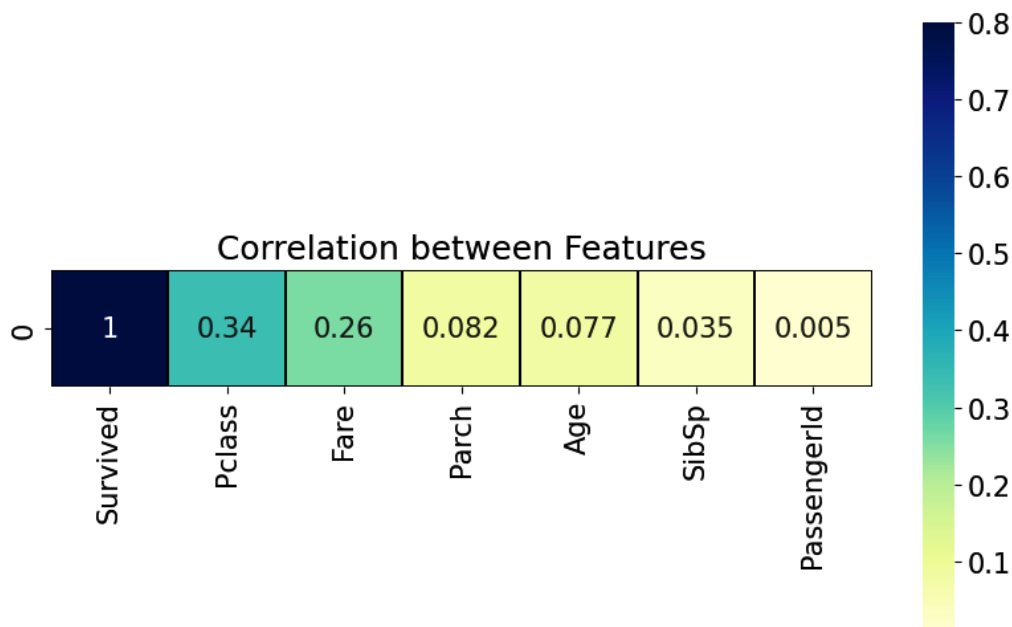
0 = Died
1 = Survived

Survival Status



Visualisation on Correlation

```
[17]: 1 numeric_df = df.select_dtypes(include=['number'])
2 corr = numeric_df.corrwith(numeric_df['Survived']).abs().sort_values(ascending=False)
3 plt.figure(figsize=(10,6))
4 sns.heatmap(
5     corr.values.reshape(1, -1),
6     vmax=0.8, linewidths=.01,
7     square=True, annot=True, cmap='YlGnBu',
8     linecolor="black", xticklabels=corr.index)
9 plt.title('Correlation between Features')
10 plt.show()
```



Seeing the missing values

```
[18]: 1 def missing_value (df):
2     missing_Number = df.isnull().sum().sort_values(ascending=False)[df.isnull().sum().sort_values(ascending=False) != 0]
3     missing_percent = round((df.isnull().sum()/df.isnull().count())*100,2)[round((df.isnull().sum()/df.isnull().count())*100,
4     missing = pd.concat([missing_Number,missing_percent],axis=1,keys=['Missing Number','Missing Percentage'])
5     return missing
```

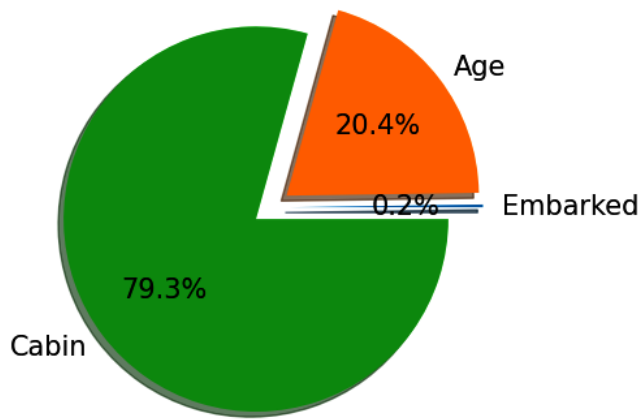
```
[19]: 1 missing_value(df).style.background_gradient(cmap='coolwarm').format(precision=2)
```

	Missing Number	Missing Percentage
Cabin	687	77.10
Age	177	19.87
Embarked	2	0.22

Pie Chart on missing values

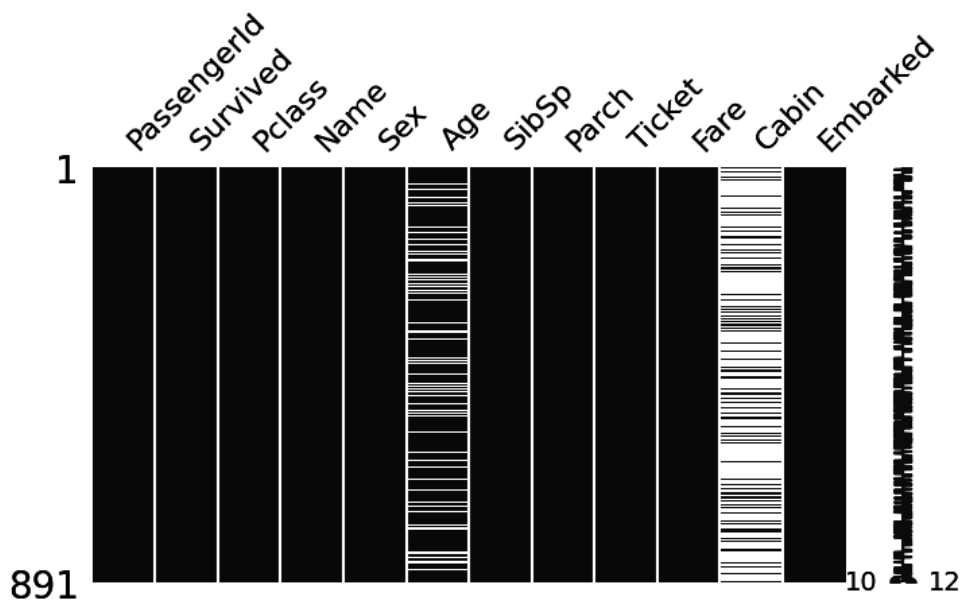
```
[20]: 1 missing_values = df.isnull().sum()
2 missing_values = missing_values[missing_values > 0]
3 missing_values.sort_values(inplace=True)
4 missing_values.plot.pie(explode=[.1,.1,.1],
5     autopct='%1.1f%%',
6     shadow=True)
7 plt.title('Missing Values',font='Lucida Calligraphy')
8
```

Missing Values



Matrix on missing values

```
[21]: 1 msno.matrix(df, figsize=(8,4))
      2 plt.show()
```



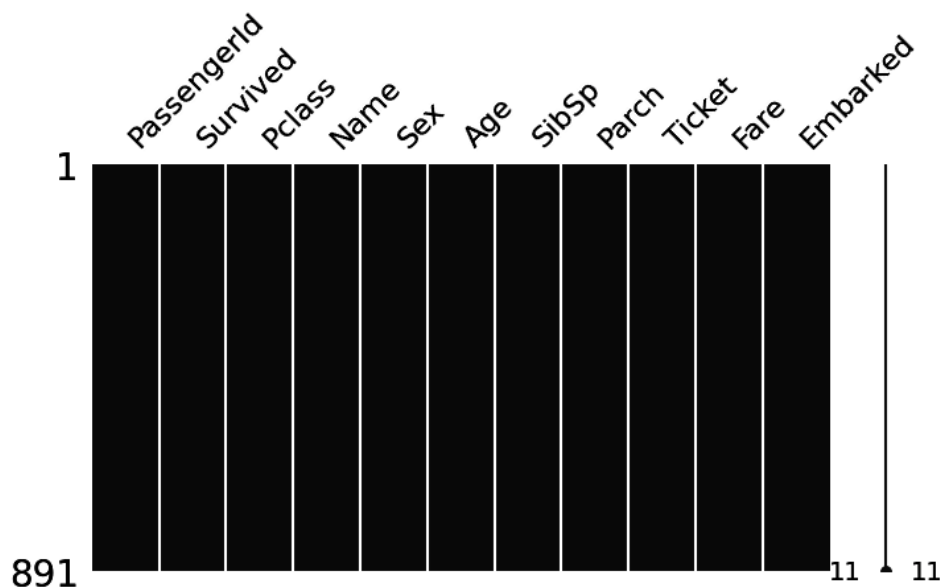
Working on missing values

```
[24]: 1 df['Age'] = df['Age'].fillna(df['Age'].mean())
```

```
[25]: 1 df['Embarked'] = df['Embarked'].fillna(method='bfill')
```

```
[26]: 1 df = df.drop(['Cabin'], axis=1)
```

```
[27]: 1 msno.matrix(df, figsize=(8,4))
      2 plt.show()
```



```
[28]: 1 df.isnull().sum().sum()
```

```
[28... 0
```

Dropping unnecessary columns

```
[29]: 1 df = df.drop(['Name', 'Ticket'], axis=1)
```

```
[30]: 1 df.head()
```

```
[30...
   PassengerId  Survived  Pclass    Sex  Age  SibSp  Parch    Fare  Embarked
0            1         0        3   male  22.0     1     0   7.2500         S
1            2         1        1  female  38.0     1     0  71.2833         C
2            3         1        3  female  26.0     0     0   7.9250         S
3            4         1        1  female  35.0     1     0  53.1000         S
4            5         0        3   male  35.0     0     0   8.0500         S
```

Transforming the categorical columns into numerical one

```
[31]: 1 df = pd.get_dummies(df, columns=['Sex', 'Embarked'], drop_first=True, dtype=int)
      2 df.head()
```

```
[31...
   PassengerId  Survived  Pclass  Age  SibSp  Parch    Fare  Sex_male  Embarked_Q  Embarked_S
0            1         0        3  22.0     1     0   7.2500         1         0         1
1            2         1        1  38.0     1     0  71.2833         0         0         0
2            3         1        3  26.0     0     0   7.9250         0         0         1
3            4         1        1  35.0     1     0  53.1000         0         0         1
4            5         0        3  35.0     0     0   8.0500         1         0         1
```

```
[34]: 1 X = df.drop('Survived',axis=1)
      2 y = df['Survived']
```

Splitting the data into Train and Test Set

```
[36]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.2, random_state=35)
```

Normalisation

```
[37]: 1 scaler = StandardScaler()
      2 X_train = scaler.fit_transform(X_train)
      3 X_test = scaler.transform(X_test)
      4
      5 X_train = pd.DataFrame(X_train, columns=X.columns)
      6 X_test = pd.DataFrame(X_test, columns=X.columns)
```

```
[38]: 1 display(X_train.head())
      2 display(X_test.head())
```

	PassengerId	Pclass	Age	SibSp	Parch	Fare	Sex_male	Embarked_Q	Embarked_S
0	0.068012	0.820268	0.655112	-0.473165	-0.470393	-0.507031	0.713074	-0.31696	0.601146
1	-0.953297	0.820268	0.340576	-0.473165	-0.470393	-0.518447	0.713074	-0.31696	0.601146
2	-0.506717	-0.386407	1.913254	-0.473165	-0.470393	-0.363866	0.713074	-0.31696	0.601146
3	0.588375	-1.593083	1.520085	0.377576	-0.470393	0.520444	0.713074	-0.31696	-1.663489
4	-0.036837	-1.593083	2.699594	0.377576	4.722159	4.765351	0.713074	-0.31696	0.601146

	PassengerId	Pclass	Age	SibSp	Parch	Fare	Sex_male	Embarked_Q	Embarked_S
0	0.203928	0.820268	-0.996201	-0.473165	-0.470393	-0.473814	0.713074	-0.31696	0.601146
1	-0.118386	-0.386407	-0.917567	-0.473165	2.125883	-0.384465	-1.402379	-0.31696	0.601146
2	-0.060137	-0.386407	0.969647	0.377576	-0.470393	-0.116675	-1.402379	-0.31696	0.601146
3	1.683468	0.820268	0.261942	-0.473165	-0.470393	-0.489608	0.713074	-0.31696	0.601146
4	-1.547443	0.820268	-0.917567	0.377576	-0.470393	-0.285589	-1.402379	-0.31696	0.601146

ML Process with Logistic Regression

```
[39]: 1 logreg = LogisticRegression()
      2 logreg.fit(X_train, y_train)
      3 y_pred = logreg.predict(X_test)
      4
      5 log_training = round(logreg.score(X_train, y_train)*100,2)
      6 log_accuracy = round(accuracy_score(y_pred,y_test)*100,2)
      7
      8 print('Training Accuracy: ',log_training)
      9 print('Model Accuracy Score: ',log_accuracy)
```

Training Accuracy: 80.06
Model Accuracy Score: 78.77

ML Process with Support Vector Machine

[40]:

```
1 svc=SVC()
2 svc.fit(X_train,y_train)
3 y_pred = svc.predict(X_test)
4
5 svc_training = round(svc.score(X_train,y_train)*100,2)
6 svc_accuracy = round(accuracy_score(y_pred, y_test)*100,2)
7
8 print('Training Accuracy: ', svc_training)
9 print('Model Accuracy Score: ', svc_accuracy)
```

Training Accuracy: 83.99
Model Accuracy Score: 83.8

ML Process with KNeighbors Classifier

[41]:

```
1 knn = KNeighborsClassifier()
2 knn.fit(X_train,y_train)
3 y_pred = knn.predict(X_test)
4
5 knn_training = round(knn.score(X_train,y_train)*100,2)
6 knn_accuracy = round(accuracy_score(y_pred, y_test)*100,2)
7
8 print('Training Accuracy: ',knn_training)
9 print('Model Accuracy Score: ', knn_accuracy)
```

Training Accuracy: 85.96
Model Accuracy Score: 81.01

ML Process with Naive Bayes

[42]:

```
1 gaussian = GaussianNB()
2 gaussian.fit(X_train, y_train)
3 y_pred = gaussian.predict(X_test)
4
5 gaussian_training = round(gaussian.score(X_train, y_train)*100,2)
6 gaussian_accuracy = round(accuracy_score(y_pred,y_test)*100,2)
7
8 print('Training Accuracy: ', gaussian_training)
9 print('Model Accuracy Score: ',gaussian_accuracy)
```

Training Accuracy: 79.78
Model Accuracy Score: 74.86

ML Process with Linear SVM

[43]:

```
1 linear_svc = LinearSVC()
2 linear_svc.fit(X_train, y_train)
3 y_pred = linear_svc.predict(X_test)
4
5 linear_training = round(linear_svc.score(X_train,y_train)*100,2)
6 linear_accuracy = round(accuracy_score(y_pred,y_test)*100,2)
7
8 print('Training Accuracy: ', linear_training)
9 print('Model Accuracy Score: ',linear_accuracy)
```

Training Accuracy: 80.2
Model Accuracy Score: 79.33

ML Process with Stochastic Gradient Descent

```
[44]: 1 sgd = SGDClassifier()
2 sgd.fit(X_train,y_train)
3 y_pred = sgd.predict(X_test)
4
5 sgd_training = round(sgd.score(X_train,y_train)*100,2)
6 sgd_accuracy = round(accuracy_score(y_pred,y_test)*100,2)
7
8 print('Training Score: ',sgd_training)
9 print('Model Accuracy Score: ',sgd_accuracy)
```

Training Score: 75.7
Model Accuracy Score: 75.98

ML Process with Decision Tree Classifier

```
[46]: 1 decision = DecisionTreeClassifier()
2 decision.fit(X_train,y_train)
3 y_pred=decision.predict(X_test)
4
5 decision_training = round(decision.score(X_train,y_train)*100,2)
6 decision_accuracy = round(accuracy_score(y_pred,y_test)*100,2)
7
8 print('Training Accuracy: ', decision_training)
9 print('Model Accuracy Score: ', decision_accuracy)
```

Training Accuracy: 100.0
Model Accuracy Score: 73.74

ML Process with Random Forest Classifier

```
[47]: 1 random_forest = RandomForestClassifier(n_estimators=100)
2 random_forest.fit(X_train,y_train)
3 y_pred = random_forest.predict(X_test)
4 random_forest.score(X_train,y_train)
5
6 random_forest_training = round(random_forest.score(X_train, y_train)*100,2)
7 random_forest_accuracy = round(accuracy_score(y_pred,y_test)*100,2)
8
9 print('Training Accuracy: ',random_forest_training)
10 print('Model Accuracy Score: ',random_forest_accuracy)
```

Training Accuracy: 100.0
Model Accuracy Score: 82.12

ML Process with XGBOOST

```
[48]: 1 xgb = Xgb.XGBClassifier()
2 xgb.fit(X_train,y_train)
3 y_pred = xgb.predict(X_test)
4 xgb.score(X_train,y_train)
5
6 xgb_training = round(xgb.score(X_train,y_train)*100,2)
7 xgb_accuracy = round(accuracy_score(y_pred,y_test)*100,2)
8
9 print('Training Accuracy: ',xgb_training)
10 print('Model Accuracy Score: ',xgb_accuracy)
```

Training Accuracy: 100.0
Model Accuracy Score: 78.21

```
[49]: 1 models = pd.DataFrame({
2     'Model': [
3         'Logistic Regression', 'Support Vector Machine', 'KNN', 'GaussianNB', 'Linear SVC', 'Stochastic Gradient Decent', 'Decision Tree',
4     ],
5     'Training Accuracy': [
6         log_training, svc_training, knn_training, gaussian_training, linear_training, sgd_training, decision_training, random_forest_training,
7     ],
8     'Model Accuracy Score': [
9         log_accuracy, svc_accuracy, knn_accuracy, gaussian_accuracy, linear_accuracy, sgd_accuracy, decision_accuracy, random_forest_accuracy,
10        xgb_accuracy
11    ]
12 })
```

Sorting the models by the accuracy

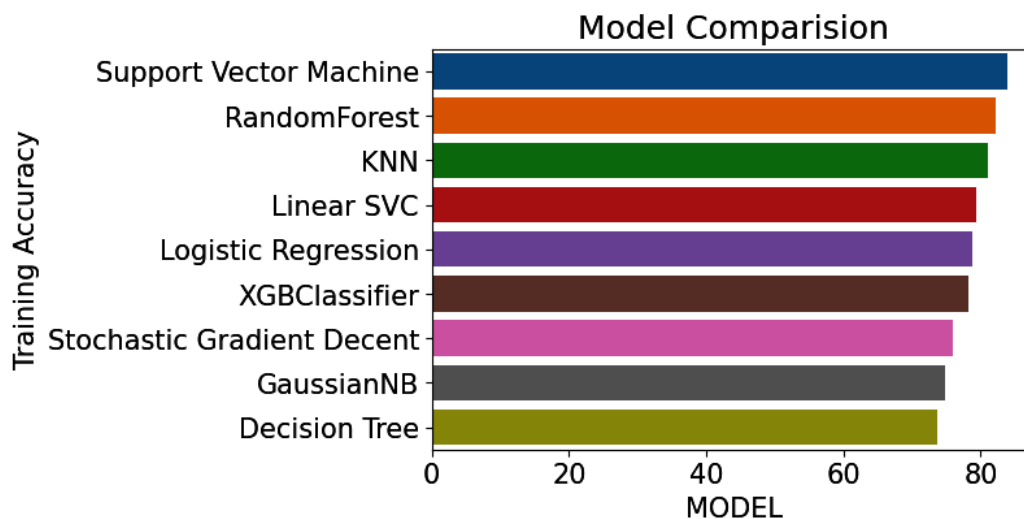
```
[50]: 1 models.sort_values(by='Training Accuracy', ascending=False)
```

```
[50...
      Model  Training Accuracy  Model Accuracy Score
6      Decision Tree          100.00             73.74
7      RandomForest          100.00             82.12
8      XGBClassifier          100.00             78.21
2              KNN           85.96             81.01
1  Support Vector Machine      83.99             83.80
4          Linear SVC          80.20             79.33
0      Logistic Regression      80.06             78.77
3          GaussianNB          79.78             74.86
5  Stochastic Gradient Decent    75.70             75.98
```

Visualisation on Model Comparison

```
[51]: 1 models = models.sort_values(by='Model Accuracy Score', ascending=False)[:20]
2 plt.figure(figsize=(6,4))
3 sns.barplot(y='Model', x='Model Accuracy Score', data=models)
4 plt.title('Model Comparision')
5 plt.xlabel('MODEL')
6 plt.ylabel('Training Accuracy')
```

```
[51... Text(0, 0.5, 'Training Accuracy')
```



```
[208]: 1 y = df['Survived']
2
3 features = ['Pclass', 'SibSp', 'Parch']
4 X = pd.get_dummies(df[features])
5 X_test = pd.get_dummies(df[features])
6
7 model = SVC()
8 model.fit(X,y)
9 predictions = model.predict(X_test)
10
11 output = pd.DataFrame({'PassangerID': df.PassengerId, 'Survived': predictions})
12 output.to_csv('submission.csv', index=False)
13 print('Your submission was successfully saved !')
```

Your submission was successfully saved !

We successfully saved our ML Module.

