

Summary

Needs to create a model that predicts which passengers survived the Titanic shipwreck.

On April 15, 1912, during her maiden voyage, the widely considered “unsinkable” RMS Titanic sank after colliding with an iceberg. Unfortunately, there weren’t enough lifeboats for everyone onboard, resulting in the death of 1502 out of 2224 passengers and crew.

In [1]:

```
# This Python 3 environment comes with many helpful analytics libraries installed
# It is defined by the kaggle/python Docker image: https://github.com/kaggle/docker-python
# For example, here's several helpful packages to load

import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib as mpl

from sklearn.ensemble import RandomForestClassifier

# Input data files are available in the read-only "../input/" directory
# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory

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

# You can write up to 20GB to the current directory (/kaggle/working/) that gets preserved as output when you create a version using "Save & Run All"
# You can also write temporary files to /kaggle/temp/, but they won't be saved outside of the current session
```

```

/kaggle/input/titanic/train.csv
/kaggle/input/titanic/test.csv
/kaggle/input/titanic/gender_submission.csv

```

In [2]:

```

train_data = pd.read_csv("/kaggle/input/titanic/train.csv")
train_data

```

Out[2]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows × 12 columns

In [3]:

```
print (train_data.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
None
```

In [4]:

```
test_data = pd.read_csv("/kaggle/input/titanic/test.csv")
test_data
```

Out[4]:

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	892	3	Kelly, Mr. James	male	34.5	0	0	330911	7.8292	NaN	Q
1	893	3	Wilkes, Mrs. James (Ellen Needs)	female	47.0	1	0	363272	7.0000	NaN	S
2	894	2	Myles, Mr. Thomas Francis	male	62.0	0	0	240276	9.6875	NaN	Q
3	895	3	Wirz, Mr. Albert	male	27.0	0	0	315154	8.6625	NaN	S
4	896	3	Hirvonen, Mrs. Alexander (Helga E Lindqvist)	female	22.0	1	1	3101298	12.2875	NaN	S
...
413	1305	3	Spector, Mr. Woolf	male	NaN	0	0	A.5. 3236	8.0500	NaN	S
414	1306	1	Oliva y Ocana, Dona. Fermina	female	39.0	0	0	PC 17758	108.9000	C105	C
415	1307	3	Saether, Mr. Simon Sivertsen	male	38.5	0	0	SOTON/O.Q. 3101262	7.2500	NaN	S
416	1308	3	Ware, Mr. Frederick	male	NaN	0	0	359309	8.0500	NaN	S
417	1309	3	Peter, Master. Michael J	male	NaN	1	1	2668	22.3583	NaN	C

418 rows × 11 columns

In [5]:

```
train_data.isnull().sum()
```

Out[5]:

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

In [6]:

```
train_data.describe()
```

Out[6]:

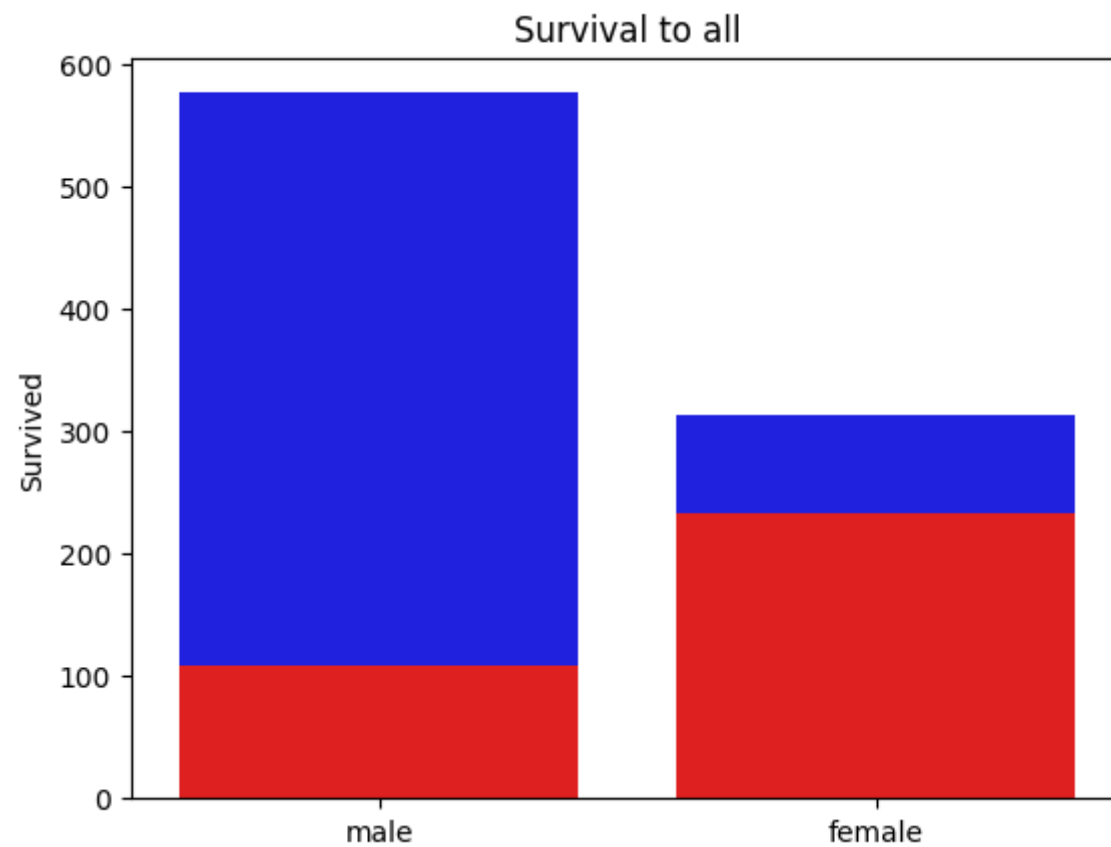
	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
std	257.353842	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
min	1.000000	0.000000	1.000000	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

Sex

In [7]:

```
mpl.style.use('default')
sns.barplot(x=train_data['Sex'].unique(), y=train_data['Sex'].value_counts(), color='blue')
sns.barplot(x=train_data['Sex'].unique(), y=train_data.groupby(['Sex'])['Survived'].sum().sort_index(ascending=False), color='red')
plt.title('Survival to all')
plt.show()

print('Survived man %s'%(train_data.groupby(['Sex'])['Survived'].sum()['male']/train_data['Sex'].value_counts()['male']))
print('Survived woman %s'%(train_data.groupby(['Sex'])['Survived'].sum()['female']/train_data['Sex'].value_counts()['female']))
```

Survived man 0.18890814558058924
Survived woman 0.7420382165605095

No comments, just left it in the dataset. May be useful in PCA with SibS or Parch.

Pclass

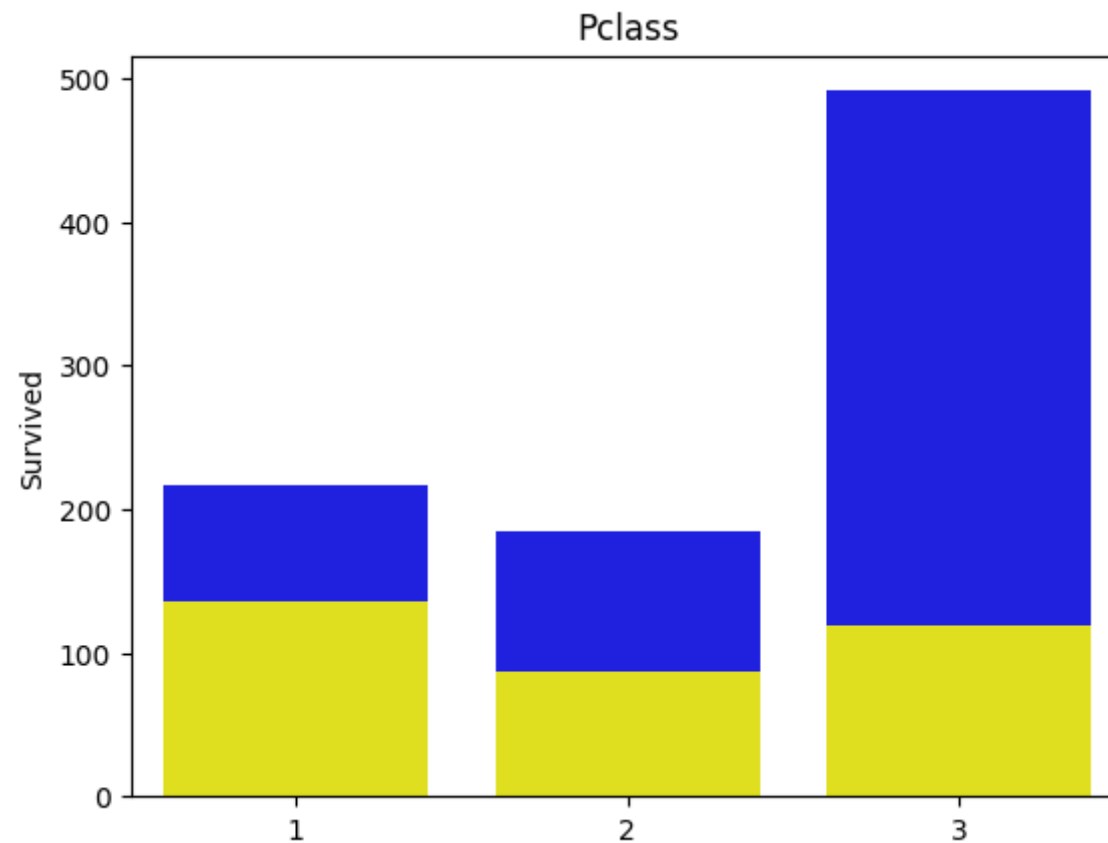
looks really useful. We need to decide which way to encode.

In [8]:

```
mpl.style.use('default')
sns.barplot(x=[1, 2, 3], y=train_data['Pclass'].value_counts().sort_index(), color='blue')
sns.barplot(x=[1, 2, 3], y=train_data.groupby(['Pclass'])['Survived'].sum().sort_index(), color='yellow')
plt.title('Pclass')
```

Out[8]:

```
Text(0.5, 1.0, 'Pclass')
```

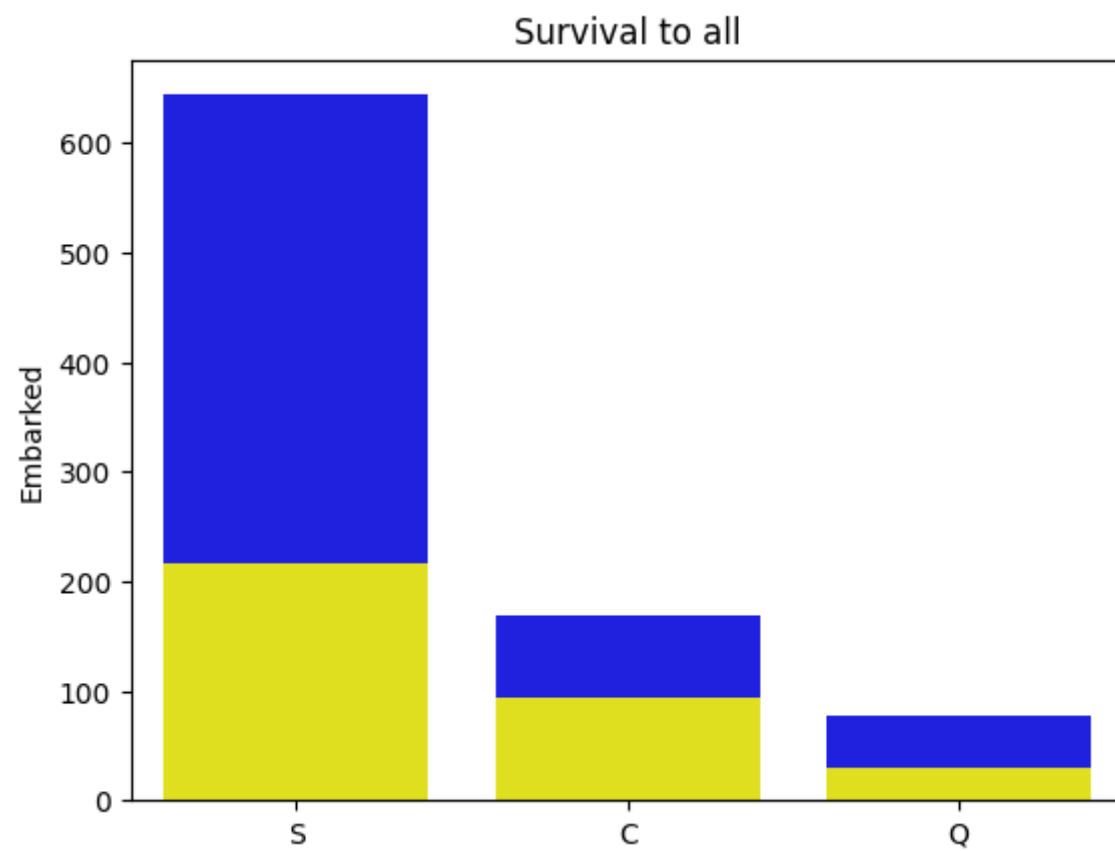


Embarked

In [9]:

```
survived_embarked = train_data.loc[train_data.Survived==1, "Embarked"].value_counts() / train_data.Embarked.value_counts()
mpl.style.use('default')
sns.barplot(x=test_data.Embarked.value_counts().index, y=train_data['Embarked'].value_counts(), color='blue')
sns.barplot(x=test_data.Embarked.value_counts().index, y=train_data.loc[train_data.Survived==1, "Embarked"].value_counts(), color='yellow')
plt.title('Survival to all')
print(survived_embarked)
```

```
S    0.336957  
C    0.553571  
Q    0.389610  
Name: Embarked, dtype: float64
```



Age

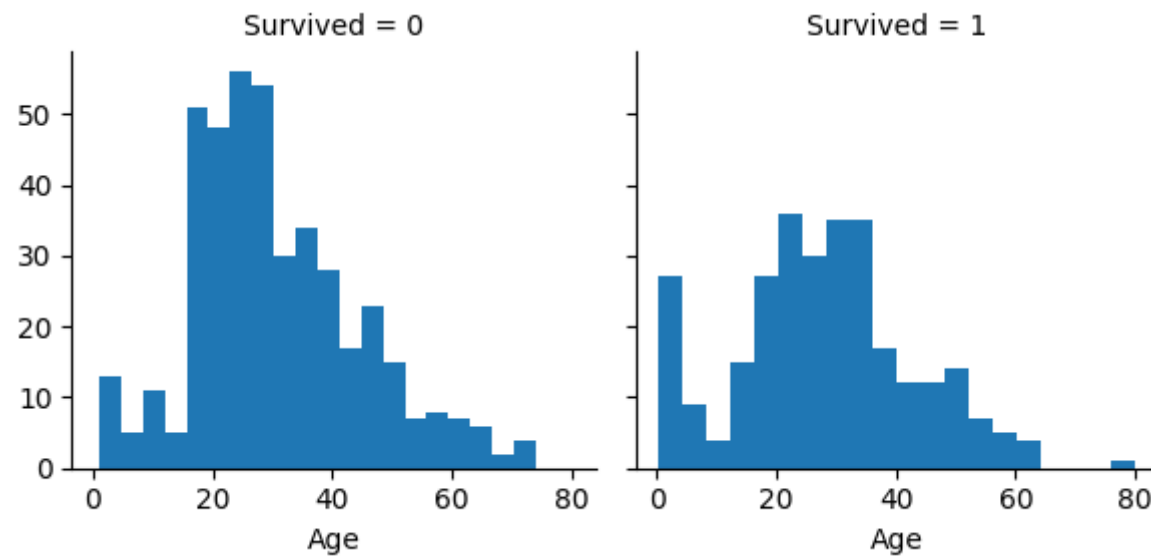
In first versions not necessary to include this to futures.

In [10]:

```
g = sns.FacetGrid(train_data, col='Survived')
g.map(plt.hist, 'Age', bins=20)
```

Out[10]:

```
<seaborn.axisgrid.FacetGrid at 0x7f9c968c4850>
```



In [11]:

#test code cell`print(train_data.columns)``print('-'*10)``print(test_data.columns)`

```
Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
      'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],
      dtype='object')
```

```
Index(['PassengerId', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp', 'Parch',
      'Ticket', 'Fare', 'Cabin', 'Embarked'],
      dtype='object')
```


In [12]:

```

for dataset in [train_data, test_data]:
    # Mapping Sex
    dataset['Sex'] = dataset['Sex'].map( {'female': 0, 'male': 1} ).astype(int)

    # Mapping Embarked
    dataset['Embarked'] = dataset['Embarked'].fillna("S")
    dataset['Embarked'] = dataset['Embarked'].map( {'C': 0, 'Q': 1, 'S': 2} ).astype(int)

    # Mapping Fare
    dataset['Fare'] = dataset['Fare'].fillna(train_data['Fare'].median())
    dataset.loc[ dataset['Fare'] <= 7.91, 'Fare'] = 0
    dataset.loc[(dataset['Fare'] > 7.91) & (dataset['Fare'] <= 14.454), 'Fare'] = 1
    dataset.loc[(dataset['Fare'] > 14.454) & (dataset['Fare'] <= 31), 'Fare'] = 2
    dataset.loc[ dataset['Fare'] > 31, 'Fare'] = 3
    dataset['Fare'] = dataset['Fare'].astype(int)

    # Mapping Age
    age_avg = dataset['Age'].mean()
    age_std = dataset['Age'].std()
    age_null_count = dataset['Age'].isnull().sum()

    age_null_random_list = np.random.randint(age_avg - age_std, age_avg + age_std, size=age_null_count)
    dataset['Age'][np.isnan(dataset['Age'])] = age_null_random_list
    dataset['Age'] = dataset['Age'].astype(int)

    dataset.loc[ dataset['Age'] <= 16, 'Age'] = 0
    dataset.loc[(dataset['Age'] > 16) & (dataset['Age'] <= 32), 'Age'] = 1
    dataset.loc[(dataset['Age'] > 32) & (dataset['Age'] <= 48), 'Age'] = 2
    dataset.loc[(dataset['Age'] > 48) & (dataset['Age'] <= 64), 'Age'] = 3
    dataset.loc[ dataset['Age'] > 64, 'Age'] = 4

```

```
# Feature Selection
```

```
drop_elements = ['PassengerId', 'Name', 'Ticket', 'Cabin', 'SibSp', \
                 'Parch']
```

```
train_data = train_data.drop(drop_elements, axis = 1)
```

```
test_data  = test_data.drop(drop_elements, axis = 1)
```

```
print (train_data.head(10))
```

```
train_data = train_data.values
```

```
test_data  = test_data.values
```

	Survived	Pclass	Sex	Age	Fare	Embarked
0	0	3	1	1	0	2
1	1	1	0	2	3	0
2	1	3	0	1	1	2
3	1	1	0	2	3	2
4	0	3	1	2	1	2
5	0	3	1	2	1	1
6	0	1	1	3	3	2
7	0	3	1	0	2	2
8	1	3	0	1	1	2
9	1	2	0	0	2	0

```
/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:23: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:23: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

TRAIN MODELS

In [13]:

```
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import StratifiedShuffleSplit
from sklearn.metrics import accuracy_score, log_loss
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis, QuadraticDiscriminantAnalysis
from sklearn.linear_model import LogisticRegression

classifiers = [
    KNeighborsClassifier(3),
    SVC(probability=True),
    DecisionTreeClassifier(),
    RandomForestClassifier(),
    AdaBoostClassifier(),
    GradientBoostingClassifier(),
    GaussianNB(),
    LinearDiscriminantAnalysis(),
    QuadraticDiscriminantAnalysis(),
    LogisticRegression()]

log_cols = ["Classifier", "Accuracy"]
log      = pd.DataFrame(columns=log_cols)

sss = StratifiedShuffleSplit(n_splits=10, test_size=0.1, random_state=0)
```

```
X = train_data[0::, 1::]
y = train_data[0::, 0]

acc_dict = {}

# saving the best classifier
best_classifier = LogisticRegression()
best_measure = 0.0

for train_index, test_index in sss.split(X, y):
    X_train, X_test = X[train_index], X[test_index]
    y_train, y_test = y[train_index], y[test_index]

    for clf in classifiers:
        name = clf.__class__.__name__
        clf.fit(X_train, y_train)
        train_predictions = clf.predict(X_test)
        acc = accuracy_score(y_test, train_predictions)
        if acc > best_measure:
            print(acc, clf.__class__.__name__)
            best_measure = acc
            best_classifier = clf
        if name in acc_dict:
            acc_dict[name] += acc
        else:
            acc_dict[name] = acc

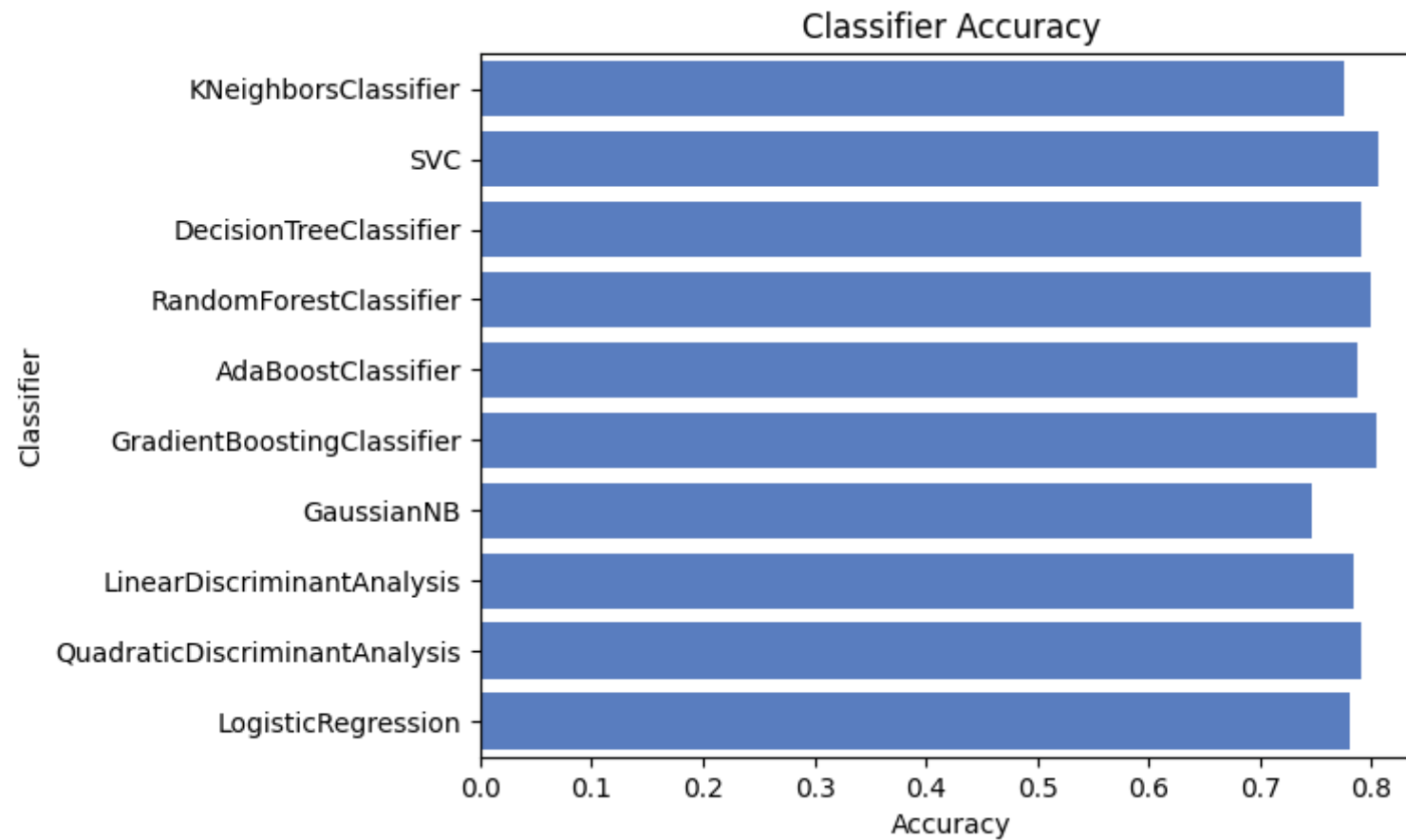
for clf in acc_dict:
    acc_dict[clf] = acc_dict[clf] / 10.0
    log_entry = pd.DataFrame([[clf, acc_dict[clf]]], columns=log_cols)
    log = log.append(log_entry)
```

```
plt.xlabel('Accuracy')  
plt.title('Classifier Accuracy')  
  
sns.set_color_codes("muted")  
sns.barplot(x='Accuracy', y='Classifier', data=log, color="b")
```

```
0.7555555555555555 KNeighborsClassifier  
0.8 SVC  
0.8555555555555555 SVC  
0.8666666666666667 AdaBoostClassifier  
0.8777777777777778 QuadraticDiscriminantAnalysis
```

Out[13]:

```
<AxesSubplot:title={'center':'Classifier Accuracy'}, xlabel='Accuracy', ylabel='Classifier'>
```



In [14]:

```
print(test_data)
```

```
[[3 1 2 0 1]
 [3 0 2 0 2]
 [2 1 3 1 1]
 ...
 [3 1 2 0 2]
 [3 1 1 1 2]
 [3 1 1 2 0]]
```

In [15]:

```
predictions = best_classifier.predict(test_data)

test_data = pd.read_csv("/kaggle/input/titanic/test.csv")
output = pd.DataFrame({'PassengerId': test_data.PassengerId, 'Survived': predictions})
output.to_csv('my_submission.csv', index=False)
print("Your submission was successfully saved!")
```

Your submission was successfully saved!

In []: