

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

```
train=pd.read_csv('/content/drive/MyDrive/InternShip/Titanic-Dataset.csv')
test=pd.read_csv('/content/drive/MyDrive/InternShip/test.csv')
```

```
print(train.shape)
print(test.shape)
```

```
(891, 12)
(418, 11)
```

```
test.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 11 columns):
 #   Column        Non-Null Count  Dtype
---  -
 0   PassengerId   418 non-null    int64
 1   Pclass        418 non-null    int64
 2   Name          418 non-null    object
 3   Sex           418 non-null    object
 4   Age           332 non-null    float64
 5   SibSp         418 non-null    int64
 6   Parch         418 non-null    int64
 7   Ticket        418 non-null    object
 8   Fare          417 non-null    float64
 9   Cabin         91 non-null     object
 10  Embarked      418 non-null    object
dtypes: float64(2), int64(4), object(5)
memory usage: 36.0+ KB
```

```
train.drop(columns=['Cabin'],inplace=True)
test.drop(columns=['Cabin'],inplace=True)
```

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

```
<ipython-input-7-6702a28d6ea0>:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assi
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col]

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

```
test['Fare'].fillna(test['Fare'].mean(), inplace=True)
```

```
<ipython-input-8-f2bd6e3cf28a>:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assi
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col]

test['Fare'].fillna(test['Fare'].mean(), inplace=True)
```

```
train.isnull().sum()
```




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

dtype: int64

```
gen_age=np.random.randint(train['Age'].mean()-train['Age'].std(),train['Age'].mean()+train['Age'].std(), size=177)
```

```
train['Age'][np.isnan(train['Age'])]=gen_age
```



<ipython-input-11-2ce148eb0d56>:1: FutureWarning: ChainedAssignmentError: behaviour will change in pandas 3.0!
You are setting values through chained assignment. Currently this works in certain cases, but when using Copy-on-Write (which will be the default behaviour in pandas 3.0), this can become an error. A typical example is when you are setting values in a column of a DataFrame, like:

```
df["col"][row_indexer] = value
```

Use `df.loc[row_indexer, "col"] = values` instead, to perform the assignment in a single step and ensure this keeps updating the original 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['Age'][np.isnan(train['Age'])]=gen_age
```


<ipython-input-11-2ce148eb0d56>:1: 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['Age'][np.isnan(train['Age'])]=gen_age
```

```
train.isnull().sum()
```




	0
PassengerId	0
Survived	0
Pclass	0
Name	0
Sex	0
Age	0
SibSp	0
Parch	0
Ticket	0
Fare	0
Embarked	0

dtype: int64

```
gen_age1=np.random.randint(test['Age'].mean()-test['Age'].std(),test['Age'].mean()+test['Age'].std(), size=86)
```

```
test['Age'][np.isnan(test['Age'])]=gen_age1
```



<ipython-input-14-0b74930e6951>:1: FutureWarning: ChainedAssignmentError: behaviour will change in pandas 3.0!
You are setting values through chained assignment. Currently this works in certain cases, but when using Copy-on-Write (which will be the default behaviour in pandas 3.0), this can become an error. A typical example is when you are setting values in a column of a DataFrame, like:

```
df["col"][row_indexer] = value
```

Use `df.loc[row_indexer, "col"] = values` instead, to perform the assignment in a single step and ensure this keeps updating the original DataFrame.


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

```
test['Age'][np.isnan(test['Age'])]=gen_age1
<ipython-input-14-0b74930e6951>:1: 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

```
test['Age'][np.isnan(test['Age'])]=gen_age1
```

```
test.isnull().sum()
```



	0
PassengerId	0
Pclass	0
Name	0
Sex	0
Age	0
SibSp	0
Parch	0
Ticket	0
Fare	0
Embarked	0

dtype: int64


```
train.isnull().sum()
```



	0
PassengerId	0
Survived	0
Pclass	0
Name	0
Sex	0
Age	0
SibSp	0
Parch	0
Ticket	0
Fare	0
Embarked	0

dtype: int64

```
train[['Pclass', 'Survived']].groupby('Pclass').mean()
```



	Survived
Pclass	
1	0.629630
2	0.472826
3	0.242363

```
train[['Sex', 'Survived']].groupby('Sex').mean()
```



	Survived
Sex	
female	0.742038
male	0.188908

```
train[['Embarked', 'Survived']].groupby('Embarked').mean()
```

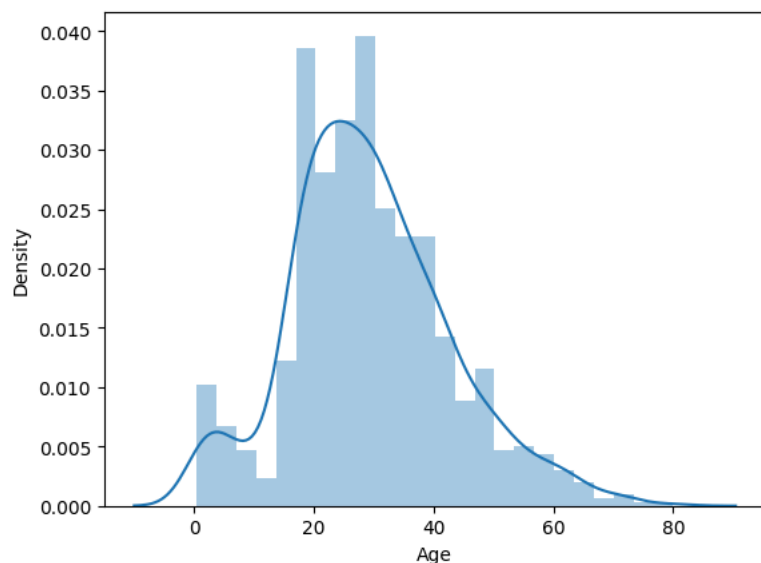
	Survived
Embarked	
C	0.553571
Q	0.389610
S	0.339009

```
sns.distplot(train['Age'])
```

<ipython-input-20-22ed6932e5e8>:1: UserWarning:
 `distplot` is a deprecated function and will be removed in seaborn v0.14.0.
 Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

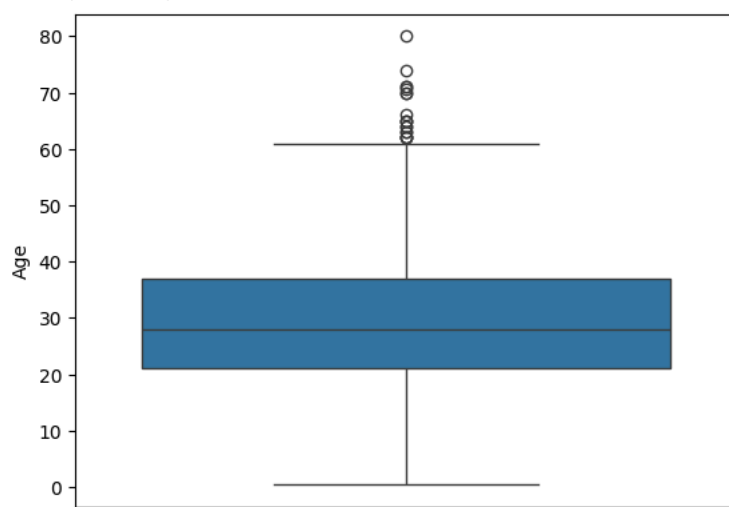
For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(train['Age'])
<Axes: xlabel='Age', ylabel='Density'>
```



```
sns.boxplot(train['Age'])
```

<Axes: ylabel='Age'>



```
train[train['Age']>75]['Survived'].value_counts()
```

```
count
Survived
1      1
```

```
dtype: int64
```

```
plt.subplots(figsize=(15,4))
sns.distplot(train[train['Survived']==0]['Age'])
sns.distplot(train[train['Survived']==1]['Age'])
```

```
<ipython-input-23-8cacb70f70fe>:2: UserWarning:
```

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(train[train['Survived']==0]['Age'])
```

```
<ipython-input-23-8cacb70f70fe>:3: UserWarning:
```

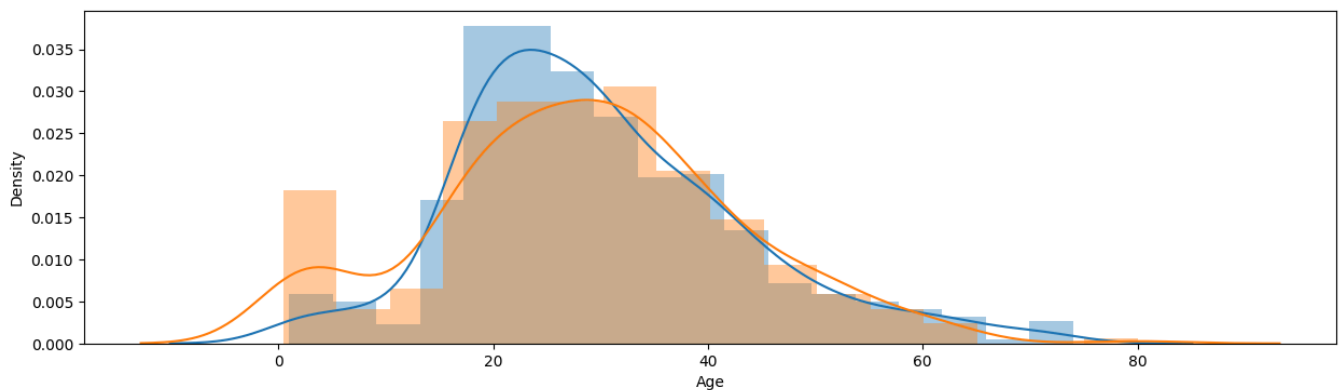
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(train[train['Survived']==1]['Age'])
```

```
<Axes: xlabel='Age', ylabel='Density'>
```



```
passengerId=test['PassengerId'].values
```

```
train.drop(columns=['PassengerId','Ticket'],inplace=True)
```

```
test.drop(columns=['PassengerId','Ticket'],inplace=True)
```

```
train.isnull().sum()
```

```
0
Survived 0
Pclass 0
Name 0
Sex 0
Age 0
SibSp 0
Parch 0
Fare 0
Embarked 0
```

```
dtype: int64
```

```
sns.distplot(train['Fare'])
```



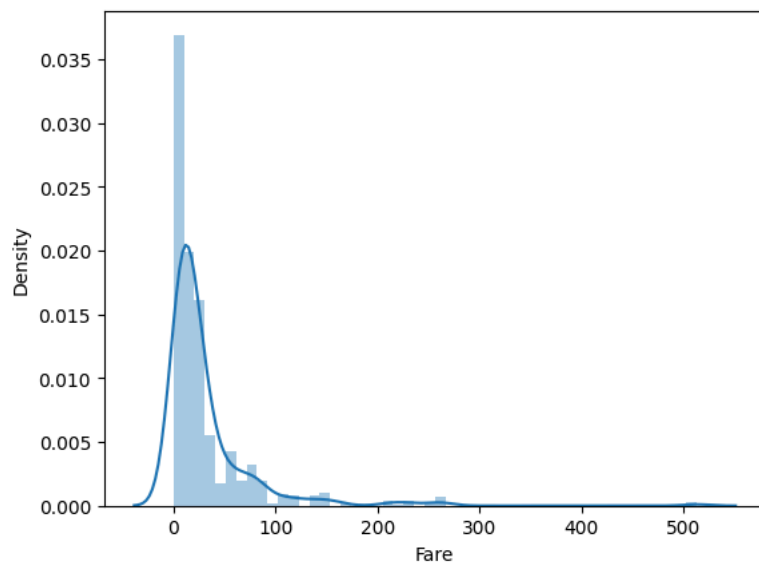
```
<ipython-input-27-98b78f01c2eb>:1: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

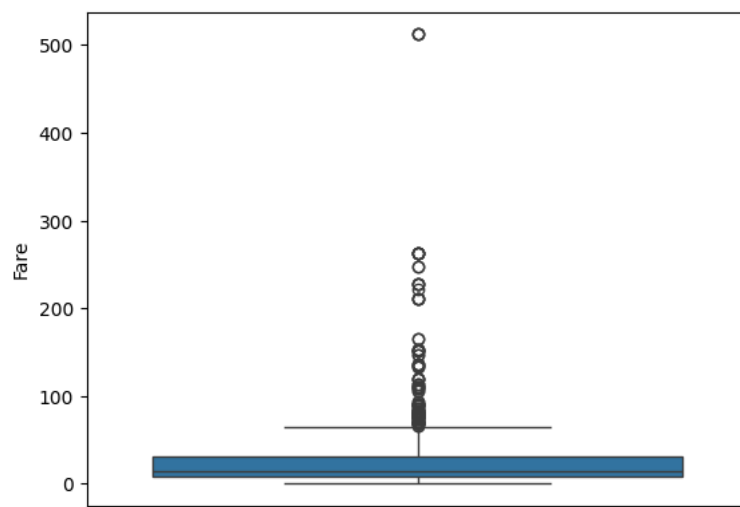
```
sns.distplot(train['Fare'])
<Axes: xlabel='Fare', ylabel='Density'>
```



```
sns.boxplot(train['Fare'])
```



```
<Axes: ylabel='Fare'>
```



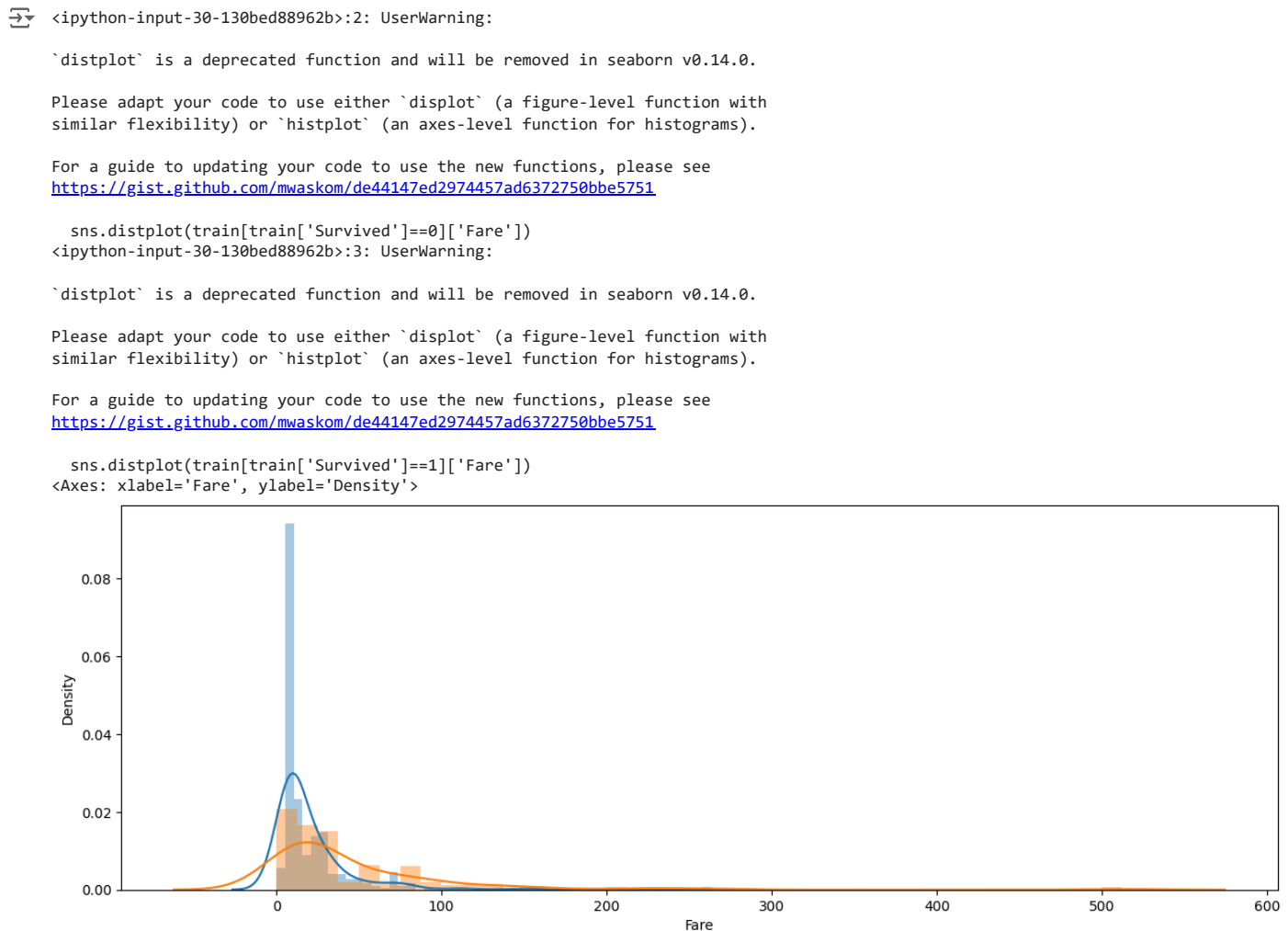
```
train[train['Fare']>400]['Survived'].value_counts()
```



```
count
Survived
1      3
```

```
dtype: int64
```

```
plt.subplots(figsize=(15,5))
sns.distplot(train[train['Survived']==0]['Fare'])
sns.distplot(train[train['Survived']==1]['Fare'])
```



```
# Don't delete this unless its 1st Jan
train['Name']
```

```
<ipython-input-30-130bed88962b>:4: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

sns.distplot(train[train['Survived']==1]['Fare'])
<Axes: xlabel='Fare', ylabel='Density'>
```

	Name
0	Braund, Mr. Owen Harris
1	Cumings, Mrs. John Bradley (Florence Briggs Th...
2	Heikkinen, Miss. Laina
3	Futelle, Mrs. Jacques Heath (Lily May Peel)
4	Allen, Mr. William Henry
...	...
886	Montvila, Rev. Juozas
887	Graham, Miss. Margaret Edith
888	Johnston, Miss. Catherine Helen "Carrie"
889	Behr, Mr. Karl Howell
890	Dooley, Mr. Patrick

891 rows × 1 columns

```
dtype: object
```

```
train.drop(columns=['Name'],inplace=True)
test.drop(columns=['Name'],inplace=True)
```

```
train['family']=train['SibSp'] + train['Parch'] + 1
test['family']=test['SibSp'] + test['Parch'] + 1
```

```
train.drop(columns=['SibSp', 'Parch'], inplace=True)
test.drop(columns=['SibSp', 'Parch'], inplace=True)
```

```
train['family'].value_counts()
```



	count
family	
1	537
2	161
3	102
4	29
6	22
5	15
7	12
11	7
8	6

dtype: int64

```
train[['family', 'Survived']].groupby('family').mean()
```



	Survived
family	
1	0.303538
2	0.552795
3	0.578431
4	0.724138
5	0.200000
6	0.136364
7	0.333333
8	0.000000
11	0.000000

```
def family_size(number):
    if number==1:
        return "Alone"
    elif number>1 and number <5:
        return "Small"
    else:
        return "Large"
```

```
family_size(5)
```



'Large'

```
train['family_size']=train['family'].apply(family_size)
```

```
test['family_size']=test['family'].apply(family_size)
```

```
train.drop(columns=['family'], inplace=True)
test.drop(columns=['family'], inplace=True)
```

```
y=train['Survived'].values
y
```



```
array([0, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1,
       1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1,
       1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1,
       1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0,
       1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0,
       0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0,
       0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0,
```



```
y_pred.shape
```

```
↔ (179,)
```

```
y_test.shape
```

```
↔ (179,)
```

```
from sklearn.metrics import accuracy_score  
accuracy_score(y_test,y_pred)
```

```
↔ 0.8100558659217877
```

```
yf=clf.predict(Xf)
```

```
yf.shape
```

```
↔ (418,)
```

```
submission=pd.DataFrame()
```

```
submission['PassengerId']=passengerId  
submission['Survived']=yf
```

```
submission.to_csv('submission.csv', index=False)
```

```
from google.colab import files
```