

```
In [2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot

In [3]: titanic_data=pd.read_csv('titanic_train.csv')

In [4]: len(titanic_data)

Out[4]: 891

In [5]: titanic_data.head()

Out[5]:
```

	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		Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0		PC 17599	71.2833	C85	C
2	3	1	3		Heikinen, Mrs. 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

```
In [6]: titanic_data.index

Out[6]: RangeIndex(start=0, stop=891, step=1)

In [7]: titanic_data.columns

Out[7]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp', 'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],
          dtype='object')

In [8]: titanic_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                  296 non-null    object
11  Embarked               889 non-null    object
dtypes: float64(5), int64(5), object(5)
memory usage: 83.7+ KB

In [9]: titanic_data.dtypes

Out[9]: PassengerId    int64
Survived          int64
Pclass            int64
Name              object
Sex               object
Age              float64
SibSp             int64
Parch             int64
Ticket            object
Fare              float64
Cabin             object
Embarked          object
dtype: object

In [10]: titanic_data.describe()

Out[10]:
```

	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

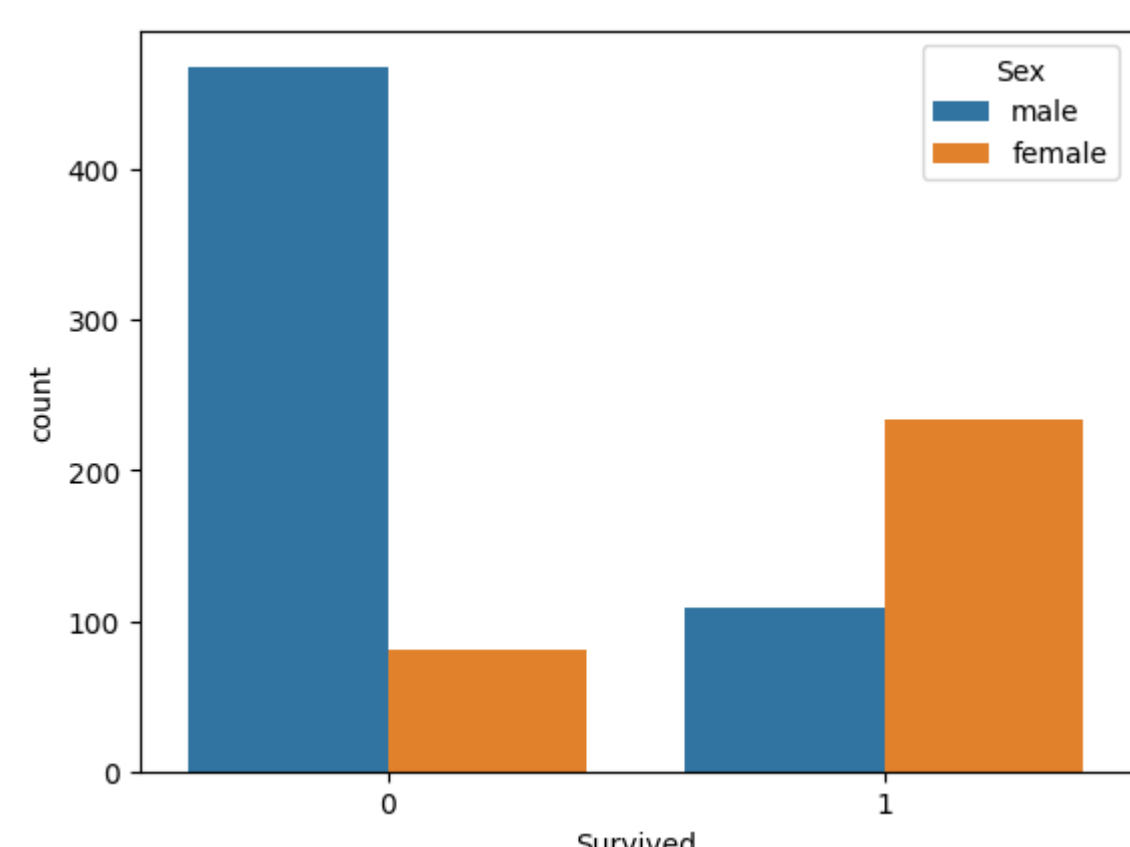
```
In [11]: sns.countplot(x='Survived',data=titanic_data)

Out[11]: <Axes: xlabel='Survived', ylabel='count'>
```



```
In [12]: sns.countplot(x='Survived',data=titanic_data,hue='Sex')

Out[12]: <Axes: xlabel='Survived', ylabel='count'>
```



```
In [13]: titanic_data.isna()

Out[13]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	False	False	False	False	False	False	False	False	False	False	True	False
1	False	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	True	False
3	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	True	False
...
886	False	False	False	False	False	False	False	False	False	False	True	False
887	False	False	False	False	False	False	False	False	False	False	False	False
888	False	False	False	False	True	False	False	False	False	True	False	False
889	False	False	False	False	False	False	False	False	False	False	False	False
890	False	False	False	False	False	False	False	False	False	False	True	False

891 rows x 12 columns

```
In [14]: titanic_data.isna().sum()

Out[14]: 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 [15]: sns.heatmap(titanic_data.isna())

Out[15]: <Axes: >
```



```
In [16]: (titanic_data['Age'].isna().sum()/len(titanic_data['Age']))*100

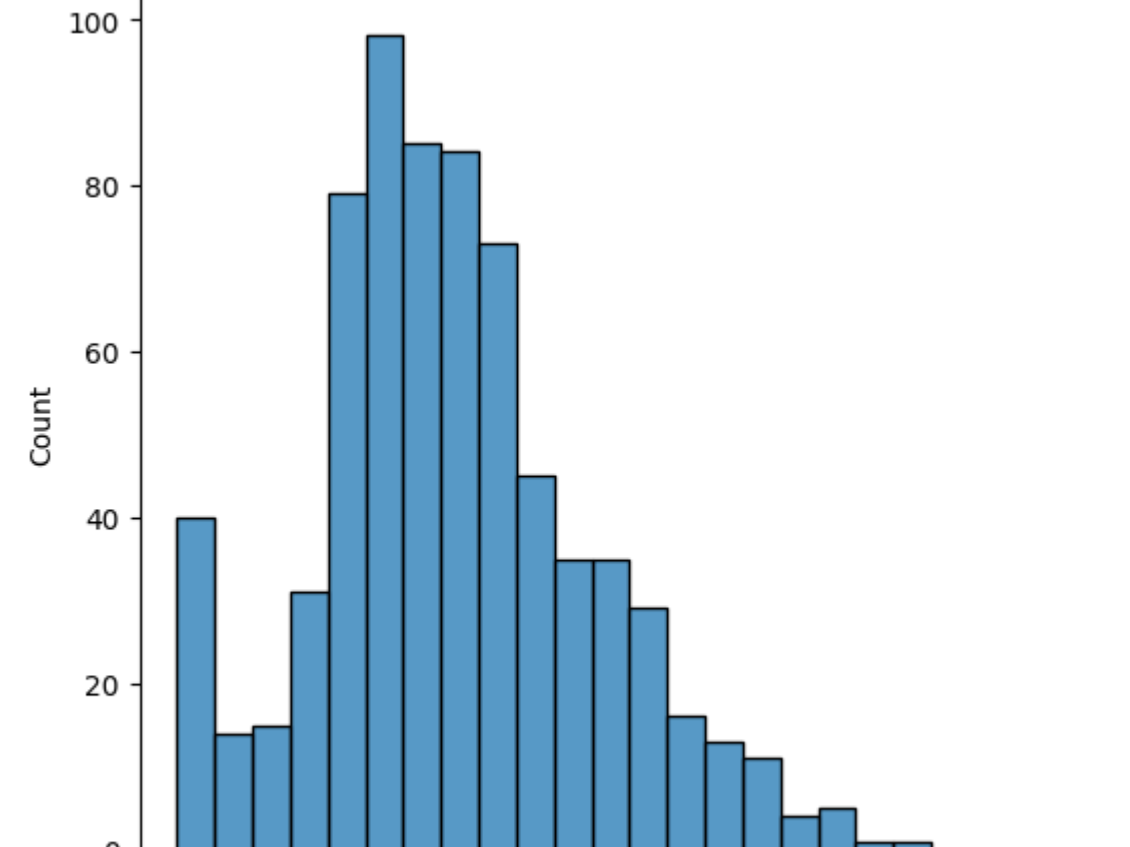
Out[16]: 19.865319865319865

In [17]: (titanic_data['Cabin'].isna().sum()/len(titanic_data['Cabin']))*100

Out[17]: 77.19437719437711

In [18]: sns.displot(x='Age',data=titanic_data)

Out[18]: <seaborn.axisgrid.FacetGrid at 8x14bad3ab898>
```



```
In [19]: # Instead of this:
# titanic_data['Age'].fillna(titanic_data['Age'].mean(), inplace=True)
# Use one of the recommended alternatives:
# Option 1: using the DataFrame's fillna method directly
titanic_data.fillna({'Age': titanic_data['Age'].mean(), inplace=True)
# Option 2: Assigning the result back to the DataFrame column
titanic_data['Age'] = titanic_data['Age'].fillna(titanic_data['Age'].mean())

In [20]: titanic_data['Age'].isna().sum()

Out[20]: 0

In [21]: sns.heatmap(titanic_data.isna())

Out[21]: <Axes: >
```



```
In [22]: titanic_data.drop('Cabin',axis=1,inplace=True)

In [23]: titanic_data.head()

Out[23]:
```

	PassengerId	Survived	Pclass		Name	Sex	Age	SibSp	Parch		Ticket	Fare	Embarked
0	1	0	3		Braund, Mr. Owen Harris	male	22.0	1	0		A/5 21171	7.2500	S
1	2	1	1		Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0		PC 17599	71.2833	C
2	3	1	3		Heikinen, Mrs. Laina	female	26.0	0	0		STON/O2. 3101282	7.9250	S
3	4	1	1		Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0		113803	53.1000	S
4	5	0	3		Allen, Mr. William Henry	male	35.0	0	0		373450	8.0500	S

```
In [24]: titanic_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 11 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                    891 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  Embarked               889 non-null    object
dtypes: float64(2), int64(5), object(4)
memory usage: 76.7+ KB

In [25]: titanic_data.dtypes

Out[25]: PassengerId    int64
Survived          int64
Pclass            int64
Name              object
Sex               object
Age              float64
SibSp             int64
Parch             int64
Ticket            object
Fare              float64
Embarked          object
dtype: object

In [26]: gender=pd.get_dummies(titanic_data['Sex'],drop_first=True)

In [27]: titanic_data['gender']=gender

In [28]: titanic_data.head()

Out[28]:
```

	PassengerId	Survived	Pclass		Name	Sex	Age	SibSp	Parch		Ticket	Fare	Embarked	Gender
0	1	0	3		Braund, Mr. Owen Harris	male	22.0	1	0		A/5 21171	7.2500	S	True
1	2	1	1		Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0		PC 17599	71.2833	C	False
2	3	1	3		Heikinen, Mrs. Laina	female	26.0	0	0		STON/O2. 3101282	7.9250	S	False
3	4	1	1		Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0		113803	53.1000	S	False
4	5	0	3		Allen, Mr. William Henry	male	35.0	0	0		373450	8.0500	S	True

```
In [29]: titanic_data.drop(['Name','Sex','Ticket','Embarked'],axis=1,inplace=True)

In [30]: titanic_data.head()

Out[30]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare	Gender
0	1	0	3	22.0	1	0	7.2500	True
1	2	1	1	38.0	1	0	71.2833	False
2	3	1	3	26.0	0	0	7.9250	False
3	4	1	1	35.0	1	0	53.1000	False
4	5	0	3	35.0	0	0	8.0500	True

```
In [31]: x=titanic_data[['PassengerId','Pclass','Age','SibSp','Parch','Fare','Gender']]
y=titanic_data['Survived']

In [32]: y

Out[32]:
```

	y
0	0
1	1
2	1
3	1
4	0
...	...
886	0
887	1
888	0
889	1
890	0

Name: Survived, Length: 891, dtype: int64

```
In [33]: from sklearn.model_selection import train_test_split

In [34]: x_train,x_test,y_train,y_test = train_test_split(x,y, test_size=0.33, random_state=42)

In [35]: from sklearn.linear_model import LogisticRegression

In [36]: lr=LogisticRegression()

In [37]: lr.fit(x_train,y_train)

In [38]: C:\Users\leah\AppData\Local\Microsoft\Python\Python32\site-packages\sklearn\linear_model\logistic.py:469: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_ = check_optimize_result(
Out[38]: LogisticRegression
LogisticRegression()

In [39]: LogisticRegression()

Out[39]: LogisticRegression
LogisticRegression()

In [40]: predict_lr.predict(x_test)

In [41]: from sklearn.metrics import confusion_matrix

In [42]: pd.DataFrame(confusion_matrix(y_test,predict),columns=['Predicted No','Predicted Yes'],index=['Actual No','Actual Yes'])

Out[42]:
```

	Predicted No	Predicted Yes
Actual No	151	24
Actual Yes	38	82

```
In [43]: from sklearn.metrics import classification_report

In [44]: print(classification_report(y_test,predict))

precision    recall    f1-score   support

0           0.88        0.86        0.83       175
1           0.77        0.68        0.73       128

accuracy          0.79        0.77        0.79       295
macro avg          0.79        0.77        0.79       295
weighted avg          0.79        0.79        0.79       295

In [45]: from sklearn.preprocessing import StandardScaler
x_train,x_test,y_train,y_test=train_test_split(x,y, test_size=0.2, random_state=42)
scaler=StandardScaler()
x_data=scaler.fit_transform(x)
y_data=y.to_numpy

In [46]: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
clf = RandomForestClassifier()
param_grid=[{"n_estimators":[10,160,200,300],"max_depth":[None,5,10],"min_samples_split":[2,3,4]}]
y_data = y.values
grid_search=GridSearchCV(clf,param_grid,cv=3,scoring='accuracy', return_train_score=True)
grid_search.fit(x_data,y_data)

In [47]: from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# 1. Prepare the Data
# Assuming x contains features and y contains the target variable
# X_train,X_test,y_train,y_test = train_test_split(X,y, test_size=0.2, random_state=42)

# 2. Create and Train the Model
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)
rf_classifier.fit(x_train,y_train)

# 3. Evaluate the Model
y_pred = rf_classifier.predict(x_test)
accuracy = accuracy_score(y_test,y_pred)
print("Accuracy:", accuracy)

# 4. (Optional) Hyperparameter Tuning
# You can perform Hyperparameter tuning to improve the model's performance

# 5. (Optional) Deploy the Model
# Once satisfied with the model's performance, you can deploy it for making predictions on new data

In [48]: # Import Random Forest Classifier
from sklearn.ensemble import RandomForestClassifier

# instantiate the classifier
rfc = RandomForestClassifier(random_state=0)

# fit the model
rfc.fit(x_train,y_train)

# Predict the Test set results
y_pred = rfc.predict(x_test)

# Check accuracy score
from sklearn.metrics import accuracy_score
print("Model accuracy score with 10 decision-trees : (0.0.4f)".format(accuracy_score(y_test,y_pred)))

Model accuracy score with 10 decision-trees : 0.8237

In [49]: rfc_100 = RandomForestClassifier(n_estimators=100, random_state=0)

# fit the model to the training set
rfc_100.fit(x_train,y_train)

# Predict on the test set results
y_pred_100 = rfc_100.predict(x_test)

# Check accuracy score
print("Model accuracy score with 100 decision-trees : (0.0.4f)".format(accuracy_score(y_test,y_pred_100)))

Model accuracy score with 100 decision-trees : 0.8237

In [50]: # create the classifier with n_estimators=100
clf = RandomForestClassifier(n_estimators=100, random_state=0)

# fit the model to the training set
clf.fit(x_train,y_train)

Out[50]:
RandomForestClassifier
RandomForestClassifier(random_state=0)

In [51]: feature_scores = pd.Series(clf.feature_importances_, index=x_train.columns).sort_values(ascending=False)
feature_scores

Out[51]:
Gender          0.245881
Fare            0.218995
PassengerId     0.206622
Pclass          0.178676
Age             0.087233
SibSp           0.043555
```

