

```
from google.colab import files
uploaded = files.upload()
```

Choose Files Titanic-Dataset.csv

Titanic-Dataset.csv(text/csv) - 61194 bytes, last modified: 6/15/2025 - 100% done

Saving Titanic-Dataset.csv to Titanic-Dataset.csv

```
import pandas as pd
df = pd.read_csv("Titanic-Dataset.csv") # or the exact filename you uploaded
df.head()
```

	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

Next steps:

Generate code with df

☒ View recommended plots

New interactive sheet

```
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
```

```
df.describe()
```



	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



```
df.isnull().sum()
```



	0
PassengerId	0
Survived	0
Pclass	0
Name	0
Sex	0
Age	177
SibSp	0
Parch	0
Ticket	0
Fare	0
Cabin	687
Embarked	2

```
df.isnull().sum()
```

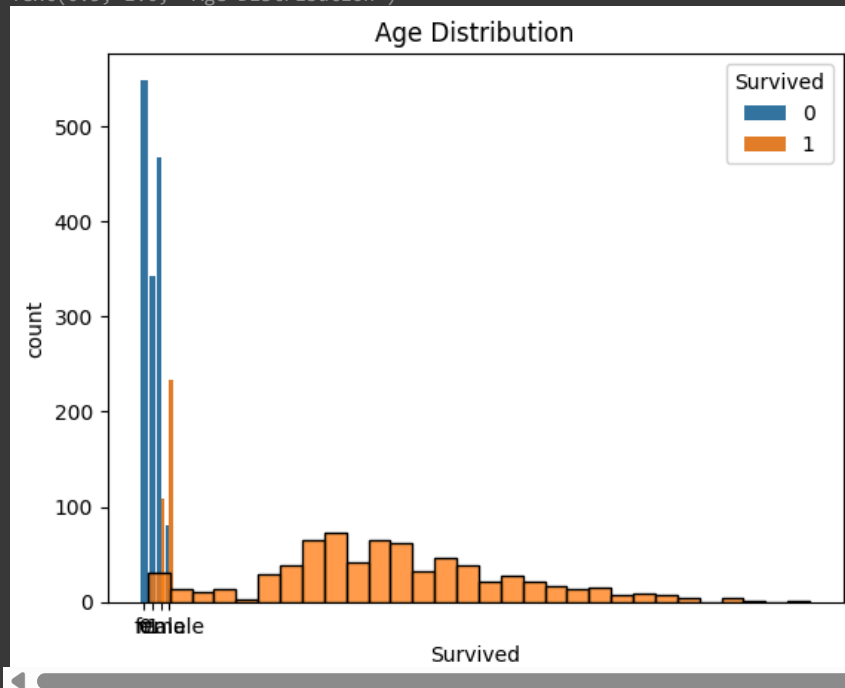
```
import seaborn as sns
import matplotlib.pyplot as plt
```

```
sns.countplot(x='Survived', data=df)
plt.title("Survival Count")

sns.countplot(x='Sex', hue='Survived', data=df)
plt.title("Survival by Gender")

sns.histplot(df['Age'].dropna(), bins=30)
plt.title("Age Distribution")
```

↻ Text(0.5, 1.0, 'Age Distribution')



```
df['Age'] = df['Age'].fillna(df['Age'].median())
df['Embarked'] = df['Embarked'].fillna(df['Embarked'].mode()[0])
```

```
# Features we'll use to predict
X = df[['Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare', 'Embarked']]

# Target variable
y = df['Survived']
```

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)
```

df.dtypes



0

PassengerId	int64
Survived	int64
Pclass	int64
Sex	object
Age	float64
SibSp	int64
Parch	int64
Fare	float64
Embarked	object

dtypes: object

```
print(df['Sex'].unique())
print(df['Embarked'].unique())
```



```
['male' 'female']
['S' 'C' 'Q']
```

```
df['Sex'] = df['Sex'].map({'male': 0, 'female': 1})
df['Embarked'] = df['Embarked'].map({'S': 0, 'C': 1, 'Q': 2})
```

```
X = df[['Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare', 'Embarked']]
y = df['Survived']
```

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)
```

```
from sklearn.linear_model import LogisticRegression
model = LogisticRegression(max_iter=200)
model.fit(X_train, y_train)
```

LogisticRegression

LogisticRegression(max_iter=200)

```
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
```

```
# --- LOGISTIC REGRESSION ---
```

```
log_model = LogisticRegression(max_iter=200)
log_model.fit(X_train, y_train)
y_pred_log = log_model.predict(X_test)
```

```
print(" ♦ Logistic Regression Results")
```

```
print("Accuracy:", accuracy_score(y_test, y_pred_log))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred_log))
print("Classification Report:\n", classification_report(y_test, y_pred_log))
```

```
# --- RANDOM FOREST ---
```

```
from sklearn.ensemble import RandomForestClassifier
rf_model = RandomForestClassifier(random_state=42)
rf_model.fit(X_train, y_train)
y_pred_rf = rf_model.predict(X_test)
```

```
print("\n ♦ Random Forest Results")
```

```
print("Accuracy:", accuracy_score(y_test, y_pred_rf))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred_rf))
print("Classification Report:\n", classification_report(y_test, y_pred_rf))
```

♦ Logistic Regression Results

Accuracy: 0.7988826815642458

Confusion Matrix:

```
[[89 16]
 [20 54]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.82	0.85	0.83	105
1	0.77	0.73	0.75	74
accuracy			0.80	179
macro avg	0.79	0.79	0.79	179
weighted avg	0.80	0.80	0.80	179

♦ Random Forest Results

Accuracy: 0.8268156424581006

Confusion Matrix:

[[92 13]

[18 56]]

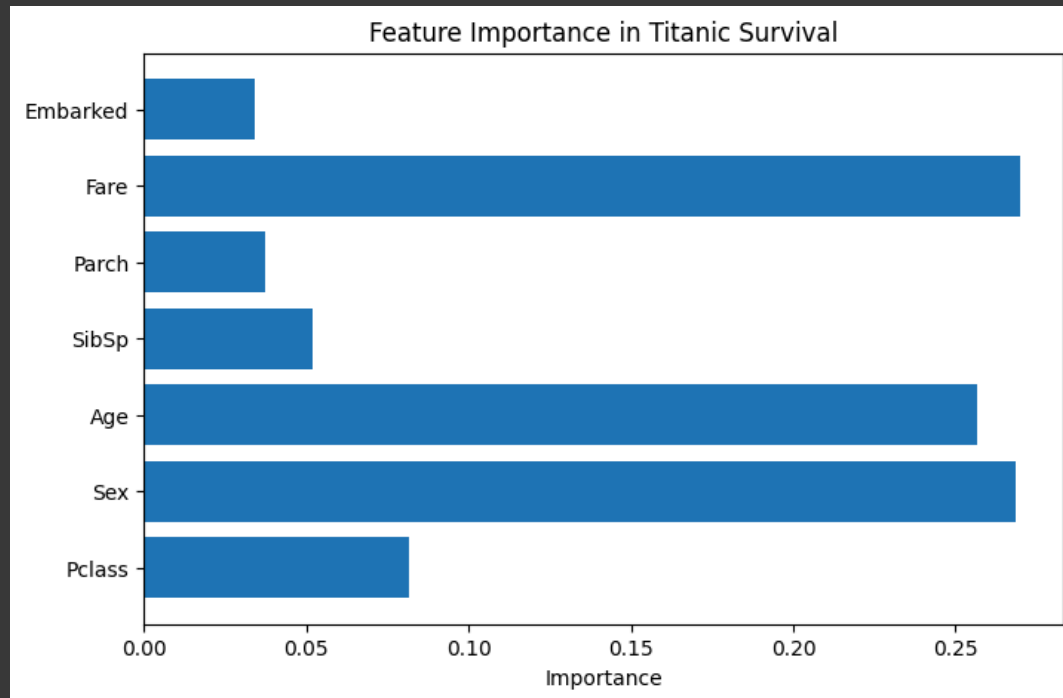
Classification Report:

	precision	recall	f1-score	support
0	0.84	0.88	0.86	105
1	0.81	0.76	0.78	74
accuracy			0.83	179
macro avg	0.82	0.82	0.82	179
weighted avg	0.83	0.83	0.83	179

```
import matplotlib.pyplot as plt
```

```
importances = rf_model.feature_importances_  
features = X.columns
```

```
plt.figure(figsize=(8,5))  
plt.barh(features, importances)  
plt.xlabel("Importance")  
plt.title("Feature Importance in Titanic Survival")  
plt.show()
```



```
import joblib
joblib.dump(rf_model, 'titanic_model.pkl')
```



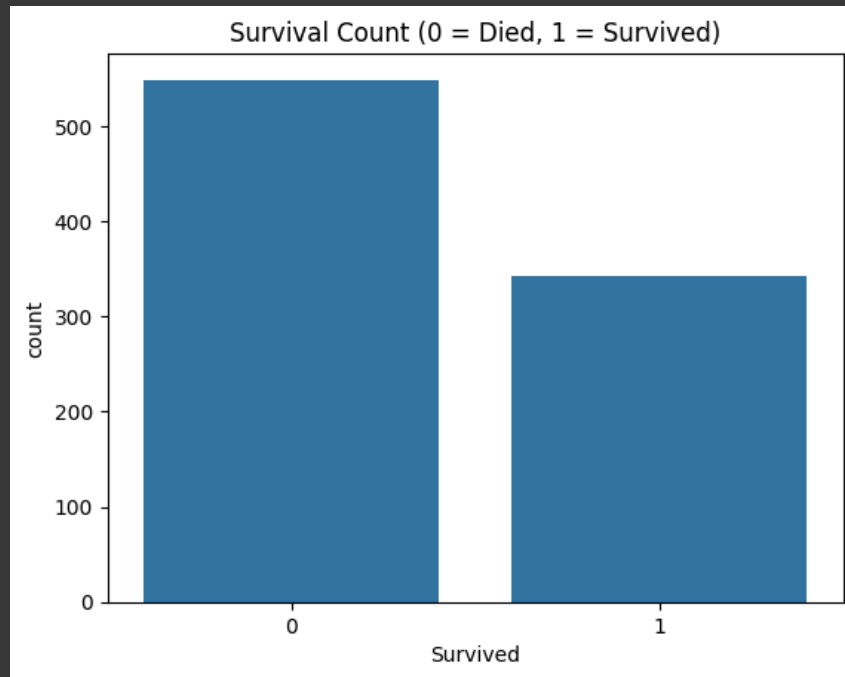
```
['titanic_model.pkl']
```

```
from google.colab import files
files.download('titanic_model.pkl')
```



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

sns.countplot(x='Survived', data=df)
plt.title('Survival Count (0 = Died, 1 = Survived)')
plt.show()
```



```
sns.histplot(data=df, x='Age', hue='Survived', kde=True, bins=30)
plt.title('Age Distribution: Survivors vs Non-Survivors')
plt.show()
```




Age Distribution: Survivors vs Non-Survivors