

DASARI.NAGAVENI

EMAIL: dasari.nagaveni2020@vitstudent.ac.in

Vellore institute of technology

Chennai

Assignment-2

ADS Assignment 2

Titanic Ship Case Study

Problem Description: On April 15, 1912, during her maiden voyage, the Titanic sank after colliding

with an iceberg, killing 1502 out of 2224 passengers and crew. Translated 32% survival rate.

☐ One of the reasons that the shipwreck led to such loss of life was that there were not

enough lifeboats for the passengers and crew.

☐ Although there was some element of luck involved in surviving the sinking, some groups of

people were more likely to survive than others, such as women, children, and the upperclass.

The problem associated with the Titanic dataset is to predict whether a passenger survived the

disaster or not. The dataset contains various features such as passenger class, age, gender,

cabin, fare, and whether the passenger had any siblings or spouses on board. These features can

be used to build a predictive model to determine the likelihood of a passenger surviving the

disaster. The dataset offers opportunities for feature engineering, data visualization, and model

selection, making it a valuable resource for developing and testing data analysis and machine

learning skills.

Perform Below Tasks to complete the assignment:-

1. Download the dataset: Dataset
2. Load the dataset.
3. Perform Below Visualizations.
 - Univariate Analysis
 - Bi - Variate Analysis
 - Multi - Variate Analysis
4. Perform descriptive statistics on the dataset.
5. Handle the Missing values.
6. Find the outliers and replace the outliers
7. Check for Categorical columns and perform encoding.
8. Split the data into dependent and independent variables.
9. Scale the independent variables
10. **Split the data into training and testing**

[1] Download the dataset: Dataset

Tools - LibreOffice Calc

FileEditViewInsertFormatStyleSheetDataToolsWindowHelp

</

[2] Load the dataset.

The screenshot shows a code editor interface. On the left is a file explorer with a 'Files' tab. It displays a directory structure with a folder named 'sample_data' and a file named 'titanic.csv'. The main area of the editor shows two code snippets. Snippet [3] contains imports for numpy, pandas, difflib, and TfIdfVectorizer from sklearn, along with cosine_similarity from sklearn.metrics.pairwise. Snippet [4] shows the loading of 'titanic.csv' into a pandas DataFrame named 'tit_data'.

[3] Perform Below Visualizations.

- Univariate Analysis
- Bi - Variate Analysis

- Multi - Variate Analysis

CODE

- Univariate Analysis

HISTOGRAM FOR AGE

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

```
# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
```

```
plt.hist(tit_data['age'].dropna(), bins=10)
```

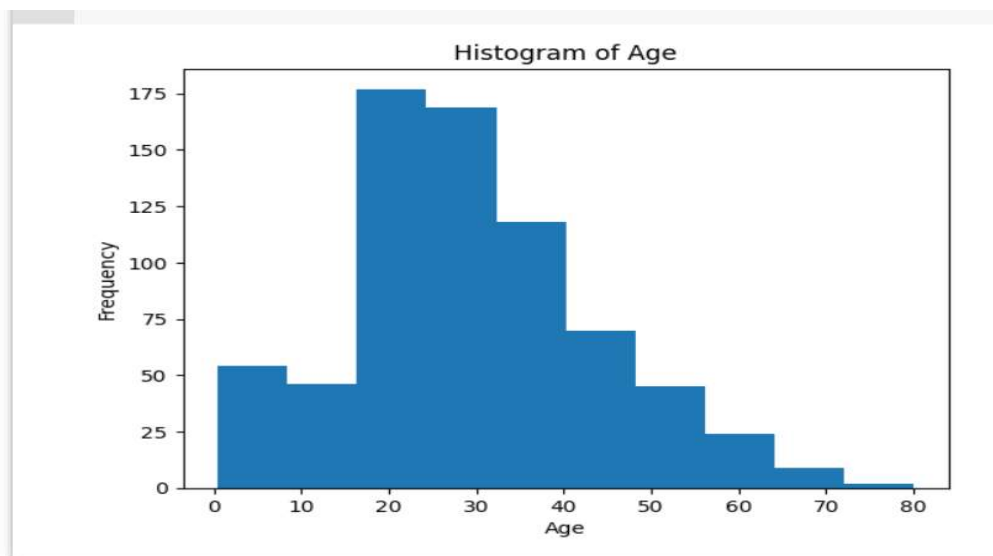
```
plt.xlabel('Age')
```

```
plt.ylabel('Frequency')
```

```
plt.title('Histogram of Age')
```

```
plt.show()
```

OUTPUT



BOXPLOT

```
import seaborn as sns
```

```
# Assuming your data is in a pandas DataFrame called 'df'
```

```
sns.boxplot(x=df['age'])
```

```
plt.xlabel('age')
```

```
plt.title('Boxplot')
```

```
plt.show()
```

OUTPUT

● Bi - Variate Analysis

SCATTER PLOT

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

```
# Assuming your data is in a pandas DataFrame called 'df'
```

```
plt.scatter(df['age'], df['fare'])
```

```
plt.xlabel('age')
```

```
plt.ylabel('fare')
```

```
plt.title('Scatter Plot')
```

```
plt.show()
```

OUTPUT

HEAT MAP

```
import seaborn as sns
```

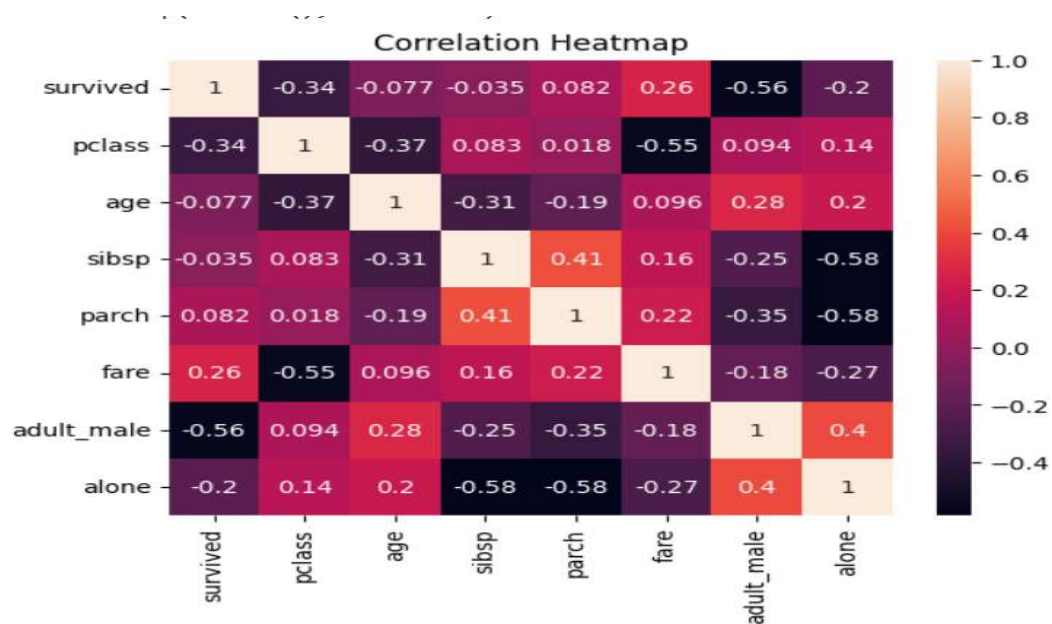
```
# Assuming your data is in a pandas DataFrame called 'df'
```

```
sns.heatmap(df.corr(), annot=True)
```

```
plt.title('Correlation Heatmap')
```

```
plt.show()
```

OUTPUT



● Multi - Variate Analysis

PAIR PLOT

```
import seaborn as sns
```

```
# Assuming your data is in a pandas DataFrame called 'df'
```

```
sns.pairplot(df)
```

```
plt.title('Pairplot')
```

```
plt.show()
```

OUTPUT



3D SCATTER PLOT

```
import matplotlib.pyplot as plt

from mpl_toolkits.mplot3d import Axes3D
```

```
# Assuming your data is in a pandas DataFrame called 'df'
```

```
fig = plt.figure()

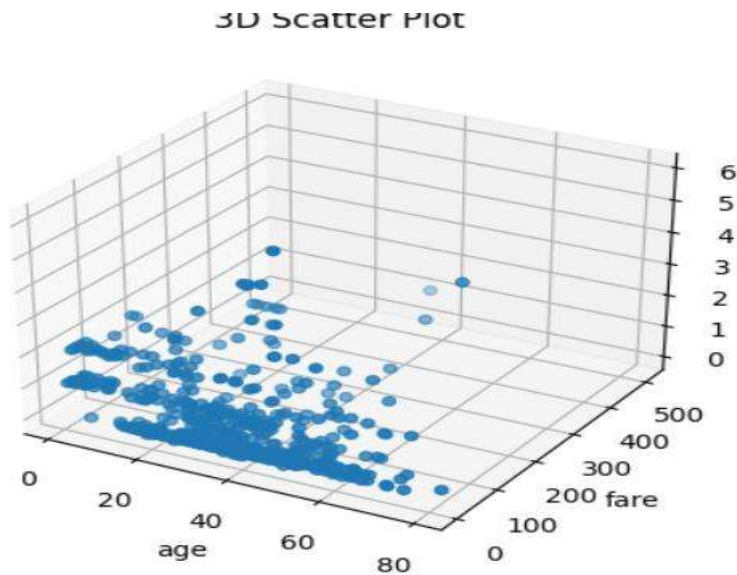
ax = fig.add_subplot(111, projection='3d')

ax.scatter(df['age'], df['fare'], df['parch'])

ax.set_xlabel('age')
ax.set_ylabel('fare')
ax.set_zlabel('parch')

plt.title('3D Scatter Plot')
```

```
plt.show()
```



[4]

Perform descriptive statistics on the dataset.

```
import pandas as pd
```

```
# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
```

```
# Select numeric columns for descriptive statistics
```

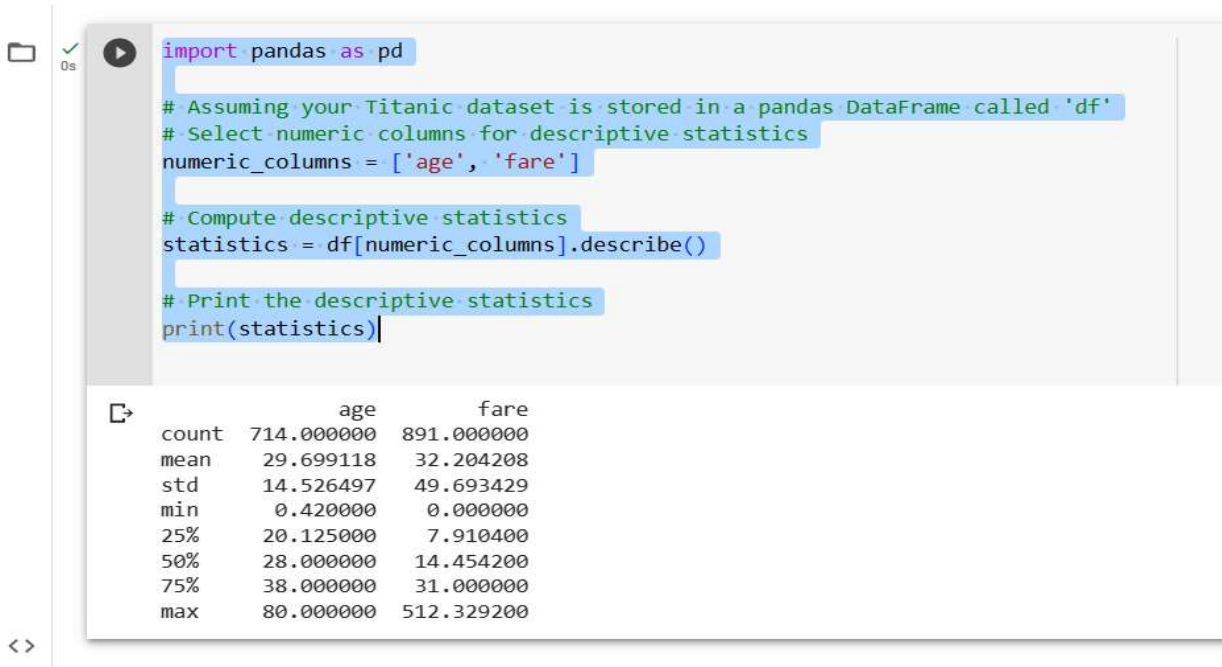
```
numeric_columns = ['age', 'fare']
```

```
# Compute descriptive statistics
```

```
statistics = df[numeric_columns].describe()
```

```
# Print the descriptive statistics
```

```
print(statistics)
```

```
import pandas as pd

# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
# Select numeric columns for descriptive statistics
numeric_columns = ['age', 'fare']

# Compute descriptive statistics
statistics = df[numeric_columns].describe()

# Print the descriptive statistics
print(statistics)
```

	age	fare
count	714.000000	891.000000
mean	29.699118	32.204208
std	14.526497	49.693429
min	0.420000	0.000000
25%	20.125000	7.910400
50%	28.000000	14.454200
75%	38.000000	31.000000
max	80.000000	512.329200

[5] Handle the Missing values.

IDENTIFY MISSING VALUES

```
import pandas as pd
```

```
# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
```

```
missing_values = df.isnull().sum()
```

```
print(missing_values)
```

```
import pandas as pd

# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
missing_values = df.isnull().sum()
print(missing_values)
```

survived	0
pclass	0
sex	0
age	177
sibsp	0
parch	0
fare	0
embarked	2
class	0
who	0
adult_male	0
deck	688
embark_town	2
alive	0
alone	0
dtype: int64	

DROP MISSING VALUES

Drop rows with any missing values

```
df.dropna(inplace=True)
```

IMPUTE MISSING VALUE

Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'

```
mean_age = df['age'].mean()
```

```
df['age'].fillna(mean_age, inplace=True)
```

CREATE INDICATOR VARIABLE

[6]

Find the outliers and replace the outliers

IDENTIFY OUTLIERS

```
import numpy as np
```

Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'

column = 'fare'

threshold = 3 # Set the threshold for outlier detection

z_scores = (df[column] - df[column].mean()) / df[column].std()

outliers = df[np.abs(z_scores) > threshold]

print(outliers)



```
import numpy as np

# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
column = 'fare'
threshold = 3 # Set the threshold for outlier detection

z_scores = (df[column] - df[column].mean()) / df[column].std()
outliers = df[np.abs(z_scores) > threshold]
print(outliers)
```

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	\
679	1	1	male	36.0	0	1	512.3292	C	First	
737	1	1	male	35.0	0	0	512.3292	C	First	
	who	adult_male	deck	embark_town	alive	alone	Age_Missing			
679	man	True	B	Cherbourg	yes	False	0			
737	man	True	B	Cherbourg	yes	True	0			

REPLACE OUTLIERS

Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'

column = 'fare'

median_value = df[column].median()

df.loc[np.abs(z_scores) > threshold, column] = median_value

[7] Check for Categorical columns and perform encoding.

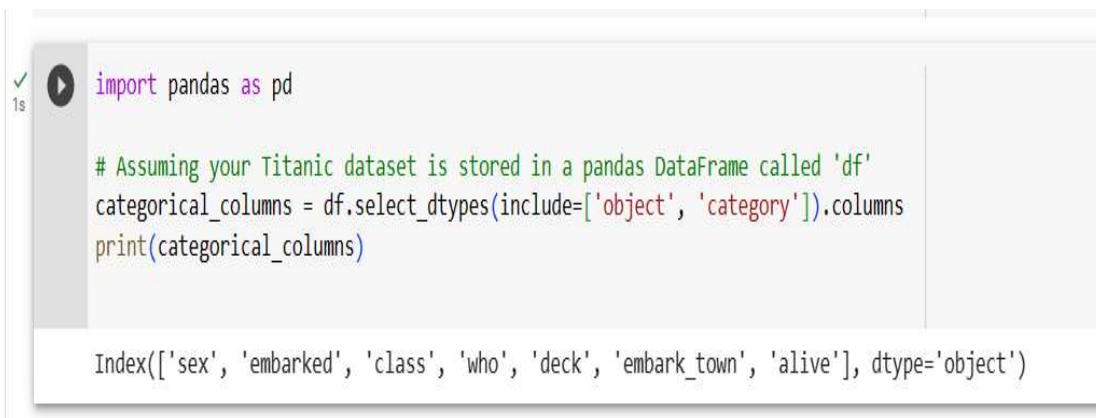
IDENTIFYING CATEGORICAL COLUMNS

```
import pandas as pd
```

```
# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
```

```
categorical_columns = df.select_dtypes(include=['object', 'category']).columns
```

```
print(categorical_columns)
```



A screenshot of a Jupyter Notebook cell. The code imports pandas as pd, assumes a Titanic dataset is in a DataFrame 'df', and uses select_dtypes to find columns of type 'object' or 'category'. The output shows the categorical columns: sex, embarked, class, who, deck, embark_town, and alive.

```
import pandas as pd

# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
categorical_columns = df.select_dtypes(include=['object', 'category']).columns
print(categorical_columns)
```

Index(['sex', 'embarked', 'class', 'who', 'deck', 'embark_town', 'alive'], dtype='object')

PERFORM ENCODING



A screenshot of a Jupyter Notebook cell. The code imports LabelEncoder from sklearn.preprocessing, creates an instance, and fits it to the 'deck' column. The output shows the encoded 'sex' column. Below this, another code block shows the creation of one-hot encoded variables for the 'sex' column and concatenating them to the original DataFrame.

```
[27] from sklearn.preprocessing import LabelEncoder

# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
label_encoder = LabelEncoder()
df['sex'] = label_encoder.fit_transform(df['deck'])
```

```
one_hot_encoded = pd.get_dummies(df['sex'], prefix='column')
df = pd.concat([df, one_hot_encoded], axis=1)
```

[8]

Split the data into dependent and independent variables.

```
import pandas as pd
```

```
# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
```

```
# Identify the dependent variable (target variable)
```

```
target_variable = 'alive'
```

```
# Identify the independent variables (features)
```

```
independent_variables = ['pclass', 'age', 'sex', 'fare']
```

```
# Split the data into dependent and independent variables
```

```
X = df[independent_variables]
```

```
y = df[target_variable]
```

```
# Print the shape of the data
```

```
print("Independent variables (X):", X.shape)
```

```
print("Dependent variable (y):", y.shape)
```

```
[28] one_hot_encoded = pd.get_dummies(df['sex'], prefix='column')
df = pd.concat([df, one_hot_encoded], axis=1)

import pandas as pd

# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
# Identify the dependent variable (target variable)
target_variable = 'alive'

# Identify the independent variables (features)
independent_variables = ['pclass', 'age', 'sex', 'fare']

# Split the data into dependent and independent variables
X = df[independent_variables]
y = df[target_variable]

# Print the shape of the data
print("Independent variables (X):", X.shape)
print("Dependent variable (y):", y.shape)
```

Independent variables (X): (182, 4)
Dependent variable (y): (182,)

[9]

Scale the independent variables

```
import pandas as pd
```

```
from sklearn.preprocessing import StandardScaler
```

```
# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
```

```
# Assuming your independent variables are stored in a DataFrame called 'X'
```

```
# Select the independent variables you want to scale
```

```
independent_variables = ['pclass', 'age', 'fare']
```

```
# Create a new DataFrame with only the selected independent variables
```

```
X_selected = df[independent_variables]
```

```
# Initialize the StandardScaler
```

```
scaler = StandardScaler()
```

```
# Fit and transform the selected independent variables
```

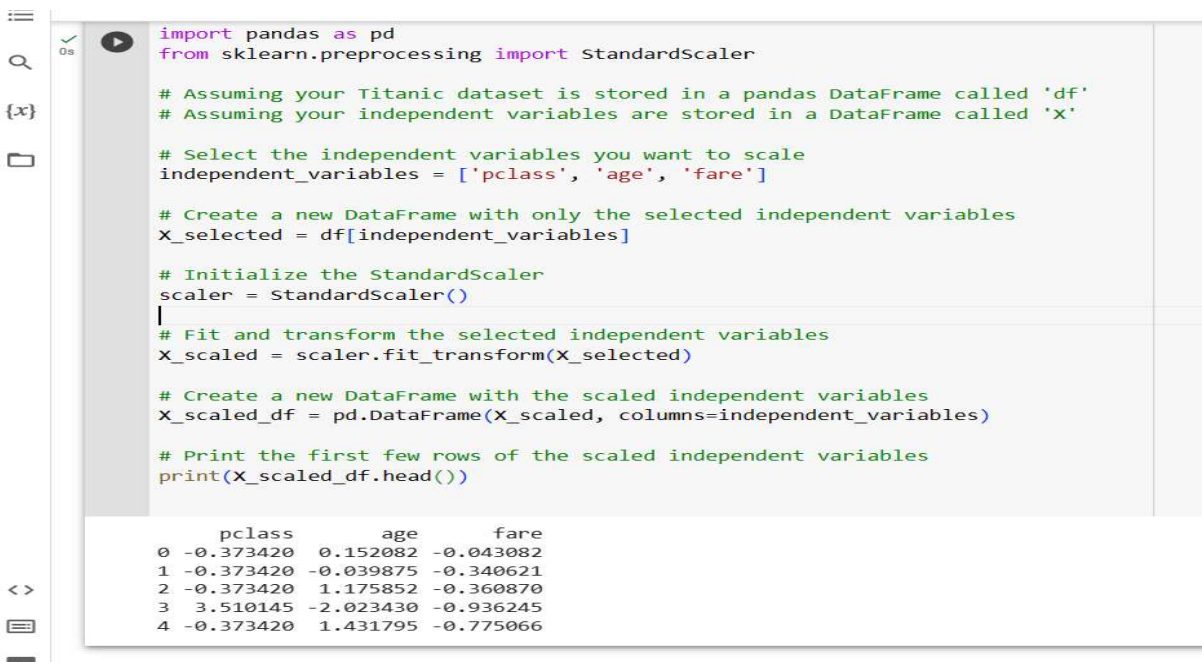
```
X_scaled = scaler.fit_transform(X_selected)
```

```
# Create a new DataFrame with the scaled independent variables
```

```
X_scaled_df = pd.DataFrame(X_scaled, columns=independent_variables)
```

```
# Print the first few rows of the scaled independent variables
```

```
print(X_scaled_df.head())
```



```
import pandas as pd
from sklearn.preprocessing import StandardScaler

# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
# Assuming your independent variables are stored in a DataFrame called 'X'

# Select the independent variables you want to scale
independent_variables = ['pclass', 'age', 'fare']

# Create a new DataFrame with only the selected independent variables
X_selected = df[independent_variables]

# Initialize the StandardScaler
scaler = StandardScaler()

# Fit and transform the selected independent variables
X_scaled = scaler.fit_transform(X_selected)

# Create a new DataFrame with the scaled independent variables
X_scaled_df = pd.DataFrame(X_scaled, columns=independent_variables)

# Print the first few rows of the scaled independent variables
print(X_scaled_df.head())
```

	pclass	age	fare
0	-0.373420	0.152082	-0.043082
1	-0.373420	-0.039875	-0.340621
2	-0.373420	1.175852	-0.360870
3	3.510145	-2.023430	-0.936245
4	-0.373420	1.431795	-0.775066

[10] Split the data into training and testing

```
import pandas as pd
```

```
from sklearn.model_selection import train_test_split
```

```
# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
```

```
# Assuming your dependent variable is stored in a Series called 'y'
```

```
# Assuming your independent variables are stored in a DataFrame called 'X'
```

```
# Split the data into training and testing sets
```

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

```
# Print the shape of the training and testing sets
```

```
print("X_train shape:", X_train.shape)
```

```
print("X_test shape:", X_test.shape)
```

```
print("y_train shape:", y_train.shape)
```

```
print("y_test shape:", y_test.shape)
```

```
import pandas as pd
from sklearn.model_selection import train_test_split

# Assuming your Titanic dataset is stored in a pandas DataFrame called 'df'
# Assuming your dependent variable is stored in a Series called 'y'
# Assuming your independent variables are stored in a DataFrame called 'X'

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Print the shape of the training and testing sets
print("X_train shape:", X_train.shape)
print("X_test shape:", X_test.shape)
print("y_train shape:", y_train.shape)
print("y_test shape:", y_test.shape)
```

```
X_train shape: (145, 4)
X_test shape: (37, 4)
y_train shape: (145,)
y_test shape: (37,)
```