

Data Analytics Lifecycle
[Data discovery and preparation]Sim

To perform data discovery and exploratory analysis on real world dataset

Code

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split.
```

```
df = pd.read_csv('titanic.csv')
print(df.head())
print(df.info())
print(df.isnull().sum())
print(df.describe())
```

```
sns.heatmap(df.isnull(), cbar=False,
             cmap='viridis')
```

```
plt.title('missing data visualization')
plt.show()
```

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

```
df[numericals_cols] = imputer.fit_transform
print(df.isnull().sum())
```

Output:

	PassengerId	survived	Pclass
0	1	0	2
1	2	1	1
2	3	1	3
3	4	1	1
4	5	0	3

	Name	Sex	Age	SibSp
0	Brund, Mr. Owen Harry	male	22.0	1
1	Cumings, Mrs. John Bradly	female	38.0	1
2	Allen, Mr. William Henry	male	35.0	0

#	columns	non-null count	DType
0	PassengerId	891 non-null	int 64
1	survived	891 non-null	int 64
2	Pclass	891 non-null	int 64
3	Name	891 non-null	object

memory usage : 83.7 + KB

passengerId : 0

survived : 0

pclass : 0

Name : 0

sex : 0

age : 177

sibsp : 0

Training Data shape : (512, 8)

Testing Data shape : (129, 8)

```
x = df[['pclass', 'Age', 'Fare', 'sibsp', 'parch']]
```

```
y = df['survived']
```

```
X = pd.get_dummies(x, drop_first = True)
```

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

```
print(f"Training Data Shape : {X_train.shape}")
```

also check for better mapping result

Result

thus the data discovery and exploratory analysis on titanic dataset has been executed successfully.