

practical-exam-02

May 23, 2023

```
[3]: from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

1 Problem Statement 2

Perform the following operations using Python on any open source dataset (e.g., data.csv) 1. Provide a clear description of the data and its source. 2. Load the Dataset into pandas dataframe. 3. Data Preprocessing: check for missing values in the data using pandas isnull(), describe() function to get some initial statistics. Provide variable descriptions. Types of variables etc. Check the dimensions of the data frame. 4. Data Formatting and Data Normalization: Summarize the types of variables by checking the data types (i.e., character, numeric, integer, factor, and logical) of the variables in the data set. If variables are not in the correct data type, apply proper type conversions.

```
[1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
[4]: df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/exam_datasets/1-2-
data.csv')
df.describe()
```

```
[4]:
```

	PassengerId	Survived	Pclass	Age	SibSp	\
count	891.000000	891.000000	891.000000	714.000000	891.000000	
mean	446.000000	0.383838	2.308642	29.699118	0.523008	
std	257.353842	0.486592	0.836071	14.526497	1.102743	
min	1.000000	0.000000	1.000000	0.420000	0.000000	
25%	223.500000	0.000000	2.000000	20.125000	0.000000	
50%	446.000000	0.000000	3.000000	28.000000	0.000000	
75%	668.500000	1.000000	3.000000	38.000000	1.000000	
max	891.000000	1.000000	3.000000	80.000000	8.000000	

	Parch	Fare
count	891.000000	891.000000
mean	0.381594	32.204208

std	0.806057	49.693429
min	0.000000	0.000000
25%	0.000000	7.910400
50%	0.000000	14.454200
75%	0.000000	31.000000
max	6.000000	512.329200

```
[5]: df.isnull().sum()
```

```
[5]: 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
```

```
[6]: 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
```

```
[7]: df.describe()
```

```
[7]:
```

	PassengerId	Survived	Pclass	Age	SibSp \
count	891.000000	891.000000	891.000000	714.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008
std	257.353842	0.486592	0.836071	14.526497	1.102743
min	1.000000	0.000000	1.000000	0.420000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000
50%	446.000000	0.000000	3.000000	28.000000	0.000000
75%	668.500000	1.000000	3.000000	38.000000	1.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000

	Parch	Fare
count	891.000000	891.000000
mean	0.381594	32.204208
std	0.806057	49.693429
min	0.000000	0.000000
25%	0.000000	7.910400
50%	0.000000	14.454200
75%	0.000000	31.000000
max	6.000000	512.329200

```
[8]: df.describe(include=['object'])
```

```
[8]:
```

	Name	Sex	Ticket	Cabin	Embarked
count	891	891	891	204	889
unique	891	2	681	147	3
top	Braund, Mr. Owen Harris	male	347082	B96 B98	S
freq	1	577	7	4	644

```
[9]: df.dtypes
```

```
[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
```

```
[10]: df.shape
```

```
[10]: (891, 12)
```

```
[14]: df = df.dropna(subset=['Age'])
df['Age'] = df['Age'].astype(int)
df.dtypes
```

<ipython-input-14-80b0b0292fc6>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using `.loc[row_indexer,col_indexer] = value` instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df['Age'] = df['Age'].astype(int)

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

```
[15]: df = df.dropna(subset=['PassengerId'])
df['PassengerId'] = df['PassengerId'].astype(str)
df.dtypes
```

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

In pandas, the object data type is used to represent strings. When you convert a column to string using the `astype(str)` method, the resulting column will have the object data type.

In Python, you can convert between several common data types using built-in functions. Here are

some examples: * int(x) * float(x) * str(x) * bool(x)