

Data Processing with Pandas

case study

Dataset used:

Loan.csv

1) Loading Data in Pandas DataFrame

```
from google.colab import drive
import pandas as pd
drive.mount('/content/drive')
file_path = '/content/drive/My Drive/case_study_dataset/LoanData.csv'
df = pd.read_csv(file_path)
print(df.head())
```

Output

```
Mounted at /content/drive
  Loan_ID  Gender  Married  Dependents  Education  Self_Employed  \
0  LP001002   Male     No           0    Graduate             No
1  LP001003   Male     Yes           1    Graduate             No
2  LP001005   Male     Yes           0    Graduate             Yes
3  LP001006   Male     Yes           0  Not Graduate             No
4  LP001008   Male     No           0    Graduate             No

  ApplicantIncome  CoapplicantIncome  LoanAmount  Loan_Amount_Term  \
0           5849              0.0         NaN         360.0
1           4583             1508.0         128.0         360.0
2           3000              0.0          66.0         360.0
3           2583             2358.0         120.0         360.0
4           6000              0.0         141.0         360.0

  Credit_History  Property_Area  Loan_Status
0             1.0           Urban            Y
1             1.0           Rural            N
2             1.0           Urban            Y
3             1.0           Urban            Y
4             1.0           Urban            Y
```

2) Printing rows of the Data

```
print('==== First 5 rows in the Dataset ====')
print(df.head()) #prints first 5 rows
print('==== Last 5 rows in the Dataset ====')
print(df.tail()) #print last 5 rows
```

Output

```

==== First 5 rows in the Dataset ====
  Loan_ID Gender Married Dependents Education Self_Employed \
0 LP001002 Male No 0 Graduate No
1 LP001003 Male Yes 1 Graduate No
2 LP001005 Male Yes 0 Graduate Yes
3 LP001006 Male Yes 0 Not Graduate No
4 LP001008 Male No 0 Graduate No

  ApplicantIncome CoapplicantIncome LoanAmount Loan_Amount_Term \
0 5849 0.0 NaN 360.0
1 4583 1508.0 128.0 360.0
2 3000 0.0 66.0 360.0
3 2583 2358.0 120.0 360.0
4 6000 0.0 141.0 360.0

  Credit_History Property_Area Loan_Status
0 1.0 Urban Y
1 1.0 Rural N
2 1.0 Urban Y
3 1.0 Urban Y
4 1.0 Urban Y

==== Last 5 rows in the Dataset ====
  Loan_ID Gender Married Dependents Education Self_Employed \
609 LP002978 Female No 0 Graduate No
610 LP002979 Male Yes 3+ Graduate No
611 LP002983 Male Yes 1 Graduate No
612 LP002984 Male Yes 2 Graduate No
613 LP002990 Female No 0 Graduate Yes

  ApplicantIncome CoapplicantIncome LoanAmount Loan_Amount_Term \
609 2900 0.0 71.0 360.0
610 4106 0.0 40.0 180.0
611 8072 240.0 253.0 360.0
612 7583 0.0 187.0 360.0
613 4583 0.0 133.0 360.0

  Credit_History Property_Area Loan_Status
609 1.0 Rural Y
610 1.0 Rural Y
611 1.0 Urban Y
612 1.0 Urban Y
613 0.0 Semiurban N

```

3) Printing the column names of the DataFrame

```
print(df.columns)
```

Output

```

Index(['Loan_ID', 'Gender', 'Married', 'Dependents', 'Education',
      'Self_Employed', 'ApplicantIncome', 'CoapplicantIncome', 'LoanAmount',
      'Loan_Amount_Term', 'Credit_History', 'Property_Area', 'Loan_Status'],
      dtype='object')

```

4) Summary of Data Frame

```
print(df.info())
```

Output

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 614 entries, 0 to 613
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Loan_ID               614 non-null   object
1   Gender                601 non-null   object
2   Married               611 non-null   object
3   Dependents            599 non-null   object
4   Education             614 non-null   object
5   Self_Employed         582 non-null   object
6   ApplicantIncome       614 non-null   int64
7   CoapplicantIncome     614 non-null   float64
8   LoanAmount            592 non-null   float64
9   Loan_Amount_Term      600 non-null   float64
10  Credit_History        564 non-null   float64
11  Property_Area         614 non-null   object
12  Loan_Status           614 non-null   object
dtypes: float64(4), int64(1), object(8)
memory usage: 62.5+ KB
None
```

5) Descriptive Statistical Measures of a DataFrame

```
print(df.describe())
```

Output

```

ApplicantIncome  CoapplicantIncome  LoanAmount  Loan_Amount_Term \
count      614.000000      614.000000  592.000000      600.000000
mean       5403.459283      1621.245798  146.412162      342.000000
std        6109.041673      2926.248369   85.587325       65.12041
min         150.000000         0.000000    9.000000       12.000000
25%        2877.500000         0.000000   100.000000      360.000000
50%        3812.500000      1188.500000   128.000000      360.000000
75%        5795.000000      2297.250000   168.000000      360.000000
max        81000.000000    41667.000000  700.000000      480.000000

Credit_History
count      564.000000
mean         0.842199
std          0.364878
min          0.000000
25%          1.000000
50%          1.000000
75%          1.000000
max          1.000000
```

6) Missing Data Handling

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

#Dropping rows with missing values
df_cleaned = df.dropna()
```

Output

```

Loan_ID      0
Gender       13
Married      3
Dependents   15
Education    0
Self_Employed 32
ApplicantIncome 0
CoapplicantIncome 0
LoanAmount   22
Loan_Amount_Term 14
Credit_History 50
Property_Area 0
Loan_Status  0
dtype: int64

```

7) Sorting DataFrame values

```

df_sorted = df.sort_values(by='LoanAmount', ascending=False)
print(df_sorted.head())

```

Output

| | Loan_ID | Gender | Married | Dependents | Education | Self_Employed | \ |
|-----|----------|--------|---------|------------|-----------|---------------|---|
| 171 | LP001585 | NaN | Yes | 3+ | Graduate | No | |
| 130 | LP001469 | Male | No | 0 | Graduate | Yes | |
| 561 | LP002813 | Female | Yes | 1 | Graduate | Yes | |
| 155 | LP001536 | Male | Yes | 3+ | Graduate | No | |
| 369 | LP002191 | Male | Yes | 0 | Graduate | No | |

| | ApplicantIncome | CoapplicantIncome | LoanAmount | Loan_Amount_Term | \ |
|-----|-----------------|-------------------|------------|------------------|---|
| 171 | 51763 | 0.0 | 700.0 | 300.0 | |
| 130 | 20166 | 0.0 | 650.0 | 480.0 | |
| 561 | 19484 | 0.0 | 600.0 | 360.0 | |
| 155 | 39999 | 0.0 | 600.0 | 180.0 | |
| 369 | 19730 | 5266.0 | 570.0 | 360.0 | |

| | Credit_History | Property_Area | Loan_Status |
|-----|----------------|---------------|-------------|
| 171 | 1.0 | Urban | Y |
| 130 | NaN | Urban | Y |
| 561 | 1.0 | Semiurban | Y |
| 155 | 0.0 | Semiurban | Y |
| 369 | 1.0 | Rural | N |

8) Merge Data Frames

```

df1 = pd.read_csv(file_path)
df2 = pd.read_csv(file_path)
#merge data
df = pd.merge(df1, df2)
print(df)

```

Output

| | Loan_ID | Gender | Married | Dependents | Education | Self_Employed | \ |
|-----|----------|--------|---------|------------|--------------|---------------|---|
| 0 | LP001002 | Male | No | 0 | Graduate | No | |
| 1 | LP001003 | Male | Yes | 1 | Graduate | No | |
| 2 | LP001005 | Male | Yes | 0 | Graduate | Yes | |
| 3 | LP001006 | Male | Yes | 0 | Not Graduate | No | |
| 4 | LP001008 | Male | No | 0 | Graduate | No | |
| ... | ... | ... | ... | ... | ... | ... | |
| 609 | LP002978 | Female | No | 0 | Graduate | No | |
| 610 | LP002979 | Male | Yes | 3+ | Graduate | No | |
| 611 | LP002983 | Male | Yes | 1 | Graduate | No | |
| 612 | LP002984 | Male | Yes | 2 | Graduate | No | |
| 613 | LP002990 | Female | No | 0 | Graduate | Yes | |

| | ApplicantIncome | CoapplicantIncome | LoanAmount | Loan_Amount_Term | \ |
|-----|-----------------|-------------------|------------|------------------|---|
| 0 | 5849 | 0.0 | NaN | 360.0 | |
| 1 | 4583 | 1500.0 | 120.0 | 360.0 | |
| 2 | 3000 | 0.0 | 66.0 | 360.0 | |
| 3 | 2503 | 2350.0 | 120.0 | 360.0 | |
| 4 | 6000 | 0.0 | 141.0 | 360.0 | |
| ... | ... | ... | ... | ... | |
| 609 | 2900 | 0.0 | 71.0 | 360.0 | |
| 610 | 4106 | 0.0 | 40.0 | 180.0 | |
| 611 | 8072 | 240.0 | 253.0 | 360.0 | |
| 612 | 7583 | 0.0 | 187.0 | 360.0 | |
| 613 | 4583 | 0.0 | 133.0 | 360.0 | |

| | Credit_History | Property_Area | Loan_Status |
|-----|----------------|---------------|-------------|
| 0 | 1.0 | Urban | Y |
| 1 | 1.0 | Rural | N |
| 2 | 1.0 | Urban | Y |
| 3 | 1.0 | Urban | Y |
| 4 | 1.0 | Urban | Y |
| ... | ... | ... | ... |
| 609 | 1.0 | Rural | Y |
| 610 | 1.0 | Rural | Y |
| 611 | 1.0 | Urban | Y |
| 612 | 1.0 | Urban | Y |
| 613 | 0.0 | Semiurban | N |

[614 rows x 13 columns]

9) Apply Function

```
def status_upper(status):
    return status.upper()

df['Loan_Status'] = df['Loan_Status'].apply(status_upper)
print(df['Loan_Status'].head())
```

Output

```
0    Y
1    N
2    Y
3    Y
4    Y
Name: Loan_Status, dtype: object
```

10) By using the lambda operator

```
df['outputcolumn'] = df['LoanAmount'].apply(lambda x:x/10)
df['appcolumn'] = df['ApplicantIncome'].apply(lambda x:x/10)
df.head()
```

Output

| | Loan_ID | Gender | Married | Dependents | Education | Self_Employed | ApplicantIncome | CoapplicantIncome | LoanAmount | Loan_Amount_Term | Credit_His1 |
|---|----------|--------|---------|------------|--------------|---------------|-----------------|-------------------|------------|------------------|-------------|
| 0 | LP001002 | Male | No | 0 | Graduate | No | 5849 | 0.0 | NaN | 360.0 | |
| 1 | LP001003 | Male | Yes | 1 | Graduate | No | 4583 | 1508.0 | 128.0 | 360.0 | |
| 2 | LP001005 | Male | Yes | 0 | Graduate | Yes | 3000 | 0.0 | 66.0 | 360.0 | |
| 3 | LP001006 | Male | Yes | 0 | Not Graduate | No | 2583 | 2358.0 | 120.0 | 360.0 | |
| 4 | LP001008 | Male | No | 0 | Graduate | No | 6000 | 0.0 | 141.0 | 360.0 | |

| Credit_History | Property_Area | Loan_Status | outputcolumn | appcolumn |
|----------------|---------------|-------------|--------------|-----------|
| 1.0 | Urban | Y | NaN | 584.9 |
| 1.0 | Rural | N | 12.8 | 458.3 |
| 1.0 | Urban | Y | 6.6 | 300.0 |
| 1.0 | Urban | Y | 12.0 | 258.3 |
| 1.0 | Urban | Y | 14.1 | 600.0 |

11) Visualizing DataFrame

```
import matplotlib.pyplot as plt
import seaborn as sns
plt.style.use('fivethirtyeight')
plt.rcParams['figure.figsize'] = (15,5)
plt.subplot(1,3,1)
sns.boxplot(df['ApplicantIncome'])
plt.subplot(1,3,2)
sns.boxplot(df['CoapplicantIncome'])
plt.subplot(1,3,3)
sns.boxplot(df['LoanAmount'])
plt.suptitle('Outliers Detection')
plt.show()
```

Output

