

```


1 # ziad aburas group A
2 # https://github.com/ziadaburas/study/tree/main/3B/DM
3 import pandas as pd
4 path="/content/sample_data/_Mall Customers toCleanAssignment (1).csv"
5 dataset = pd.read_csv(path,encoding_errors='ignore')
6 dataset.rename(columns={"Spending Score (1-100)": "Spending_Score_1_100", "Annual
7

```



```

1 df1 = dataset.drop('CustomerID',axis=1)
2 df1.head()

```



| | Gender | Age | Annual_Income_k_Dollar | Spending_Score_1_100 |
|---|--------|------|------------------------|----------------------|
| 0 | Male | 19.0 | 15.0 | 39.0 |
| 1 | 1 | 19.0 | NaN | 39.0 |
| 2 | Female | 20.0 | 16.0 | 6.0 |
| 3 | Male | NaN | 15.0 | 39.0 |
| 4 | Female | NaN | 16.0 | 77.0 |

Next steps:

[Generate code with df1](#)[View recommended plots](#)[New interactive sheet](#)

```

1 df2=df1.copy()
2 df2.Gender=df2.Gender.astype(str)
3 df2.Gender.unique()
4
5 df2.loc[df2.Gender=='1','Gender']='Male'
6 df2.loc[df2.Gender=='0','Gender']='Female'
7 df2.Gender.unique()
8
9 df2.where((df2.Gender=='Male')| (df2.Gender=='Female'),other=None,inplace=True)
10 df2.Gender.unique()
11
12 df2.drop(index=df2[df2.Gender.isnull()].index,inplace=True)
13 df2.Gender.unique()
14

```



```
array(['Male', 'Female'], dtype=object)
```

```

1 df3=df2.copy()
2 df3.describe()
3
4 # df3.Spending_Score_1_100.fillna(df3.Spending_Score_1_100.mean(),inplace=True)
5 df3.Spending_Score_1_100=df3.Spending_Score_1_100.fillna(df3.Spending_Score_1_
6 df3.isnull().sum()

```

```

7
8 df3.Spending_Score_1_100[df3.Spending_Score_1_100 <1]
9 df3.Spending_Score_1_100[df3.Spending_Score_1_100 >100]
10 #
11 # df3.Spending_Score_1_100[df3.Spending_Score_1_100 <1]=abs(df3.Spending_Score_1_100)
12 # df3.Spending_Score_1_100[df3.Spending_Score_1_100 >100]=df3.Spending_Score_1_100
13 df3.loc[df3.Spending_Score_1_100 <1, 'Spending_Score_1_100']=abs(df3.Spending_Score_1_100)
14 df3.loc[df3.Spending_Score_1_100 >100, 'Spending_Score_1_100']=df3.Spending_Score_1_100
15 df3.describe()
16
17

```



| | Age | Annual_Income_k_Dollar | Spending_Score_1_100 |
|--------------|------------|------------------------|----------------------|
| count | 841.000000 | 835.000000 | 989.000000 |
| mean | 37.353151 | 68.415569 | 49.691770 |
| std | 18.468252 | 34.212392 | 23.042389 |
| min | -53.000000 | -43.000000 | 3.000000 |
| 25% | 23.000000 | 40.000000 | 36.000000 |
| 50% | 35.000000 | 65.000000 | 49.798949 |
| 75% | 54.000000 | 97.000000 | 65.000000 |
| max | 70.000000 | 137.000000 | 99.000000 |

```

1 df4=df3.copy()
2 df4.describe()
3 # df4.Annual_Income_k_Dollar.fillna(df4.Annual_Income_k_Dollar.mean(),inplace=True)
4 df4.Annual_Income_k_Dollar=df4.Annual_Income_k_Dollar.fillna(df4.Annual_Income_k_Dollar.mean())
5 df4.isnull().sum()
6
7 df4.Annual_Income_k_Dollar[df4.Annual_Income_k_Dollar <1]
8 # df4.Annual_Income_k_Dollar[df4.Annual_Income_k_Dollar <1]=abs(df4.Annual_Income_k_Dollar)
9 df4.loc[df4.Annual_Income_k_Dollar <1, 'Annual_Income_k_Dollar']=abs(df4.Annual_Income_k_Dollar)
10 df4.describe()

```



| | Age | Annual_Income_k_Dollar | Spending_Score_1_100 |
|--------------|------------|------------------------|----------------------|
| count | 841.000000 | 989.000000 | 989.000000 |
| mean | 37.353151 | 68.587460 | 49.691770 |
| std | 18.468252 | 31.055916 | 23.042389 |
| min | -53.000000 | 15.000000 | 3.000000 |
| 25% | 23.000000 | 45.000000 | 36.000000 |
| 50% | 35.000000 | 68.415569 | 49.798949 |

| | | | |
|------------|-----------|------------|-----------|
| 50% | 55.000000 | 55.713333 | 75.750000 |
| 75% | 54.000000 | 89.000000 | 65.000000 |
| max | 70.000000 | 137.000000 | 99.000000 |

```

1 df5=df4.copy()
2 df5.isnull().sum()
3 df5 = df5.dropna()
4 df5.isnull().sum()
5
6 df5.Age[df5.Age <1]
7 df5.loc[df5.Age <1, 'Age']=abs(df5.Age[df5.Age <1])
8
9 df5.describe()
10 df5[df5.Age < 15].count()
11
12 # Drop if age <15 because they are children not working
13 df5.drop(index=df5[df5.Age < 15].index,inplace=True)
14 df5[df5.Age < 15].count()
15
16 df5.describe()
17
18 df5.Age=df5.Age.astype(int)
19 print(df5.Age.dtype)
20 df5.info()

```

```

int64
<class 'pandas.core.frame.DataFrame'>
Index: 752 entries, 0 to 997
Data columns (total 4 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Gender                                752 non-null    object
1   Age                                   752 non-null    int64
2   Annual_Income_k_Dollar               752 non-null    float64
3   Spending_Score_1_100                 752 non-null    float64
dtypes: float64(2), int64(1), object(1)
memory usage: 29.4+ KB

```

```

1 df5.head(20)
2 df5.to_csv("/content/drive/MyDrive/Colab Notebooks/Mall_Customers_toCleanAssignm

```

