

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: df=pd.read_excel(r"C:\Users\UP\Desktop\Q1walli Sales Data project.xlsx")

Out[2]:
   User_ID  Cust_name  Product_ID  Gender  Age Group  Age  Marital_Status  State  Zone  Occupation  Product_Category  Orders  Amount  Status  unnamed1
0  1002903  Sarskriti  P00125942  F  26-35  28  0  Maharashtra  Western  Healthcare  Auto  1  23952.0  NaN  NaN
1  1000732  Karik  P00110942  F  26-35  35  1  Andhra Pradesh  Southern  Govt  Auto  3  23934.0  NaN  NaN
2  1001990  Bindu  P00118542  F  26-35  35  1  Uttar Pradesh  Central  Automobile  Auto  3  23924.0  NaN  NaN
3  1001425  Sudevi  P00237842  M  0-17  16  0  Karnataka  Southern  Construction  Auto  2  23912.0  NaN  NaN
4  1000568  Joni  P00057942  M  26-35  28  1  Gujarat  Western  Food Processing  Auto  2  23877.0  NaN  NaN
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
11246  1000695  Manning  P00226942  M  18-25  19  1  Maharashtra  Western  Chemical  Office  4  370.0  NaN  NaN
11247  1004089  Reichenbach  P00171342  M  26-35  33  0  Haryana  Northern  Healthcare  Veterinary  3  367.0  NaN  NaN
11248  1001209  Oshin  P00201342  F  36-45  40  0  Madhya Pradesh  Central  Textile  Office  4  213.0  NaN  NaN
11249  1004023  Nooran  P00208442  M  36-45  37  0  Karnataka  Southern  Agriculture  Office  3  206.0  NaN  NaN
11250  1002744  Brunney  P00281742  F  18-25  19  0  Maharashtra  Western  Healthcare  Office  3  188.0  NaN  NaN

11251 rows x 15 columns

In [3]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11251 entries, 0 to 11250
Data columns (total 15 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   User_ID               11251 non-null  int64
 1   Cust_name            11251 non-null  object
 2   Product_ID           11251 non-null  object
 3   Gender               11251 non-null  object
 4   Age Group            11251 non-null  object
 5   Age                  11251 non-null  int64
 6   Marital_Status       11251 non-null  object
 7   State                11251 non-null  object
 8   Zone                 11251 non-null  object
 9   Occupation            11251 non-null  object
10  Product_Category     11251 non-null  object
11  Orders               11251 non-null  int64
12  Amount               11239 non-null  float64
13  Status               1189 non-null  float64
14  unnamed1             0 non-null     float64
dtypes: float64(3), int64(4), object(8)
memory usage: 3.3+ MB

In [4]: df.shape

(11251, 15)

Out[4]: df.head()

In [5]: df.info()

Out[5]:
   User_ID  Cust_name  Product_ID  Gender  Age Group  Age  Marital_Status  State  Zone  Occupation  Product_Category  Orders  Amount  Status  unnamed1
0  1002903  Sarskriti  P00125942  F  26-35  28  0  Maharashtra  Western  Healthcare  Auto  1  23952.0  NaN  NaN
1  1000732  Karik  P00110942  F  26-35  35  1  Uttar Pradesh  Southern  Govt  Auto  3  23934.0  NaN  NaN
2  1001990  Bindu  P00118542  F  26-35  35  1  Uttar Pradesh  Central  Automobile  Auto  3  23924.0  NaN  NaN
3  1001425  Sudevi  P00237842  M  0-17  16  0  Karnataka  Southern  Construction  Auto  2  23912.0  NaN  NaN
4  1000568  Joni  P00057942  M  26-35  28  1  Gujarat  Western  Food Processing  Auto  2  23877.0  NaN  NaN
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
11246  1000695  Manning  P00226942  M  18-25  19  1  Maharashtra  Western  Chemical  Office  4  370.0  NaN  NaN
11247  1004089  Reichenbach  P00171342  M  26-35  33  0  Haryana  Northern  Healthcare  Veterinary  3  367.0  NaN  NaN
11248  1001209  Oshin  P00201342  F  36-45  40  0  Madhya Pradesh  Central  Textile  Office  4  213.0  NaN  NaN
11249  1004023  Nooran  P00208442  M  36-45  37  0  Karnataka  Southern  Agriculture  Office  3  206.0  NaN  NaN
11250  1002744  Brunney  P00281742  F  18-25  19  0  Maharashtra  Western  Healthcare  Office  3  188.0  NaN  NaN

11251 rows x 13 columns

In [6]: #DATA CLEANING
#check null value
df.isnull().sum()

Out[6]:
User_ID      0
Cust_name    0
Product_ID    0
Gender        0
Age Group     0
Age           0
Marital_Status  0
State         0
Zone          0
Occupation    0
Product_Category  0
Orders        0
Amount        0
Status        0
dtype: int64

In [8]: #fill null value by mean
df["Amount"]=df["Amount"].fillna(df["Amount"].mean())

In [9]: #checking null value is fill
df.isnull().sum()

Out[9]:
User_ID      0
Cust_name    0
Product_ID    0
Gender        0
Age Group     0
Age           0
Marital_Status  0
State         0
Zone          0
Occupation    0
Product_Category  0
Orders        0
Amount        0
dtype: int64

In [10]: df

Out[11]:
   User_ID  Cust_name  Product_ID  Gender  Age Group  Age  Marital_Status  State  Zone  Occupation  Product_Category  Orders  Amount
0  1002903  Sarskriti  P00125942  F  26-35  28  0  Maharashtra  Western  Healthcare  Auto  1  23952.0
1  1000732  Karik  P00110942  F  26-35  35  1  Andhra Pradesh  Southern  Govt  Auto  3  23934.0
2  1001990  Bindu  P00118542  F  26-35  35  1  Uttar Pradesh  Central  Automobile  Auto  3  23924.0
3  1001425  Sudevi  P00237842  M  0-17  16  0  Karnataka  Southern  Construction  Auto  2  23912.0
4  1000568  Joni  P00057942  M  26-35  28  1  Gujarat  Western  Food Processing  Auto  2  23877.0
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
11246  1000695  Manning  P00226942  M  18-25  19  1  Maharashtra  Western  Chemical  Office  4  370.0
11247  1004089  Reichenbach  P00171342  M  26-35  33  0  Haryana  Northern  Healthcare  Veterinary  3  367.0
11248  1001209  Oshin  P00201342  F  36-45  40  0  Madhya Pradesh  Central  Textile  Office  4  213.0
11249  1004023  Nooran  P00208442  M  36-45  37  0  Karnataka  Southern  Agriculture  Office  3  206.0
11250  1002744  Brunney  P00281742  F  18-25  19  0  Maharashtra  Western  Healthcare  Office  3  188.0

11251 rows x 13 columns

In [12]: #Change data type
df["Amount"]=df["Amount"].astype('int')

In [13]: df["Amount"].dtype

dtype('int32')

Out[13]:

In [14]: #to check column name
df.columns

Out[14]:
Index(['User_ID', 'Cust_name', 'Product_ID', 'Gender', 'Age Group', 'Age', 'Marital_Status', 'State', 'Zone', 'Occupation', 'Product_Category', 'Orders', 'Amount'],
      dtype='object')

In [15]: df=df.rename(columns={'Gender':'Sex'})
df

Out[15]:
   User_ID  Cust_name  Product_ID  Sex  Age Group  Age  Marital_Status  State  Zone  Occupation  Product_Category  Orders  Amount
0  1002903  Sarskriti  P00125942  F  26-35  28  0  Maharashtra  Western  Healthcare  Auto  1  23952
1  1000732  Karik  P00110942  F  26-35  35  1  Andhra Pradesh  Southern  Govt  Auto  3  23934.0
2  1001990  Bindu  P00118542  F  26-35  35  1  Uttar Pradesh  Central  Automobile  Auto  3  23924
3  1001425  Sudevi  P00237842  M  0-17  16  0  Karnataka  Southern  Construction  Auto  2  23912
4  1000568  Joni  P00057942  M  26-35  28  1  Gujarat  Western  Food Processing  Auto  2  23877
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
11246  1000695  Manning  P00226942  M  18-25  19  1  Maharashtra  Western  Chemical  Office  4  370
11247  1004089  Reichenbach  P00171342  M  26-35  33  0  Haryana  Northern  Healthcare  Veterinary  3  367
11248  1001209  Oshin  P00201342  F  36-45  40  0  Madhya Pradesh  Central  Textile  Office  4  213
11249  1004023  Nooran  P00208442  M  36-45  37  0  Karnataka  Southern  Agriculture  Office  3  206
11250  1002744  Brunney  P00281742  F  18-25  19  0  Maharashtra  Western  Healthcare  Office  3  188

11251 rows x 13 columns

In [16]: #describe function return the description of the data in Data frame(i.e count,mean,std etc)
df.describe()

Out[16]:
   User_ID      Age  Marital_Status      Orders      Amount
count  1125100e+04  11251.000000  11251.000000  11251.000000
mean    1.003004e+06  35.421207  0.420318  2.489290  9453.609901
std    1.718125e+03  12.754122  0.493632  1.110047  5219.569169
min    1.000001e+06  12.000000  0.000000  1.000000  388.000000
25%    1.001400e+06  21.000000  0.000000  1.500000  5443.500000
50%    1.002656e+06  33.000000  0.000000  2.000000  8110.000000
75%    1.004430e+06  43.000000  1.000000  3.000000  12671.000000
max    1.006040e+06  92.000000  1.000000  4.000000  23952.000000

In [17]: #describe function for specific column
df[["Age", "Orders", "Amount"]].describe()

Out[17]:
   Age      Orders      Amount
count  11251.000000  11251.000000  11251.000000
mean    35.421207  0.489290  9453.609901
std    12.754122  0.115047  5219.569169
min    12.000000  1.000000  388.000000
25%    21.000000  1.500000  5443.500000
50%    33.000000  2.000000  8110.000000
75%    43.000000  3.000000  12671.000000
max    92.000000  4.000000  23952.000000

In [18]: #sns.countplot(x='Sex',data=df1)
#for data label
for bar in ax.containers:
    ax.bar_label(bar)

Out[18]:
count
8000
7000
6000
5000
4000
3000
2000
1000
0
F M
Sex
7842 3409

In [19]: data=df1[["Sex","Amount"]]
data

Out[19]:
   Sex  Amount
0  F  23952
1  F  23934
2  F  23924
3  M  23912
4  M  23877
...  ...
11246  M  370
11247  M  367
11248  F  213
11249  M  206
11250  F  188

11251 rows x 2 columns

In [20]: data.groupby('Sex').sum()

Out[20]:
   Sex  Amount
F  7443393
M  3193282

In [21]: sns=df1.groupby(['Sex'],as_index=False)[['Amount']].sum().sort_values(by='Amount',ascending=False)
ss

Out[21]:
   Sex  Amount
0  F  7443393
1  M  3193282

In [22]: sns.barplot(x='Sex',y='Amount',data=ss)

Out[22]:
<Axes: xlabel='Sex', ylabel='Amount'>

In [23]:
fig,ax=plt.subplots()
sns.barplot(x='Sex',y='Amount',data=ss)
ax.bar_label(ax.containers[0].bars)

Out[23]:
count
3000
2500
2000
1500
1000
500
0
26-35 0-17 18-25 46-50 51-55 36-45
Age Group
Sex
F M

In [24]: #Total Amount vs Age Group
data=sales_df1.groupby(['Age Group'],as_index=False)[['Amount']].sum().sort_values(by='Amount',ascending=False)

Out[24]:
   Age Group  Amount
0  26-35  4263248
1  36-45  2217353
3  18-25  1274732
4  46-50  824566
5  51-55  628283
6  55+  499440
0  0-17  269963

In [25]: sns.barplot(x='Age Group',y='Amount',data=sales_age)

Out[25]:
<Axes: xlabel='Age Group', ylabel='Amount'>

In [26]:
fig,ax=plt.subplots()
sns.barplot(x='Age Group',y='Amount',data=sales_age)
ax.bar_label(ax.containers[0].bars)

Out[26]:
count
4.0
3.5
3.0
2.5
2.0
1.5
1.0
0.5
0.0
26-35 36-45 18-25 46-50 51-55 55+ 0-17
Age Group
Amount

In [27]: #Total no. of order from 10 states
statewise=df1[["State","Orders"]]
statewise

Out[27]:
   State  Orders
0  Maharashtra  1
1  Andhra Pradesh  3
2  Uttar Pradesh  3
3  Karnataka  2
4  Gujarat  2
...  ...
11246  Maharashtra  4
11247  Haryana  3
11248  Madhya Pradesh  4
11249  Karnataka  3
11250  Maharashtra  3

11251 rows x 2 columns

In [27]: statewise_order=df1.groupby(['State'],as_index=False)[['Orders']].sum().sort_values(by='Orders',ascending=False)[0:10]

Out[27]:
   State  Orders
10  Uttar Pradesh  4613
14  Maharashtra  3811
7  Karnataka  2741
2  Delhi  2744
9  Madhya Pradesh  2259
0  Andhra Pradesh  2054
5  Himachal Pradesh  1568
8  Kerala  1137
4  Haryana  1109
3  Gujarat  1070

In [28]: sns.set(rc={'figure.figsize':(15,5)})
sns.barplot(x='State',y='Orders',data=statewise_order)

Out[28]:
<Axes: xlabel='State', ylabel='Orders'>

In [29]:
fig,ax=plt.subplots()
sns.barplot(x='State',y='Orders',data=statewise_order)
ax.bar_label(ax.containers[0].bars)

Out[29]:
count
5000
4000
3000
2000
1000
0
Uttar Pradesh Maharashtra Karnataka Delhi Madhya Pradesh Andhra Pradesh Himachal Pradesh Kerala Haryana Gujarat
State
Orders

In [30]: #Sales from top 10 states
top_state_sales=df1.groupby(['State'],as_index=False)[['Amount']].sum().sort_values(by='Amount',ascending=False)[0:10]

Out[30]:
   State  Amount
14  Uttar Pradesh  10392874
10  Maharashtra  1453996
7  Karnataka  13532993
2  Delhi  11632177
9  Madhya Pradesh  8120048
0  Andhra Pradesh  8046599
5  Himachal Pradesh  4963368
4  Haryana  4220175
1  Bihar  4027757
3  Gujarat  3964888

In [30]: sns.barplot(x='Amount',y='State',data=top_state_sales)

Out[30]:
<Axes: xlabel='Amount', ylabel='State'>

In [31]:
fig,ax=plt.subplots()
sns.barplot(x='Amount',y='State',data=top_state_sales)
ax.bar_label(ax.containers[0].bars)

Out[31]:
count
6000
5000
4000
3000
2000
1000
0
0 1
Marital_Status
Amount

In [32]: df1

Out[32]:
   User_ID  Cust_name  Product_ID  Sex  Age Group  Age  Marital_Status  State  Zone  Occupation  Product_Category  Orders  Amount
0  1002903  Sarskriti  P00125942  F  26-35  28  0  Maharashtra  Western  Healthcare  Auto  1  23952
1  1000732  Karik  P00110942  F  26-35  35  1  Andhra Pradesh  Southern  Govt  Auto  3  23934
2  1001990  Bindu  P00118542  F  26-35  35  1  Uttar Pradesh  Central  Automobile  Auto  3  23924
3  1001425  Sudevi  P00237842  M  0-17  16  0  Karnataka  Southern  Construction  Auto  2  23912
4  1000568  Joni  P00057942  M  26-35  28  1  Gujarat  Western  Food Processing  Auto  2  23877
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
11246  1000695  Manning  P00226942  M  18-25  19  1  Maharashtra  Western  Chemical  Office  4  370
11247  1004089  Reichenbach  P00171342  M  26-35  33  0  Haryana  Northern  Healthcare  Veterinary  3  367
11248  1001209  Oshin  P00201342  F  36-45  40  0  Madhya Pradesh  Central  Textile  Office  4  213
11249  1004023  Nooran  P00208442  M  36-45  37  0  Karnataka  Southern  Agriculture  Office  3  206
11250  1002744  Brunney  P00281742  F  18-25  19  0  Maharashtra  Western  Healthcare  Office  3  188

11251 rows x 13 columns

In [33]: #Amount according to marital status with sex
sns=df1.groupby(['Marital_Status','Sex'],as_index=False)[['Amount']].sum().sort_values(by='Amount',ascending=False)
sss

Out[33]:
   Marital_Status  Sex  Amount
0  0  F  43815005
2  1  F  30615378
1  0  M  18348191
3  1  M  13683991

In [33]: sns.barplot(x='Marital_Status',data=sss,y='Amount',hue='Sex')

Out[33]:
<Axes: xlabel='Marital_Status', ylabel='Amount'>

In [34]:
fig,ax=plt.subplots()
sns.barplot(x='Marital_Status',data=sss,y='Amount',hue='Sex')
ax.bar_label(ax.containers[0].bars)

Out[34]:
count
4
3
2
1
0
0 1
Marital_Status
Amount
Sex
F M

In [35]: sns.set(rc={'figure.figsize':(20,5)})
sns.countplot(x='Occupation',data=df1)
fig,ax=plt.subplots()
sns.countplot(x='Occupation',data=df1)
ax.bar_label(ax.containers[0].bars)

Out[35]:
count
1800
1400
1000
800
600
400
200
0
Healthcare Govt Automobile Construction Food Processing Lawyer IT Sector Media Banking Retail Hospitality Aviation Agriculture Textile Chemical
Occupation

In [36]: occupation_wise_amount=df1.groupby(['Occupation'],as_index=False)[['Amount']].sum().sort_values(by='Amount',ascending=False)

Out[36]:
   Occupation  Amount
10  IT Sector  14802344
8  Healthcare  13034086
2  Aviation  12602298
3  Banking  10789516
9  Govt  8502226
6  Hospitality  696321
12  Media  629532
1  Automobile  5378049
4  Chemical  5306889
11  Lawyer  4981665
13  Retail  4783170
5  Food Processing  4070670
14  Textile  3234425
0  Agriculture  2992067

In [37]: sns.barplot(x='Occupation',y='Amount',data=occupation_wise_amount)

Out[37]:
<Axes: xlabel='Occupation', ylabel='Amount'>

In [38]:
fig,ax=plt.subplots()
sns.barplot(x='Occupation',y='Amount',data=occupation_wise_amount)
ax.bar_label(ax.containers[0].bars)

Out[38]:
count
1.4
1.2
1.0
0.8
0.6
0.4
0.2
0.0
IT Sector Healthcare Aviation Banking Govt Hospitality Media Automobile Occupation Chemical Lawyer Retail Food Processing Construction Textile Agriculture
Amount

In [39]: sns.countplot(x='Product_Category',data=df1)

Out[39]:
count
2500
2000
1500
1000
500
0
AutoHand & Power ToolBathroomry TupperwareFood&Shoe&ShoeFurniture Food Games & ToysSports Products Book&Electronics & GadgetsClothing & ApparelBeauty Household ItemsPet Care Veterinary Office
Product_Category

In [ ]:
```