

```

1 # Filename:Lab2.py
2 # Date: 07/20/2022
3 # Qiang Liu
4 # NCSU-GTI Summer 2022 Data Science
5 #-----import library-----
6 import sys
7 import math
8 import os,io, csv
9
10 #-----Data Science library-----
11 import numpy as np
12 import matplotlib.pyplot as plt
13 import pandas as pd
14 import seaborn as sn
15
16 #-----global variables-----
17 dash = "-" * 80
18 print(dash)
19
20 def main():
21     print("Start Program....")
22     #a. Import data(VandelaySales2015.csv) into a Pandas dataframe
23     FILENAME = "VandelaySales2015.csv"
24     df = pd.read_csv(FILENAME,encoding="utf-8")
25     print('VandelaySales2015 data are following:\n',df)
26     print(dash)
27     # check if file has missing values
28     print("\n Any null values?")
29     print(df.isnull().any())
30     print(dash)
31     print(df.isnull().sum())
32     print(dash)
33     print(df.info())
34     print(dash)
35     #a. print head for top 10 record only for the dataframe call it "df"
36     print('A. Top 10 record of VandelaySales2015 are following:\n',df.head(10))
37     print(dash)
38
39     #b. Find the sum, mean, max, min value of total_product_price (COLUMN I) of VandelaySales2015.csv file.
40     print('B. The sum, mean, max, min value of total_product_price are following:\n',df['total_product_price'].describe())
41     print(dash)
42
43     #c. Show correlation
44     print('C. The correlations between variables are following:\n',df.corr())
45     print(dash)
46     # show corr in heatmap
47     dataplot = sn.heatmap(df.corr(),
48                           cmap='YlGnBu',
49                           annot=True)
50     plt.show()
51
52     #d. Show distribution analysis Histogram for item_product_price(COLUMN H) of VandelaySales2015.csv file.
53     plt.hist(df['item_product_price'],
54             bins=15,
55             color='steelblue',
56             edgecolor='k',
57             label='Histogram')
58     plt.title("D. Distribution analysis Histogram for item_product_price")
59     plt.ylabel("Frequency for item_product_price")
60     plt.xlabel("Bins")
61     plt.show()
62
63 #starting point... launch
64 # __name__ is predefined class attribute
65 if __name__ == '__main__':
66     main()

```

```
C:\Users\Liuqiang\AppData\Local\Programs\Python\Python310\python.exe C:/Users/Liuqiang/PycharmProjects/GTISummer2022/Lab2.py
```

```
-----  
Start Program.....
```

```
VandelaySales2015 data are following:
```

	order_id	customer_id	...	category_name	product_color
0	59929	11914	...	T-Shirts	Blue
1	59929	11914	...	Sweatshirts	Purple
2	59966	14644	...	Sweatshirts	Yellow
3	59973	23186	...	Jeans	Green
4	59973	23186	...	Visors	Purple
...	...	...	...	...	...
47296	31061	12976	...	T-Shirts	Yellow
47297	31061	12976	...	Hats	Green
47298	31033	1759	...	Socks	Orange
47299	31033	1759	...	T-Shirts	Yellow
47300	31033	1759	...	Sweatshirts	Purple

```
[47301 rows x 13 columns]  
-----
```

```
Any null values?
```

order_id	False
customer_id	False
order_date	False
order_short_date	False
promo_code	False
referral_source	False
product_quantity	False
item_product_price	False
total_product_price	False
product_name	False
vendor_name	False
category_name	False
product_color	False

```
dtype: bool  
-----
```

order_id	0
customer_id	0
order_date	0
order_short_date	0
promo_code	0
referral_source	0
product_quantity	0
item_product_price	0
total_product_price	0
product_name	0
vendor_name	0
category_name	0
product_color	0

```
dtype: int64  
-----
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 47301 entries, 0 to 47300
```

```
Data columns (total 13 columns):
```

#	Column	Non-Null Count	Dtype
0	order_id	47301 non-null	int64
1	customer_id	47301 non-null	int64
2	order_date	47301 non-null	object
3	order_short_date	47301 non-null	object
4	promo_code	47301 non-null	object
5	referral_source	47301 non-null	object
6	product_quantity	47301 non-null	int64
7	item_product_price	47301 non-null	float64
8	total_product_price	47301 non-null	float64
9	product_name	47301 non-null	object
10	vendor_name	47301 non-null	object
11	category_name	47301 non-null	object
12	product_color	47301 non-null	object

```
dtypes: float64(2), int64(3), object(8)
```

```
memory usage: 4.7+ MB
```

```
None  
-----
```

-----  
A. Top 10 record of VandelaySales2015 are following:

	order_id	customer_id	...	category_name	product_color
0	59929	11914	...	T-Shirts	Blue
1	59929	11914	...	Sweatshirts	Purple
2	59966	14644	...	Sweatshirts	Yellow
3	59973	23186	...	Jeans	Green
4	59973	23186	...	Visors	Purple
5	59935	7437	...	T-Shirts	Orange
6	59935	7437	...	T-Shirts	Blue
7	59980	23193	...	T-Shirts	Blue
8	59916	360	...	T-Shirts	Red
9	59916	360	...	Boxers	Orange

[10 rows x 13 columns]

-----  
B. The sum, mean, max, min value of total\_product\_price are following:

count	47301.000000
mean	91.612787
std	100.521693
min	13.530000
25%	25.720000
50%	57.390000
75%	124.320000
max	893.280000

Name: total\_product\_price, dtype: float64

-----  
C. The correlations between variables are following:

	order_id	...	total_product_price
order_id	1.000000	...	-0.013055
customer_id	0.218983	...	0.360153
product_quantity	-0.010340	...	0.792602
item_product_price	-0.003488	...	0.465881
total_product_price	-0.013055	...	1.000000

[5 rows x 5 columns]

