



Darshan
UNIVERSITY

(<https://www.darshan.ac.in/>)

Data Mining

Lab - 4

Step 1. Import the necessary libraries

```
In [4]: import pandas as pd  
import numpy as np
```

Step 2. Import the dataset from this [address](#)

(<https://raw.githubusercontent.com/justmarkham/DAT8/master/data/ch>



Step 3. Assign it to a variable called chipo.

```
In [10]: df= pd.read_csv("https://raw.githubusercontent.com/justmarkham/DAT8/master/dat
```

Step 4. See the first 10 entries

In [11]: `df.head(10)`

Out[11]:

	order_id	quantity	item_name	choice_description	item_price
0	1	1	Chips and Fresh Tomato Salsa	NaN	\$2.39
1	1	1	Izze	[Clementine]	\$3.39
2	1	1	Nantucket Nectar	[Apple]	\$3.39
3	1	1	Chips and Tomatillo-Green Chili Salsa	NaN	\$2.39
4	2	2	Chicken Bowl	[Tomatillo-Red Chili Salsa (Hot), [Black Beans...	\$16.98
5	3	1	Chicken Bowl	[Fresh Tomato Salsa (Mild), [Rice, Cheese, Sou...	\$10.98
6	3	1	Side of Chips	NaN	\$1.69
7	4	1	Steak Burrito	[Tomatillo Red Chili Salsa, [Fajita Vegetables...	\$11.75
8	4	1	Steak Soft Tacos	[Tomatillo Green Chili Salsa, [Pinto Beans, Ch...	\$9.25
9	5	1	Steak Burrito	[Fresh Tomato Salsa, [Rice, Black Beans, Pinto...	\$9.25

Step 5. What is the number of observations in the dataset?

In [14]: `# Solution 1
len(df)`

Out[14]: 4622

In [15]: `# Solution 2
df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4622 entries, 0 to 4621
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  -
0   order_id              4622 non-null   int64
1   quantity              4622 non-null   int64
2   item_name             4622 non-null   object
3   choice_description     3376 non-null   object
4   item_price            4622 non-null   object
dtypes: int64(2), object(3)
memory usage: 180.7+ KB
```

Step 6. What is the number of columns in the dataset?

In [19]: `df.shape[1]`

Out[19]: 5

Step 7. Print the name of all the columns.

In [22]: `df.columns`

Out[22]: Index(['order_id', 'quantity', 'item_name', 'choice_description',
 'item_price'],
 dtype='object')

Step 8. How is the dataset indexed?

In [23]: `df.index`

Out[23]: RangeIndex(start=0, stop=4622, step=1)

Step 9. Number of Unique Items ?

In [34]: `df["item_name"].nunique()`

Out[34]: 50

Step 10. Which was the most-ordered item?

In [41]: `c = df.groupby('item_name');
c1=c.sum(numeric_only = True)
c2=c1.sort_values('quantity').tail(1)
c2`

Out[41]:

	order_id	quantity
item_name		
Chicken Bowl	713926	761

Step 11. How many items were ordered in total?

In [47]: `df['quantity'].sum()`

Out[47]: 4972

Step 12. Turn the item price into a float

Step 12.a. Check the item price type

```
In [48]: df['item_price'].dtypes
```

```
Out[48]: dtype('O')
```

Step 12.b. Create a lambda function and change the type of item price

```
In [51]: l1 = lambda x: float(x[1:])  
df['item_price'] = df['item_price'].apply(l1)
```

Step 12.c. Check the item price type

```
In [52]: df['item_price'].dtypes
```

```
Out[52]: dtype('float64')
```

Step 14. How much was the revenue for the period in the dataset?

```
In [63]: df['revenue'] = df['item_price'] * df['quantity']  
d1 = df['revenue'].sum()  
d1
```

```
Out[63]: 39237.02
```

Step 15. How many orders were made ?

```
In [57]: df['order_id'].nunique()
```

```
Out[57]: 1834
```

Step 17. How many different choice descriptions are there?

```
In [54]: df['choice_description'].nunique()
```

```
Out[54]: 1043
```

Step 18. What items have been ordered more than 100 times?

```
In [60]: c6=df.groupby('item_name')
c7 =c6['quantity'].sum()
c7[c7>100]
```

```
Out[60]: item_name
Bottled Water                211
Canned Soda                  126
Canned Soft Drink           351
Chicken Bowl                 761
Chicken Burrito              591
Chicken Salad Bowl          123
Chicken Soft Tacos           120
Chips                        230
Chips and Fresh Tomato Salsa 130
Chips and Guacamole          506
Side of Chips                 110
Steak Bowl                   221
Steak Burrito                386
Name: quantity, dtype: int64
```

Step 19. What is the average revenue amount per order?

```
In [66]: c6=df['order_id'].nunique()
c7 = df['revenue'].sum()
print(c7/c6)
```

```
21.39423118865867
```

In [64]: df

Out[64]:

	order_id	quantity	item_name	choice_description	item_price	revenue
0	1	1	Chips and Fresh Tomato Salsa	NaN	2.39	2.39
1	1	1	Izze	[Clementine]	3.39	3.39
2	1	1	Nantucket Nectar	[Apple]	3.39	3.39
3	1	1	Chips and Tomatillo-Green Chili Salsa	NaN	2.39	2.39
4	2	2	Chicken Bowl	[Tomatillo-Red Chili Salsa (Hot), [Black Beans...	16.98	33.96
...
4617	1833	1	Steak Burrito	[Fresh Tomato Salsa, [Rice, Black Beans, Sour ...	11.75	11.75
4618	1833	1	Steak Burrito	[Fresh Tomato Salsa, [Rice, Sour Cream, Cheese...	11.75	11.75
4619	1834	1	Chicken Salad Bowl	[Fresh Tomato Salsa, [Fajita Vegetables, Pinto...	11.25	11.25
4620	1834	1	Chicken Salad Bowl	[Fresh Tomato Salsa, [Fajita Vegetables, Lettu...	8.75	8.75
4621	1834	1	Chicken Salad Bowl	[Fresh Tomato Salsa, [Fajita Vegetables, Pinto...	8.75	8.75

4622 rows × 6 columns

In []: *# Solution 2*
P

Out[32]: 21.394231188658654

In []:

In []:

In []:

In []: