

Instacart EDA 4 Assignment [ANSWER]

In this assignment you will answer to four individual questions (Questions 0,1,2,3). You can complete any question in any order. **The produced results from one question will not be used on another.**

When you complete this assignment keep this kernel as private, make a succesful commit and add as collaborator the [ISLAB](#).

Note that you will be assessed based on the final version submitted before the deadline.

To answer the questions of this assignment you may create as many code blocks as you wish. In addition, you can use `head()` or any other method to display your results. **You can use comments to describe the rationale of your solution.** You can also use comments when you use new methods (that we haven't used in our notebooks), new arguments or alternative approaches for the same problem.

In the case where a code block takes more than 3 minutes to execute, hit the stop button on the left bottom corner and check for any mistakes on your code or your solution in general.

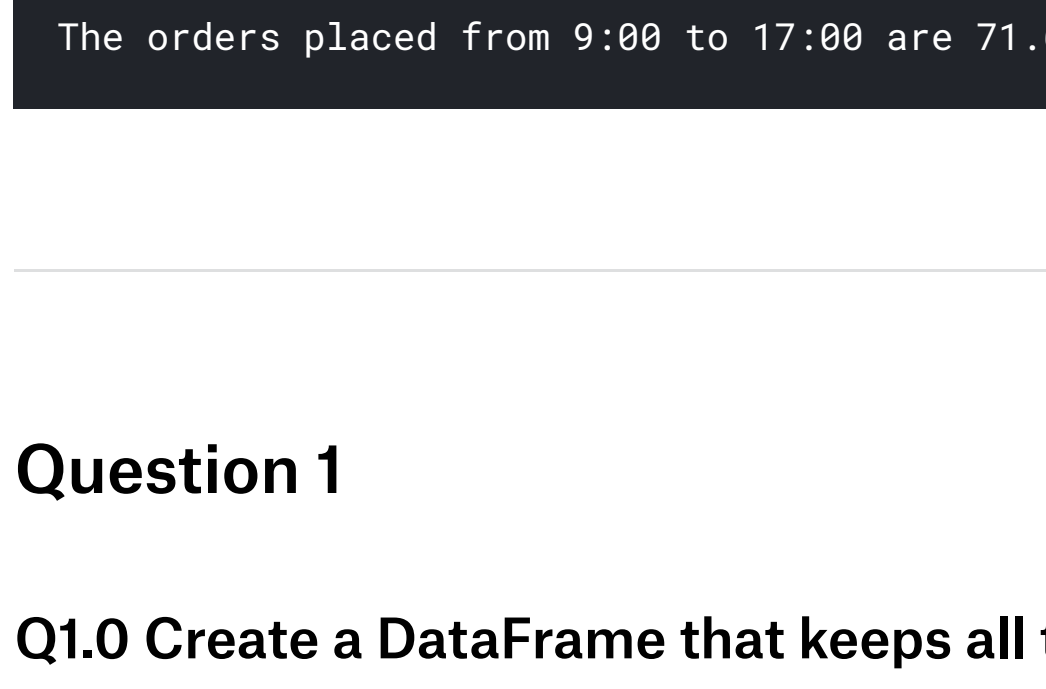
First load the requested packages and data files for this assignment:

```
In [1]:  
  
#load packages  
import pandas as pd          # for data manipulation  
import matplotlib.pyplot as plt # for plotting  
import seaborn as sns        # an extension of matplotlib for statistical graphics  
  
#load data  
orders = pd.read_csv('../input/orders.csv')  
products = pd.read_csv('../input/products.csv')  
order_products_prior = pd.read_csv('../input/order_products__prior.csv')  
aisles = pd.read_csv('../input/aisles.csv')
```

Question 0

Q0.0 Create a plot which shows how many orders were placed in each hour of the day.

- use `groupby()` method

```
In [2]:  
  
orders.groupby('order_hour_of_day')[['order_id']].count().plot.bar()  
  
Out[2]:  
  
<matplotlib.axes._subplots.AxesSubplot at 0x7fe08ecf01d0>  
  

```

Q0.1 What is the percentage allocation of the orders placed from 9 to 5 (09:00-17:00 inclusive) compared to all orders?

- store the final percentage on a variable with name 'pct'
- print the result with the appropriate text

```
In [3]:  
  
orders_hour = orders.groupby('order_hour_of_day')[['order_id']].count()  
orders_hour.columns = ['volume']  
orders_hour = orders_hour.reset_index()  
pct = orders_hour[(orders_hour.order_hour_of_day >= 9) & (orders_hour.order_hour_of_day <= 17)].volume.sum() / orders_hour.volume.sum() * 100  
print('The orders placed from 9:00 to 17:00 are ' + str(round(pct,2)) + '% of the total orders')  
  
The orders placed from 9:00 to 17:00 are 71.61% of the total orders
```

Question 1

Q1.0 Create a DataFrame that keeps all the prior orders and the products that have been purchased

```
In [4]:  
  
prd = pd.merge(orders, order_products_prior, on='order_id', how='inner')
```

Q1.1 For each product find its average position in the cart (in the orders of all customers).

```
In [5]:  
  
item_avg = prd.groupby('product_id')[['add_to_cart_order']].mean()  
item_avg.head()  
  
Out[5]:  
  


|            | add_to_cart_order |
|------------|-------------------|
| product_id |                   |
| 1          | 5.801836          |
| 2          | 9.888889          |
| 3          | 6.415162          |
| 4          | 9.507599          |
| 5          | 6.466667          |


```

Q1.2 For a given customer, find the average position of a product on its orders

```
In [6]:  
  
position = prd.groupby(['user_id', 'product_id'])[['add_to_cart_order']].mean()  
position.columns = ['mean_pos_cart']  
position.head()  
  
Out[6]:  
  


|         |            | mean_pos_cart |
|---------|------------|---------------|
| user_id | product_id |               |
| 1       | 196        | 1.400000      |
|         | 10258      | 3.333333      |
|         | 10326      | 5.000000      |
|         | 12427      | 3.300000      |
|         | 13032      | 6.333333      |


```

Q1.3 Select from user with id==35 the product with id==4942. Is the average position of this product for this user greater than the average position from all users?

```
In [7]:  
  
position = position.reset_index()  
user_pr = position[(position.user_id==35) & (position.product_id==4942)]  
user_pr.head()  
  
Out[7]:  
  


|      | user_id | product_id | mean_pos_cart |
|------|---------|------------|---------------|
| 2044 | 35      | 4942       | 12.625        |


```

```
In [8]:  
  
user_pr_avg = user_pr.mean_pos_cart.values  
user_pr_avg  
  
Out[8]:  
  
array([12.625])  
  
In [9]:  
  
avg_all = item_avg[item_avg.index==4942].add_to_cart_order.values  
avg_all  
  
Out[9]:  
  
array([6.21933842])  
  
In [10]:  
  
user_pr_avg > avg_all  
  
Out[10]:  
  
array([ True])
```

Question 2

Q2.0 Create a DataFrame that keeps only the prior orders

- You need to keep only the orders (from orders.csv) and not the products that have been purchased in each order.
- Save the DataFrame as 'orders_prior'

```
In [11]:  
  
orders_prior = orders[orders.eval_set== 'prior']  
orders_prior.head(20)  
  
Out[11]:  
  


|    | order_id | user_id | eval_set | order_number | order_dow | order_hour_of_day | days_since_prior_order |
|----|----------|---------|----------|--------------|-----------|-------------------|------------------------|
| 0  | 2539329  | 1       | prior    | 1            | 2         | 8                 | NaN                    |
| 1  | 2398795  | 1       | prior    | 2            | 3         | 7                 | 15.0                   |
| 2  | 473747   | 1       | prior    | 3            | 3         | 12                | 21.0                   |
| 3  | 2254736  | 1       | prior    | 4            | 4         | 7                 | 29.0                   |
| 4  | 431534   | 1       | prior    | 5            | 4         | 15                | 28.0                   |
| 5  | 3367565  | 1       | prior    | 6            | 2         | 7                 | 19.0                   |
| 6  | 550135   | 1       | prior    | 7            | 1         | 9                 | 20.0                   |
| 7  | 3108588  | 1       | prior    | 8            | 1         | 14                | 14.0                   |
| 8  | 2295261  | 1       | prior    | 9            | 1         | 16                | 0.0                    |
| 9  | 2550362  | 1       | prior    | 10           | 4         | 8                 | 30.0                   |
| 11 | 2168274  | 2       | prior    | 1            | 2         | 11                | NaN                    |
| 12 | 1501582  | 2       | prior    | 2            | 5         | 10                | 10.0                   |
| 13 | 1901567  | 2       | prior    | 3            | 1         | 10                | 3.0                    |
| 14 | 738281   | 2       | prior    | 4            | 2         | 10                | 8.0                    |
| 15 | 1673511  | 2       | prior    | 5            | 3         | 11                | 8.0                    |
| 16 | 1199898  | 2       | prior    | 6            | 2         | 9                 | 13.0                   |
| 17 | 3194192  | 2       | prior    | 7            | 2         | 12                | 14.0                   |
| 18 | 788338   | 2       | prior    | 8            | 1         | 15                | 27.0                   |
| 19 | 1718559  | 2       | prior    | 9            | 2         | 9                 | 8.0                    |
| 20 | 1447487  | 2       | prior    | 10           | 1         | 11                | 6.0                    |


```

Q2.1 Keep the orders from the customers that have made at least 25 orders

- Use 'orders_prior' DataFrame that you have created on the previous step
- Save the DataFrame as 'orders_25'

```
In [12]:  
  
orders_25 = orders_prior.groupby('user_id').filter(lambda x: x.order_number.max()  
>= 25)  
orders_25.head()  
  
Out[12]:  
  


|     | order_id | user_id | eval_set | order_number | order_dow | order_hour_of_day | days_since_prior_order |
|-----|----------|---------|----------|--------------|-----------|-------------------|------------------------|
| 160 | 1737705  | 17      | prior    | 1            | 2         | 13                | NaN                    |
| 161 | 1681401  | 17      | prior    | 2            | 5         | 10                | 3.0                    |
| 162 | 2680214  | 17      | prior    | 3            | 3         | 10                | 5.0                    |
| 163 | 3197376  | 17      | prior    | 4            | 1         | 14                | 5.0                    |
| 164 | 3237467  | 17      | prior    | 5            | 6         | 17                | 5.0                    |


```

```
In [13]:  
  
# Sanity check on user_id 28  
orders_25[orders_25.user_id==28]  
  
Out[13]:  
  


|  | order_id | user_id | eval_set | order_number | order_dow | order_hour_of_day | days_since_prior_order |
|--|----------|---------|----------|--------------|-----------|-------------------|------------------------|
|--|----------|---------|----------|--------------|-----------|-------------------|------------------------|


```

Q2.2 How many days on average pass for these customers to place an order?

```
In [14]:  
  
orders_25.days_since_prior_order.mean()  
  
Out[14]:  
  
7.025531251863962
```

Q2.3 Create a histogram for the days that pass since a prior order for the customers with at least 25 orders

- use `bins=100`

```
In [15]:  
  
orders_25.days_since_prior_order.plot.hist(bins=100)  
  
Out[15]:  
  
<matplotlib.axes._subplots.AxesSubplot at 0x7fe08b3ce400>  
  

```

Q2.4 Calculate the rate of change for the average days since prior order between customers with at least 25 orders & all customers

- Create the following formula:

$$ROC = \frac{avg_days_since_prior(\geq 25) - avg_days_since_prior(ALL)}{avg_days_since_prior(ALL)} \times 100$$

```
In [16]:  
  
((orders_25.days_since_prior_order.mean() - orders[orders.eval_set== 'prior'].days_since_prior_order.mean()) / orders[orders.eval_set== 'prior'].days_since_prior_order.mean()) * 100  
  
Out[16]:  
  
-34.385248487509314
```

Question 3

Q3.0 Create a DataFrame that keeps all the prior orders and the products that have been purchased

```
In [17]:  
  
prd = pd.merge(orders, order_products_prior, on='order_id', how='inner')
```

Q3.1 Create a new column which keeps the order_number in reverse order

```
In [18]:  
  
prd['order_number_back'] = prd.groupby('user_id')['order_number'].transform(max) - prd.order_number + 1  
  
Q3.2 Keep only the last 10 orders from each user
```

```
In [19]:  
  
prd10 = prd[prd.order_number_back <= 10]  
prd10.head()  
  
Out[19]:  
  


|   | order_id | user_id | eval_set | order_number | order_dow | order_hour_of_day | days_since_prior_order | product_id |
|---|----------|---------|----------|--------------|-----------|-------------------|------------------------|------------|
| 0 | 2539329  | 1       | prior    | 1            | 2         | 8                 | NaN                    | 196        |
| 1 | 2539329  | 1       | prior    | 1            | 2         | 8                 | NaN                    | 14084      |
| 2 | 2539329  | 1       | prior    | 1            | 2         | 8                 | NaN                    | 12427      |
| 3 | 2539329  | 1       | prior    | 1            | 2         | 8                 | NaN                    | 26088      |
| 4 | 2539329  | 1       | prior    | 1            | 2         | 8                 | NaN                    | 26405      |


```

Q3.3 Exclude all the customers who have less than 10 orders in total

```
In [20]:  
  
prd10 = prd10.groupby('user_id').filter(lambda x: x.order_id.nunique() >= 10)  
  
In [21]:  
  
#sanity test  
prd10[prd10.user_id==19]  
  
Out[21]:  
  


|  | order_id | user_id | eval_set | order_number | order_dow | order_hour_of_day | days_since_prior_order | product_id |
|--|----------|---------|----------|--------------|-----------|-------------------|------------------------|------------|
|--|----------|---------|----------|--------------|-----------|-------------------|------------------------|------------|


```

Q3.4 For the last 10 orders, find the average basket size of each user.

```
In [22]:  
  
basket_size = prd10.groupby(['user_id', 'order_id'])[['product_id']].count()  
basket_size.columns = ['volume']  
basket_size.head()  
  
Out[22]:  
  


|         |          | volume |
|---------|----------|--------|
| user_id | order_id |        |
| 1       | 431534   | 8      |
|         | 473747   | 5      |
|         | 550135   | 5      |
|         | 2254736  | 5      |
|         | 2295261  | 6      |


```

```
In [23]:  
  
basket_size = basket_size.reset_index()  
basket_size_user = basket_size.groupby('user_id')['volume'].mean()  
basket_size_user.head()  
  
Out[23]:  
  


|         | volume |
|---------|--------|
| user_id |        |
| 1       | 5.9    |
| 2       | 15.8   |
| 3       | 6.9    |
| 7       | 9.4    |
| 13      | 7.2    |


```

Q3.5 How many customers have average basket size = 1 for their last 10 orders?

```
In [24]:  
  
basket_size_user[basket_size_user.volume==1].count()  
  
Out[24]:  
  
volume    132  
dtype: int64
```