

In [1]:

```
1 import pandas as pd
2 import seaborn as sns
3
4 import matplotlib.pyplot as plt
5 import numpy as np
6 import datetime as dt
7
8 import missingno as msno
9 from textwrap import wrap
```

In [2]:

```
1 transaction_df = pd.read_excel('transaction.xlsx')
2 print(transaction_df.shape)
3 transaction_df.head().T
```

(20000, 13)

Out[2]:

|                                | 0                      | 1                      | 2                      | 3                      | 4                      |
|--------------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| <b>transaction_id</b>          | 1                      | 2                      | 3                      | 4                      | 5                      |
| <b>product_id</b>              | 2                      | 3                      | 37                     | 88                     | 78                     |
| <b>customer_id</b>             | 2950                   | 3120                   | 402                    | 3135                   | 787                    |
| <b>transaction_date</b>        | 2017-02-25<br>00:00:00 | 2017-05-21<br>00:00:00 | 2017-10-16<br>00:00:00 | 2017-08-31<br>00:00:00 | 2017-10-01<br>00:00:00 |
| <b>online_order</b>            | 0.0                    | 1.0                    | 0.0                    | 0.0                    | 1.0                    |
| <b>order_status</b>            | Approved               | Approved               | Approved               | Approved               | Approved               |
| <b>brand</b>                   | Solex                  | Trek Bicycles          | OHM Cycles             | Norco<br>Bicycles      | Giant Bicycles         |
| <b>product_line</b>            | Standard               | Standard               | Standard               | Standard               | Standard               |
| <b>product_class</b>           | medium                 | medium                 | low                    | medium                 | medium                 |
| <b>product_size</b>            | medium                 | large                  | medium                 | medium                 | large                  |
| <b>list_price</b>              | 71.49                  | 2091.47                | 1793.43                | 1198.46                | 1765.3                 |
| <b>standard_cost</b>           | 53.62                  | 388.92                 | 248.82                 | 381.1                  | 709.48                 |
| <b>product_first_sold_date</b> | 41245.0                | 41701.0                | 36361.0                | 36145.0                | 42226.0                |

In [3]:

```
1
2 for col in ['order_status', 'brand', 'product_line', 'product_class', 'product_size']
3     print(col, transaction_df[col].unique())
4
```

order\_status ['Approved' 'Cancelled']

brand ['Solex' 'Trek Bicycles' 'OHM Cycles' 'Norco Bicycles' 'Giant Bicycles'

'WeareA2B' nan]

product\_line ['Standard' 'Road' 'Mountain' 'Touring' nan]

product\_class ['medium' 'low' 'high' nan]

product\_size ['medium' 'large' 'small' nan]

In [4]:

```
1 df = transaction_df.copy()
2 stats = [df.isna().sum(), df.count(), df.nunique(), df.dtypes, df.min(), df.
3 stats_name = {0: 'missing',
4               1: 'total',
5               2: 'unique',
6               3: 'dtypes',
7               4: 'minimum',
8               5: 'median',
9               6: 'maximum',
10              7: 'skewness'}
11
12 df_info = pd.DataFrame(stats).T.rename(columns=stats_name)
13 df_info
```

/opt/anaconda3/lib/python3.7/site-packages/ipykernel\_launcher.py:2: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric\_only=None') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

/opt/anaconda3/lib/python3.7/site-packages/ipykernel\_launcher.py:2: FutureWarning: DataFrame.mean and DataFrame.median with numeric\_only=None will include datetime64 and datetime64tz columns in a future version.

Out[4]:

|                         | missing | total | unique | dtypes         | minimum             | median  | maximum             | sk |
|-------------------------|---------|-------|--------|----------------|---------------------|---------|---------------------|----|
| transaction_id          | 0       | 20000 | 20000  | int64          | 1                   | 10000.5 | 20000               |    |
| product_id              | 0       | 20000 | 101    | int64          | 0                   | 44.0    | 100                 | C  |
| customer_id             | 0       | 20000 | 3494   | int64          | 1                   | 1736.0  | 5034                |    |
| transaction_date        | 0       | 20000 | 364    | datetime64[ns] | 2017-01-01 00:00:00 | NaN     | 2017-12-30 00:00:00 |    |
| online_order            | 360     | 19640 | 2      | float64        | 0.0                 | 1.0     | 1.0                 | -C |
| order_status            | 0       | 20000 | 2      | object         | Approved            | NaN     | Cancelled           |    |
| brand                   | 197     | 19803 | 6      | object         | NaN                 | NaN     | NaN                 |    |
| product_line            | 197     | 19803 | 4      | object         | NaN                 | NaN     | NaN                 |    |
| product_class           | 197     | 19803 | 3      | object         | NaN                 | NaN     | NaN                 |    |
| product_size            | 197     | 19803 | 3      | object         | NaN                 | NaN     | NaN                 |    |
| list_price              | 0       | 20000 | 296    | float64        | 12.01               | 1163.89 | 2091.47             | -C |
| standard_cost           | 197     | 19803 | 103    | float64        | 7.21                | 507.58  | 1759.85             | C  |
| product_first_sold_date | 197     | 19803 | 100    | float64        | 33259.0             | 38216.0 | 42710.0             | -C |

In [5]:

```
1 import pandas as pd
2 import seaborn as sns
3 from tqdm.notebook import tqdm
4 import warnings
5 warnings.simplefilter('ignore')
6 tqdm.pandas()
```

In [6]:

```
1
2 def prepare_data(df, select_columns):
3     df = df[select_columns]
4     df['margin'] = df['list_price'] - df['standard_cost']
5     return df
6
7 df = pd.read_excel('transaction.xlsx')
8 select_columns = ['transaction_date', 'customer_id', 'list_price', 'standard_cost']
9 df = prepare_data(df, select_columns)
10 df.head(2)
11
```

Out[6]:

|   | transaction_date | customer_id | list_price | standard_cost | margin  |
|---|------------------|-------------|------------|---------------|---------|
| 0 | 2017-02-25       | 2950        | 71.49      | 53.62         | 17.87   |
| 1 | 2017-05-21       | 3120        | 2091.47    | 388.92        | 1702.55 |

In [7]:

```
1
2 def get_cohort_timestamp(df, col_name, new_col_name):
3     transaction_year = df[col_name].dt.year
4     transaction_month = df[col_name].dt.month
5     df[new_col_name] = [dt.datetime(i,j,1) for i,j in zip(transaction_year, transaction_month)]
6
7     return df
8
9 df = get_cohort_timestamp(df, 'transaction_date', 'cohort_timestamp')
10 df.head()
11
```

Out[7]:

|   | transaction_date | customer_id | list_price | standard_cost | margin  | cohort_timestamp |
|---|------------------|-------------|------------|---------------|---------|------------------|
| 0 | 2017-02-25       | 2950        | 71.49      | 53.62         | 17.87   | 2017-02-01       |
| 1 | 2017-05-21       | 3120        | 2091.47    | 388.92        | 1702.55 | 2017-05-01       |
| 2 | 2017-10-16       | 402         | 1793.43    | 248.82        | 1544.61 | 2017-10-01       |
| 3 | 2017-08-31       | 3135        | 1198.46    | 381.10        | 817.36  | 2017-08-01       |
| 4 | 2017-10-01       | 787         | 1765.30    | 709.48        | 1055.82 | 2017-10-01       |

In [8]:

```
1
2 def add_cohort_index(df, transaction_column, cohort_column, focus_column):
3     gpy_info = df.groupby(by = focus_column).min()[cohort_column].reset_index()
4     df = df.drop(columns = cohort_column).merge(gpy_info, on=focus_column)
5     year_diff = df[transaction_column].dt.year - df[cohort_column].dt.year
6     month_diff = df[transaction_column].dt.month - df[cohort_column].dt.month
7     df['cohort_index'] = [12*i + j + 1 for i, j in zip(year_diff, month_diff)]
8     return df
9
10
11 df = add_cohort_index(df, 'transaction_date', 'cohort_timestamp', 'customer_id')
12 df.head()
13
```

Out[8]:

|   | transaction_date | customer_id | list_price | standard_cost | margin  | cohort_timestamp | cohort_index |
|---|------------------|-------------|------------|---------------|---------|------------------|--------------|
| 0 | 2017-02-25       | 2950        | 71.49      | 53.62         | 17.87   | 2017-02-01       | 1            |
| 1 | 2017-10-16       | 2950        | 1403.50    | 954.82        | 448.68  | 2017-02-01       | 1            |
| 2 | 2017-04-26       | 2950        | 478.16     | 298.72        | 179.44  | 2017-02-01       | 1            |
| 3 | 2017-05-21       | 3120        | 2091.47    | 388.92        | 1702.55 | 2017-01-01       | 1            |
| 4 | 2017-10-05       | 3120        | 1129.13    | 677.48        | 451.65  | 2017-01-01       | 1            |

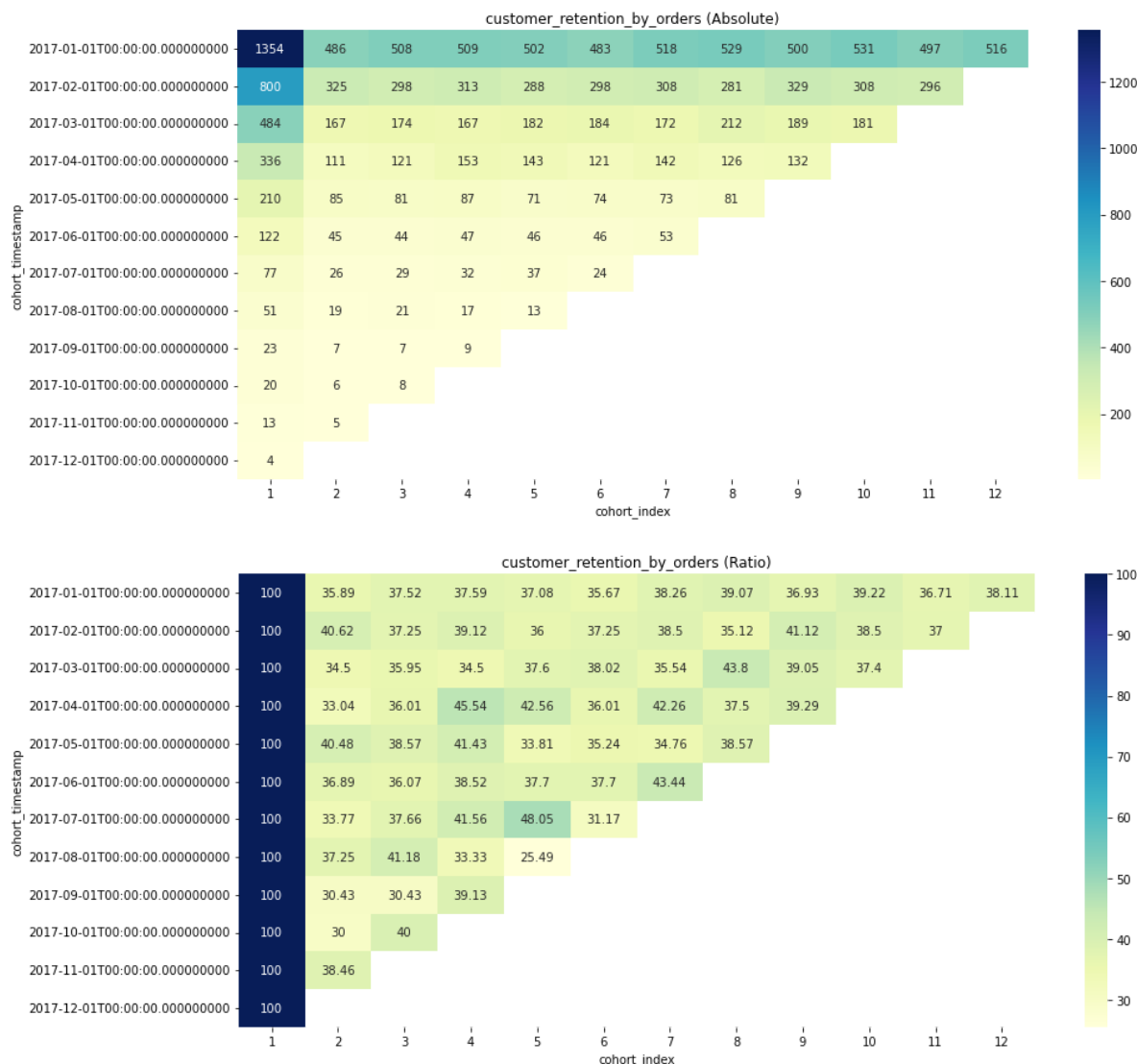
In [9]:

```
1
2 def plot_heatmap(df, title=''):
3     plt.figure(figsize=(15,7))
4     plt.title(title)
5     sns.heatmap(df, cmap='YlGnBu', annot=True, fmt='g')
6     plt.show()
7
```

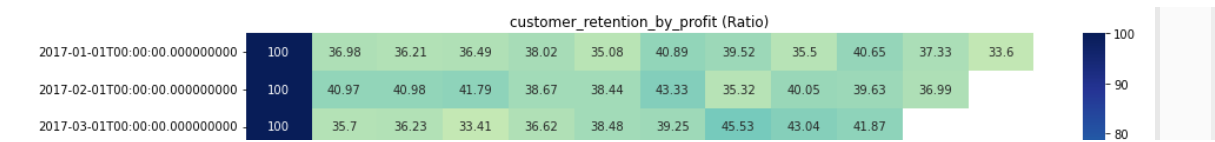
```

1
2 def get_chort_chart(df, row, column, value, function):
3
4     if function == 'customer_retention_by_orders':
5         table = df.groupby(by = [row, column]).nunique()[value]
6
7     elif function in ['customer_retention_by_revenue', 'customer_retention_by_pr
8         table = df.groupby(by = [row, column]).sum()[value]
9
10    data      = pd.DataFrame(table).reset_index()
11    absolute  = data.pivot(index = row, columns = column, values=value)
12    ratio     = round(absolute.divide(absolute.iloc[:,0], axis=0)*100,2)
13
14    plot_heatmap(absolute, function+' (Absolute)')
15    plot_heatmap(ratio, function+' (Ratio)')
16
17
18 get_chort_chart(df, 'cohort_timestamp', 'cohort_index', 'customer_id', 'customer
19 get_chort_chart(df, 'cohort_timestamp', 'cohort_index', 'list_price' , 'customer
20 get_chort_chart(df, 'cohort_timestamp', 'cohort_index', 'margin' , 'customer
21

```







In [ ]:

|   |  |
|---|--|
| 1 |  |
|---|--|

In [ ]:

|   |  |
|---|--|
| 1 |  |
|---|--|