

## Практическое задание по клиентской аналитике

(итоговые данные можно предоставить в документе ворд или pdf, сильно расписывать не надо, просто выводы по пунктам)

Действия:

1. Залить в свою БД данные по продажам (часть таблицы Orders в csv, исходник здесь <https://drive.google.com/drive/folders/1C3HqIJcABbIKM2tz8vPGiXTFT7MisrML?usp=sharing> (<https://drive.google.com/drive/folders/1C3HqIJcABbIKM2tz8vPGiXTFT7MisrML?usp=sharing>) )
2. Проанализировать, какой период данных выгружен
3. Посчитать кол-во строк, кол-во заказов и кол-во уникальных пользователей, кот совершали заказы.
4. По годам посчитать средний чек, среднее кол-во заказов на пользователя, сделать вывод , как изменялись это показатели Год от года.
5. Найти кол-во пользователей, кот покупали в одном году и перестали покупать в следующем.
6. Найти ID самого активного по кол-ву покупок пользователя.

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 %matplotlib inline
7 %config InlineBackend.figure_format = 'svg'
```

In [2]:

```
1 orders = pd.read_csv('orders_20190822.csv', sep=';')
2 a = orders.shape[0]
3 print(f'Первоначальное количество строк = {a}')
```

Первоначальное количество строк = 2002804

Конвертация данных

In [3]:

```
1 def convert_data(data):
2     new_data = data.replace(',', '.')
3     return float(new_data)
```

In [4]:

```
1 orders['price'] = orders['price'].apply(convert_data)
```

In [5]:

```
1 orders['o_date'] = pd.to_datetime(pd.Series(orders['o_date']), format="%d.%m.%Y")
```

Убираем строки с неположительной и слишком большой ценой

In [6]:

```
1 negative_list = orders.loc[(orders['price'] <= 0)|(orders['price'] > 100000)]
2 orders = orders.drop(negative_list.index, axis = 0)
```

### 1. Проанализировать, какой период данных выгружен

In [7]:

```
1 orders['o_date'].max(), orders['o_date'].min()
```

Out[7]:

```
(Timestamp('2017-12-31 00:00:00'), Timestamp('2016-01-01 00:00:00'))
```

### 2. Посчитать кол-во строк, кол-во заказов и кол-во уникальных пользователей, кот совершали заказы.

In [8]:

```
1 b = orders.shape[0]
2 print(f'Количество строк после удаления некорректных = {b}')
```

Количество строк после удаления некорректных = 2002696

In [9]:

```
1 print(f'После того, как были убраны строки с неположительной и слишком большой ценой к
```

После того, как были убраны строки с неположительной и слишком большой ценой количество строк сократилось на 108

In [10]:

```
1 number_of_orders = orders['id_o'].unique()
2 len(number_of_orders)
```

Out[10]:

2002696

In [11]:

```
1 number_of_users = orders['user_id'].unique()
2 len(number_of_users)
```

Out[11]:

1015088

### 3. По годам посчитать средний чек, среднее кол-во заказов на пользователя, сделать вывод, как изменялись эти показатели Год от года.

In [12]:

```
1 mean_order = orders.groupby(orders['o_date'].dt.year, as_index=False)[['price']].mean(
2 mean_order
```

Out[12]:

|   | mean_price  |
|---|-------------|
| 0 | 2090.294426 |
| 1 | 2392.783860 |

In [13]:

```
1 users_2016 = orders.loc[orders['o_date'].dt.year == 2016, 'user_id'].values
2 unique_users_2016 = orders.loc[orders['o_date'].dt.year == 2016, 'user_id'].unique()
3 order_per_user_2016 = len(users_2016) / len(unique_users_2016)
4
5 users_2017 = orders.loc[orders['o_date'].dt.year == 2017, 'user_id'].values
6 unique_users_2017 = orders.loc[orders['o_date'].dt.year == 2017, 'user_id'].unique()
7 order_per_user_2017 = len(users_2017) / len(unique_users_2017)
8
9 order_per_user_2016, order_per_user_2017
```

Out[13]:

```
(1.935104238984639, 1.742957402362599)
```

**Вывод:**

средний чек вырос с 2090 руб. в 2016 г. до 2393 руб. в 2017 г., среднее количество заказов на покупателя упало с 1,93 в 2016 г. до 1,74 в 2017 г.

**4. Найти кол-во пользователей, кот покупали в одном году и перестали покупать в следующем.**

In [14]:

```
1 ex_users = []
2 for i in unique_users_2016:
3     if i not in unique_users_2017:
4         ex_users.append(i)
5 len(ex_users)
```

Out[14]:

```
360216
```

**5. Найти ID самого активного по кол-ву покупок пользователя.**

**К сожалению, на всей БД из более чем 2 млн.строк мой компьютер исполнение кода не потянул**

даже за несколько часов, поэтому специально вырезаю из исходного более короткий датафрейм на 100 тыс. строк и отрабатываю на нём.

In [18]:

```
1 orders = orders[0:100000]
```

In [19]:

```
1 max_value = 1
2 max_user_id = 1
3
4 for i in orders['user_id']:
5     if (orders['user_id'] == i).sum() > max_value:
6         max_value = (orders['user_id'] == i).sum()
7         max_user_id = i
8         a = orders[orders['user_id'] == i].index
9         orders['user_id'].drop(a)
10
11 max_value, max_user_id
```

Out[19]:

```
(65, 765861)
```

In [20]:

```
1 print(f'На выборке из первых 100000 строк победил user #{max_user_id}, который произвёл')
```

На выборке из первых 100000 строк победил user #765861, который произвёл 65 покупок.

In [ ]:

```
1
```