

# 1 Описание проекта

Вы — аналитик крупного интернет-магазина. Вместе с отделом маркетинга вы подготовили список гипотез для увеличения выручки. Приоритизируйте гипотезы, запустите A/B-тест и проанализируйте результаты.

## 2 Часть 1. Приоритизация гипотез.

### 2.1 Описание данных

Файл `/datasets/hypothesis.csv`.

- `Hypothesis` — краткое описание гипотезы;
- `Reach` — охват пользователей по 10-балльной шкале;
- `Impact` — влияние на пользователей по 10-балльной шкале;
- `Confidence` — уверенность в гипотезе по 10-балльной шкале;
- `Efforts` — затраты ресурсов на проверку гипотезы по 10-балльной шкале. Чем больше значение `Efforts`, тем дороже проверка гипотезы.

### 2.2 Загрузите данные и подготовьте их к анализу

In [1]:



```
1 import pandas as pd
2
3 data=pd.read_csv('/datasets/hypothesis.csv')
4 data
5
```

Out[1]:

	Hypothesis	Reach	Impact	Confidence	Efforts
0	Добавить два новых канала привлечения трафика,...	3	10	8	6
1	Запустить собственную службу доставки, что сок...	2	5	4	10
2	Добавить блоки рекомендаций товаров на сайт ин...	8	3	7	3
3	Изменить структура категорий, что увеличит кон...	8	3	3	8
4	Изменить цвет фона главной страницы, чтобы уве...	3	1	1	1
5	Добавить страницу отзывов клиентов о магазине,...	3	2	2	3
6	Показать на главной странице баннеры с актуаль...	5	3	8	3
7	Добавить форму подписки на все основные страни...	10	7	8	5
8	Запустить акцию, дающую скидку на товар в день...	1	9	9	5

Приведем название столбцов к общему формату. Необходимо сделать переименование колонок. (предобработка данных)

In [2]:



```
1 data.columns = [x.lower() for x in data.columns.values]
2 data.columns
```

Out[2]:

```
Index(['hypothesis', 'reach', 'impact', 'confidence', 'efforts'], dtype='object')
```

## 2.3 Примените фреймворк ICE для приоритизации гипотез. Отсортируйте их по убыванию приоритета.

In [3]:



```
1 data['ICE'] = (data['impact'] * data['confidence']) / data['efforts']
2 data[['hypothesis', 'ICE']].sort_values(by='ICE', ascending=False)
```

Out[3]:

	hypothesis	ICE
8	Запустить акцию, дающую скидку на товар в день...	16.200000
0	Добавить два новых канала привлечения трафика,...	13.333333
7	Добавить форму подписки на все основные страни...	11.200000
6	Показать на главной странице баннеры с актуаль...	8.000000
2	Добавить блоки рекомендаций товаров на сайт ин...	7.000000
1	Запустить собственную службу доставки, что сок...	2.000000
5	Добавить страницу отзывов клиентов о магазине,...	1.333333
3	Изменить структура категорий, что увеличит кон...	1.125000
4	Изменить цвет фона главной страницы, чтобы уве...	1.000000

## 2.4 Примените фреймворк RICE для приоритизации гипотез. Отсортируйте их по убыванию приоритета.

In [4]:



```
1 data['RICE'] = (data['reach'] * data['impact'] * data['confidence']) / data['effort']
2 data[['hypothesis', 'RICE']].sort_values(by='RICE', ascending=False)
```

Out[4]:

	hypothesis	RICE
7	Добавить форму подписки на все основные страни...	112.0
2	Добавить блоки рекомендаций товаров на сайт ин...	56.0
0	Добавить два новых канала привлечения трафика,...	40.0
6	Показать на главной странице баннеры с актуаль...	40.0
8	Запустить акцию, дающую скидку на товар в день...	16.2
3	Изменить структура категорий, что увеличит кон...	9.0
1	Запустить собственную службу доставки, что сок...	4.0
5	Добавить страницу отзывов клиентов о магазине,...	4.0
4	Изменить цвет фона главной страницы, чтобы уве...	3.0

## 2.5 Укажите, как изменилась приоритизация гипотез при применении RICE вместо ICE. Объясните, почему так произошло.

Метод RICE учитывает больше параметров чем ICE. При расчете RICE мы учитываем такой параметр как Reach — сколько пользователей затронет изменение, которое вы хотите внести. Из-за этого он более точно позволяет приоритезировать гипотезы. Данный метод целесообразно применять, когда оцениваемые гипотезы затрагивают разные по объемам группы пользователей. Гипотеза №7 затронет всех пользователей интернет магазина, в то время как гипотеза №4 коснется меньшего числа пользователей (не все пользователи начинают свой путь по сайту с главной страницы). В нашем случае использование метода RICE более предпочтительнее, и приоритезация по данному методу более точно отражает значимость гипотез. Гипотеза №7 имеет максимальный приоритет.

## 3 Часть 2. Анализ A/B-теста

### 3.1 Описание данных

Вы провели A/B-тест и получили результаты, которые описаны в файлах `/datasets/orders.csv` и `/datasets/visitors.csv`

Файл `/datasets/orders.csv`.

- `transactionId` — идентификатор заказа;
- `visitorId` — идентификатор пользователя, совершившего заказ;
- `date` — дата, когда был совершён заказ;
- `revenue` — выручка заказа;
- `group` — группа A/B-теста, в которую попал заказ.

Файл `/datasets/visitors.csv`.

- `date` — дата;

- `group` — группа A/B-теста;
- `visitors` — количество пользователей в указанную дату в указанной группе A/B-теста

## 3.2 Постройте график кумулятивной выручки по группам. Сделайте выводы и предположения.

In [5]:



```
1 import pandas as pd
2 import datetime as dt
3 import numpy as np
4
5 import warnings
6
7 warnings.simplefilter('ignore')
8
9 orders = pd.read_csv('//datasets/orders.csv')
10 orders['date'] = orders['date'].map(lambda x: dt.datetime.strptime(x, '%Y-%m-%d'))
11
12 visitors = pd.read_csv('/datasets/visitors.csv')
13 visitors['date'] = visitors['date'].map(lambda x: dt.datetime.strptime(x, '%Y-%m-%d'))
14
15 display(orders.head())
16 display(visitors.head())
17
```

	transactionId	visitorId	date	revenue	group
0	3667963787	3312258926	2019-08-15	1650	B
1	2804400009	3642806036	2019-08-15	730	B
2	2961555356	4069496402	2019-08-15	400	A
3	3797467345	1196621759	2019-08-15	9759	B
4	2282983706	2322279887	2019-08-15	2308	B

	date	group	visitors
0	2019-08-01	A	719
1	2019-08-02	A	619
2	2019-08-03	A	507
3	2019-08-04	A	717
4	2019-08-05	A	756

In [6]:



```
1 # создаем массив уникальных пар значений дат и групп теста
2 datesGroups = orders[['date', 'group']].drop_duplicates()
3
4 # получаем агрегированные кумулятивные по дням данные о заказах
5 ordersAggregated = datesGroups.apply(
6     lambda x: orders[np.logical_and(orders['date'] <= x['date'],
7         orders['group'] == x['group'])].agg({'date' : 'max',
8         'group' : 'max',
9         'transactionId' : 'nunique',
10        'visitorId' : 'nunique',
11        'revenue' : 'sum'}), axis=1).sort_v
12
13 # получаем агрегированные кумулятивные по дням данные о посетителях интернет-ма
14 visitorsAggregated = datesGroups.apply(
15     lambda x: visitors[np.logical_and(visitors['date'] <= x['date'],
16         visitors['group'] == x['group'])].agg({'date' : 'max',
17         'group' : 'max',
18         'visitors' : 'sum'}), axis=1).so
19
20 # объединяем кумулятивные данные в одной таблице и присваиваем ее столбцам поня
21 cumulativeData = ordersAggregated.merge(visitorsAggregated, left_on=['date', 'g
22     right_on=['date', 'group'])
23 cumulativeData.columns = ['date', 'group', 'orders', 'buyers', 'revenue', 'visi
24
25 cumulativeData.head(5)
```

Out[6]:

	date	group	orders	buyers	revenue	visitors
0	2019-08-01	A	24	20	148579	719
1	2019-08-01	B	21	20	101217	713
2	2019-08-02	A	44	38	242401	1338
3	2019-08-02	B	45	43	266748	1294
4	2019-08-03	A	68	62	354874	1845

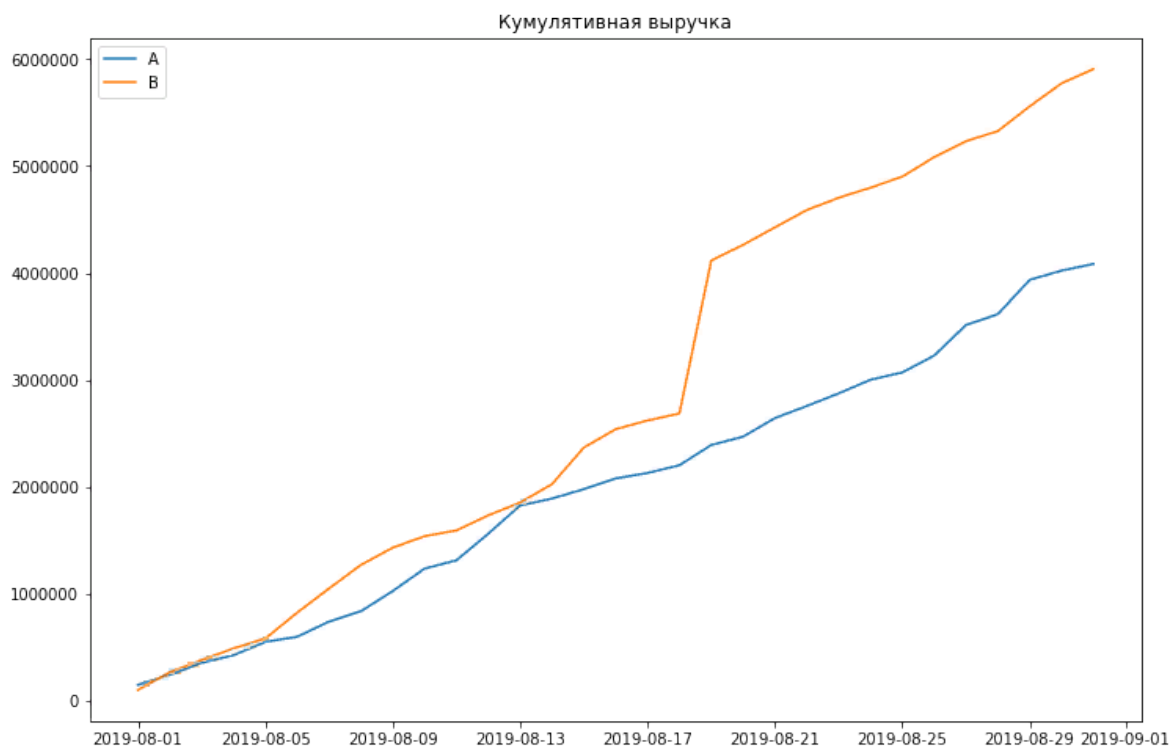
In [7]:



```
1 import matplotlib.pyplot as plt
2
3 # датафрейм с кумулятивным количеством заказов и кумулятивной выручкой по дням
4 cumulativeRevenueA = cumulativeData[cumulativeData['group']=='A']['date', 'revenue']
5
6 # датафрейм с кумулятивным количеством заказов и кумулятивной выручкой по дням
7 cumulativeRevenueB = cumulativeData[cumulativeData['group']=='B']['date', 'revenue']
8
9 # Строим график выручки группы A
10 plt.figure(figsize=(12,8))
11 plt.plot(cumulativeRevenueA['date'], cumulativeRevenueA['revenue'], label='A',)
12
13 # Строим график выручки группы B
14 plt.plot(cumulativeRevenueB['date'], cumulativeRevenueB['revenue'], label='B')
15 plt.title('Кумулятивная выручка')
16 plt.legend()
```

Out[7]:

<matplotlib.legend.Legend at 0x7f2c78240810>



**Вывод:** Выручка почти равномерно увеличивается в течение всего теста. Хороший знак. Однако график выручки группы В в нескольких точках резко растет. Это может сигнализировать о всплесках числа заказов, либо о появлении очень дорогих заказов в выборке.

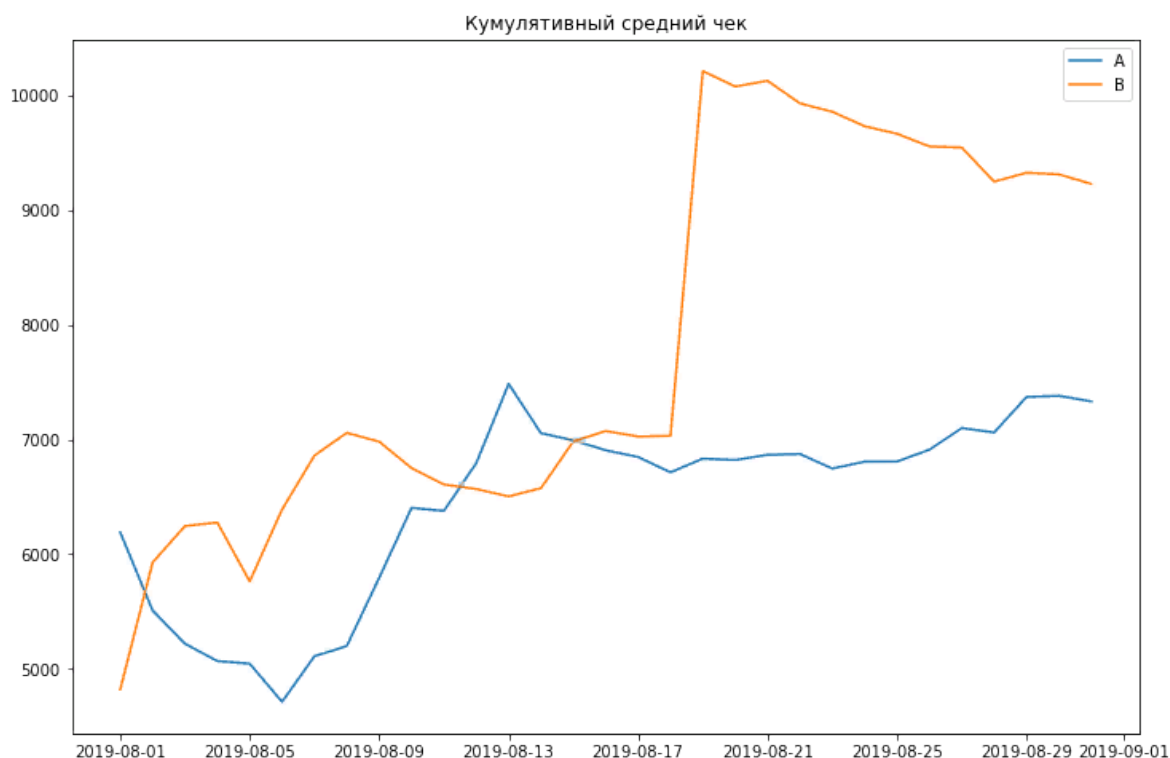
### 3.3 Постройте график кумулятивного среднего чека по группам. Сделайте выводы и предположения.

In [8]:

```
1 # Построим графики среднего чека по группам – разделим кумулятивную выручку на
2 plt.figure(figsize=(12,8))
3 plt.plot(cumulativeRevenueA['date'], cumulativeRevenueA['revenue']/cumulativeRe
4 plt.plot(cumulativeRevenueB['date'], cumulativeRevenueB['revenue']/cumulativeRe
5 plt.title('Кумулятивный средний чек')
6 plt.legend()
```

Out[8]:

<matplotlib.legend.Legend at 0x7f2c6c1cf8d0>



**Вывод:** За месяц средний чек по группе А становится более равномерным, а вот по группе В мы наблюдаем резкий всплеск в районе 18-19 числа. Скорее всего была совершена аномально большая покупка (такое бывает). Данный заказ исказил данные по группе. Для получения корректных данных на графике необходимо большее число наблюдений за группой В. Иными словами, опираясь на собранные данные мы можем сделать предположение о среднем чеке по группе А на уровне 7000 руб. Но для группы В данных недостаточно.

### 3.4 Постройте график относительного изменения кумулятивного среднего чека группы В к группе А.

In [21]:

```
1 # собираем данные в одном датафрейме
2
3 mergedCumulativeRevenue = cumulativeRevenueA.merge(cumulativeRevenueB,
4                                                     left_on='date',
5                                                     right_on='date',
6                                                     how='left', suffixes=['A', 'B'])
7
8 # строим отношение средних чеков
9 plt.figure(figsize=(12,8))
10 plt.plot(mergedCumulativeRevenue['date'], (mergedCumulativeRevenue['revenueB']/mergedCumulativeRevenue['revenueA']))
11 plt.title('Относительное изменение кумулятивного среднего чека группы В к группе А')
12
13 # добавляем ось X
14 plt.axhline(y=0, color='black', linestyle='--')
```

Out[21]:

<matplotlib.lines.Line2D at 0x7f2c6bee85d0>



**Вывод:** В нескольких точках график различия между сегментами резко «скачет». Это говорит о крупных заказах.

### 3.5 Постройте график кумулятивной конверсии по группам.



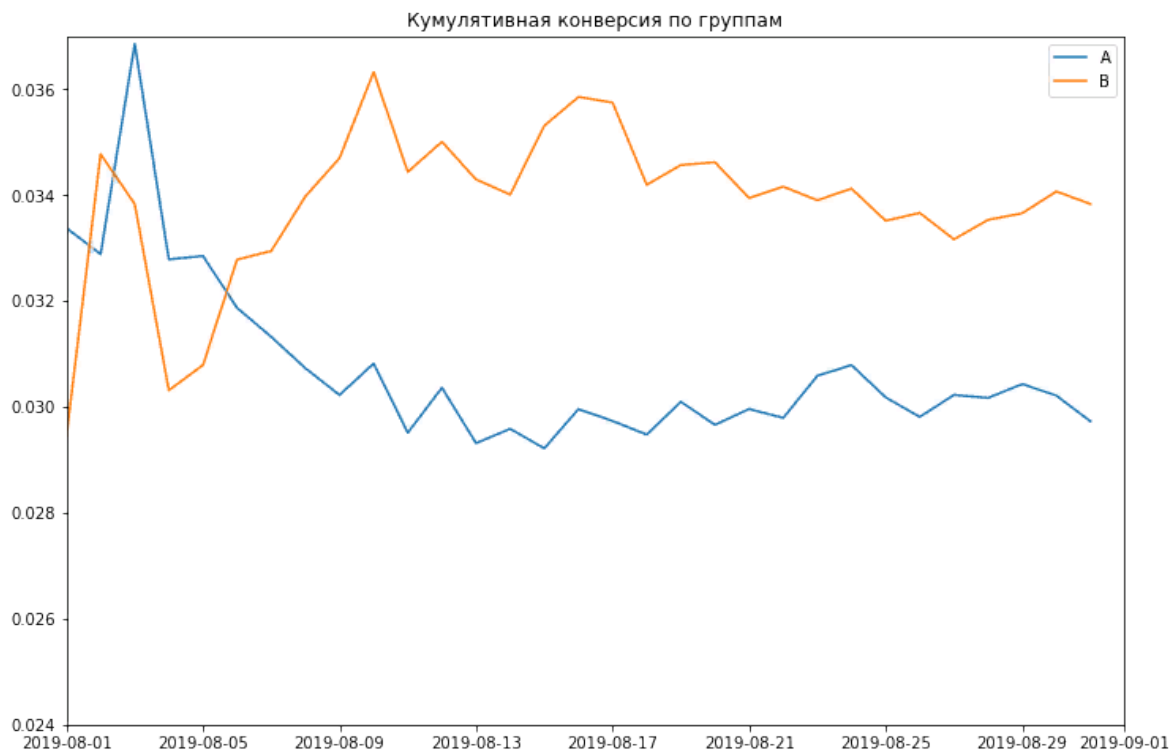
In [22]:



```
1 # считаем кумулятивную конверсию
2 cumulativeData['conversion'] = cumulativeData['orders']/cumulativeData['visitor']
3
4 # отделяем данные по группе A
5 cumulativeDataA = cumulativeData[cumulativeData['group']=='A']
6
7 # отделяем данные по группе B
8 cumulativeDataB = cumulativeData[cumulativeData['group']=='B']
9
10 # строим графики
11 plt.figure(figsize=(12,8))
12 plt.plot(cumulativeDataA['date'], cumulativeDataA['conversion'], label='A')
13 plt.plot(cumulativeDataB['date'], cumulativeDataB['conversion'], label='B')
14 plt.legend()
15 plt.title('Кумулятивная конверсия по группам')
16
17 # задаем масштаб осей
18 plt.axis(["2019-08-01", '2019-09-01', 0.024, 0.037])
```

Out[22]:

```
['2019-08-01', '2019-09-01', 0.024, 0.037]
```



**Вывод:** Конверсия группы В стабильно выше. Мы наблюдаем стандартный график. В начале идут большие скачки в связи с небольшим промежутком анализа, и далее они сглаживаются. К концу месяца можно говорить об установившемся значении кумулятивной конверсии для обеих групп.

### 3.6 Постройте график относительного изменения кумулятивной конверсии

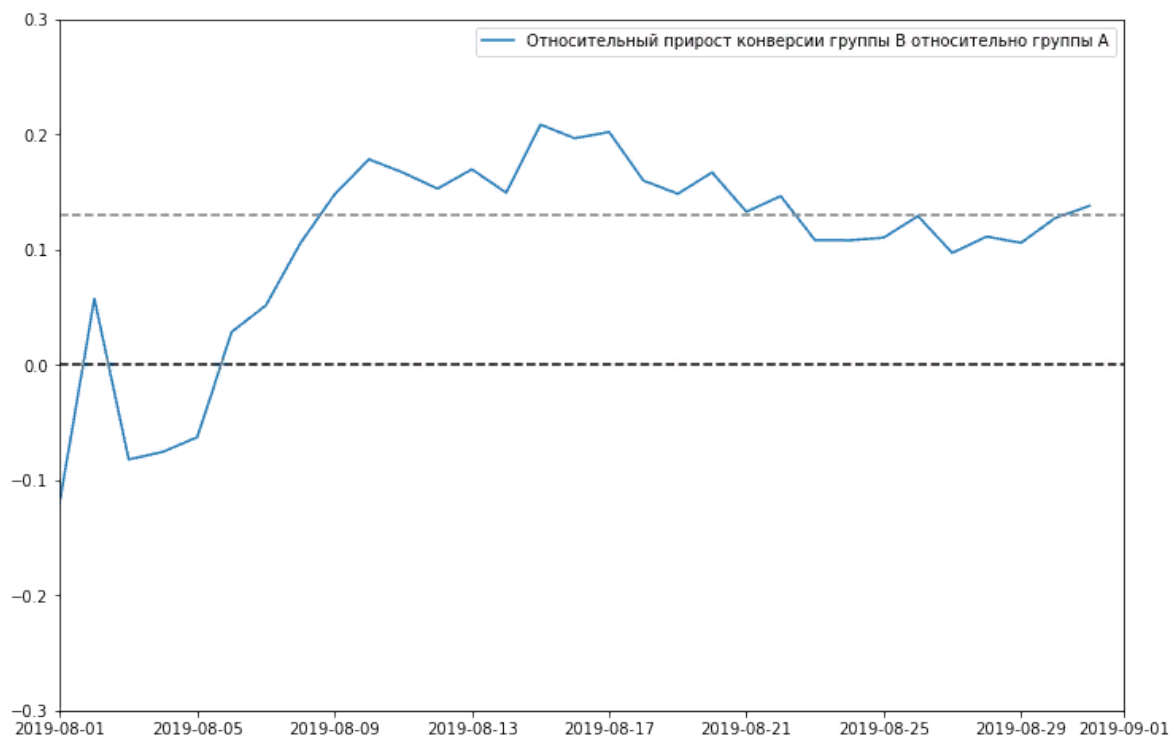
группы В к группе А.

In [11]:

```
1 mergedCumulativeConversions = cumulativeDataA[['date', 'conversion']].merge(cumu
2                                     left
3                                     right
4                                     how=
5
6 # строим график
7 plt.figure(figsize=(12,8))
8 plt.plot(mergedCumulativeConversions['date'],
9          mergedCumulativeConversions['conversionB']/mergedCumulativeConversions
10         label="Относительный прирост конверсии группы В относительно группы А"
11 plt.legend()
12
13 plt.axhline(y=0, color='black', linestyle='--')
14 plt.axhline(y=0.13, color='grey', linestyle='--')
15 plt.axis(["2019-08-01", '2019-09-01', -0.3, 0.3])
```

Out[11]:

['2019-08-01', '2019-09-01', -0.3, 0.3]



**Вывод:** Мы наблюдаем постоянный рост относительного изменения конверсии группы В к группе А. В целом отношение конверсии ещё не установилось, и сейчас делать какие-либо выводы по тесту нельзя. Впрочем, сперва стоит проанализировать аномалии, возможно, они изменят картину

### 3.7 Постройте точечный график количества заказов по пользователям

In [12]:



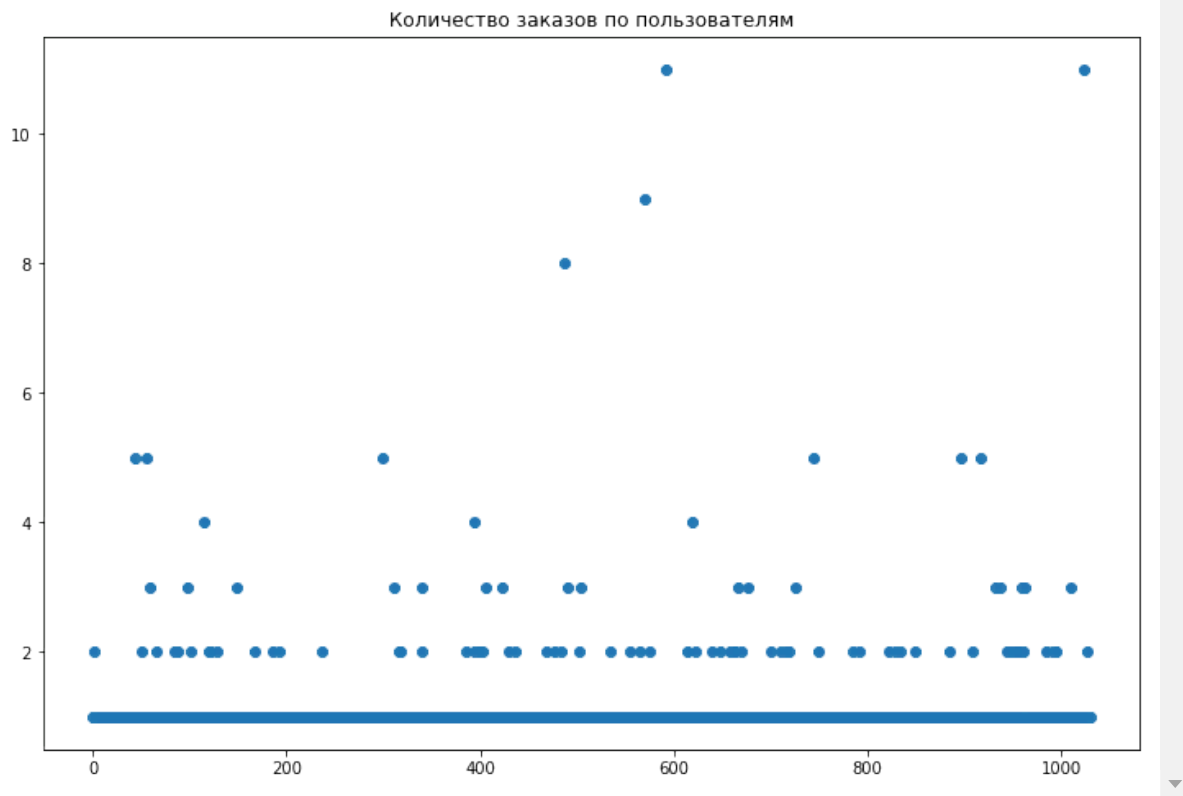
```
1 # получаю данные о количестве заказов
2 ordersByUsers = (
3     orders.groupby('visitorId', as_index=False)
4     .agg({'transactionId': pd.Series.nunique})
5 )
6 ordersByUsers.columns = ['userId', 'orders']
7
8 display(ordersByUsers.sort_values(by='orders', ascending=False).head(10))
9
10 # серия из чисел от 0 до количества наблюдений в ordersByUsers
11 plt.figure(figsize=(12,8))
12 x_values = pd.Series(range(0, len(ordersByUsers)))
13 plt.title('Количество заказов по пользователям')
14 plt.scatter(x_values, ordersByUsers['orders'])
```

	userId	orders
1023	4256040402	11
591	2458001652	11
569	2378935119	9
487	2038680547	8
44	199603092	5
744	3062433592	5
55	237748145	5
917	3803269165	5
299	1230306981	5
897	3717692402	5

Out[12]:

<matplotlib.collections.PathCollection at 0x7f2c75300f50>





**Вывод:** Большинство пользователей делали 1-3 заказа, но есть и "уникальные" которые совершили 11 заказов. Их точная доля не ясна — непонятно, считать их аномалиями или нет.

**3.8 Посчитайте 95-й и 99-й перцентили количества заказов на пользователя. Выберите границу для определения аномальных пользователей.**

In [13]:



```
1 display(ordersByUsers.sort_values(by='orders', ascending=False).head(10))
2
3 # Метод np.percentile('column', [percentile1, percentile2, percentile3]) находи
4 np.percentile(ordersByUsers['orders'], [95, 99])
```

	userId	orders
1023	4256040402	11
591	2458001652	11
569	2378935119	9
487	2038680547	8
44	199603092	5
744	3062433592	5
55	237748145	5
917	3803269165	5
299	1230306981	5
897	3717692402	5

Out[13]:

```
array([2., 4.])
```

**Вывод:** Всего 5% покупателей оформили более 2 заказов и всего 1% оформили более 4. Разумно выбрать 2 заказа на одного пользователя за нижнюю границу числа заказов, и отсеять аномальных пользователей по ней. Пользователи, совершившие много заказов, влияют на числитель формулы конверсии. Скорее всего, их поведение отличается от нормального. Если речь не идёт об интернет-магазине с регулярным спросом, «обычный» пользователь совершает не больше двух заказов за короткий срок.

### 3.9 Постройте точечный график стоимостей заказов. Сделайте выводы и предположения.

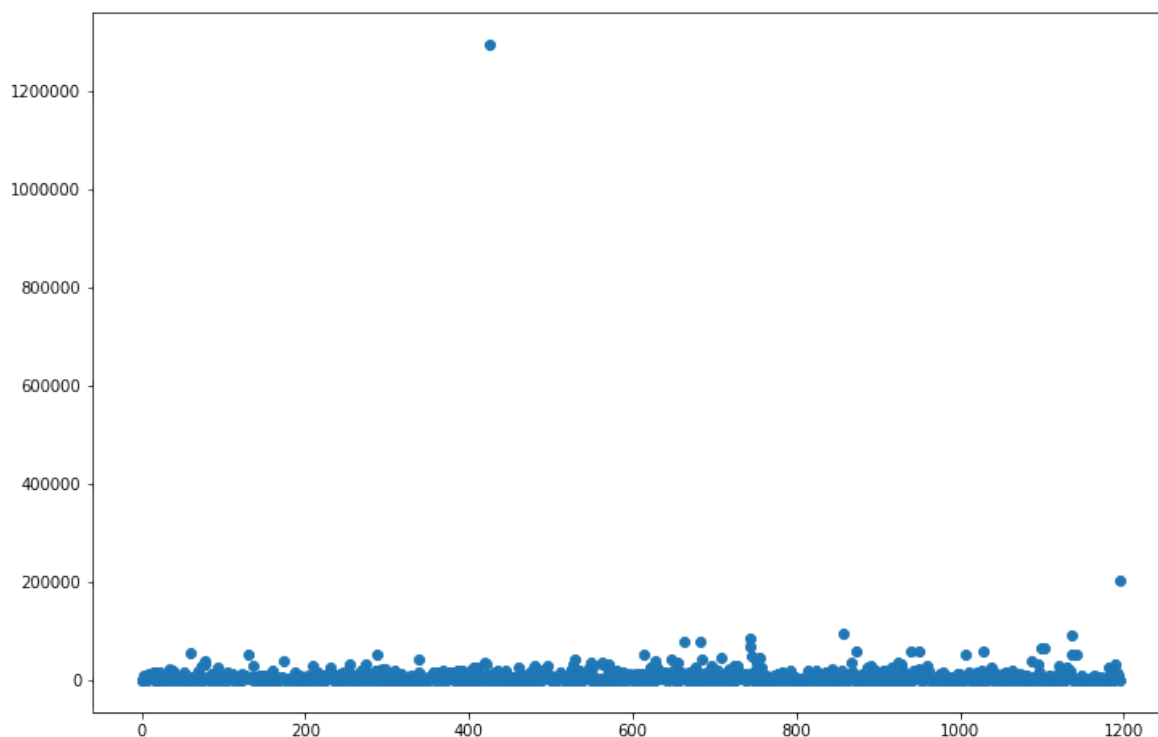
In [14]:



```
1 x_values = pd.Series(range(0, len(orders['revenue'])))
2 plt.figure(figsize=(12, 8))
3 plt.title('Точечный график стоимостей заказов')
4 plt.scatter(x_values, orders['revenue'])
```

Out[14]:

<matplotlib.collections.PathCollection at 0x7f2c6bf5b9d0>



**Вывод:** Мы видим один аномальный заказ более чем на 1 млн и еще один чуть больше чем 200 тыс. Эти заказы явно оказывают негативное влияние на нашу выборку. Вполне возможно, что в дальнейших расчетах мы не будем их учитывать.

**3.10 Посчитайте 95-й и 99-й перцентили стоимости заказов. Выберите**

границу для определения аномальных заказов.

In [15]:



```
1 np.percentile(orders['revenue'], [95, 99])
```

Out[15]:

```
array([28000. , 58233.2])
```

**Вывод:** 5% пользователей оформили заказ на сумму свыше 28 тыс. Если такие дорогие заказы попадут в одну из групп теста, они сразу же исказят результаты и сделают победителем ту группу, где оказались. При этом такие покупки — редкость, исключение, которое проявилось не из-за тестирования гипотезы, а случайно.

Такие аномально дорогие заказы следует удалять из теста. Можно провести и анализ самых дешёвых заказов. Однако из-за маленькой стоимости они слабо влияют на средний чек.

### 3.11 Посчитайте статистическую значимость различий в конверсии между группами по «сырым» данным.

In [16]:



```
1 import scipy.stats as stats
2
3 # пользователи группы A по дням
4 visitorsADaily = visitors[visitors['group'] == 'A'][['date', 'visitors']]
5 visitorsADaily.columns = ['date', 'visitorsPerDateA']
6
7 # кумулятивная суума пользователей A по дням
8 visitorsACummulative = visitorsADaily.apply(
9     lambda x: visitorsADaily[visitorsADaily['date'] <= x['date']].agg(
10         {'date': 'max', 'visitorsPerDateA': 'sum'}
11     ),
12     axis=1,
13 )
14 visitorsACummulative.columns = ['date', 'visitorsCummulativeA']
15
16 # пользователи группы B по дням
17 visitorsBDaily = visitors[visitors['group'] == 'B'][['date', 'visitors']]
18 visitorsBDaily.columns = ['date', 'visitorsPerDateB']
19
20 # кумулятивная суума пользователей B по дням
21 visitorsBCummulative = visitorsBDaily.apply(
22     lambda x: visitorsBDaily[visitorsBDaily['date'] <= x['date']].agg(
23         {'date': 'max', 'visitorsPerDateB': 'sum'}
24     ),
25     axis=1,
26 )
27 visitorsBCummulative.columns = ['date', 'visitorsCummulativeB']
28
29 # пользователи группы A совершившие хотя бы 1 заказ
30 ordersADaily = (
31     orders[orders['group'] == 'A'][['date', 'transactionId', 'visitorId', 'revenue']]
32     .groupby('date', as_index=False)
33     .agg({'transactionId': pd.Series.nunique, 'revenue': 'sum'})
34 )
35 ordersADaily.columns = ['date', 'ordersPerDateA', 'revenuePerDateA']
36
37 # кумулятивные данные по группе A
38 ordersACummulative = ordersADaily.apply(
39     lambda x: ordersADaily[ordersADaily['date'] <= x['date']].agg(
40         {'date': 'max', 'ordersPerDateA': 'sum', 'revenuePerDateA': 'sum'}
41     ),
42     axis=1,
43 ).sort_values(by=['date'])
44 ordersACummulative.columns = [
45     'date',
46     'ordersCummulativeA',
47     'revenueCummulativeA',
48 ]
49
50 # пользователи группы B совершившие хотя бы 1 заказ
51 ordersBDaily = (
52     orders[orders['group'] == 'B'][['date', 'transactionId', 'visitorId', 'revenue']]
53     .groupby('date', as_index=False)
54     .agg({'transactionId': pd.Series.nunique, 'revenue': 'sum'})
55 )
56 ordersBDaily.columns = ['date', 'ordersPerDateB', 'revenuePerDateB']
57
58 ordersBCummulative = ordersBDaily.apply(
59     lambda x: ordersBDaily[ordersBDaily['date'] <= x['date']].agg(
```



```

60         {'date': 'max', 'ordersPerDateB': 'sum', 'revenuePerDateB': 'sum'}
61     ),
62     axis=1,
63 ).sort_values(by=['date'])
64 ordersBCummulative.columns = [
65     'date',
66     'ordersCummulativeB',
67     'revenueCummulativeB',
68 ]
69
70 # общая таблица с кумулятивными данными
71 data = (
72     ordersADaily.merge(
73         ordersBDaily, left_on='date', right_on='date', how='left'
74     )
75     .merge(ordersACummulative, left_on='date', right_on='date', how='left')
76     .merge(ordersBCummulative, left_on='date', right_on='date', how='left')
77     .merge(visitorsADaily, left_on='date', right_on='date', how='left')
78     .merge(visitorsBDaily, left_on='date', right_on='date', how='left')
79     .merge(visitorsACummulative, left_on='date', right_on='date', how='left')
80     .merge(visitorsBCummulative, left_on='date', right_on='date', how='left')
81 )
82
83 display(data.head(5))
84 ordersByUsersA = (
85     orders[orders['group'] == 'A']
86     .groupby('visitorId', as_index=False)
87     .agg({'transactionId': pd.Series.nunique})
88 )
89 ordersByUsersA.columns = ['userId', 'orders']
90
91 ordersByUsersB = (
92     orders[orders['group'] == 'B']
93     .groupby('visitorId', as_index=False)
94     .agg({'transactionId': pd.Series.nunique})
95 )
96 ordersByUsersB.columns = ['userId', 'orders']
97 sampleA = pd.concat([ordersByUsersA['orders'],
98                     pd.Series(0, index=np.arange(data['visitorsPerDateA'].sum(
99
100 sampleB = pd.concat([ordersByUsersB['orders'],
101                     pd.Series(0, index=np.arange(data['visitorsPerDateB'].sum(
102 ordersByUsersA = (
103     orders[orders['group'] == 'A']
104     .groupby('visitorId', as_index=False)
105     .agg({'transactionId': pd.Series.nunique})
106 )
107 ordersByUsersA.columns = ['userId', 'orders']
108
109 ordersByUsersB = (
110     orders[orders['group'] == 'B']
111     .groupby('visitorId', as_index=False)
112     .agg({'transactionId': pd.Series.nunique})
113 )
114 ordersByUsersB.columns = ['userId', 'orders']
115
116 sampleA = pd.concat(
117     [
118         ordersByUsersA['orders'],
119         pd.Series(
120             0,

```

```

121         index=np.arange(
122             data['visitorsPerDateA'].sum() - len(ordersByUsersA['orders'])
123         ),
124         name='orders',
125     ),
126 ],
127 axis=0,
128 )
129
130 sampleB = pd.concat(
131     [
132         ordersByUsersB['orders'],
133         pd.Series(
134             0,
135             index=np.arange(
136                 data['visitorsPerDateB'].sum() - len(ordersByUsersB['orders'])
137             ),
138             name='orders',
139         ),
140     ],
141     axis=0,
142 )
143 display("{0:.5f}".format(stats.mannwhitneyu(sampleA, sampleB)[1]))
144
145 display("{0:.3f}".format(sampleB.mean() / sampleA.mean() - 1))

```

	date	ordersPerDateA	revenuePerDateA	ordersPerDateB	revenuePerDateB	ordersCumulative
0	2019-08-01	24	148579	21	101217	
1	2019-08-02	20	93822	24	165531	
2	2019-08-03	24	112473	16	114248	
3	2019-08-04	16	70825	17	108571	
4	2019-08-05	25	124218	23	92428	

'0.00840'

'0.138'

**Вывод** P-value значительно меньше 0.05, потому нулевую гипотезу отвергаем (различий между группами нет). Анализ сырых данных показывает, что в конверсии между группами есть статистически значимые различия. Относительный прирост конверсии группы В к конверсии группы А равен 13,8%

### 3.12 Посчитайте статистическую значимость различий в среднем чеке заказа между группами по «сырым» данным.

In [17]:



```
1 display('{0:.3f}'.format(stats.mannwhitneyu(orders[orders['group']=='A']['revenue'],
2                                             orders[orders['group']=='B']['revenue'])))
3 display('{0:.3f}'.format(orders[orders['group']=='B']['revenue'].mean()/
4                           orders[orders['group']=='A']['revenue'].mean()-1))
```

'0.365'

'0.259'

**Вывод** P-value больше 0.05 - статистически значимых отличий в среднем чеке между группами нет. Хотя разница между группами значительна почти 26%. Это говорит нам о том что ввыборке есть выбросы, аномально дорогие заказы.

### 3.13 Посчитайте статистическую значимость различий в конверсии между группами по «очищенным» данным.

In [18]:



```
1 # выборка по количеству заказов
2 usersWithManyOrders = pd.concat(
3     [
4         ordersByUsersA[ordersByUsersA['orders'] > 2]['userId'],
5         ordersByUsersB[ordersByUsersB['orders'] > 2]['userId'],
6     ],
7     axis=0,
8 )
9
10 # выборка по сумме заказов
11 usersWithExpensiveOrders = orders[orders['revenue'] > 28000]['visitorId']
12 abnormalUsers = (
13     pd.concat([usersWithManyOrders, usersWithExpensiveOrders], axis=0)
14     .drop_duplicates()
15     .sort_values()
16 )
17 display(abnormalUsers.head(5))
```

```
1099    148427295
18      199603092
928     204675465
23      237748145
37      249864742
dtype: int64
```

In [19]:

```
1 sampleAFiltered = pd.concat([ordersByUsersA[np.logical_not(ordersByUsersA['user
2     pd.Series(0,
3     index=np.arange(data['visitorsPerDateA'].sum() - len(ordersByUsersA
4     name='orders')], axis=0,)
5
6 sampleBFiltered = pd.concat([ordersByUsersB[np.logical_not(ordersByUsersB['user
7     ]['orders'],
8     pd.Series(0, index=np.arange(data['visitorsPerDateB'].sum() - len(orders
9     ),
10    name='orders')], axis=0,)
11 print('{0:.5f}'.format(stats.mannwhitneyu(sampleAFiltered, sampleBFiltered)[1]))
12 print('{0:.3f}'.format(sampleBFiltered.mean()/sampleAFiltered.mean()-1))
```

0.00652

0.173

**Вывод:** P-value меньше 0.05 поэтому отвергаем нулевую гипотезу. Как и сырыми данными статистическая значимость достигнута. Сегмент В значительно лучше сегмента А. (на 17%)

### 3.14 Посчитайте статистическую значимость различий в среднем чеке заказа между группами по «очищенным» данным.

In [20]:

```
1 print('{0:.3f}'.format(stats.mannwhitneyu(orders[np.logical_and(orders['group']
2     np.logical_not(orders['visitorId'].isin(abnormalUsers))][['revenue']][1]))
3     orders[np.logical_and(orders['group'] == 'B', np.logical_not(orders
4     )['revenue'])[1]))
5
6 print(
7     "{0:.3f}".format(
8         orders[
9             np.logical_and(
10                orders['group'] == 'B',
11                np.logical_not(orders['visitorId'].isin(abnormalUsers)),
12            )
13        ]['revenue'].mean()
14        / orders[
15            np.logical_and(
16                orders['group'] == 'A',
17                np.logical_not(orders['visitorId'].isin(abnormalUsers)),
18            )
19        ]['revenue'].mean()
20        - 1
21    )
22 )
```

0.369

-0.020

**Выводы:** По все видимости изменения которые проверяем не оказывают существенного влияния средний чек. Это хорошо показывают очищенные данные. P-value значительно больше 0.05

следовательно статистически значимых различий в среднем чеке нет. Да и относительная разница всего 2%

## 4 Общий вывод и решение по результатам теста

В рамках предоставленных данным мы не видим какие изменения в работе интернет магазина мы тестируем. Если изменения направлены на изменения среднего чека, то результат не достигнут и требуется дальнейшее продолжения теста. Визуально видно, что средний чек по группе В вырос. Правда сновной рост достигнут благодаря аномальным заказам случившимся в период тестирования. Но тем не менен изменения среднего чека по "сырым" и по "очищенным" данным не показали стаистической значимости

Если же изменения направлены на изменения конверсии, то результат достигнут. Конверсия по группе В лучше, что и показали статистические тесты по "сырым" и очищенным данным. Как видно из графика кумулятивной конверсии по группам, показатели конверсии в конце месяца достаточно стабильна.