

# Проект: "Принятие решений в бизнесе"

## Описание проекта

**Интернет-магазин:**  
Вместе с отделом маркетинга подготовлен список гипотез для увеличения выручки, были получены результаты A/B-теста.

- Задачи:**
- приоритизировать гипотезы;
  - запустить A/B-тест;
  - проанализировать результаты.

## Оглавление

- ▼ [1 Часть 1. Приоритизация гипотез](#)
  - [1.1 Метод ICE](#)
  - [1.2 Метод RICE](#)
- ▼ [2 Часть 2. Анализ A/B-теста](#)
  - [2.1 График кумулятивной выручки по группам](#)
  - [2.2 График кумулятивного среднего чека по группам](#)
  - [2.3 График относительного изменения кумулятивного среднего чека группы В к группе А](#)
  - [2.4 График кумулятивного среднего количества заказов на посетителя по группам](#)
  - [2.5 График относительного изменения кумулятивного среднего количества заказов на посетителя группы В к группе А](#)
  - [2.6 Точечный график количества заказов по пользователям](#)
  - [2.7 95-й и 99-й перцентили количества заказов на пользователя. Выбор границы для определения аномальных пользователей](#)
  - [2.8 Точечный график стоимостей заказов](#)
  - [2.9 95-й и 99-й перцентили стоимости заказов. Выбор границ для определения аномальных заказов](#)
  - [2.10 Статистическая значимость различий в среднем количестве заказов на посетителя между группами по «сырым» данным](#)
  - [2.11 Статистическая значимость различий в среднем чеке заказа между группами по «сырым» данным](#)
  - [2.12 Статистическая значимость различий в среднем количестве заказов на посетителя между группами по «очищенным» данным](#)
  - [2.13 Статистическая значимость различий в среднем чеке заказа между группами по «очищенным» данным](#)
  - [2.14 Решение по результатам теста](#)

## 1 Часть 1. Приоритизация гипотез

- Задачи:**
1. Применить фреймворк ICE для приоритизации гипотез.
  2. Применить фреймворк RICE для приоритизации гипотез.
  3. Указать, как изменилась приоритизация гипотез при применении RICE вместо ICE.

- Файл 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.

In [1]:

▼

```
1 # Импорт библиотек
2 import pandas as pd
3 import scipy.stats as stats
4 import datetime as dt
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import seaborn as sns
```

In [2]:

▼

▼

▼

```
1 # Чтение файлов
2 try:
3     df = pd.read_csv('/datasets/hypothesis.csv')
4     # если не получилось прочитать файл из локальной папки, то загружаем данные из сети
5 except:
6     df = pd.read_csv('datasets/hypothesis.csv')
```

In [3]:

▼

```
1 # Снимаем ограничение на вывод всех символов в записях
2 pd.set_option('display.max_colwidth', None)
```

```
In [4]: 1 df
```

Out[4]:

		Hypothesis	Reach	Impact	Confidence	Efforts
0		Добавить два новых канала привлечения трафика, что позволит привлекать на 30% больше пользователей	3	10	8	6
1		Запустить собственную службу доставки, что сократит срок доставки заказов	2	5	4	10
2		Добавить блоки рекомендаций товаров на сайт интернет магазина, чтобы повысить конверсию и средний чек заказа	8	3	7	3
3		Изменить структура категорий, что увеличит конверсию, т.к. пользователи быстрее найдут нужный товар	8	3	3	8
4		Изменить цвет фона главной страницы, чтобы увеличить вовлеченность пользователей	3	1	1	1
5		Добавить страницу отзывов клиентов о магазине, что позволит увеличить количество заказов	3	2	2	3
6		Показать на главной странице баннеры с актуальными акциями и распродажами, чтобы увеличить конверсию	5	3	8	3
7		Добавить форму подписки на все основные страницы, чтобы собрать базу клиентов для email-рассылок	10	7	8	5
8		Запустить акцию, дающую скидку на товар в день рождения	1	9	9	5

1.1 Метод ICE

ICE  
SCORE

=

Impact × Confidence

Efforts

```
In [5]: 1 # Создаем копию df для работы
2 df_ice = df[['Hypothesis', 'Impact', 'Confidence', 'Efforts']].copy()
```

```
In [6]: 1 # Расчёт ICE
2 df_ice['ICE'] = round(df_ice['Impact'] * df_ice['Confidence'] / df_ice['Efforts'], 2)
3
4 # Сортировка по ICE
5 df_ice = df_ice.sort_values(by='ICE', ascending=False)
6
7 # Вывод df_ice
8 df_ice
```

Out[6]:

		Hypothesis	Impact	Confidence	Efforts	ICE
8		Запустить акцию, дающую скидку на товар в день рождения	9	9	5	16.20
0		Добавить два новых канала привлечения трафика, что позволит привлекать на 30% больше пользователей	10	8	6	13.33
7		Добавить форму подписки на все основные страницы, чтобы собрать базу клиентов для email-рассылок	7	8	5	11.20
6		Показать на главной странице баннеры с актуальными акциями и распродажами, чтобы увеличить конверсию	3	8	3	8.00
2		Добавить блоки рекомендаций товаров на сайт интернет магазина, чтобы повысить конверсию и средний чек заказа	3	7	3	7.00
1		Запустить собственную службу доставки, что сократит срок доставки заказов	5	4	10	2.00
5		Добавить страницу отзывов клиентов о магазине, что позволит увеличить количество заказов	2	2	3	1.33
3		Изменить структура категорий, что увеличит конверсию, т.к. пользователи быстрее найдут нужный товар	3	3	8	1.12
4		Изменить цвет фона главной страницы, чтобы увеличить вовлеченность пользователей	1	1	1	1.00

**ВЫВОД:**  
Фреймворк ICE рекомендует начинать с 8-ой гипотезы. Последней рассматривается 4-ая.

1.2 Метод RICE

RICE  
SCORE

=

Reach × Impact × Confidence

Efforts

```
In [7]: 1 # Создаем копию df для работы
2 df_rice = df.copy()
```

```
In [8]: 1 # Расчёт RICE↔
```

Out[8]:

	Hypothesis	Reach	Impact	Confidence	Efforts	RICE
7	Добавить форму подписки на все основные страницы, чтобы собрать базу клиентов для email-рассылок	10	7	8	5	112.0
2	Добавить блоки рекомендаций товаров на сайт интернет магазина, чтобы повысить конверсию и средний чек заказа	8	3	7	3	56.0
0	Добавить два новых канала привлечения трафика, что позволит привлекать на 30% больше пользователей	3	10	8	6	40.0
6	Показать на главной странице баннеры с актуальными акциями и распродажами, чтобы увеличить конверсию	5	3	8	3	40.0
8	Запустить акцию, дающую скидку на товар в день рождения	1	9	9	5	16.2
3	Изменить структура категорий, что увеличит конверсию, т.к. пользователи быстрее найдут нужный товар	8	3	3	8	9.0
1	Запустить собственную службу доставки, что сократит срок доставки заказов	2	5	4	10	4.0
5	Добавить страницу отзывов клиентов о магазине, что позволит увеличить количество заказов	3	2	2	3	4.0
4	Изменить цвет фона главной страницы, чтобы увеличить вовлеченность пользователей	3	1	1	1	3.0

**ВЫВОД:**  
Фреймворк RICE рекомендует начинать с 7-ой гипотезы. Последней рассматривается 4-ая.

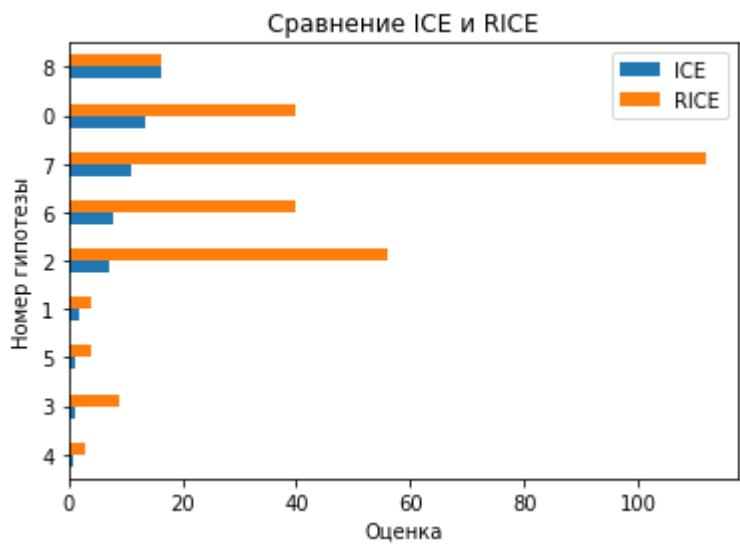
Соединим таблицы

```
In [9]: 1 # Объединение таблиц
2 df_ice.merge(df_rice)[['Hypothesis', 'ICE', 'RICE']]\
3 .set_axis(df_ice.index)\
4 .sort_values(by=['ICE', 'RICE'], ascending=False)
```

Out[9]:

	Hypothesis	ICE	RICE
8	Запустить акцию, дающую скидку на товар в день рождения	16.20	16.2
0	Добавить два новых канала привлечения трафика, что позволит привлекать на 30% больше пользователей	13.33	40.0
7	Добавить форму подписки на все основные страницы, чтобы собрать базу клиентов для email-рассылок	11.20	112.0
6	Показать на главной странице баннеры с актуальными акциями и распродажами, чтобы увеличить конверсию	8.00	40.0
2	Добавить блоки рекомендаций товаров на сайт интернет магазина, чтобы повысить конверсию и средний чек заказа	7.00	56.0
1	Запустить собственную службу доставки, что сократит срок доставки заказов	2.00	4.0
5	Добавить страницу отзывов клиентов о магазине, что позволит увеличить количество заказов	1.33	4.0
3	Изменить структура категорий, что увеличит конверсию, т.к. пользователи быстрее найдут нужный товар	1.12	9.0
4	Изменить цвет фона главной страницы, чтобы увеличить вовлеченность пользователей	1.00	3.0

```
In [10]: 1 df_ice.merge(df_rice)[['Hypothesis', 'ICE', 'RICE']]\
2 .set_axis(df_ice.index)\
3 .sort_values(by=['ICE', 'RICE'], ascending=True).plot(kind='barh', title='Сравнение ICE и RICE');
4 plt.xlabel("Оценка")
5 plt.ylabel("Номер гипотезы");
```



**Выводы:**

Два фреймворка производят приоритизацию по-разному. Вариант RICE учитывает охват пользователей (REACH) - данный показатель поднимает оценки показателей (в данном случае).

Если есть возможность оценить охват пользователей, то оптимальнее использовать RICE, однако, при одинаковом охвате достаточно ICE.

## 2 Часть 2. Анализ A/B-теста

Задача:

Проанализировать A/B-тест.

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.

```
In [11]: 1 # Чтение файлов
        2 try:
        3     orders = pd.read_csv('/datasets/orders.csv')
        4     visitors = pd.read_csv('/datasets/visitors.csv')
        5 # если не получилось прочитать файл из локальной папки, то загружаем данные из сети
        6 except:
        7     orders = pd.read_csv('datasets/orders.csv')
        8     visitors = pd.read_csv('datasets/visitors.csv')
```

Выполним предобработку данных:

```
In [12]: 1 orders.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1197 entries, 0 to 1196
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0   transactionId    1197 non-null  int64
1   visitorId        1197 non-null  int64
2   date             1197 non-null  object
3   revenue          1197 non-null  int64
4   group            1197 non-null  object
dtypes: int64(3), object(2)
memory usage: 46.9+ KB
```

```
In [13]: 1 visitors.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 62 entries, 0 to 61
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        62 non-null    object
1   group       62 non-null    object
2   visitors    62 non-null    int64
dtypes: int64(1), object(2)
memory usage: 1.6+ KB
```

```
In [14]: 1 orders.duplicated().sum()

Out[14]: 0
```

```
In [15]: 1 visitors.duplicated().sum()

Out[15]: 0
```

Полных дубликатов нет

```
In [16]: 1 orders.nunique()

Out[16]: transactionId    1197
visitorId      1031
date           31
revenue        713
group           2
dtype: int64
```

```
In [17]: 1 visitors.nunique()

Out[17]: date           31
group           2
visitors       58
dtype: int64
```

```
In [18]: 1 print('Начало теста:', orders['date'].min())
        2 print('Конец теста:', orders['date'].max())

Начало теста: 2019-08-01
Конец теста: 2019-08-31
```

```
In [19]: 1 # Число ошибочных распределений
2 orders\
3 .pivot_table(index=['visitorId'],
4             values=['group','revenue'],
5             aggfunc='nunique'
6             )\
7 .pivot_table(index=['group'],
8             values=['revenue'],
9             aggfunc='count'
10            )
```

Out[19]:

	revenue
group	
1	973
2	58

```
In [20]: 1 # id "дефектных" пользователей
2 index_ord = orders.pivot_table(index=['visitorId'],
3                               values=['group'],
4                               aggfunc='nunique'
5                               )\
6 .query('group == 2').index
```

```
In [21]: 1 orders.query('visitorId in @index_ord').sort_values(by='visitorId')
```

Out[21]:

	transactionId	visitorId	date	revenue	group
703	4293855558	8300375	2019-08-07	1790	A
71	3679129301	8300375	2019-08-01	10510	B
823	2971973105	199603092	2019-08-27	2790	A
246	437656952	199603092	2019-08-02	3488	B
26	2223239646	199603092	2019-08-15	3488	A
...	...	...	...	...	...
187	2048878902	4256040402	2019-08-17	1550	A
114	1120327437	4256040402	2019-08-01	5800	A
60	1421016313	4256040402	2019-08-16	56650	B
662	1811671147	4266935830	2019-08-29	78990	A
682	1216533772	4266935830	2019-08-29	78990	B

181 rows × 5 columns

Итоги:

- есть повторяющиеся клиенты;
- групп в тесте 2;
- 31 день: с 01/08/2019 по 31/08/2019;
- в двух группах находятся 58 покупателей;
- "дефектные" покупатели занимают 181 запись.

Так как в данных не должно быть покупателей, которые попадают в обе группы, то удалим их.

```
In [22]: 1 print('Кол-во "дефектных" записей в группе A:',
2       orders.query('visitorId in @index_ord & group == "A"')['visitorId'].count()
3       )
4
5 print('Кол-во "дефектных" записей в группе B:',
6       orders.query('visitorId in @index_ord & group == "B"')['visitorId'].count()
7       )
```

Кол-во "дефектных" записей в группе A: 89  
Кол-во "дефектных" записей в группе B: 92

```
In [23]: 1 print('Кол-во записей в группе A:',
2       orders.query('group == "A"')['visitorId'].count()
3       )
4
5 print('Кол-во записей в группе B:',
6       orders.query('group == "B"')['visitorId'].count()
7       )
```

Кол-во записей в группе A: 557  
Кол-во записей в группе B: 640

```
In [24]: 1 print('Потери записей в группе A после очистки:',
2         round((1 - (orders.query('group == "A"')['visitorId'].count()
3         - orders.query('visitorId in @index_ord & group == "A"')['visitorId'].count())
4         /orders.query('group == "A"')['visitorId'].count()) * 100, 2),
5         '%')
6     )
7
8 print('Потери записей в группе B после очистки:',
9       round((1 - (orders.query('group == "B"')['visitorId'].count()
10      - orders.query('visitorId in @index_ord & group == "B"')['visitorId'].count())
11      / orders.query('group == "B"')['visitorId'].count()) * 100, 2),
12      '%')
13  )
```

Потери записей в группе A после очистки: 15.98 %  
Потери записей в группе B после очистки: 14.38 %

```
In [25]: 1 # Данные для вычитания из visitors
2 delete_visit = orders.query('visitorId in @index_ord').pivot_table(index=['date', 'group'], values='revenue', aggfunc='count')
```

```
In [26]: 1 orders.head()
```

Out[26]:

	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

```
In [27]: 1 delete_visit.head()
```

Out[27]:

revenue		
date	group	
2019-08-01	A	1
	B	4
2019-08-02	A	1
	B	1
2019-08-03	B	2

```
In [28]: 1 visitors.head()
```

Out[28]:

	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 [29]: 1 # Вычитаем из общего кол-ва посетителей число "дефектных"
2 for visitor_ind in visitors.index:
3     # Попытка вычитания
4     try:
5         dt = visitors.loc[visitor_ind, 'date']
6         gr = visitors.loc[visitor_ind, 'group']
7
8         visitors.loc[visitor_ind, 'visitors'] = visitors.loc[visitor_ind, 'visitors'] - delete_visit.loc[dt, gr]['revenue']
9         # если не получилось, выводи дату
10    except:
11        print(visitors.loc[visitor_ind, 'date'])
```

2019-08-03  
2019-08-13  
2019-08-16  
2019-08-26  
2019-08-19  
2019-08-20

In [30]:

1

visitors.head()

Out[30]:

	date	group	visitors
0	2019-08-01	A	718
1	2019-08-02	A	618
2	2019-08-03	A	507
3	2019-08-04	A	712
4	2019-08-05	A	753

In [31]:

▼

1

# Удаляем проблемный записи

2

orders = orders.query('visitorId not in @index\_ord').sort\_values(by='visitorId')

3

print('Кол-во записей в orders после очистки:', orders.shape[0])

Кол-во записей в orders после очистки: 1016

2.1 График кумулятивной выручки по группам

In [32]:

▼

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\

▼

6

.apply(lambda x: orders[np.logical\_and(orders['date'] <= x['date'], orders['group'] == x['group'])]\

7

.agg({'date': 'max', 'group': 'max', 'transactionId': 'nunique', 'visitorId': 'nunique', 'revenue': 'sum'}), axis=1\

8

.sort\_values(by=['date','group'])

9

10

# Получаем агрегированные кумулятивные по дням данные о посетителях интернет-магазина

11

visitorsAggregated = datesGroups\

▼

12

.apply(lambda x: visitors[np.logical\_and(visitors['date'] <= x['date'], visitors['group'] == x['group'])]\

13

.agg({'date': 'max', 'group': 'max', 'visitors': 'sum'}), axis=1)\

14

.sort\_values(by=['date','group'])

15

16

# Объединяем кумулятивные данные в одной таблице и присваиваем ее столбцам понятные названия

17

cumulativeData = ordersAggregated.merge(visitorsAggregated, left\_on=['date', 'group'], right\_on=['date', 'group'])

18

cumulativeData.columns = ['date', 'group', 'orders', 'buyers', 'revenue', 'visitors']

19

20

cumulativeData.head(5)

◀

▶

Out[32]:

	date	group	orders	buyers	revenue	visitors
0	2019-08-01	A	23	19	142779	718
1	2019-08-01	B	17	17	59758	709
2	2019-08-02	A	42	36	234381	1336
3	2019-08-02	B	40	39	221801	1289
4	2019-08-03	A	66	60	346854	1843

In [33]:

▼

1

# Перевод в мин даты

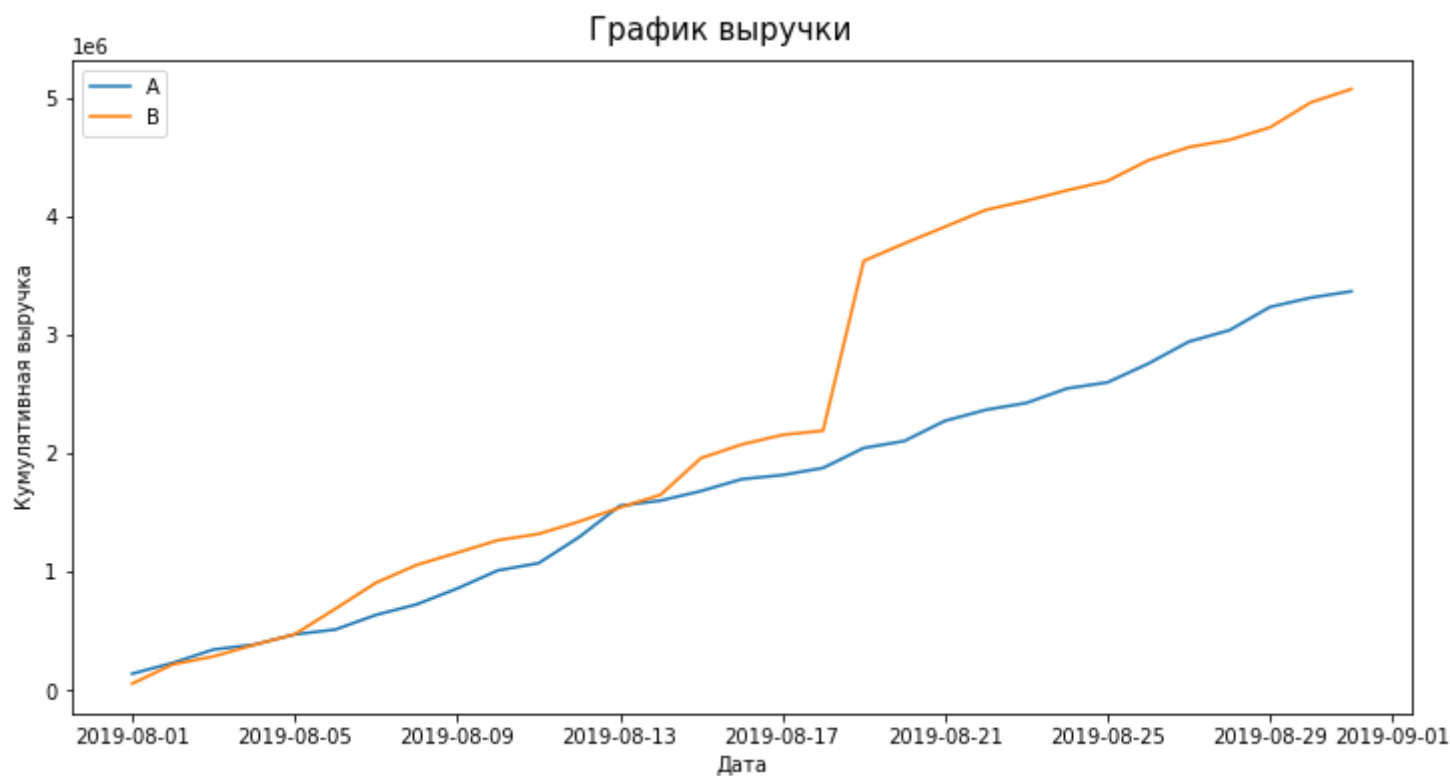
2

cumulativeData['date'] = pd.to\_datetime(cumulativeData['date'])



In [34]:

```
1 # Датафрейм с кумулятивным количеством заказов и кумулятивной выручкой по дням в группе A
2 cumulativeRevenueA = cumulativeData[cumulativeData['group']=='A'][['date', 'revenue', 'orders']]
3
4 # Датафрейм с кумулятивным количеством заказов и кумулятивной выручкой по дням в группе B
5 cumulativeRevenueB = cumulativeData[cumulativeData['group']=='B'][['date', 'revenue', 'orders']]
6
7 # Размер поля построения
8 plt.figure(figsize=(12, 6))
9
10 # Строим график выручки группы A
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.legend()
16 plt.ylabel("Кумулятивная выручка")
17 plt.xlabel("Дата")
18 plt.suptitle('График выручки', x=0.5, y=0.93, fontsize=15);
```



## Выводы:

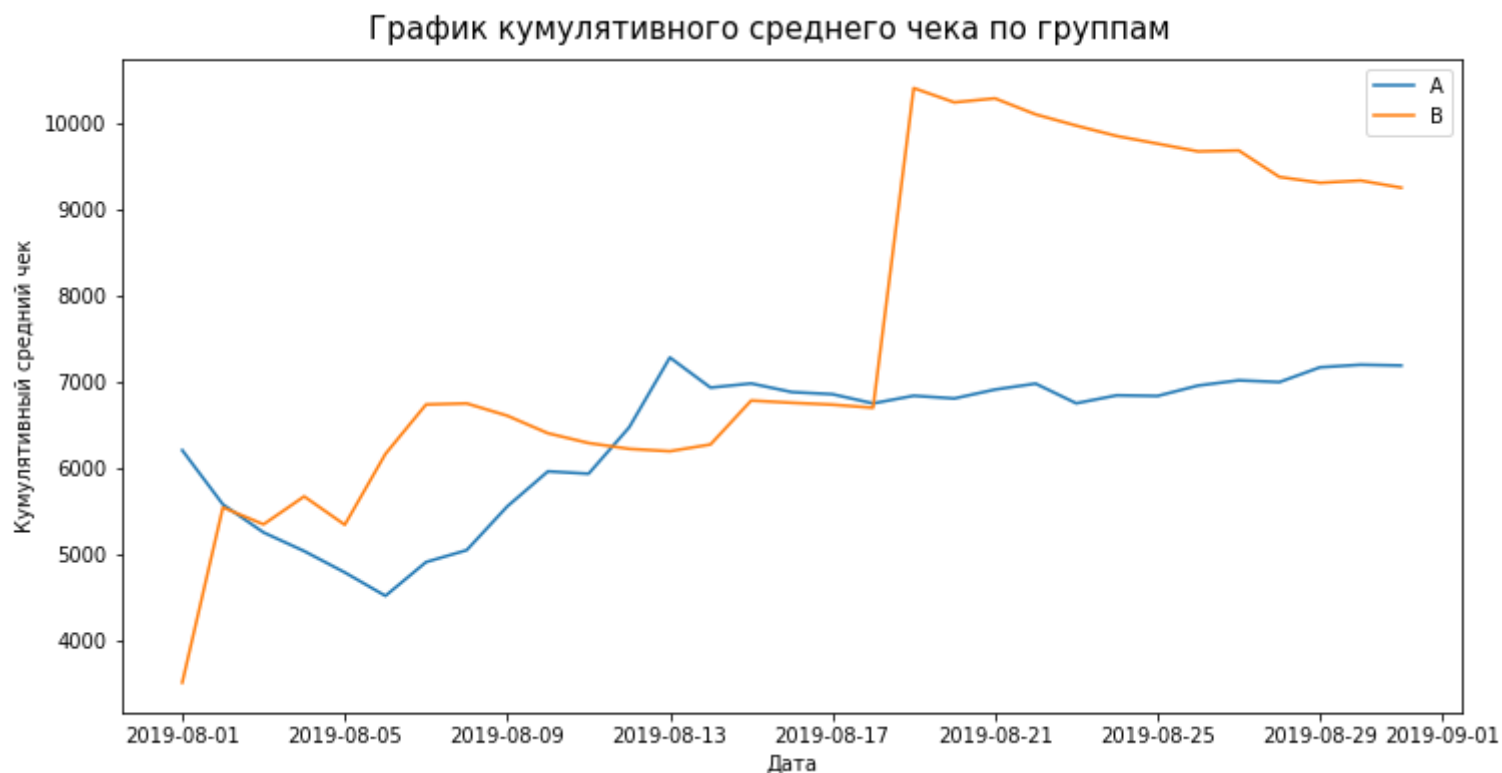
Изначально результаты групп были близкими, однако с 18-го августа наблюдается скачок группы B, после чего графики идут практически параллельно.

**Заметка:** необходимо понять причину скачка

## 2.2 График кумулятивного среднего чека по группам

In [35]:

```
1 # Размер поля построения
2 plt.figure(figsize=(12, 6))
3
4 plt.plot(cumulativeRevenueA['date'], cumulativeRevenueA['revenue']/cumulativeRevenueA['orders'], label='A')
5 plt.plot(cumulativeRevenueB['date'], cumulativeRevenueB['revenue']/cumulativeRevenueB['orders'], label='B')
6 plt.legend()
7 plt.ylabel("Кумулятивный средний чек")
8 plt.xlabel("Дата")
9 plt.suptitle('График кумулятивного среднего чека по группам', x=0.5, y=0.93, fontsize=15);
```



## Выводы:



Наблюдается скачок в величине среднего чека около 18-го числа. Вероятно, кто-то произвел крупную покупку, что отобразилось на среднем чеке - **аномалия?**

In [36]:

```
1 # Создадим таблицу
2 combo_rev = cumulativeRevenueA.merge(cumulativeRevenueB, on='date')\
3 .rename(columns={'revenue_x':'revenue_A',
4                  'revenue_y':'revenue_B',
5                  'orders_x':'orders_A',
6                  'orders_y':'orders_B'})
7
```

In [37]:

```
1 # Доп. колонки = delta
2 combo_rev['d_rev'] = combo_rev['revenue_B'] - combo_rev['revenue_A']
3 combo_rev['d_ord'] = combo_rev['orders_B'] - combo_rev['orders_A']
```

In [38]:

```
1 # Выведем таблицу
2 combo_rev.style.background_gradient(cmap='RdYlGn')
```

Out[38]:

	date	revenue_A	orders_A	revenue_B	orders_B	d_rev	d_ord
0	2019-08-01 00:00:00	142779	23	59758	17	-83021	-6
1	2019-08-02 00:00:00	234381	42	221801	40	-12580	-2
2	2019-08-03 00:00:00	346854	66	288850	54	-58004	-12
3	2019-08-04 00:00:00	388030	77	385740	68	-2290	-9
4	2019-08-05 00:00:00	474413	99	475648	89	1235	-10
5	2019-08-06 00:00:00	515332	114	690490	112	175158	-2
6	2019-08-07 00:00:00	638580	130	909654	135	271074	5
7	2019-08-08 00:00:00	727219	144	1059795	157	332576	13
8	2019-08-09 00:00:00	861456	155	1162961	176	301505	21
9	2019-08-10 00:00:00	1013731	170	1268123	198	254392	28
10	2019-08-11 00:00:00	1074396	181	1321183	210	246787	29
11	2019-08-12 00:00:00	1294788	200	1425237	229	130449	29
12	2019-08-13 00:00:00	1558426	214	1542928	249	-15498	35
13	2019-08-14 00:00:00	1601692	231	1650268	263	48576	32
14	2019-08-15 00:00:00	1682569	241	1960427	289	277858	48
15	2019-08-16 00:00:00	1782420	259	2074677	307	292257	48
16	2019-08-17 00:00:00	1817160	265	2155542	320	338382	55
17	2019-08-18 00:00:00	1876741	278	2190865	327	314124	49
18	2019-08-19 00:00:00	2044934	299	3620785	348	1575851	49
19	2019-08-20 00:00:00	2103613	309	3768059	368	1664446	59
20	2019-08-21 00:00:00	2273782	329	3908406	380	1634624	51
21	2019-08-22 00:00:00	2366147	339	4050134	401	1683987	62
22	2019-08-23 00:00:00	2424010	359	4127403	414	1703393	55
23	2019-08-24 00:00:00	2546023	372	4215269	428	1669246	56
24	2019-08-25 00:00:00	2597907	380	4294526	440	1696619	60
25	2019-08-26 00:00:00	2755121	396	4467965	462	1712844	66
26	2019-08-27 00:00:00	2940471	419	4579312	473	1638841	54
27	2019-08-28 00:00:00	3036933	434	4640510	495	1603577	61
28	2019-08-29 00:00:00	3233233	451	4746610	510	1513377	59
29	2019-08-30 00:00:00	3311413	460	4955833	531	1644420	71
30	2019-08-31 00:00:00	3364656	468	5068972	548	1704316	80

**Выводы:**

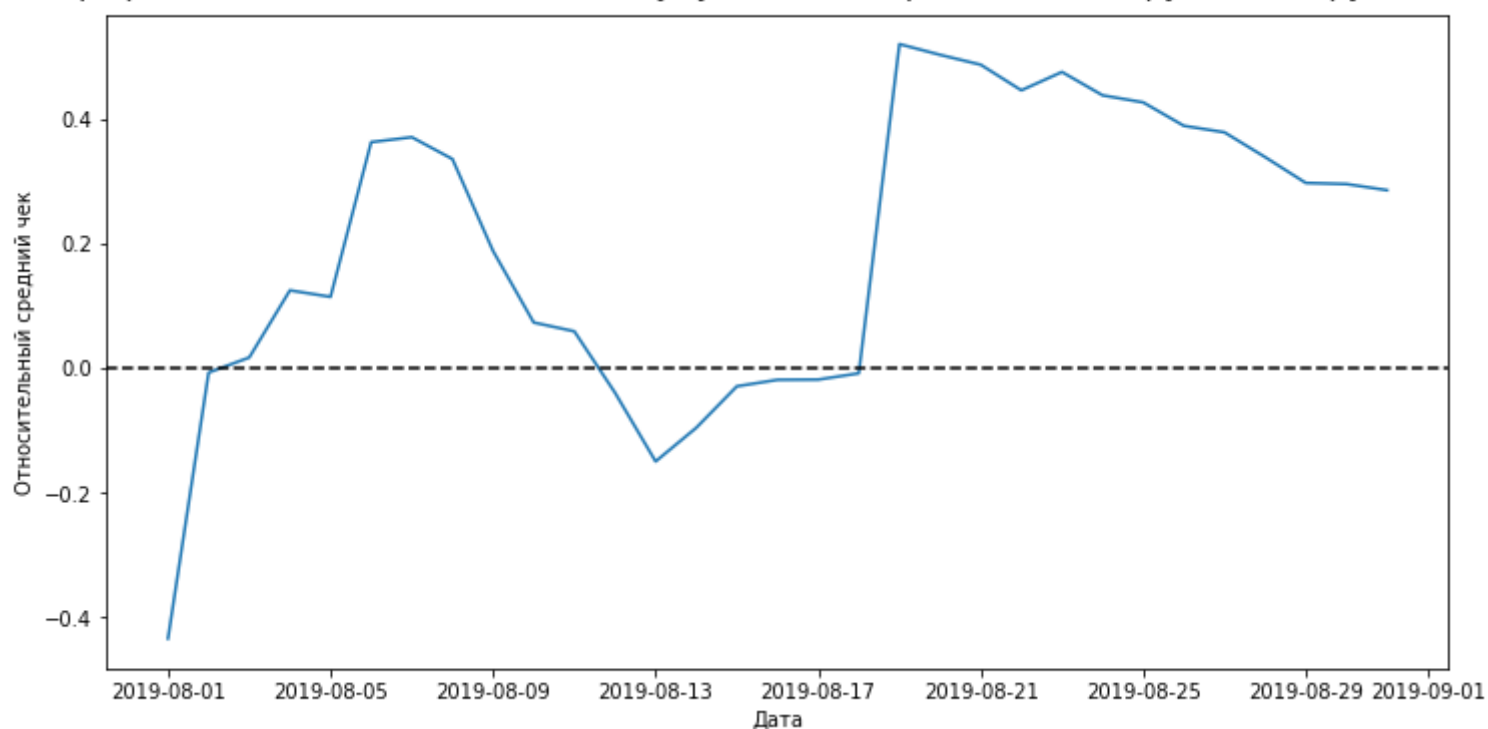
В целом, наблюдаются лучшие показатели у группы В, но нужно внимательнее рассмотреть 19-ое число, так как они могут изменить картину.

## 2.3 График относительного изменения кумулятивного среднего чека группы В к группе А

In [39]:

```
1 # Размер поля построения
2 plt.figure(figsize=(12, 6))
3
4 # Собираем данные в одном датафрейме
5 mergedCumulativeRevenue = cumulativeRevenueA.merge(cumulativeRevenueB,
6                                                     left_on='date',
7                                                     right_on='date',
8                                                     how='left',
9                                                     suffixes=['A', 'B'])
10
11
12 # Строим отношение средних чеков
13 plt.plot(mergedCumulativeRevenue['date'],
14          (mergedCumulativeRevenue['revenueB'] / mergedCumulativeRevenue['ordersB']) \
15          / (mergedCumulativeRevenue['revenueA'] / mergedCumulativeRevenue['ordersA']) - 1)
16
17 # Добавляем ось X
18 plt.axhline(y=0, color='black', linestyle='--')
19
20 plt.suptitle('График относительного изменения кумулятивного среднего чека группы В к группе А',
21             x=0.5,
22             y=0.93,
23             fontsize=15)
24
25
26 plt.ylabel("Относительный средний чек")
27 plt.xlabel("Дата");
```

График относительного изменения кумулятивного среднего чека группы В к группе А



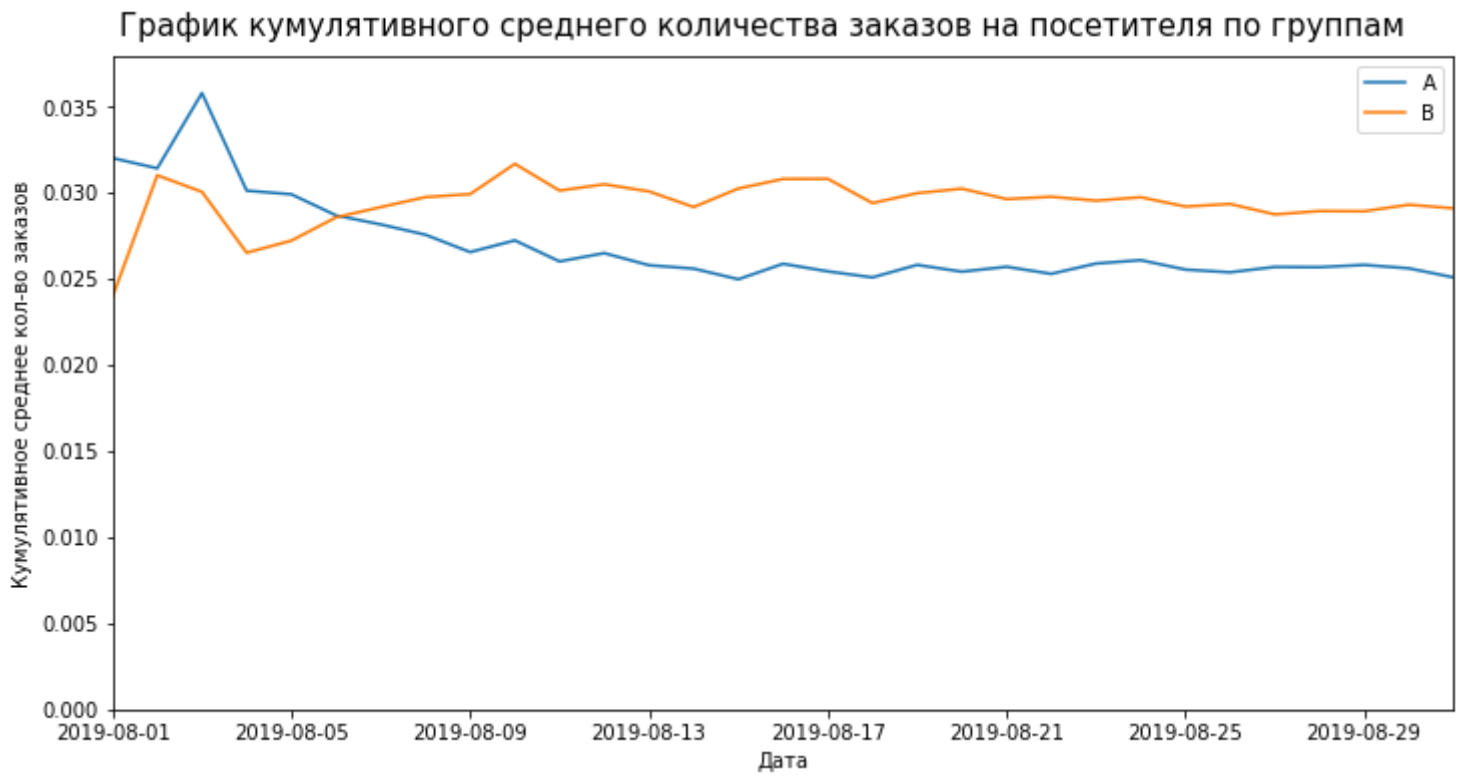
### Выводы:

Рост, падение, малый рост, скачок, плавный спуск - скачок 19-го числа. Установившегося нет.

## 2.4 График кумулятивного среднего количества заказов на посетителя по группам

In [40]:

```
1 import datetime as dt
2
3 # Размер поля построения
4 plt.figure(figsize=(12, 6))
5
6 # Считаем кумулятивную конверсию
7 cumulativeData['conversion'] = cumulativeData['orders'] / cumulativeData['visitors']
8 # Отделяем данные по группе A
9 cumulativeDataA = cumulativeData[cumulativeData['group'] == 'A']
10 # Отделяем данные по группе B
11 cumulativeDataB = cumulativeData[cumulativeData['group'] == 'B']
12 # Строим графики
13 plt.plot(cumulativeDataA['date'], cumulativeDataA['conversion'], label = 'A')
14 plt.plot(cumulativeDataB['date'], cumulativeDataB['conversion'], label = 'B')
15 plt.legend()
16
17 # Задаем масштаб осей
18 plt.axis([dt.date(2019, 8, 1), dt.date(2019, 8, 31), 0, 0.038])
19
20 # Название
21 plt.suptitle('График кумулятивного среднего количества заказов на посетителя по группам',
22             x=0.5,
23             y=0.93,
24             fontsize=15
25             )
26
27 plt.ylabel("Кумулятивное среднее кол-во заказов")
28 plt.xlabel("Дата");
```



### Выводы:

Количество заказов в группе В стабильно больше, чем в группе А, начиная с 9-го числа.

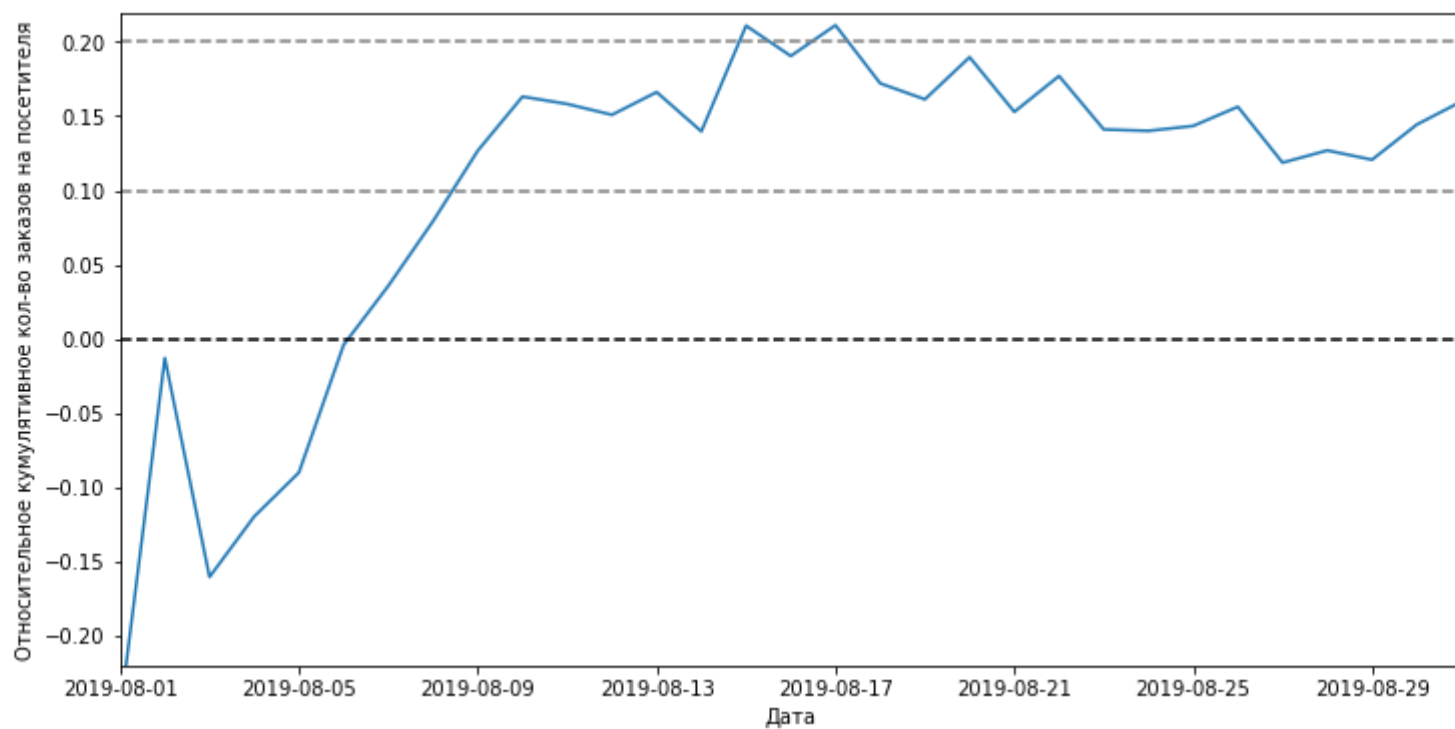
График количества заказов уже больше половины месяца похож на прямую, не имеет резких изменений. Количество заказов можно считать стабилизировавшимся.

(Конечно, дальше сентябрь, а это август - предшкольная пора, вероятно, дальше будет *синхронный* спуск).

## 2.5 График относительного изменения кумулятивного среднего количества заказов на посетителя группы В к группе А

```
In [41]: 1 # Размер поля построения
2 plt.figure(figsize = (12, 6))
3
4 mergedCumulativeConversions = cumulativeDataA[['date', 'conversion']]\
5 .merge(cumulativeDataB[['date', 'conversion']],
6       left_on='date',
7       right_on='date',
8       how='left',
9       suffixes=['A', 'B'])
10
11
12 plt.plot(mergedCumulativeConversions['date'],
13         mergedCumulativeConversions['conversionB']/mergedCumulativeConversions['conversionA']-1,
14         label="Относительный прирост конверсии группы В относительно группы А"
15         )
16
17 plt.axhline(y=0, color='black', linestyle='--')
18 plt.axhline(y=0.1, color='grey', linestyle='--')
19 plt.axhline(y=0.2, color='grey', linestyle='--')
20
21 plt.axis([dt.date(2019, 8, 1), dt.date(2019, 8, 31), -0.22, 0.22])
22
23 # Название
24 plt.suptitle('График относительного изменения кумулятивного среднего количества заказов на посетителя группы В к группе А'
25            x=0.5,
26            y=0.93,
27            fontsize=15
28            )
29
30 plt.ylabel("Относительное кумулятивное кол-во заказов на посетителя")
31 plt.xlabel("Дата");
```

График относительного изменения кумулятивного среднего количества заказов на посетителя группы В к группе А



### Выводы:

Подтверждает предыдущий вывод:

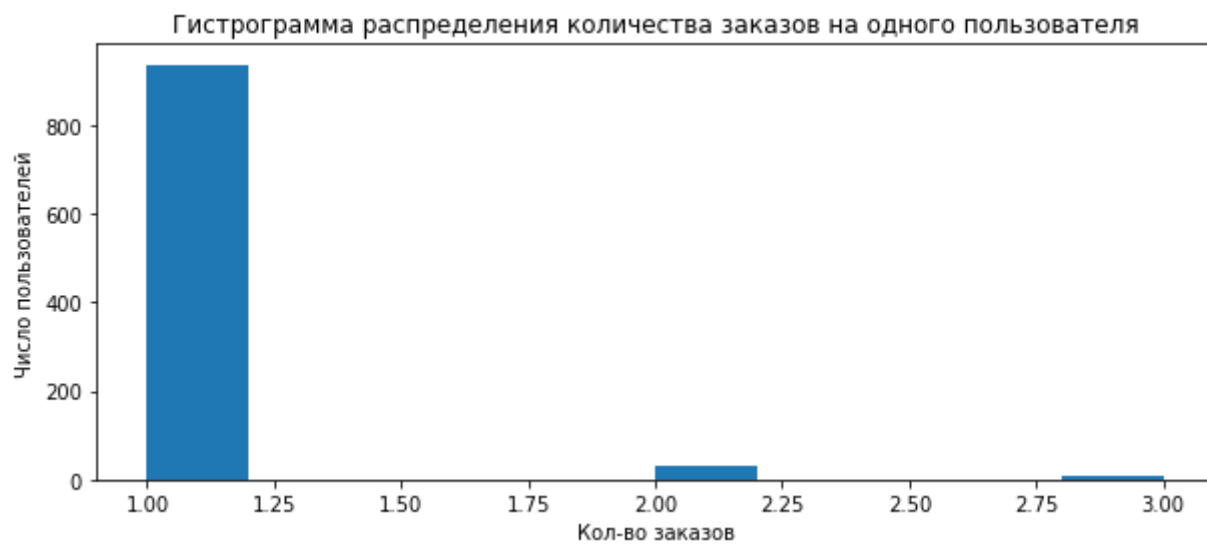
Количество заказов в группе В стабильно больше, чем в группе А, начиная с 9-го числа.

## 2.6 Точечный график количества заказов по пользователям

```
In [42]: 1 # Массив для анализа
2 ordersByUsers = orders.groupby('visitorId', as_index = False)\
3 .agg({'transactionId': 'nunique'})\
4 .sort_values(by='transactionId', ascending = False) # для удобства
5
6 ordersByUsers.columns = ['visitorId', 'orders']
```

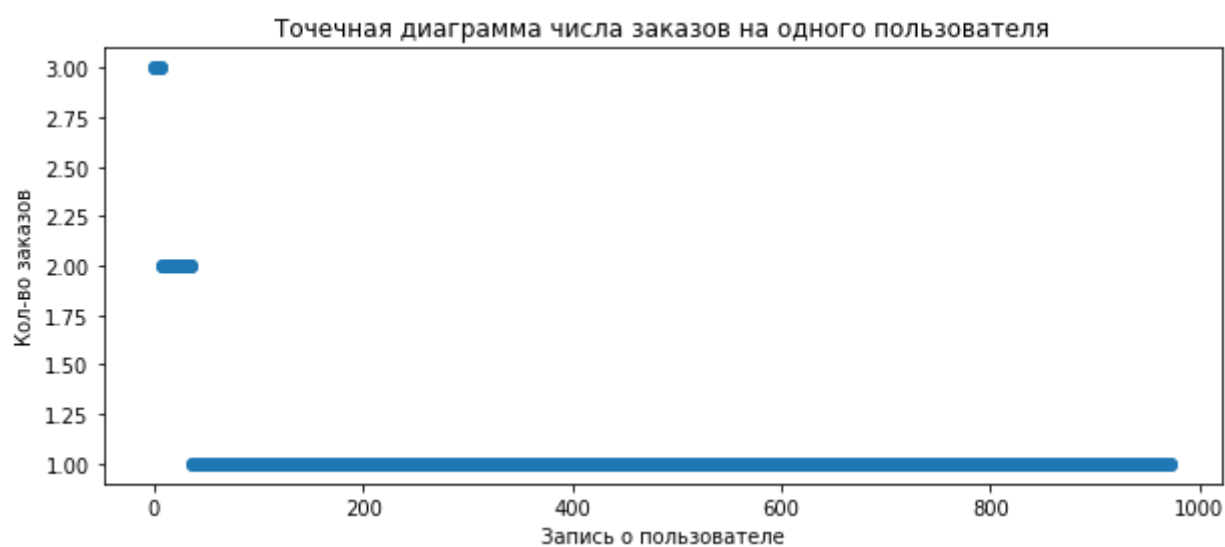
In [43]:

```
1 # Размер поля построения
2 plt.figure(figsize=(10,4))
3
4 plt.hist(ordersByUsers['orders'], bins = 10)
5 plt.title('Гистограмма распределения количества заказов на одного пользователя')
6 plt.ylabel("Число пользователей")
7 plt.xlabel("Кол-во заказов")
8 plt.show();
```



In [44]:

```
1 # Размер поля построения
2 plt.figure(figsize=(10,4))
3
4 x_values = pd.Series(range(0, len(ordersByUsers)))
5 plt.scatter(x_values, ordersByUsers['orders'])
6 plt.title('Точечная диаграмма числа заказов на одного пользователя')
7 plt.ylabel("Кол-во заказов")
8 plt.xlabel("Запись о пользователе")
9 plt.show()
```



## Выводы:

Большинство совершило 1 или 2 заказа.

## 2.7 95-й и 99-й перцентили количества заказов на пользователя. Выбор границы для определения аномальных пользователей

In [45]:

```
1 print('Перцентили количества заказов на пользователя [95, 99]:',
2       np.percentile(ordersByUsers['orders'], [95, 99])
3       )
```

Перцентили количества заказов на пользователя [95, 99]: [1. 2.]

## Выводы:

Не более 1% пользователей оформляли больше 2-ух заказов - выберем это значение в качестве границы аномальных пользователей.

## 2.8 Точечный график стоимостей заказов

```
In [46]: 1 # Размер поля построения
2 plt.figure(figsize=(10,4))
3
4 plt.hist(orders['revenue'], range=(0,100000), bins = 100)
5 plt.title('Гистограмма распределения сумм заказов на одного пользователя')
6 plt.ylabel("Число пользователей")
7 plt.xlabel("Сумма заказа")
8 plt.show()
```



```
In [47]: 1 x_values = pd.Series(range(0, len(orders)))
2
3 # Размер поля построения
4 plt.figure(figsize=(10,4))
5
6 plt.scatter(x_values, orders['revenue'])
7 plt.title('Точечная диаграмма сумм заказов на одного пользователя')
8 plt.axis([0, 1050, 0, 110000])
9 plt.ylabel("Сумма заказа")
10 plt.xlabel("Запись о пользователе")
11 plt.show()
```



### Выводы:

В подавляющем большинстве суммы заказов не превышали 30 000 Р.

## 2.9 95-й и 99-й перцентили стоимости заказов. Выбор границ для определения аномальных заказов

```
In [48]: 1 print('Перцентили количества заказов на пользователя [95, 99]:',
2       np.percentile(orders['revenue'], [95, 99])
3       )
```

Перцентили количества заказов на пользователя [95, 99]: [26785. 53904.]

### Выводы:

У не более 5% пользователей стоимость заказа была более 26 785 Р - будем считать границей.

## 2.10 Статистическая значимость различий в среднем количестве заказов на посетителя между группами по «сырым» данным

```
In [49]: 1 # Информация о группах
2 ordersByUsersA = orders[orders['group']=='A'].groupby('visitorId', as_index=False).agg({'transactionId': 'nunique'})
3 ordersByUsersA.columns = ['visitorId', 'orders']
4 ordersByUsersB = orders[orders['group']=='B'].groupby('visitorId', as_index=False).agg({'transactionId': 'nunique'})
5 ordersByUsersB.columns = ['visitorId', 'orders']
6
7 print('Кол-во покупателей в группе A: {}'.format(len/ordersByUsersA)))
8 print('Кол-во покупок в группе A: {}'.format/ordersByUsersA['orders'].sum()))
9 print('Кол-во посетителей в группе A: {}'.format(visitors[visitors['group']=='A']['visitors'].sum()))
10
11 print('Кол-во покупателей в группе B: {}'.format(len/ordersByUsersB)))
12 print('Кол-во покупок в группе B: {}'.format/ordersByUsersB['orders'].sum()))
13 print('Кол-во посетителей в группе B: {}'.format(visitors[visitors['group']=='B']['visitors'].sum()))
```

Кол-во покупателей в группе A: 445  
Кол-во покупок в группе A: 468  
Кол-во посетителей в группе A: 18647

Кол-во покупателей в группе B: 528  
Кол-во покупок в группе B: 548  
Кол-во посетителей в группе B: 18824

### Гипотезы

Согласно рейтингу ICE была принята к проверке гипотеза:

Запустить акцию, дающую скидку на товар в день рождения

Для проверки применим статистический критерий Манна-Уитни к полученным выборкам,  $\alpha = 5\%$ .

### Проверка №1

Относительный прирост конверсии группы B по отношению к группе A:

$H_0$ : Запуск акции, дающей скидку на товар в день рождения вызовет прирост конверсии.  
 $H_1$ : Запуск акции, дающей скидку на товар в день рождения не изменит конверсию.

```
In [50]: 1 # Составим списки кол-ва заказов sampleA и sampleB со всеми пользователями (+ без покупок) по группам
2 list_orders_1 = []
3 for i in range(0, (visitors[visitors['group']=='A']['visitors'].sum() - len/ordersByUsersA)):
4     list_orders_1.append(0)
5 orders_by_non_purchased_users_A = pd.Series(data = list_orders_1, name = 'orders')
6
7 list_orders_2 = []
8 for i in range(0, (visitors[visitors['group']=='B']['visitors'].sum() - len/ordersByUsersB)):
9     list_orders_2.append(0)
10 orders_by_non_purchased_users_B = pd.Series(data = list_orders_2, name = 'orders')
11
12 sampleA = pd.concat([ordersByUsersA['orders'], orders_by_non_purchased_users_A], axis=0)
13 sampleB = pd.concat([ordersByUsersB['orders'], orders_by_non_purchased_users_B], axis=0)
14 print('Относительный прирост конверсии группы B по отношению к группе A: {:.3f}'.format(sampleB.mean()/sampleA.mean()-1))
15
16 alpha = 0.05 # выбранный уровень
17 results = stats.mannwhitneyu(sampleA, sampleB) # тест
18
19 print('P-value: {}'.format(results.pvalue), '\n')
20
21 if results.pvalue < alpha:
22     print('Вывод:')
23     print('Различие в количестве заказов СТАТИСТИЧЕСКИ ЗНАЧИМО')
24 else:
25     print('Вывод:')
26     print('Различие в количестве заказов нельзя считать значимым')
```

Относительный прирост конверсии группы B по отношению к группе A: 0.160  
P-value: 0.0109567830835148

Вывод:  
Различие в количестве заказов СТАТИСТИЧЕСКИ ЗНАЧИМО

### Выводы:

С учетом пользователей без покупок различие в среднем количестве заказов по «сырым» данным является **статистически значимым**.

## 2.11 Статистическая значимость различий в среднем чеке заказа между группами по «сырым» данным

### Проверка №2

Относительное изменение среднего чека в группе B по отношению группы A:



$H_0$ : Запуск акции, дающей скидку на товар в день рождения увеличит средний чек заказа.  
 $H_1$ : Запуск акции, дающей скидку на товар в день рождения не изменит средний чек заказа.

```
In [51]: 1 print(
2     'Относительное изменение среднего чека в группе В по отношению группы А: {:.3f}'
3     .format(
4         orders[orders['group'] == 'B']['revenue'].mean() / orders[orders['group'] == 'A']['revenue'].mean()-1
5     )
6 )
7
8 results = stats.mannwhitneyu(
9     orders[orders['group'] == 'A']['revenue'], orders[orders['group'] == 'B']['revenue']
10 )
11
12 print('P-value: {}'.format(results.pvalue), '\n')
13
14 if results.pvalue < alpha:
15     print('Вывод:')
16     print('Различия в среднем чеке заказа СТАТИСТИЧЕСКИ ЗНАЧИМЫ')
17 else:
18     print('Вывод:')
19     print('Различия в среднем чеке заказа нельзя считать статистически значимыми')
```

Относительное изменение среднего чека в группе В по отношению группы А: 0.287  
P-value: 0.8294908998149533

Вывод:  
Различия в среднем чеке заказа нельзя считать статистически значимыми

**Выводы:**

Статистическую значимость **различий в среднем чеке** заказа между группами по «сырым» **нельзя считать значимой**.

2.12 Статистическая значимость различий в среднем количестве заказов на посетителя между группами по «очищенным» данным

```
In [52]: 1 # Пределы для аномалий
2 limit_orders = 2
3 limit_revenue = 26785
```

```
In [53]: 1 # Кол-во пользователей с аномалиями
2 usersWithManyOrders = pd.concat([ordersByUsersA[ordersByUsersA['orders']>limit_orders]['visitorId'],
3     ordersByUsersB[ordersByUsersB['orders']>limit_orders]['visitorId']],
4     axis = 0
5 )
6
7 usersWithExpensiveOrders = orders[orders['revenue'] > limit_revenue]['visitorId']
8
9 abnormalUsers = pd.concat([usersWithManyOrders,
10     usersWithExpensiveOrders],
11     axis = 0
12 ).drop_duplicates().sort_values()
13
14 print('Кол-во пользователей с аномалиями: {}'.format(len(abnormalUsers)),
15     'шт.'
16 )
17
```

Кол-во пользователей с аномалиями: 58 шт.

**Провека №3**

Относительный прирост конверсии группы В по отношению к группе А:

$H_0$ : Запуск акции, дающей скидку на товар в день рождения вызовет прирост конверсии.  
 $H_1$ : Запуск акции, дающей скидку на товар в день рождения не изменит конверсию.

```
In [54]: 1 sampleAFiltered = pd.concat([
2     ordersByUsersA[np.logical_not(ordersByUsersA['visitorId'].isin(abnormalUsers))]['orders'],
3     orders_by_non_purchased_users_A
4 ], axis = 0
5 )
6
7 sampleBFiltered = pd.concat([
8     ordersByUsersB[np.logical_not(ordersByUsersB['visitorId'].isin(abnormalUsers))]['orders'],
9     orders_by_non_purchased_users_B
10 ], axis = 0
11 )
12
13 print('Относительный прирост конверсии группы В по отношению к группе А после очистки данных: {:.3f}'
14       .format(sampleBFiltered.mean() / sampleAFiltered.mean() - 1))
15
16 results = stats.mannwhitneyu(sampleAFiltered,
17                               sampleBFiltered
18                               )
19
20 print('P-value: {}'.format(results.pvalue), '\n')
21
22 if results.pvalue < alpha:
23     print('Вывод:')
24     print('В среднем количестве заказов на посетителя РАЗНИЦА СТАТИСТИЧЕСКИ ЗНАЧИМА')
25 else:
26     print('Вывод:')
27     print('В среднем количестве заказов на посетителя разницу нельзя считать значимой')
```

Относительный прирост конверсии группы В по отношению к группе А после очистки данных: 0.182  
P-value: 0.012285273602771688

Вывод:  
В среднем количестве заказов на посетителя РАЗНИЦА СТАТИСТИЧЕСКИ ЗНАЧИМА

Выводы:

Статистическая значимость различий в среднем количестве заказов на посетителя между группами по «очищенным» данным СТАТИСТИЧЕСКИ ЗНАЧИМА.

Запуск акции, дающей скидку на товар в день рождения вызвал прирост конверсии. Группа В на 18,9% лучше А.

2.13 Статистическая значимость различий в среднем чеке заказа между группами по «очищенным» данным

Проверка №4

Относительное изменение среднего чека в группе В по отношению группы А:

Н<sub>0</sub>: Запуск акции, дающей скидку на товар в день рождения увеличит средний чек заказа.  
Н<sub>1</sub>: Запуск акции, дающей скидку на товар в день рождения не изменит средний чек заказа.

```
In [55]: 1 print('Относительное изменение среднего чека в группе В по отношению группы А после очистки данных: {:.3f}'
2       .format(
3     orders[np.logical_and(orders['group'] == 'B', np.logical_not(orders['visitorId'].isin(abnormalUsers)
4                                                                    )
5                                                                    )]['revenue'].mean()\
6     / orders[np.logical_and(orders['group'] == 'A',
7                             np.logical_not(orders['visitorId'].isin(abnormalUsers))
8                             )]['revenue'].mean() - 1
9     )
10    )
11
12
13
14 results = stats.mannwhitneyu(orders[np.logical_and(orders['group'] == 'A',
15                                                     np.logical_not(orders['visitorId'].isin(abnormalUsers))
16                                                     )]['revenue'],
17                               orders[np.logical_and(orders['group'] == 'B',
18                                                       np.logical_not(orders['visitorId'].isin(abnormalUsers))
19                                                       )]['revenue'])
20
21 print('P-value: {}'.format(results.pvalue), '\n')
22
23 if results.pvalue < alpha:
24     print('Вывод:')
25     print('Различия в среднем чеке заказа СТАТИСТИЧЕСКИ ЗНАЧИМЫ')
26 else:
27     print('Вывод:')
28     print('Различия в среднем чеке заказа нельзя считать статистически значимыми')
```

Относительное изменение среднего чека в группе В по отношению группы А после очистки данных: -0.048  
P-value: 0.6458964038091206

Вывод:  
Различия в среднем чеке заказа нельзя считать статистически значимыми

Выводы:

Статистическую значимость **различий в среднем чеке** заказа между группами по «очищенным» данным **нельзя считать статистически значимой**.

запуск акции, дающей скидку на товар в день рождения не изменил средний чек заказа.

До очистки 0.287, после -0.022. Это говорит о значительном влиянии аномалий.

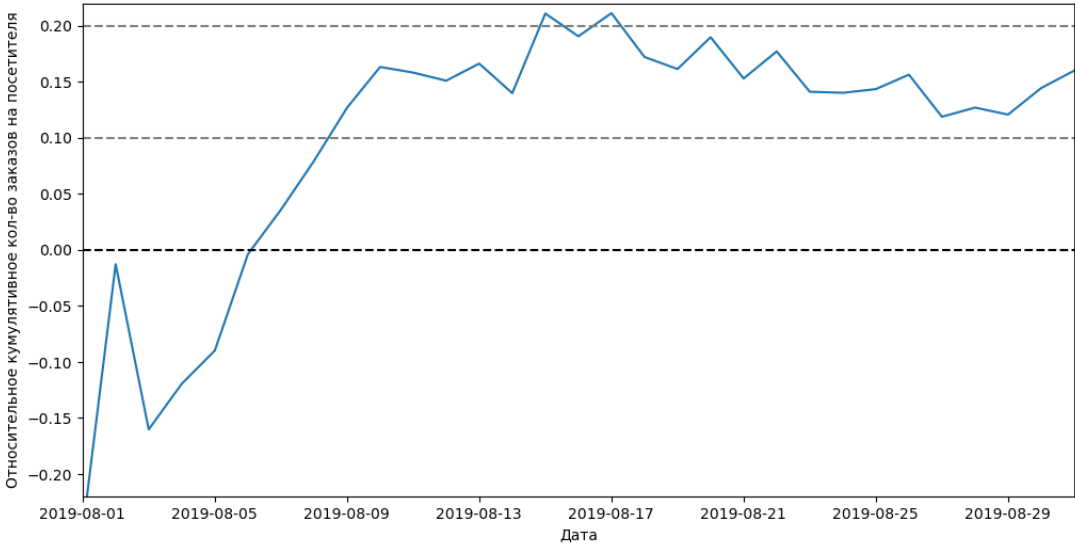
2.14 Решение по результатам теста

Выводы

- Различия в количестве заказов между группами по «сырым» и по «очищенным» отсутствуют. Между группами А и В они **статистически значимы**.
- Различия в среднем чеке заказа между группами по «сырым» и по «очищенным» отсутствуют. Между группами А и В их нельзя считать статистически значимыми.
- График различия среднего количества заказов между группами сообщает, что результаты группы В лучше группы А, кол-во заказов продолжает расти.

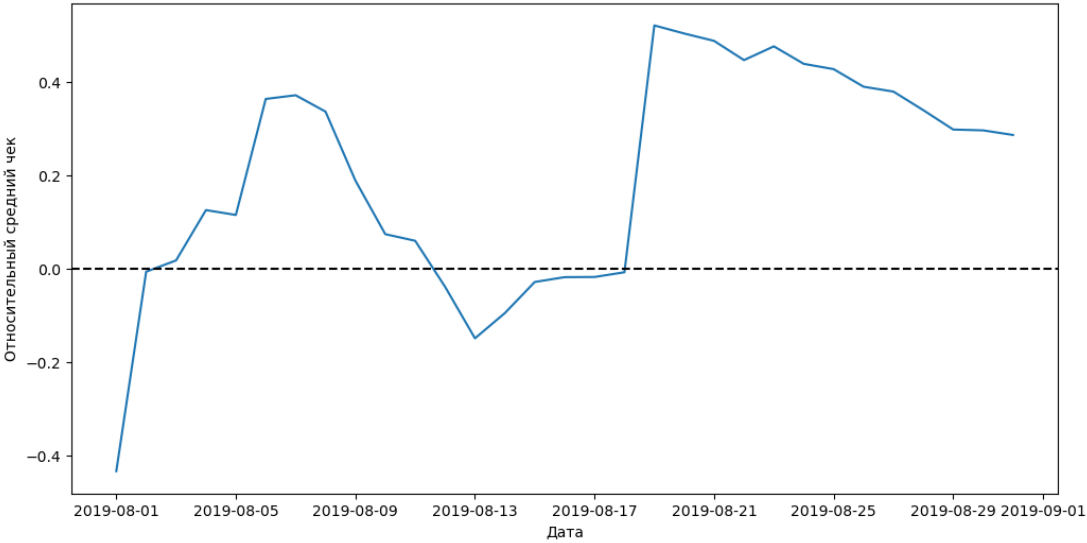
"Количество заказов в группе В стабильно больше, чем в группе А, начиная с 9-го числа"

График относительного изменения кумулятивного среднего количества заказов на посетителя группы В к группе А



- График различия среднего чека говорит о том, что результаты группы В лучше результатов группы А:

График относительного изменения кумулятивного среднего чека группы В к группе А



В кумулятивных метриках наблюдается лидерство группы В, однако были обнаружены "всплески" на графиках выше. Была проведена очистка аномалий, после которой различия между группами уменьшились.

До очистки показатель конверсии 16,0%, а после 18,9% - разница в 2,9%.

Исходя из обнаруженных фактов, тест *можно* остановить и стоит **признать победу группы В**.

Заметки для других отделов:

Также стоит отметить, что в данных были обнаружены пользователи, которые попали в обе группы тестирования - они были удалены, однако стоит рассмотреть механизм распределения, так как он может в дальнейшем повлиять на другие тесты. Ниже предствалены даты, в которых нет "проблемных" пользователей.

2019-08-03  
2019-08-13  
2019-08-16  
2019-08-26  
2019-08-19  
2019-08-20

