

Проект по А/В-тестированию

Задача: провести оценку результатов А/В-теста.

Цель: оценить корректность проведения теста и проанализировать его результаты.

Исходные данные:

Датасет с действиями пользователей, техническое задание и несколько вспомогательных датасетов.

Техническое задание:

- Название теста: `recommender_system_test` ;
- Группы: А (контрольная), В (новая платёжная воронка);
- Дата запуска: 2020-12-07 ;
- Дата остановки набора новых пользователей: 2020-12-21 ;
- Дата остановки: 2021-01-04 ;
- Аудитория: 15% новых пользователей из региона EU ;
- Назначение теста: тестирование изменений, связанных с внедрением улучшенной рекомендательной системы;
- Ожидаемое количество участников теста: 6000 .
- Ожидаемый эффект: за 14 дней с момента регистрации в системе пользователи покажут улучшение каждой метрики не менее, чем на 10% :
 - конверсии в просмотр карточек товаров — событие `product_page` ;
 - просмотры корзины — `product_cart` ;
 - покупки — `purchase` .

Данные:

- `ab_project_marketing_events.csv`
- `final_ab_new_users.csv`
- `final_ab_events.csv`
- `final_ab_participants.csv`

`/datasets/ab_project_marketing_events.csv` — календарь маркетинговых событий на 2020 год;

Структура файла:

- `name` — название маркетингового события;
- `regions` — регионы, в которых будет проводиться рекламная кампания;
- `start_dt` — дата начала кампании;
- `finish_dt` — дата завершения кампании.

`/datasets/final_ab_new_users.csv` — все пользователи, зарегистрировавшиеся в интернет-магазине в период с 7 по 21 декабря 2020 года;

Структура файла:

- `user_id` — идентификатор пользователя;
- `first_date` — дата регистрации;
- `region` — регион пользователя;
- `device` — устройство, с которого происходила регистрация.

`/datasets/final_ab_events.csv` — все события новых пользователей в период с 7 декабря 2020 по 4 января 2021 года;

Структура файла:

- `user_id` — идентификатор пользователя;
- `event_dt` — дата и время события;
- `event_name` — тип события;
- `details` — дополнительные данные о событии. Например, для покупок, `purchase`, в этом поле хранится стоимость покупки в долларах.

`/datasets/final_ab_participants.csv` — таблица участников тестов.

Структура файла:

- `user_id` — идентификатор пользователя;
- `ab_test` — название теста;
- `group` — группа пользователя.

Содержание

- ▼ [1 Исследование данных](#)
 - [1.1 Пропуски и типы данных](#)
 - [1.2 Дубликаты](#)
- [2 Оценка корректности проведения теста](#)
- [3 Исследовательский анализ данных](#)
- [4 Оценка результатов А/В-тестирования](#)

```
In [1]: 1 # Импорт библиотек
2 import pandas as pd
3 import plotly.express as px
4 from plotly import graph_objects as go
5 import matplotlib.pyplot as plt
6 import math
7 from scipy import stats as st
```

```
In [2]: 1 # Чтение файлов
2 try:
3     df_ab_project_marketing_events = pd.read_csv('datasets/ab_project_marketing_events.csv',
4                                                    parse_dates=['start_dt', 'finish_dt'])
5
6     df_final_ab_new_users = pd.read_csv('datasets/final_ab_new_users.csv', parse_dates=['first_date'])
7
8     df_final_ab_events = pd.read_csv('datasets/final_ab_events.csv', parse_dates=['event_dt'])
9
10    df_final_ab_participants = pd.read_csv('datasets/final_ab_participants.csv')
11
12 except:
13     df_ab_project_marketing_events = pd.read_csv('/datasets/ab_project_marketing_events.csv',
14                                                    parse_dates=['start_dt', 'finish_dt'])
15
16     df_final_ab_new_users = pd.read_csv('/datasets/final_ab_new_users.csv',
17                                           parse_dates=['first_date'])
18
19     df_final_ab_events = pd.read_csv('/datasets/final_ab_events.csv',
20                                       parse_dates=['event_dt'])
21
22     df_final_ab_participants = pd.read_csv('/datasets/final_ab_participants.csv')
```

1 Исследование данных

1.1 Пропуски и типы данных

```
In [3]: 1 # Информация о данных, пропусках
2 for df, name in zip([df_ab_project_marketing_events,
3                     df_final_ab_new_users,
4                     df_final_ab_events,
5                     df_final_ab_participants],
6                     ['df_ab_project_marketing_events',
7                     'df_final_ab_new_users',
8                     'df_final_ab_events',
9                     'df_final_ab_participants']):
10     print('Датасет', name)
11     display(df.head())
12     display(df.info())
13     print('Уникальные значения:')
14     display(df.nunique())
15     print()
```

Датасет df_ab_project_marketing_events

	name	regions	start_dt	finish_dt
0	Christmas&New Year Promo	EU, N.America	2020-12-25	2021-01-03
1	St. Valentine's Day Giveaway	EU, CIS, APAC, N.America	2020-02-14	2020-02-16
2	St. Patric's Day Promo	EU, N.America	2020-03-17	2020-03-19
3	Easter Promo	EU, CIS, APAC, N.America	2020-04-12	2020-04-19
4	4th of July Promo	N.America	2020-07-04	2020-07-11

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14 entries, 0 to 13
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  ---
0    name         14 non-null     object
1    regions      14 non-null     object
2    start_dt     14 non-null     datetime64[ns]
3    finish_dt    14 non-null     datetime64[ns]
dtypes: datetime64[ns](2), object(2)
memory usage: 576.0+ bytes
```

None

Уникальные значения:

```
name         14
regions       6
start_dt     14
finish_dt    14
dtype: int64
```

Датасет df_final_ab_new_users

	user_id	first_date	region	device
0	D72A72121175D8BE	2020-12-07	EU	PC
1	F1C668619DFE6E65	2020-12-07	N.America	Android
2	2E1BF1D4C37EA01F	2020-12-07	EU	PC
3	50734A22C0C63768	2020-12-07	EU	iPhone
4	E1BDDCE0DAFA2679	2020-12-07	N.America	iPhone

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 61733 entries, 0 to 61732
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  ---
0    user_id     61733 non-null  object
1    first_date  61733 non-null  datetime64[ns]
2    region      61733 non-null  object
3    device      61733 non-null  object
dtypes: datetime64[ns](1), object(3)
memory usage: 1.9+ MB
```

None

Уникальные значения:

```
user_id      61733
first_date    17
region        4
device        4
dtype: int64
```

Датасет df_final_ab_events

	user_id	event_dt	event_name	details
0	E1BDDCE0DAFA2679	2020-12-07 20:22:03	purchase	99.99
1	7B6452F081F49504	2020-12-07 09:22:53	purchase	9.99
2	9CD9F34546DF254C	2020-12-07 12:59:29	purchase	4.99
3	96F27A054B191457	2020-12-07 04:02:40	purchase	4.99
4	1FD7660FDF94CA1F	2020-12-07 10:15:09	purchase	4.99

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 440317 entries, 0 to 440316
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  -
0   user_id     440317 non-null   object
1   event_dt    440317 non-null   datetime64[ns]
2   event_name  440317 non-null   object
3   details     62740 non-null    float64
dtypes: datetime64[ns](1), float64(1), object(2)
memory usage: 13.4+ MB
```

None

Уникальные значения:

```
user_id      58703
event_dt     267268
event_name    4
details       4
dtype: int64
```

Датасет df_final_ab_participants

	user_id	group	ab_test
0	D1ABA3E2887B6A73	A	recommender_system_test
1	A7A3664BD6242119	A	recommender_system_test
2	DABC14FDDFADD29E	A	recommender_system_test
3	04988C5DF189632E	A	recommender_system_test
4	482F14783456D21B	B	recommender_system_test

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18268 entries, 0 to 18267
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   user_id     18268 non-null   object
1   group       18268 non-null   object
2   ab_test     18268 non-null   object
dtypes: object(3)
memory usage: 428.3+ KB
```

None

Уникальные значения:

```
user_id      16666
group        2
ab_test      2
dtype: int64
```

Заметка:

Датасет *df_final_ab_events* имеет пропуски в графе *details*, но вполне вероятно, что там пусто вследствие отсутствия покупок. Даты были преобразованы на этапе чтения файлов.

1.2 Дубликаты

```
In [4]: 1 # Явные дубликаты
2 print('df_ab_project_marketing_events:', df_ab_project_marketing_events.duplicated().sum())
3 print('df_final_ab_events:', df_final_ab_events.duplicated().sum())
4 print('df_final_ab_new_users:', df_final_ab_new_users.duplicated().sum())
5 print('df_final_ab_participants:', df_final_ab_participants.duplicated().sum())
```

```
df_ab_project_marketing_events: 0
df_final_ab_events: 0
df_final_ab_new_users: 0
df_final_ab_participants: 0
```

```
In [5]: 1 # df_ab_project_marketing_events
2 for col in ['name', 'regions']:
3     print('Оригинал', col, ':', df_ab_project_marketing_events[col].nunique())
4     print('Нижний регистр', col, ':', df_ab_project_marketing_events[col].str.lower().nunique())
5     print()
```

Оригинал name : 14
Нижний регистр name : 14

Оригинал regions : 6
Нижний регистр regions : 6

```
In [6]: 1 # df_final_ab_events
2 for col in ['user_id', 'region', 'device']:
3     print('Оригинал', col, ':', df_final_ab_new_users[col].nunique())
4     print('Нижний регистр', col, ':', df_final_ab_new_users[col].str.lower().nunique())
5     print()
```

Оригинал user_id : 61733
Нижний регистр user_id : 61733

Оригинал region : 4
Нижний регистр region : 4

Оригинал device : 4
Нижний регистр device : 4

```
In [7]: 1 # df_final_ab_new_users
2 for col in ['user_id', 'event_name']:
3     print('Оригинал', col, ':', df_final_ab_events[col].nunique())
4     print('Нижний регистр', col, ':', df_final_ab_events[col].str.lower().nunique())
5     print()
```

Оригинал user_id : 58703
Нижний регистр user_id : 58703

Оригинал event_name : 4
Нижний регистр event_name : 4

```
In [8]: 1 # df_final_ab_participants
2 for col in ['user_id', 'group', 'ab_test']:
3     print('Оригинал', col, ':', df_final_ab_participants[col].nunique())
4     print('Нижний регистр', col, ':', df_final_ab_participants[col].str.lower().nunique())
5     print()
```

Оригинал user_id : 16666
Нижний регистр user_id : 16666

Оригинал group : 2
Нижний регистр group : 2

Оригинал ab_test : 2
Нижний регистр ab_test : 2

Заметка:

Дубликатов нет

2 Оценка корректности проведения теста

```
In [9]: 1 test_start = pd.to_datetime('2020-12-07')
2 test_end = pd.to_datetime('2020-12-21')
```

```
In [10]: 1 df_final_ab_participants.group.unique()
```

```
Out[10]: array(['A', 'B'], dtype=object)
```

```
In [11]: 1 display(df_final_ab_events.event_dt.min())
2 display(df_final_ab_events.event_dt.max())
```

Timestamp('2020-12-07 00:00:33')

Timestamp('2020-12-30 23:36:33')

```
In [12]: 1 # Получение данных о датах
2 def get_info(data, col):
3     df = data.merge(df_final_ab_participants,
4                     how='inner',
5                     left_on=['user_id'],
6                     right_on=['user_id']
7                     )
8     df = df.query('(ab_test == "recommender_system_test")') # выбираем нужные данные (group == "A") &
9
10    display(df.head())
11
12    display(df[col].min())
13    display(df[col].max())
```

```
In [13]: 1 get_info(df_final_ab_new_users, 'first_date')
```

	user_id	first_date	region	device	group	ab_test
0	D72A72121175D8BE	2020-12-07	EU	PC	A	recommender_system_test
3	E6DE857AFBDC6102	2020-12-07	EU	PC	B	recommender_system_test
7	DD4352CDCF8C3D57	2020-12-07	EU	Android	B	recommender_system_test
10	831887FE7F2D6CBA	2020-12-07	EU	Android	A	recommender_system_test
12	4CB179C7F847320B	2020-12-07	EU	iPhone	B	recommender_system_test

Timestamp('2020-12-07 00:00:00')

Timestamp('2020-12-21 00:00:00')

Заметка:

все пользователи, зарегистрировавшиеся в интернет-магазине в период с 7 по 21 декабря 2020 года

Дата запуска: 2020-12-07;

Дата остановки набора новых пользователей: 2020-12-21;

С датами все в порядке

```
In [14]: 1 get_info(df_final_ab_events, 'event_dt')
```

	user_id	event_dt	event_name	details	group	ab_test
12	831887FE7F2D6CBA	2020-12-07 06:50:29	purchase	4.99	A	recommender_system_test
13	831887FE7F2D6CBA	2020-12-09 02:19:17	purchase	99.99	A	recommender_system_test
14	831887FE7F2D6CBA	2020-12-07 06:50:30	product_cart	NaN	A	recommender_system_test
15	831887FE7F2D6CBA	2020-12-08 10:52:27	product_cart	NaN	A	recommender_system_test
16	831887FE7F2D6CBA	2020-12-09 02:19:17	product_cart	NaN	A	recommender_system_test

Timestamp('2020-12-07 00:05:57')

Timestamp('2020-12-30 12:42:57')

Заметка:

Все события новых пользователей в период с 7 декабря 2020 по 4 января 2021 года

Дата остановки: 2021-01-04;

В рамки укладываются

```
In [15]: 1 # Дополнительный фрейм для работы
2 df = (df_final_ab_new_users
3       .merge(df_final_ab_participants,
4             how='inner',
5             left_on=['user_id'],
6             right_on=['user_id'])
7       .query('(ab_test == "recommender_system_test")'))
```

```
In [16]: 1 print('Охват аудитории новых пользователей по региону EU: {:.2%}'
2       .format(df
3               .query('region == "EU"')
4               .user_id.nunique()
5               /
6               df_final_ab_new_users.query('(region == "EU")&(first_date >= @test_start)&(first_date <= @test_end)')
7               .user_id
8               .nunique()
9               )
10       )
```

Охват аудитории новых пользователей по региону EU: 15.00%

Заметка:

Аудитория: 15% новых пользователей из региона EU;

У нас 15,00%

```
In [17]: 1 # Количество участников
        2 df.user_id.nunique()
```

Out[17]: 6701

Заметка:

Ожидаемое количество участников теста: 6000

Количество участников чуть больше (6701)

Проверим:

Время проведения теста. Убедимся, что оно не совпадает с маркетинговыми и другими активностями.

Все события новых пользователей в период с 7 декабря 2020 по 4 января 2021 года

```
In [18]: 1 # Совпадение с маркетинговыми и другими активностями
        2 df_ab_project_marketing_events.query('(start_dt >= "2020-12-07")|(finish_dt >= "2020-12-07")')
```

Out[18]:

	name	regions	start_dt	finish_dt
0	Christmas&New Year Promo	EU, N.America	2020-12-25	2021-01-03
10	CIS New Year Gift Lottery	CIS	2020-12-30	2021-01-07

```
In [ ]: 1 marketing_events.query('finish_dt >= "2020-12-07" and start_dt <="2021-01-04"')
```

Заметка:

Последние действия пользователей отмечаются 2020-12-30 12:42:57 .

В наш тест вошла одна промо акция Christmas&New Year Promo

Проверим:

Аудитория теста: удостовериться, что нет пересечений с конкурирующим тестом и нет пользователей, участвующих в двух группах теста одновременно. Проверим равномерность распределения пользователей по тестовым группам и правильность их формирования.

```
In [19]: 1 # Временный датафрейм
        2 df_temp = (df_final_ab_events
        3             .merge(df_final_ab_participants,
        4                   how='inner',
        5                   left_on=['user_id'],
        6                   right_on=['user_id']))
        7
        8 df_temp['date'] = df_temp.event_dt.dt.date
        9
        10 df_temp.head()
```

Out[19]:

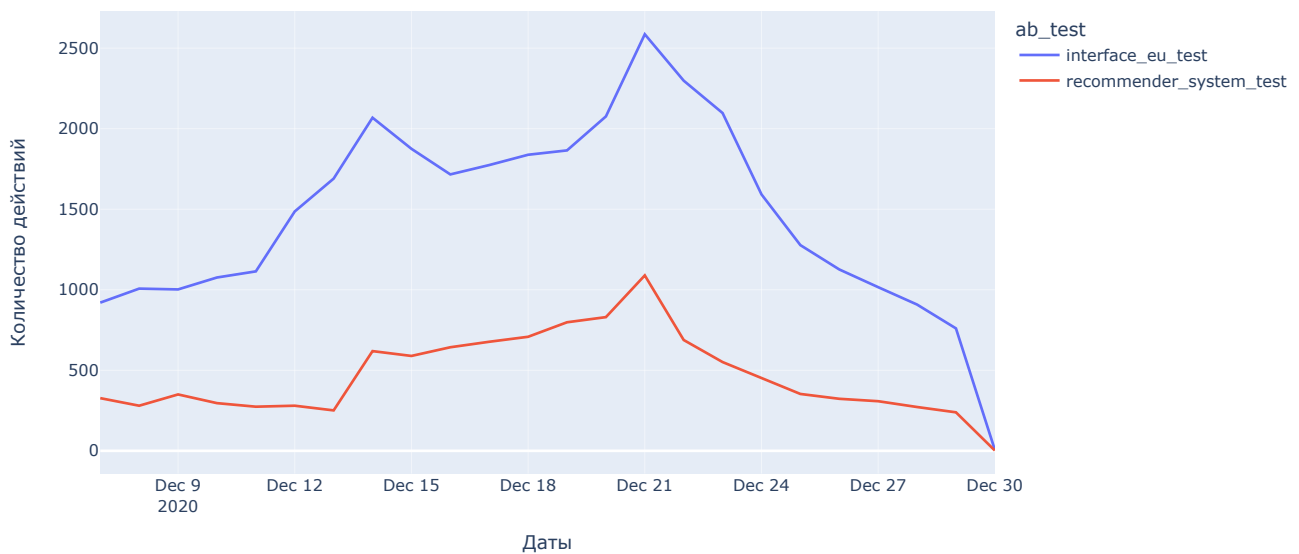
	user_id	event_dt	event_name	details	group	ab_test	date
0	96F27A054B191457	2020-12-07 04:02:40	purchase	4.99	B	interface_eu_test	2020-12-07
1	96F27A054B191457	2020-12-08 09:43:14	purchase	4.99	B	interface_eu_test	2020-12-08
2	96F27A054B191457	2020-12-09 00:44:10	purchase	4.99	B	interface_eu_test	2020-12-09
3	96F27A054B191457	2020-12-26 00:33:57	purchase	9.99	B	interface_eu_test	2020-12-26
4	96F27A054B191457	2020-12-07 04:02:41	product_page	NaN	B	interface_eu_test	2020-12-07

```

In [20]: 1 # Отрисовываем график
2 fig = px.line(df_temp.pivot_table(index=['ab_test', 'date'], values='user_id', aggfunc='nunique').reset_index(),
3             x='date',
4             y='user_id',
5             title='Динамика количества уникальных пользователей за день',
6             color='ab_test'
7             )
8
9 #fig.update_xaxes(dtick="M1")
10
11 fig.update_layout(yaxis_title='Количество действий',
12                  xaxis_title='Даты'
13                  )
14
15 fig.show()

```

Динамика количества уникальных пользователей за день



Заметка:

Два теста идут одновременно

```

In [21]: 1 # Проверим есть ли пересечение между типами тестов
2 if df_final_ab_participants.groupby('user_id')['ab_test'].nunique().max() > 1:
3     print('Есть пересечение между тестами')
4 else:
5     print('Нет пересечения между тестами')

```

Есть пересечение между тестами

```

In [22]: 1 # Проверим есть ли пользователи, которые участвовали в двух группах
2 if df_final_ab_participants.groupby('user_id')['group'].nunique().max() > 1:
3     print('Некоторые пользователи участвовали в двух группах')
4 else:
5     print('Нет пользователей, участвовавших в двух группах')

```

Некоторые пользователи участвовали в двух группах

```

In [23]: 1 # Создадим df в котором будут содержаться информация об участниках двух тестов
2 two_tests = df_final_ab_participants\
3             .groupby('user_id')\
4             .agg({'ab_test': 'nunique'})\
5             .reset_index()\
6             .query('ab_test == 2')
7
8 print ('Количество пользователей участвовавших в обоих тестах:', len(two_tests))

```

Количество пользователей участвовавших в обоих тестах: 1602

Заметка:

Необходимо исключить этих пользователей, так как на них влияют два источника, а не один - выводы неоднозначны.


```
In [24]: 1 print('Количество до удаления:', df_final_ab_participants.shape[0])
2 df_final_ab_participants = df_final_ab_participants.query('user_id not in @two_tests["user_id"]')
3 print('Количество после удаления:', df_final_ab_participants.shape[0])
```

Количество до удаления: 18268

Количество после удаления: 15064

Заметка:

Повторим проверки после удаления данных

```
In [25]: 1 def get_info1():↔
```

```
In [26]: 1 get_info1()
```

Группы: ['A' 'B']
Мин. дата: 2020-12-07 00:00:33
Макс. дата: 2020-12-30 23:36:33

Информация о df_final_ab_new_users

	user_id	first_date	region	device	group	ab_test
0	D72A72121175D8BE	2020-12-07	EU	PC	A	recommender_system_test
6	831887FE7F2D6CBA	2020-12-07	EU	Android	A	recommender_system_test
8	4CB179C7F847320B	2020-12-07	EU	iPhone	B	recommender_system_test
9	29C92313A98B1176	2020-12-07	APAC	Android	B	recommender_system_test
11	7D1BFB181017EB46	2020-12-07	CIS	PC	B	recommender_system_test

Мин. дата: 2020-12-07 00:00:00
Макс. дата: 2020-12-21 00:00:00
None

Информация о df_final_ab_events

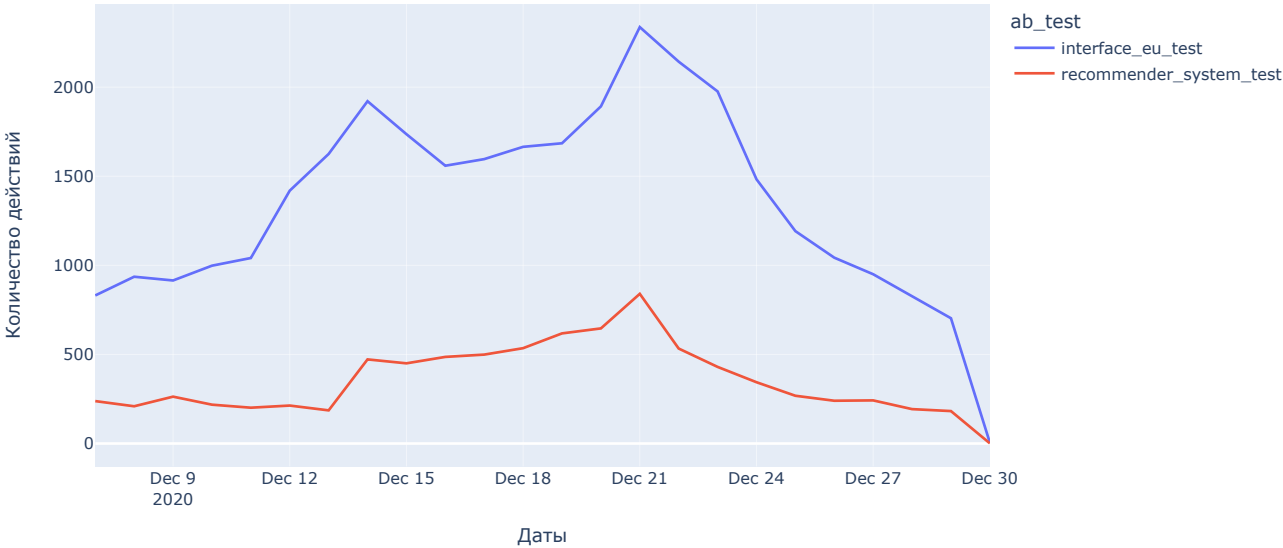
	user_id	event_dt	event_name	details	group	ab_test
12	831887FE7F2D6CBA	2020-12-07 06:50:29	purchase	4.99	A	recommender_system_test
13	831887FE7F2D6CBA	2020-12-09 02:19:17	purchase	99.99	A	recommender_system_test
14	831887FE7F2D6CBA	2020-12-07 06:50:30	product_cart	NaN	A	recommender_system_test
15	831887FE7F2D6CBA	2020-12-08 10:52:27	product_cart	NaN	A	recommender_system_test
16	831887FE7F2D6CBA	2020-12-09 02:19:17	product_cart	NaN	A	recommender_system_test

Мин. дата: 2020-12-07 00:16:00
Макс. дата: 2020-12-30 06:42:52
None

Охват аудитории новых пользователей по региону EU: 11.22%

Количество участников 5099

Динамика количества уникальных пользователей за день



Нет пересечения между тестами
Нет пользователей, участвовавших в двух группах

Количество пользователей участвовавших в обоих тестах: 0

Доли участников

	user_id	%
group		
A	2903	56.93
B	2196	43.07

Всего пользователей во фрейме: 15064

Заметка:

- охват аудитории уменьшился;

- пересечений больше нет;
- участников в двух группах теста нет;
- доли различны ~14%;
- даты остались в порядке.

3 Исследовательский анализ данных

Создадим общий фрейм

```
In [27]: 1 # Дополнительный фрейм для работы
2 df_for_work = (df_final_ab_new_users.query('(first_date >= @test_start)&(first_date <= @test_end)')
3             .merge(df_final_ab_events,
4                   how='left',
5                   left_on=['user_id'],
6                   right_on=['user_id'])
7             )
8
9 df_for_work.head()
```

Out[27]:

	user_id	first_date	region	device	event_dt	event_name	details
0	D72A72121175D8BE	2020-12-07	EU	PC	2020-12-07 21:52:10	product_page	NaN
1	D72A72121175D8BE	2020-12-07	EU	PC	2020-12-07 21:52:07	login	NaN
2	F1C668619DFE6E65	2020-12-07	N.America	Android	2020-12-07 16:38:09	product_page	NaN
3	F1C668619DFE6E65	2020-12-07	N.America	Android	2020-12-08 02:02:34	product_page	NaN
4	F1C668619DFE6E65	2020-12-07	N.America	Android	2020-12-23 14:35:41	product_page	NaN

```
In [28]: 1 # Дополнительный фрейм для работы
2 df_for_work = (df_for_work
3             .merge(df_final_ab_participants,
4                   how='left',
5                   left_on=['user_id'],
6                   right_on=['user_id'])
7             )
8
9 df_for_work = df_for_work.query('ab_test == "recommender_system_test"').reset_index(drop=True)
10
11 display(df_for_work.head())
12 print('Записей:', df_for_work.shape[0])
13 print('Количество пользователей:', df_for_work.user_id.nunique())
```

	user_id	first_date	region	device	event_dt	event_name	details	group	ab_test
0	D72A72121175D8BE	2020-12-07	EU	PC	2020-12-07 21:52:10	product_page	NaN	A	recommender_system_test
1	D72A72121175D8BE	2020-12-07	EU	PC	2020-12-07 21:52:07	login	NaN	A	recommender_system_test
2	831887FE7F2D6CBA	2020-12-07	EU	Android	2020-12-07 06:50:29	purchase	4.99	A	recommender_system_test
3	831887FE7F2D6CBA	2020-12-07	EU	Android	2020-12-09 02:19:17	purchase	99.99	A	recommender_system_test
4	831887FE7F2D6CBA	2020-12-07	EU	Android	2020-12-07 06:50:30	product_cart	NaN	A	recommender_system_test

Записей: 21115
Количество пользователей: 5099

```
In [29]: 1 df_for_work['date_event'] = df_for_work.event_dt.dt.date
2 df_for_work['date_event'] = pd.to_datetime(df_for_work['date_event'])
3 df_for_work['dt_delta'] = df_for_work['date_event'] - df_for_work['first_date']
4
5 display(df_for_work.head())
6 print('Количество записей:', df_for_work.shape[0])
```

	user_id	first_date	region	device	event_dt	event_name	details	group	ab_test	date_event	dt_delta
0	D72A72121175D8BE	2020-12-07	EU	PC	2020-12-07 21:52:10	product_page	NaN	A	recommender_system_test	2020-12-07	0 days
1	D72A72121175D8BE	2020-12-07	EU	PC	2020-12-07 21:52:07	login	NaN	A	recommender_system_test	2020-12-07	0 days
2	831887FE7F2D6CBA	2020-12-07	EU	Android	2020-12-07 06:50:29	purchase	4.99	A	recommender_system_test	2020-12-07	0 days
3	831887FE7F2D6CBA	2020-12-07	EU	Android	2020-12-09 02:19:17	purchase	99.99	A	recommender_system_test	2020-12-09	2 days
4	831887FE7F2D6CBA	2020-12-07	EU	Android	2020-12-07 06:50:30	product_cart	NaN	A	recommender_system_test	2020-12-07	0 days

Количество записей: 21115

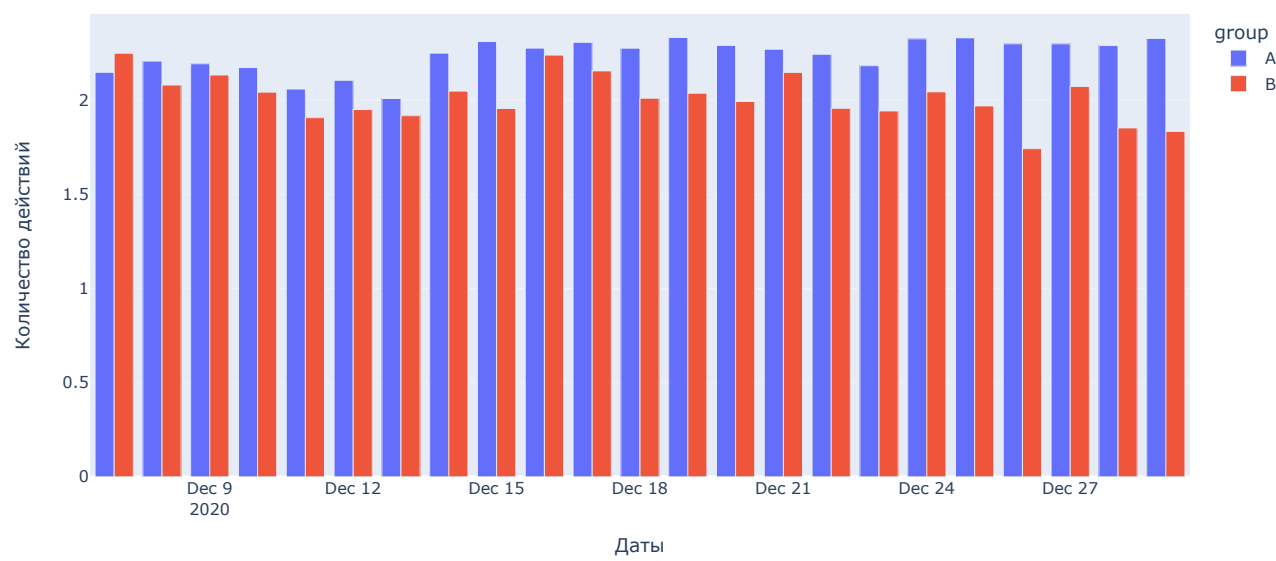
```
In [30]: 1 # Оставим записи с действиями за две недели после регистрации
2 df_for_work = df_for_work[df_for_work['dt_delta'] <= pd.Timedelta(days=14)]
3
4 display(df_for_work.head())
5 print('Количество записей:', df_for_work.shape[0])
```

	user_id	first_date	region	device	event_dt	event_name	details	group	ab_test	date_event	dt_delta
0	D72A72121175D8BE	2020-12-07	EU	PC	2020-12-07 21:52:10	product_page	NaN	A	recommender_system_test	2020-12-07	0 days
1	D72A72121175D8BE	2020-12-07	EU	PC	2020-12-07 21:52:07	login	NaN	A	recommender_system_test	2020-12-07	0 days
2	831887FE7F2D6CBA	2020-12-07	EU	Android	2020-12-07 06:50:29	purchase	4.99	A	recommender_system_test	2020-12-07	0 days
3	831887FE7F2D6CBA	2020-12-07	EU	Android	2020-12-09 02:19:17	purchase	99.99	A	recommender_system_test	2020-12-09	2 days
4	831887FE7F2D6CBA	2020-12-07	EU	Android	2020-12-07 06:50:30	product_cart	NaN	A	recommender_system_test	2020-12-07	0 days

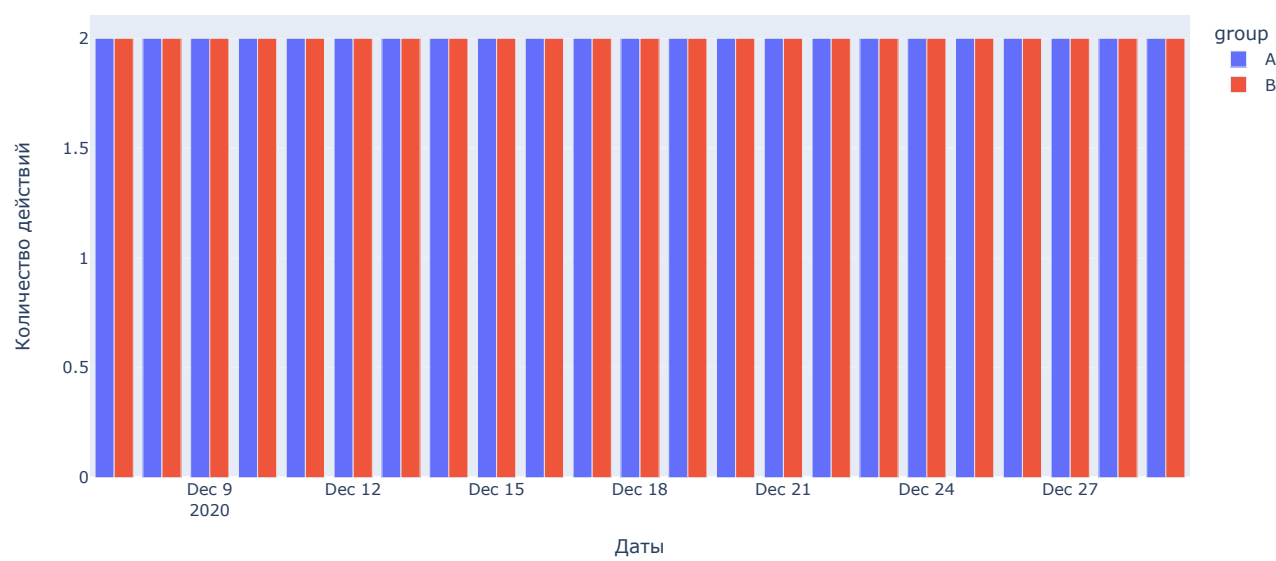
Количество записей: 18329

```
In [31]: 1 # Распределение в выборках количества событий на пользователя↔
```

Динамика среднего количества действий на пользователя по группам



Динамика медианного количества действий на пользователя по группам

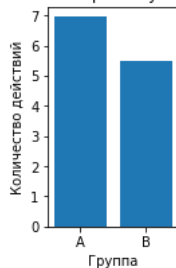


```

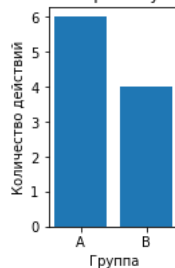
In [32]: 1 # Среднее за весь промежуток по группам
2 plt.figure(figsize=(8,3))
3 ax1 = plt.subplot(1, 4, 1)
4 df_temp = df_for_work.pivot_table(index=['group', 'user_id'], values='event_dt', aggfunc='count')
5 df_temp = df_temp.pivot_table(index=['group'],
6                               values='event_dt',
7                               aggfunc='mean').reset_index()
8 plt.bar(data=df_temp,
9         x='group',
10        height='event_dt'
11        )
12
13 plt.ylabel('Количество действий')
14 plt.xlabel('Группа')
15 plt.title("Среднее за весь промежуток по группам")
16
17 ax2 = plt.subplot(1, 4, 4)
18 df_temp = df_for_work.pivot_table(index=['group', 'user_id'], values='event_dt', aggfunc='count')
19 df_temp = df_temp.pivot_table(index=['group'],
20                               values='event_dt',
21                               aggfunc='median').reset_index()
22 plt.bar(data=df_temp,
23         x='group',
24         height='event_dt'
25         )
26
27 plt.ylabel('Количество действий')
28 plt.xlabel('Группа')
29 plt.title("Медиана за весь промежуток по группам")
30
31 plt.show()

```

Среднее за весь промежуток по группам



Медиана за весь промежуток по группам



Заметка:

Среднее количество действий находится в пределах 2 шт., медианное = 2 шт.

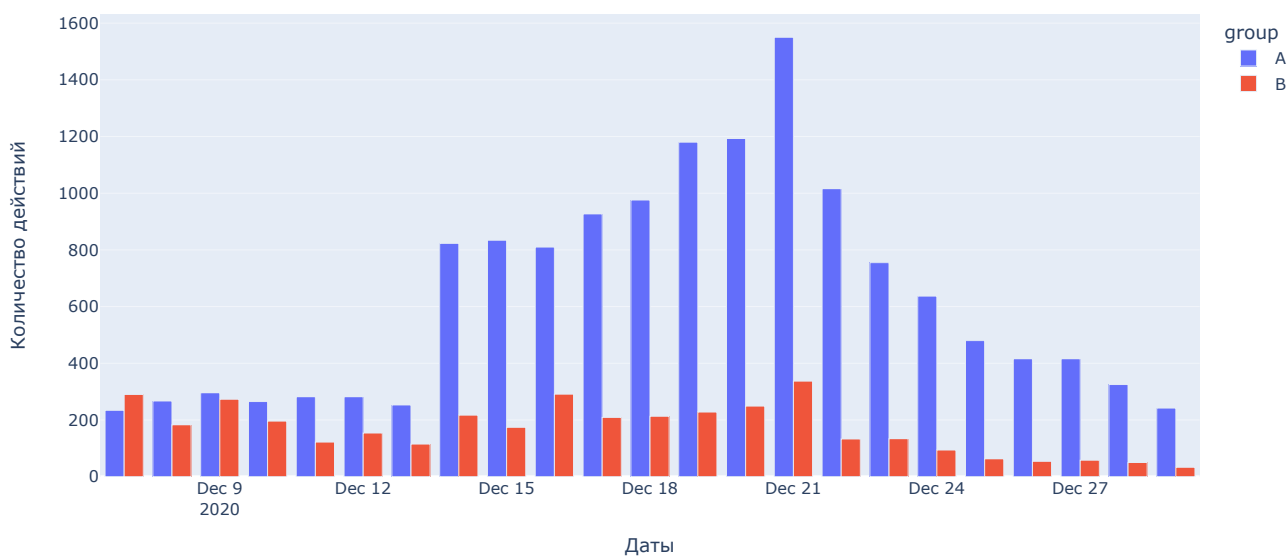
Если посмотреть на весь промежуток, то количество вырастет и станет сильнее разниться.

```

In [33]: 1 # Распределение числа событий в выборках по дням↔

```

Динамика количества действий пользователей по группам



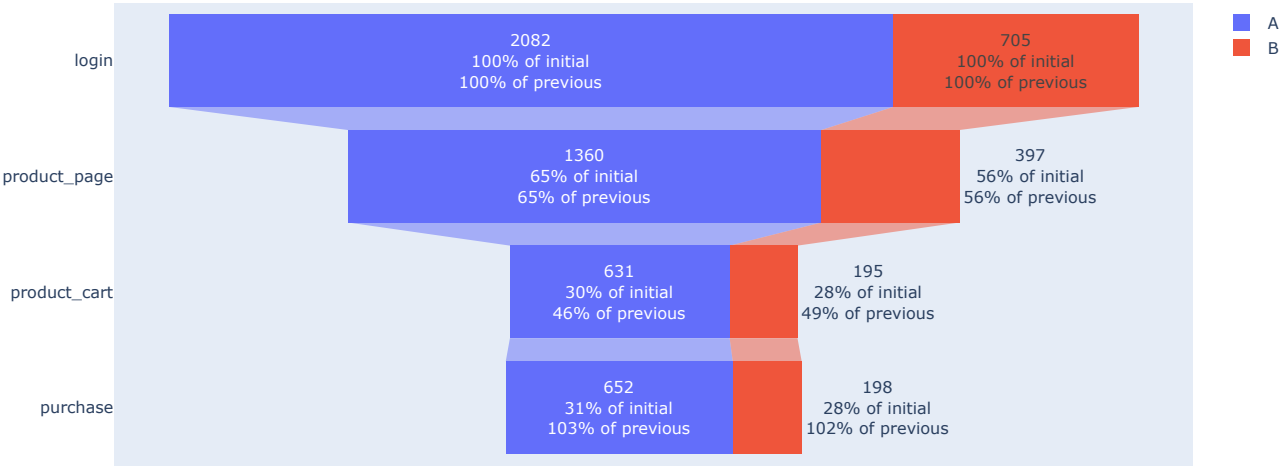
```
In [34]: 1 # Для упорядочивания воронки событий создадим датафрейм
2 group_events_count = df_for_work.groupby(['event_name', 'group']).agg({'user_id': 'nunique'}).reset_index()
3
4 # Изменим индексы для порядка действий
5 new_index = [0, 1, 4, 5, 2, 3, 6, 7]
6
7 group_events_count = group_events_count.reindex(new_index)
8
9 group_events_count
```

Out[34]:

	event_name	group	user_id
0	login	A	2082
1	login	B	705
4	product_page	A	1360
5	product_page	B	397
2	product_cart	A	631
3	product_cart	B	195
6	purchase	A	652
7	purchase	B	198

```
In [35]: 1 # Строим воронку событий по уникальным пользователям↔
```

Воронка событий



Заметка:

Часть заказов минует product_cart , переходя сразу к purchase .
Наибольшие потери конверсии наблюдаются на шаге 3 - переход с product_page к product_cart .

Особенности данных:

- (для A/B- тестирования)
- отсутствуют данные о действиях пользователей с 31.12.2020 по 04.01.2021;
 - в период проведения теста для интересующих нас пользователей параллельно проходила акция Christmas&New Year Promo ;
 - во время исследования проходило другое тестирование. Все пользователи попавшие в оба теста были удалены, т.к. их действия не будут репрезентативными;
 - без учета исключенных пользователей аудитория теста равна 15%, что соответствует ТЗ, однако после исключения не репрезентативных пользователей доля снизилась до 11.22%, что является нарушением ТЗ;
 - медианы по количеству совершенных действий у пользователей в группах различаются (A=6, B=4), однако в разрезе времени они примерно одинаковы.

? Заметка ?

Возможно, стоило исключать пользователей только группы B?

4 Оценка результатов A/B-тестирования

Подготовим данные к тесту.

Для удобства работы по строкам, создадим датафрейм с группами в виде столбцов и строк - событий.

```
In [36]: 1 # Датафрейм воронки для теста - различия в одинаковых событиях
2 event_group_test = df_for_work.pivot_table(index='event_name',
3                                             columns='group',
4                                             values='user_id',
5                                             aggfunc='nunique'
6                                             ).reset_index()
7
8 event_group_test
```

```
Out[36]:
```

	group	event_name	A	B
0		login	2082	705
1		product_cart	631	195
2		product_page	1360	397
3		purchase	652	198

```
In [37]: 1 # Упорядочим события (как в воронке)
2 new_index = [0, 2, 1, 3]
3
4 event_group_test = event_group_test.reindex(new_index)
5
6 event_group_test
```

```
Out[37]:
```

	group	event_name	A	B
0		login	2082	705
2		product_page	1360	397
1		product_cart	631	195
3		purchase	652	198

При проведении Z-теста в знаменателях пропорций успеха указывается размер группы. Для более короткого обращения запишем в датафрейм.

```
In [38]: 1 users_bygroup = df_for_work.groupby('group')['user_id'].nunique()
2 users_bygroup
```

```
Out[38]: group
A      2082
B       706
Name: user_id, dtype: int64
```

Гипотезы:

Наличие изменения конверсии группы В по отношению к группе А:

H_0 : Среднее количество пользователей, совершивших значимое событие в группах А и В, равно.

H_1 : Среднее количество пользователей, совершивших значимое событие в группах А и В, различается.

Функция Z-теста:

Z-тест будет проходить по этапам (событиям) воронки, в следствие чего будет выполняться несколько раз (4 теста одновременно), что вводит дополнительную погрешность

Заметка: применим метод Шидака для расчёта требуемого уровня значимости

```
In [39]: 1 # Задаём параметры
2 alpha = 0.05
3 print('Заданный уровень значимости:  $\alpha$  = ', alpha)
4
5 m = event_group_test.shape[0]
6 print('Число тестов: m = ', m)
7
8 alpha_by_shidok = 1 - (1 - alpha) ** (1 / m)
9 print('Требуемый уровень значимости по методу Шидака:', alpha_by_shidok)
```

Заданный уровень значимости: α = 0.05

Число тестов: m = 4

Требуемый уровень значимости по методу Шидака: 0.012741455098566168

```
In [40]: 1 # Функция для проведения Z-теста
2 def z_test(group1, group2, alpha):
3     for i in event_group_test.index:
4
5         print('Действие: {}'.format(event_group_test['event_name'][i]))
6
7         # Пропорция успехов в первой группе:
8         p1 = event_group_test[group1][i] / users_bygroup[group1]
9
10        # Пропорция успехов во второй группе:
11        p2 = event_group_test[group2][i] / users_bygroup[group2]
12
13        print('Группа', group1, ':',
14              event_group_test[group1][i],
15              users_bygroup[group1])
16        print('Группа', group2, ':',
17              event_group_test[group2][i],
18              users_bygroup[group2])
19
20
21        # Пропорция успехов в комбинированном датасете:
22        p_combined = (
23            (event_group_test[group1][i] + event_group_test[group2][i]) /
24            (users_bygroup[group1] + users_bygroup[group2]))
25
26
27        # Разница пропорций в датасетах
28        difference = p1 - p2
29
30        # Считаем статистику в ст.отклонениях стандартного нормального распределения
31        z_value = difference / \
32            math.sqrt(p_combined * (1 - p_combined) * (1/users_bygroup[group1] + 1/users_bygroup[group2]))
33
34        # Задаем стандартное нормальное распределение (среднее 0, ст.отклонение 1)
35        distr = st.norm(0, 1)
36
37        p_value = (1 - distr.cdf(abs(z_value))) * 2
38
39        print('{} p-значение: {}'.format(event_group_test['event_name'][i], p_value))
40
41        if (p_value < alpha):
42            print("ОТВЕРГАЕМ нулевую гипотезу: между группами есть значимая разница")
43        else:
44            print("НЕ получилось отвергнуть нулевую гипотезу, нет оснований считать группы разными")
45
46        print('')
```

```
In [41]: 1 # Запускаем Z-тест
2 z_test("A", "B", alpha_by_shidok)
```

Действие: login
Группа A : 2082 2082
Группа B : 705 706
login p-значение: 0.08587401754779211
НЕ получилось отвергнуть нулевую гипотезу, нет оснований считать группы разными

Действие: product_page
Группа A : 1360 2082
Группа B : 397 706
product_page p-значение: 1.5371909704686715e-05
ОТВЕРГАЕМ нулевую гипотезу: между группами есть значимая разница

Действие: product_cart
Группа A : 631 2082
Группа B : 195 706
product_cart p-значение: 0.1766337419130104
НЕ получилось отвергнуть нулевую гипотезу, нет оснований считать группы разными

Действие: purchase
Группа A : 652 2082
Группа B : 198 706
purchase p-значение: 0.10281767567786759
НЕ получилось отвергнуть нулевую гипотезу, нет оснований считать группы разными

Итоги:

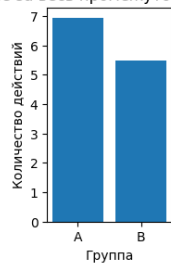
- отсутствуют данные о действиях пользователей с 31.12.2020 по 04.01.2021;
- в период проведения теста проходила акция Christmas&New Year Promo;
- параллельно проходило другое тестирование. Пользователи, попавшие в обе группы были удалены т.к. их действия не будут репрезентативными;
- после удаления осталось количество пользователей сократилось (требуется 6000) - нарушение ТЗ;
- с исключенными пользователями аудитория теста равна 15% (EU) - соответствует ТЗ, однако после исключения не репрезентативных пользователей доля падает до 11,2% - нарушение ТЗ;
- доли количества участников теста в группах А и В;

Доли участников

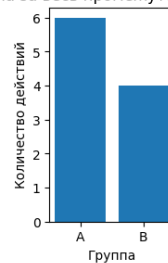
	user_id	%
group		
A	2903	56.93
B	2196	43.07

- количество совершенных действий у пользователей в группах различаются, если брать медиану за все дни;

Среднее за весь промежуток по группам

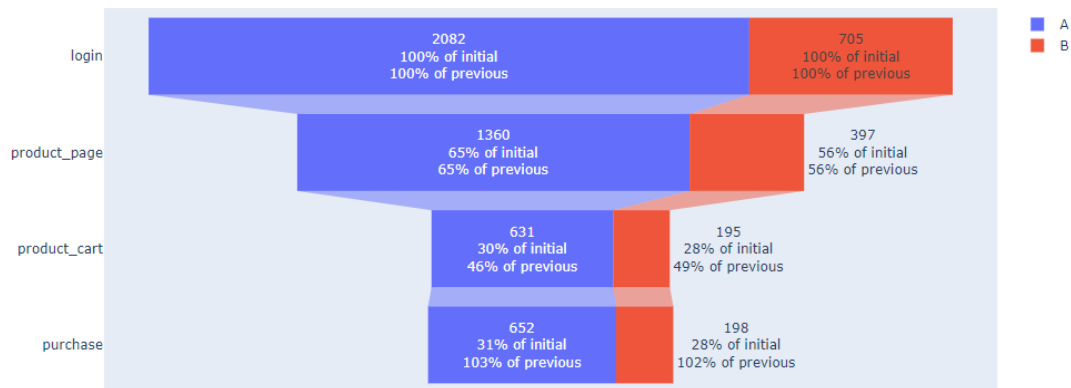


Медиана за весь промежуток по группам



- в группе А до конца воронки (совершения целевого действия - покупки) доходят на 3% больше пользователей, чем в группе В;

Воронка событий



Результаты множественного тестирования

Был применён метод Шидака для расчёта требуемого уровня значимости.

По результатам множественного тестирования, по 3-м действиям (login , product_cart , purchase) не удалось отвергнуть нулевую гипотезу (нет оснований считать, что доли пользователей разные). По действию product_page отвергнули гипотезу (доли разные).

Исходя из произведенного тестирования изменения, связанные с внедрением улучшенной рекомендательной системы, не дали нужного результата, однако результаты **нельзя назвать корректными**.

В ходе проверки данных было выявлено множество ошибок, нарушений ТЗ, что не дает уверенности в проведенном тестировании.

Рекомендации:

- подготовить данные и провести тест повторно, с учетом сделанных замечаний;
- стоит проверить механизм распределения по группам.