

Ускоряем Pandas и строим правильные графики

Александр Ольферук,
наставник

Яндекс Практикум

Ускорение pandas

Разбираемся с датами

| | | date_time |
|---|------------|-----------|
| 0 | 2013-01-01 | 00:00:00 |
| 1 | 2013-01-01 | 01:00:00 |
| 2 | 2013-01-01 | 02:00:00 |
| 3 | 2013-01-01 | 03:00:00 |
| 4 | 2013-01-01 | 04:00:00 |

Разбираемся с датами

```
def convert(df, column_name):  
    return pd.to_datetime(df[column_name])
```

Best of 3 trials with 10 function calls per trial:
Function `convert` ran in average of 1.610 seconds.

```
def convert(df, column_name):  
    return pd.to_datetime(df[column_name],  
                           format='%d/%m/%y %H:%M'  
                           )
```

Best of 3 trials with 100 function calls per trial:
Function `convert_with_format` ran in average of 0.032 seconds.

Разбираемся с датами

Если дата не в ISO 8601, то явно указываем ее формат, чтобы не делать лишние операции (которые, к слову, выполняются через **dateutil**)

Примеры:

- 2020-09-18
- 2020-09-18T16:29:08+00:00
- 2020-W38
- 2020-W38-5

https://en.wikipedia.org/wiki/ISO_8601

Никогда не делайте так!

```
for i in range(len(df)):  
    df.iloc[i, 0] = ...
```

```
for index, row in df.iterrows():  
    something = row['column_name']
```

Только так

```
df.apply(lambda row: row[ 'column' ]*123, axis=1)
```

Нормально, в принципе

```
peak_hours = df.index.hour.isin(range(17, 24))
shoulder_hours = df.index.hour.isin(range(7, 17))
off_peak_hours = df.index.hour.isin(range(0, 7))

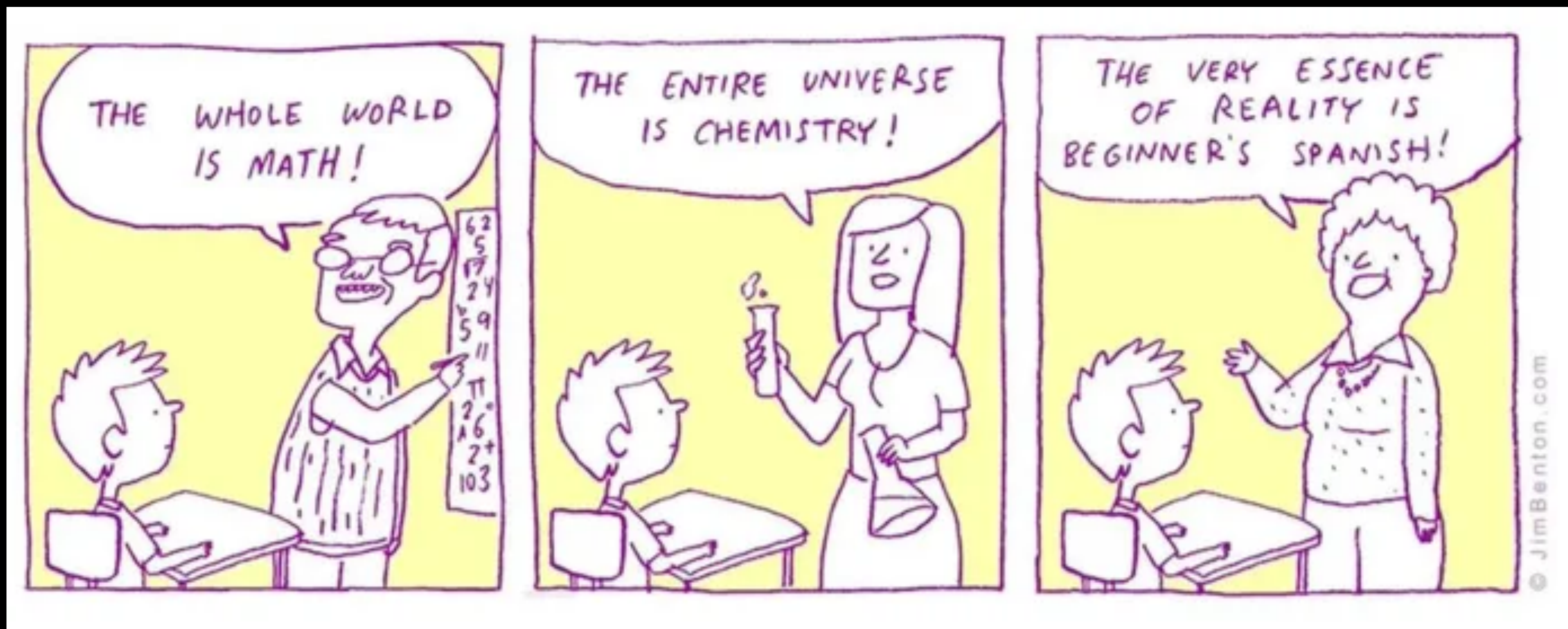
df.loc[peak_hours, 'cost_cents'] = df.loc[peak_hours, 'energy_kwh'] * 28
df.loc[shoulder_hours, 'cost_cents'] = df.loc[shoulder_hours, 'energy_kwh'] * 20
df.loc[off_peak_hours, 'cost_cents'] = df.loc[off_peak_hours, 'energy_kwh'] * 12
```


Еще лучше

Если задача позволяет, то лучше воспользоваться **pd.cut**

```
cents_per_kwh = pd.cut(x=df.index.hour,  
                        bins=[0, 7, 17, 24],  
                        include_lowest=True,  
                        labels=[12, 20, 28]).astype(int)  
df['cost_cents'] = cents_per_kwh * df['energy_kwh']
```

Сначала шютка



Суть

Все, что вы видите в pandas - **векторы** и **матрицы**.

Нужно потратить некоторое количество времени, чтобы научиться мыслить **векторами**, а не отдельными записями, пусть по первому времени и неудобно.

Любая **векторная** функция будет гарантированно быстрее.

Но groupby – не вектор!

Да, и поэтому можно костылить самому с помощью пакета **multiprocessing**

А можно не изобретать, что уже существует: [Pandarallel](#)

Можно почитать [на Хабре](#)

Numba

Переходим [на сайт](#), смотрим

```
from numba import jit
import numpy as np

x = np.arange(100).reshape(10, 10)

@jit(nopython=True) # Set "nopython" mode for best performance, equivalent to @njit
def go_fast(a): # Function is compiled to machine code when called the first time
    trace = 0.0
    for i in range(a.shape[0]): # Numba likes loops
        trace += np.tanh(a[i, i]) # Numba likes NumPy functions
    return a + trace # Numba likes NumPy broadcasting

print(go_fast(x))
```


Numba

```
from numba import jit
import pandas as pd

x = {'a': [1, 2, 3], 'b': [20, 30, 40]}

@jit
def use_pandas(a): # Function will not benefit from Numba jit
    df = pd.DataFrame.from_dict(a) # Numba doesn't know about pd.DataFrame
    df += 1 # Numba doesn't understand what this is
    return df.cov() # or this!

print(use_pandas(x))
```

Numba

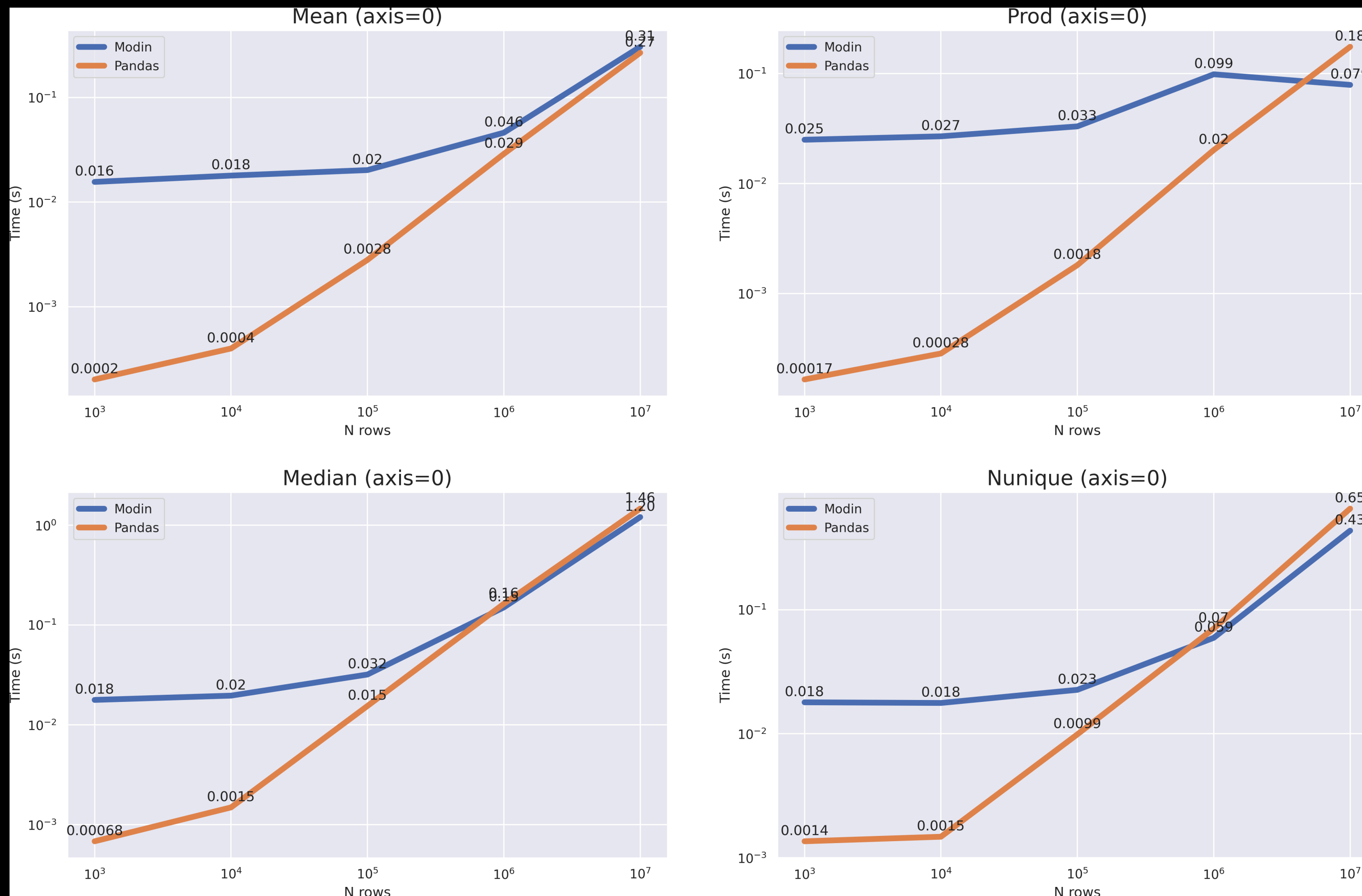
- Возможно добиться ускорения в тысячи раз
- Однако, если ускоряем не чистые математические функции – результат может быть не такой эффектный
- Иногда нужно переделать свой код под Numba

Modin

[Ссылка на него](#)

- Аналог pandas, **полностью повторяющий** API pandas
- Эффективен при чтении больших файлов
- Однако с табличками меньшего размера у него проблемы (см. след. слайд)
- Киллер-фича – возможность использовать дисковое пространство, если объема оперативной памяти не хватает
- Норм, если файл 1Gb+

Modin



Общие советы и что почитать

- [Статья на Medium](#), из которой вы узнаете, что numpy работает быстрее pandas, и возможно, вызов **.values** вам не повредит
- [Статья](#) по сравнению pandas и Modin
- [Тетрадка](#) по измерению query и eval (которые я, к слову, ненавижу)
- Неплохая [статья из документации](#) pandas по теме ускорения вычислений
- Какая-то невероятная вандервафля про работе с табличками на GPU! [Статья](#) и [репозиторий](#)

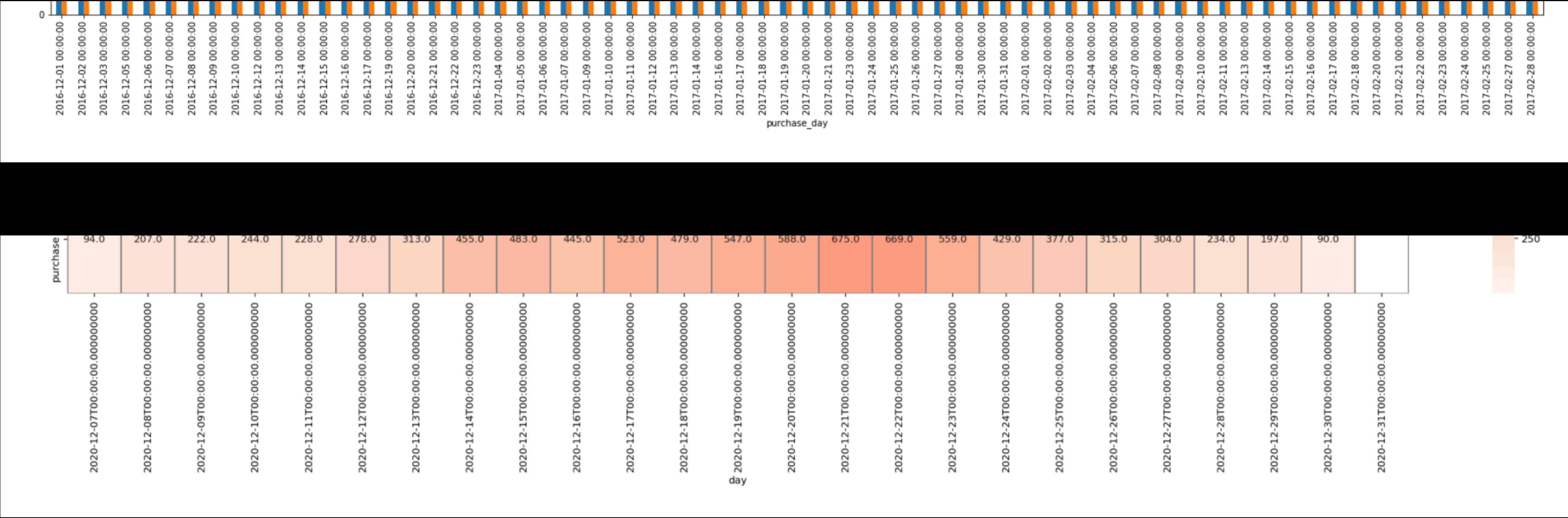
Перерыв



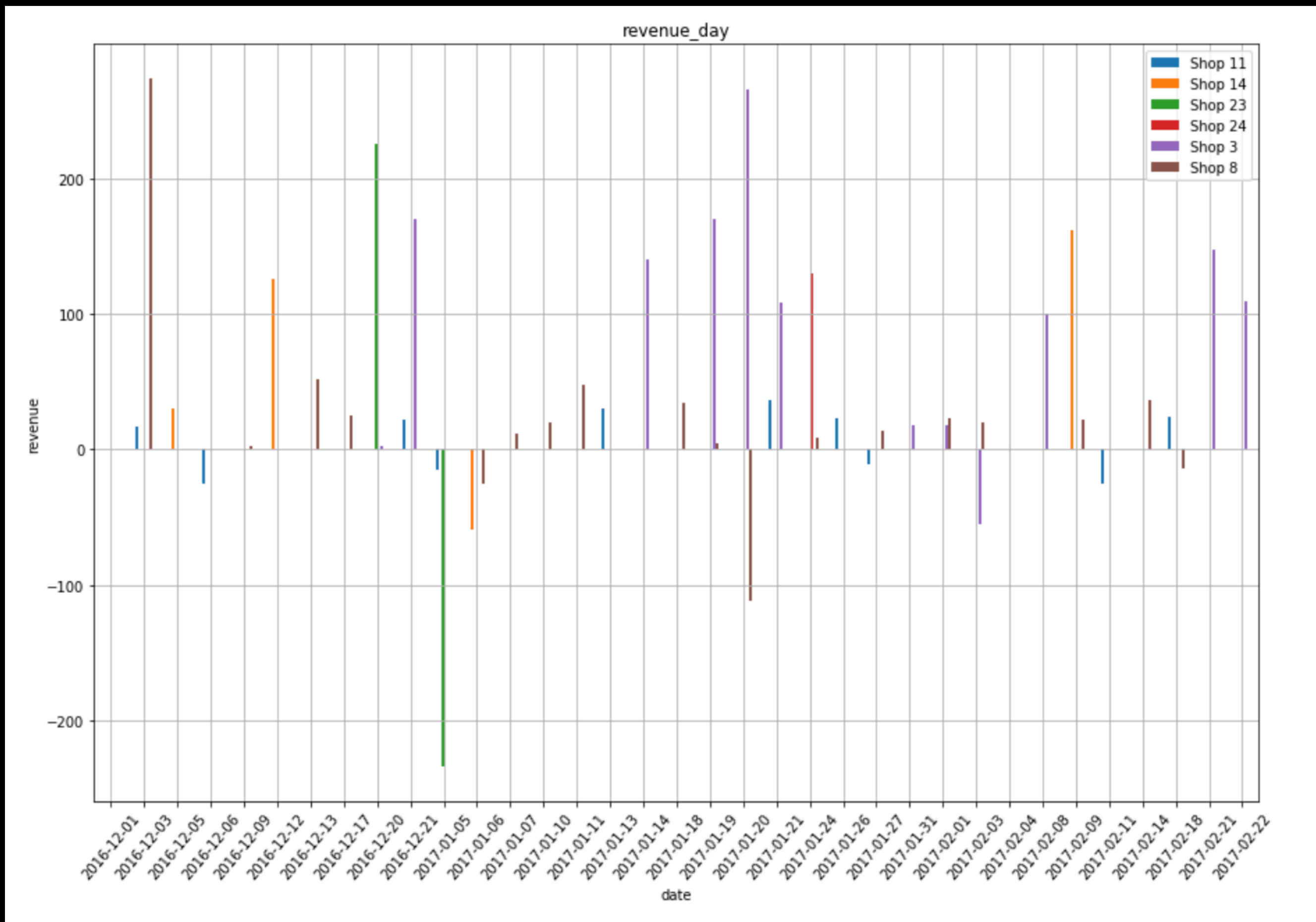
5 минут

Ошибки визуализации

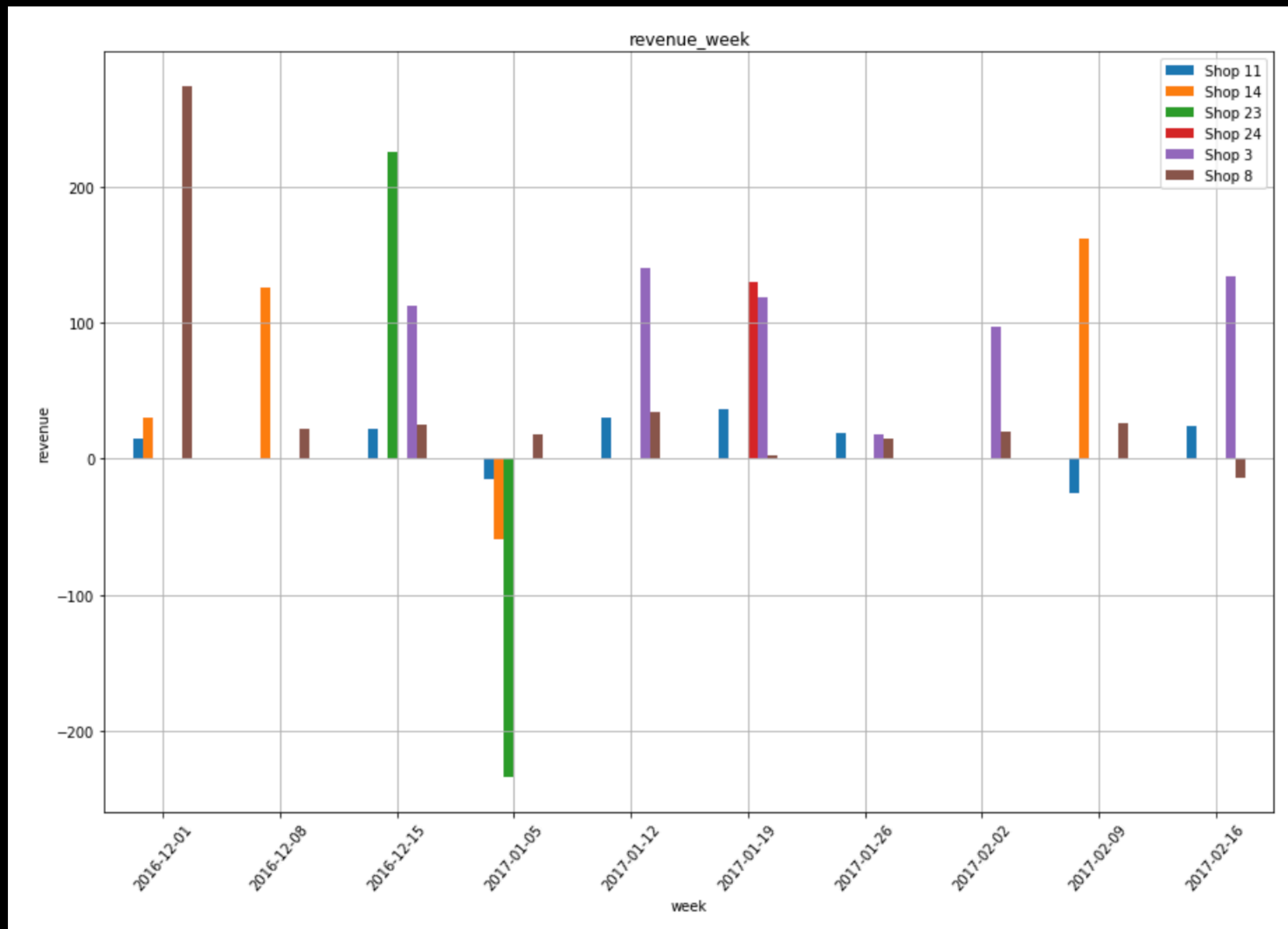
Форматирование дат и времени



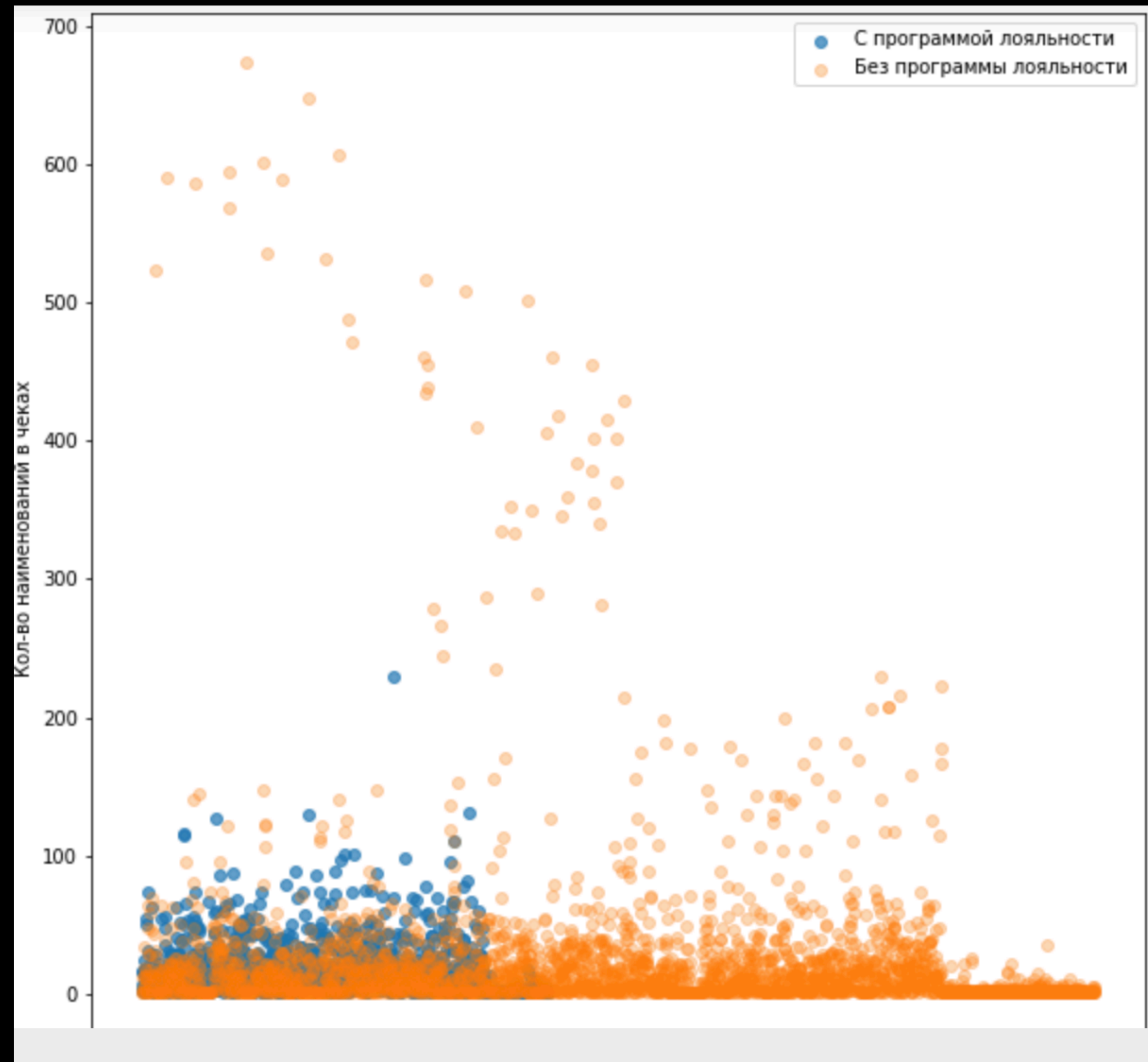
Неверный выбор графика



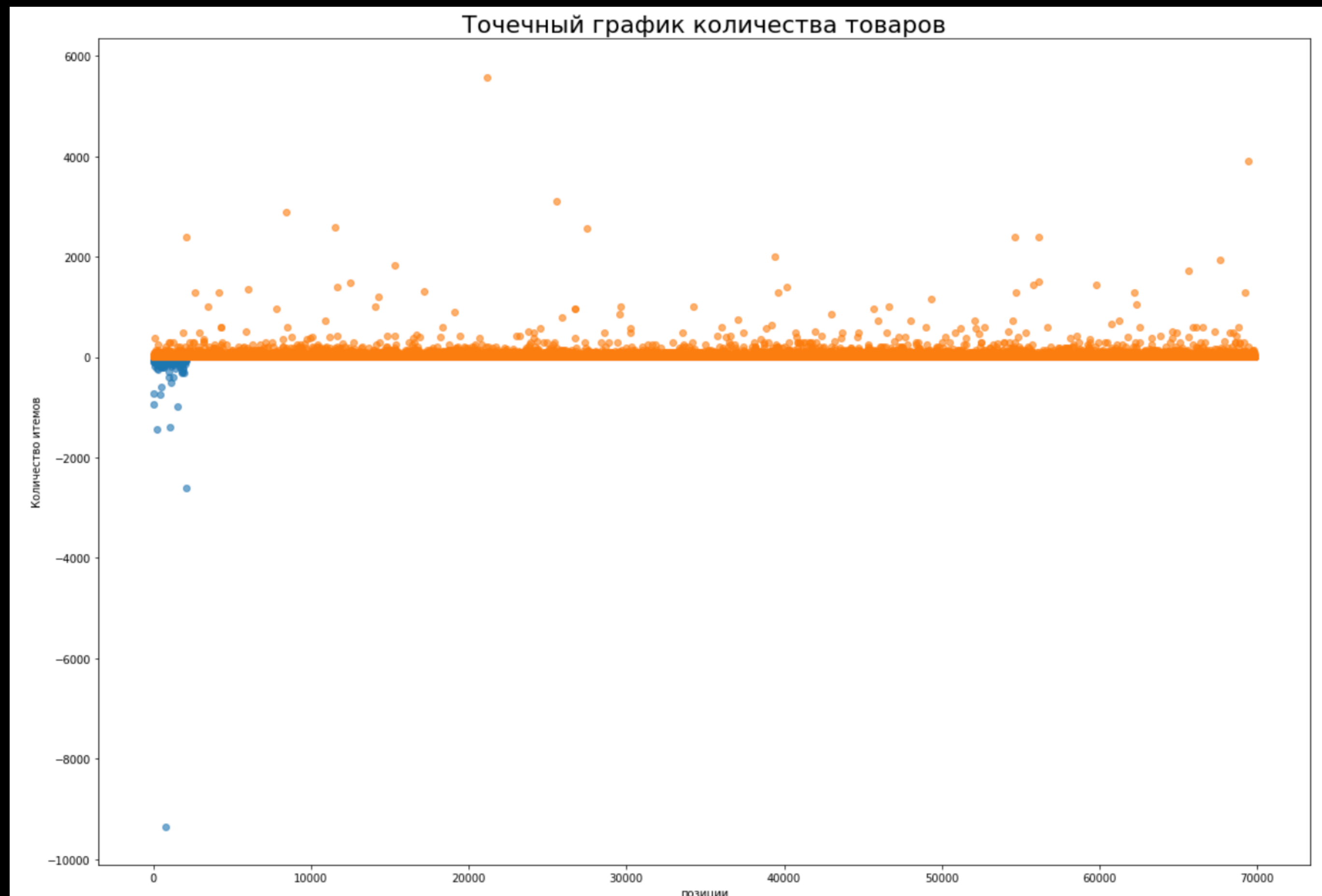
Неверный выбор графика



Неверный выбор графика №2

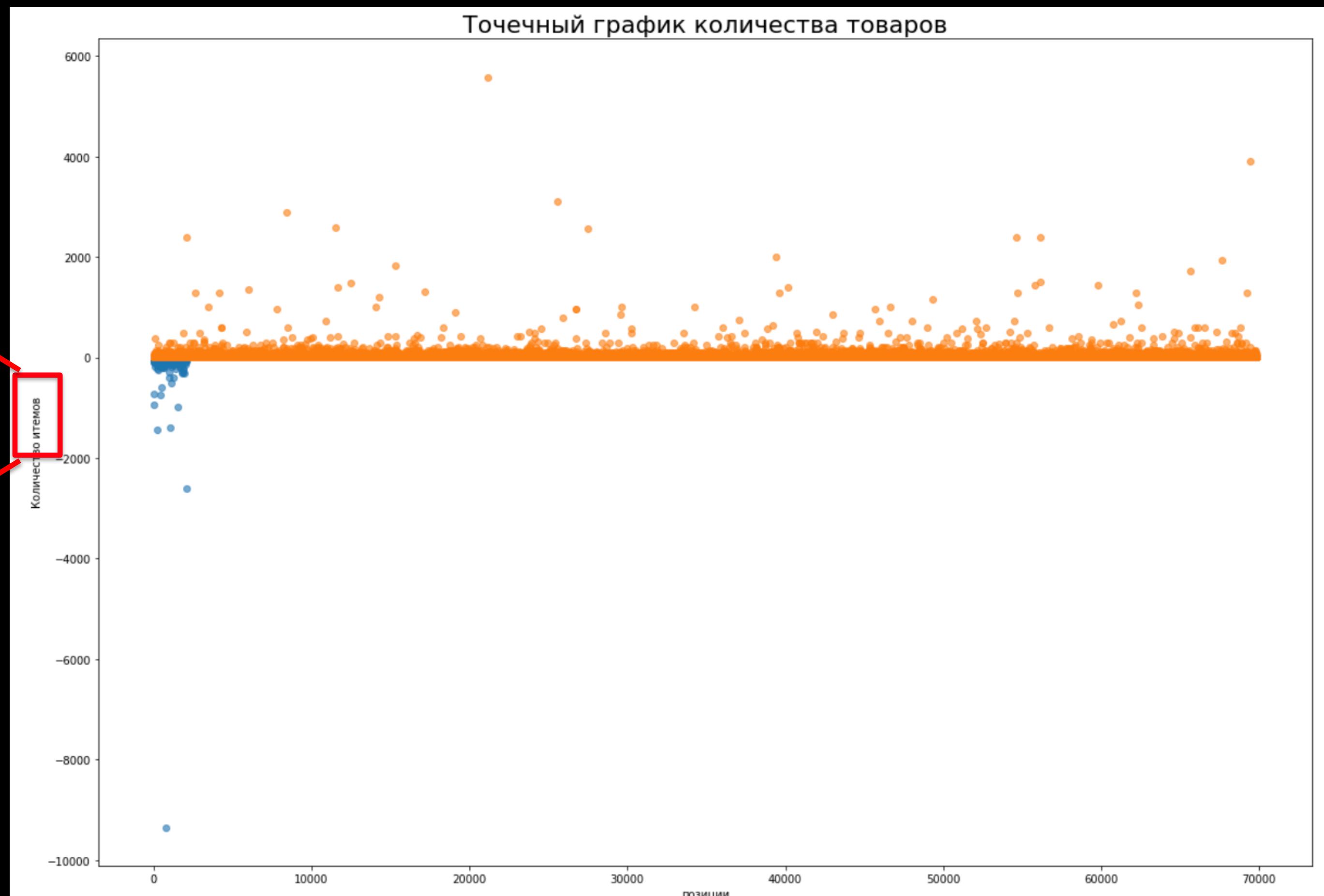


Неверный выбор графика №2

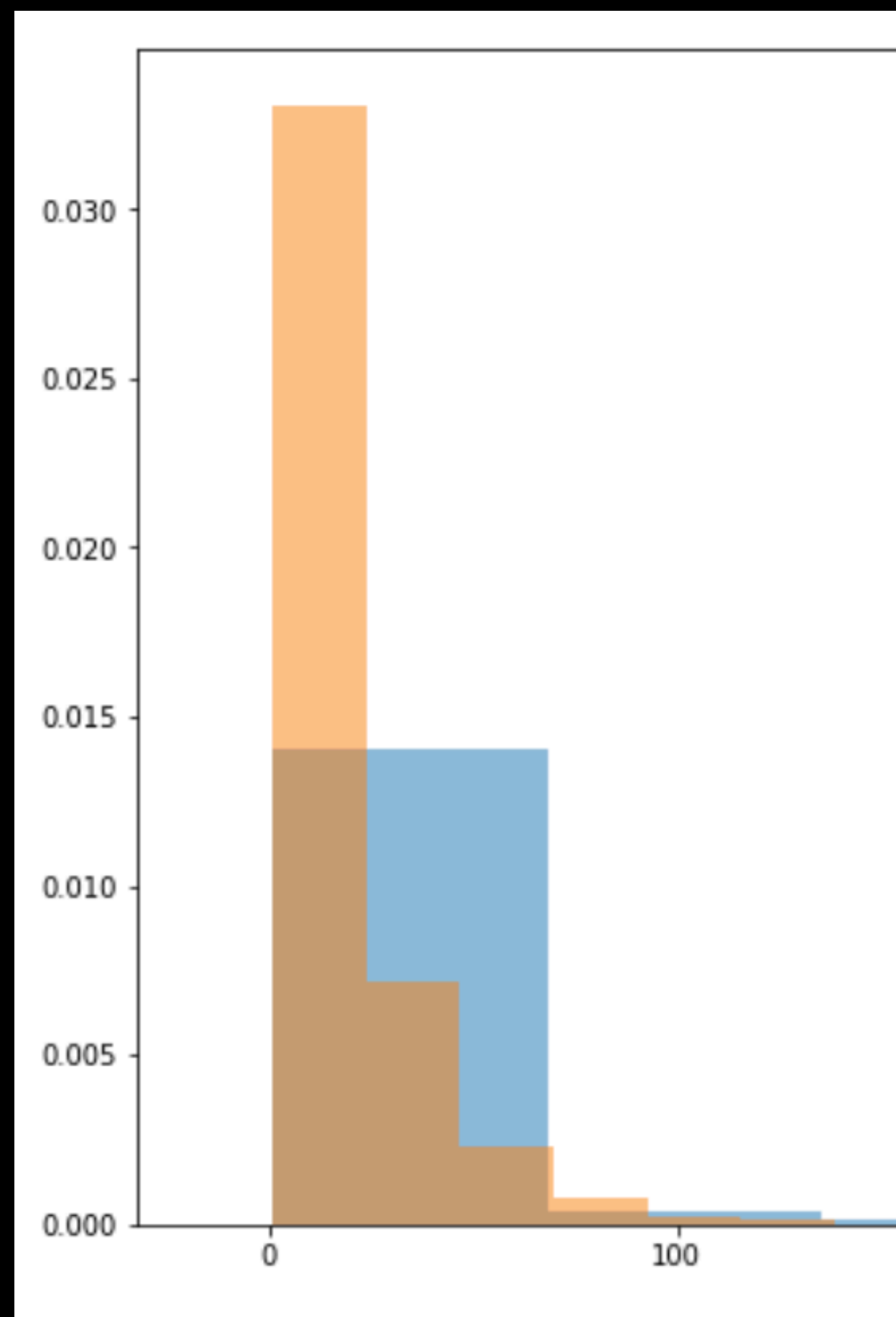


Неверный выбор графика №2

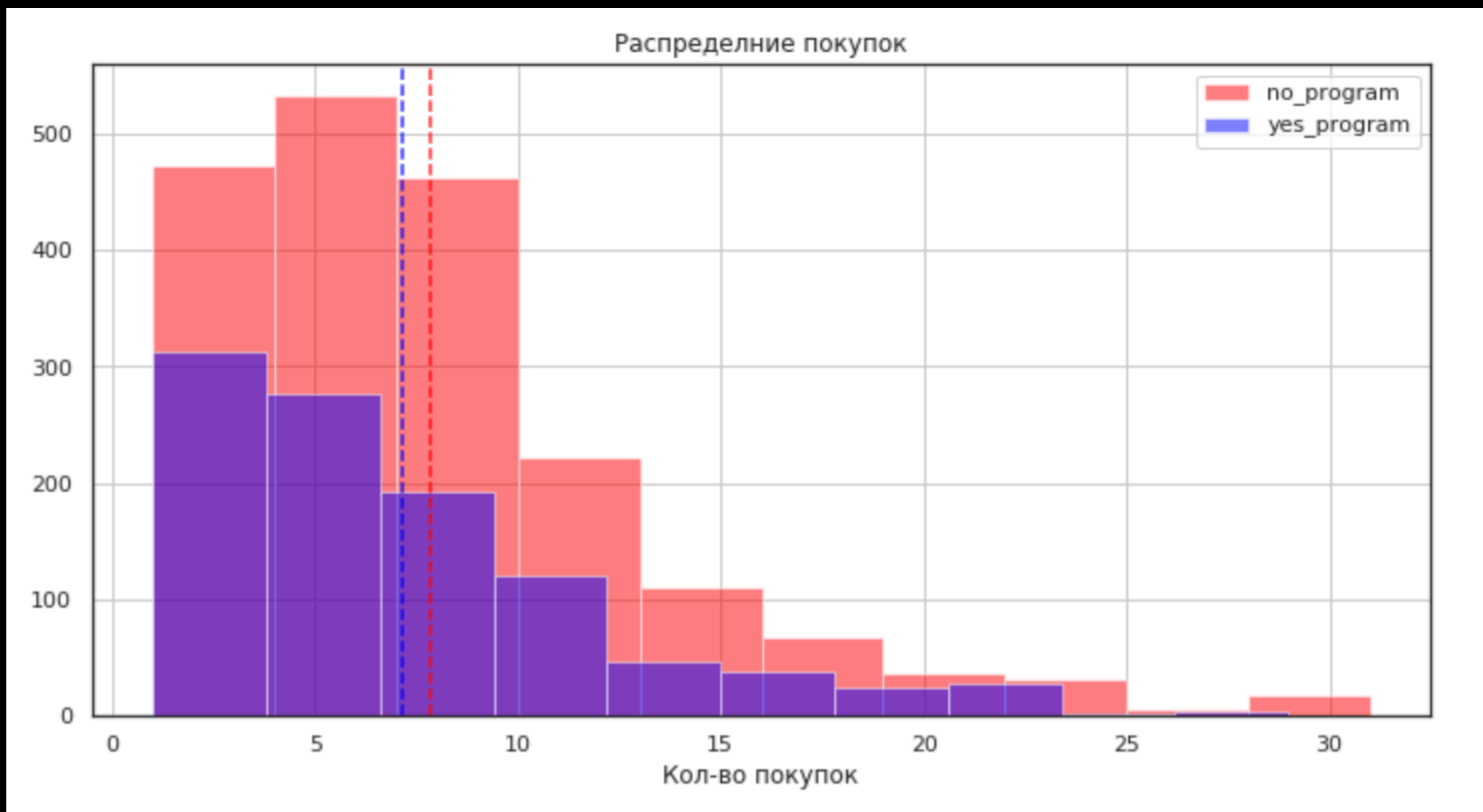
ТВО ИТЕМОВ



Неравномерные бины гистограммы



Неравномерные бины гистограммы



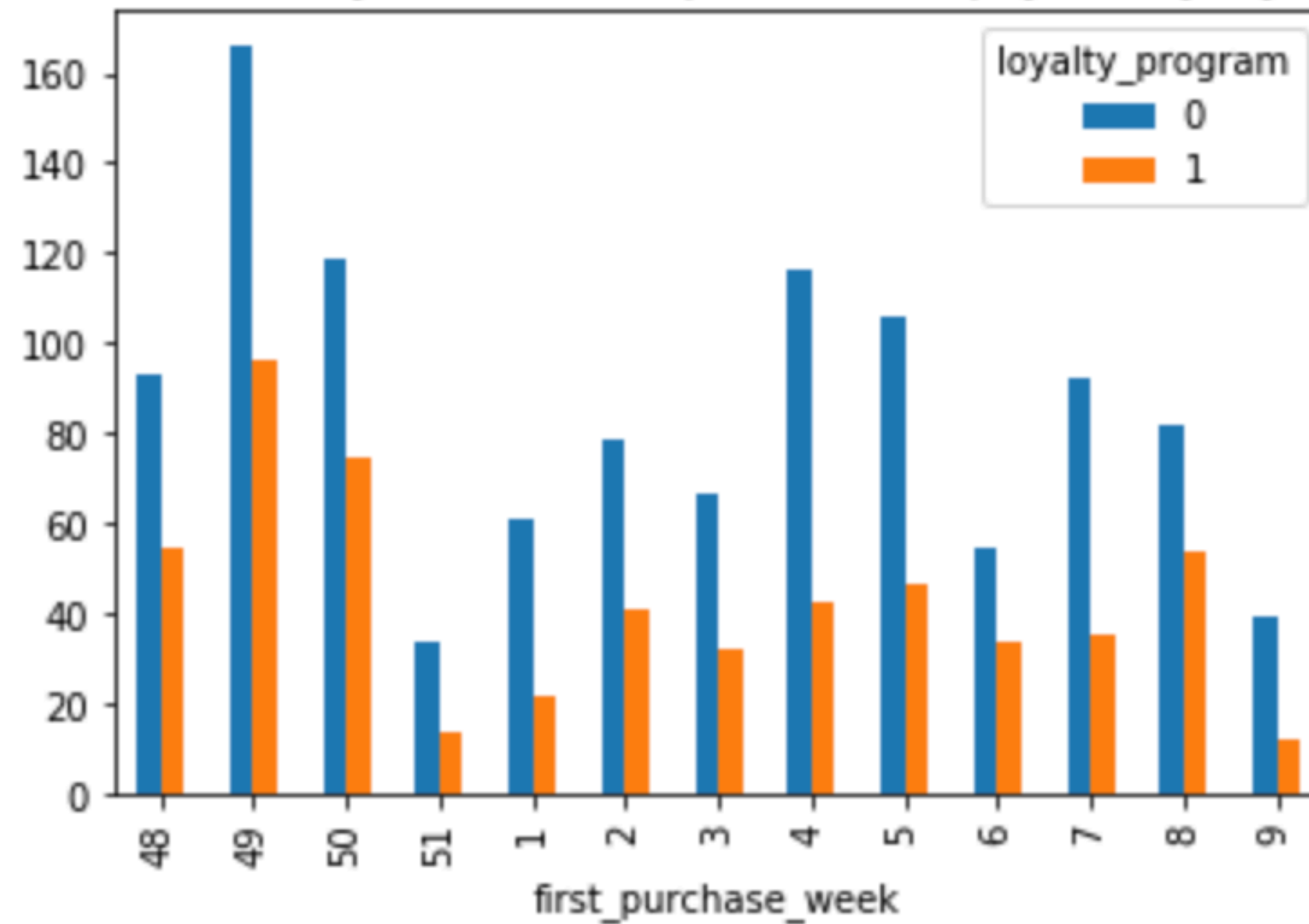
Грамматика и правописание

| Распределние покупок | | |
|----------------------|--|--|
| | | |

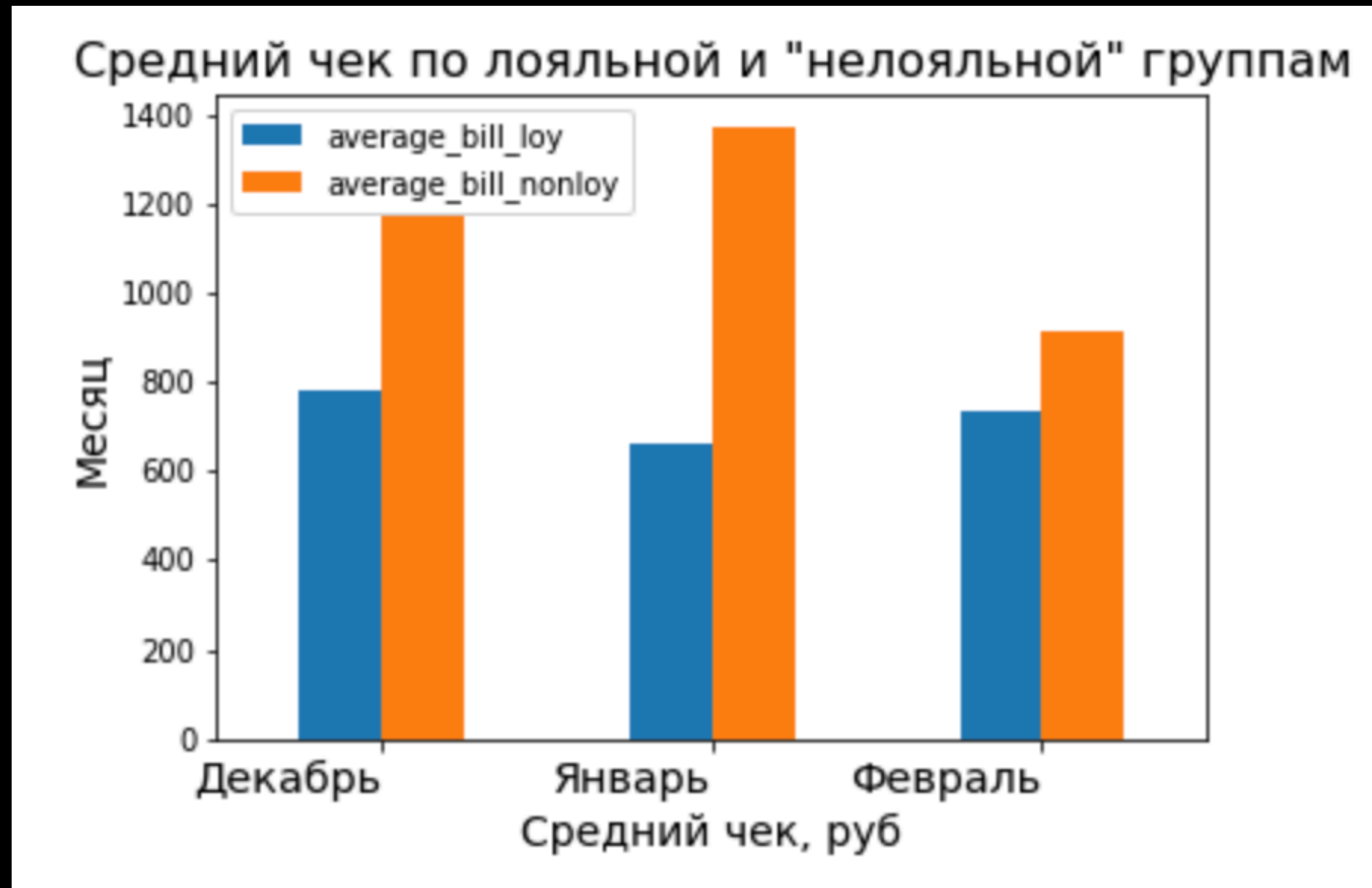
| Распределние средней выручки | | | |
|------------------------------|--|--|--|
| | | | <div><div></div> no_program</div> <div><div></div> yes_program</div> |
| | | | |

Смешивание языков

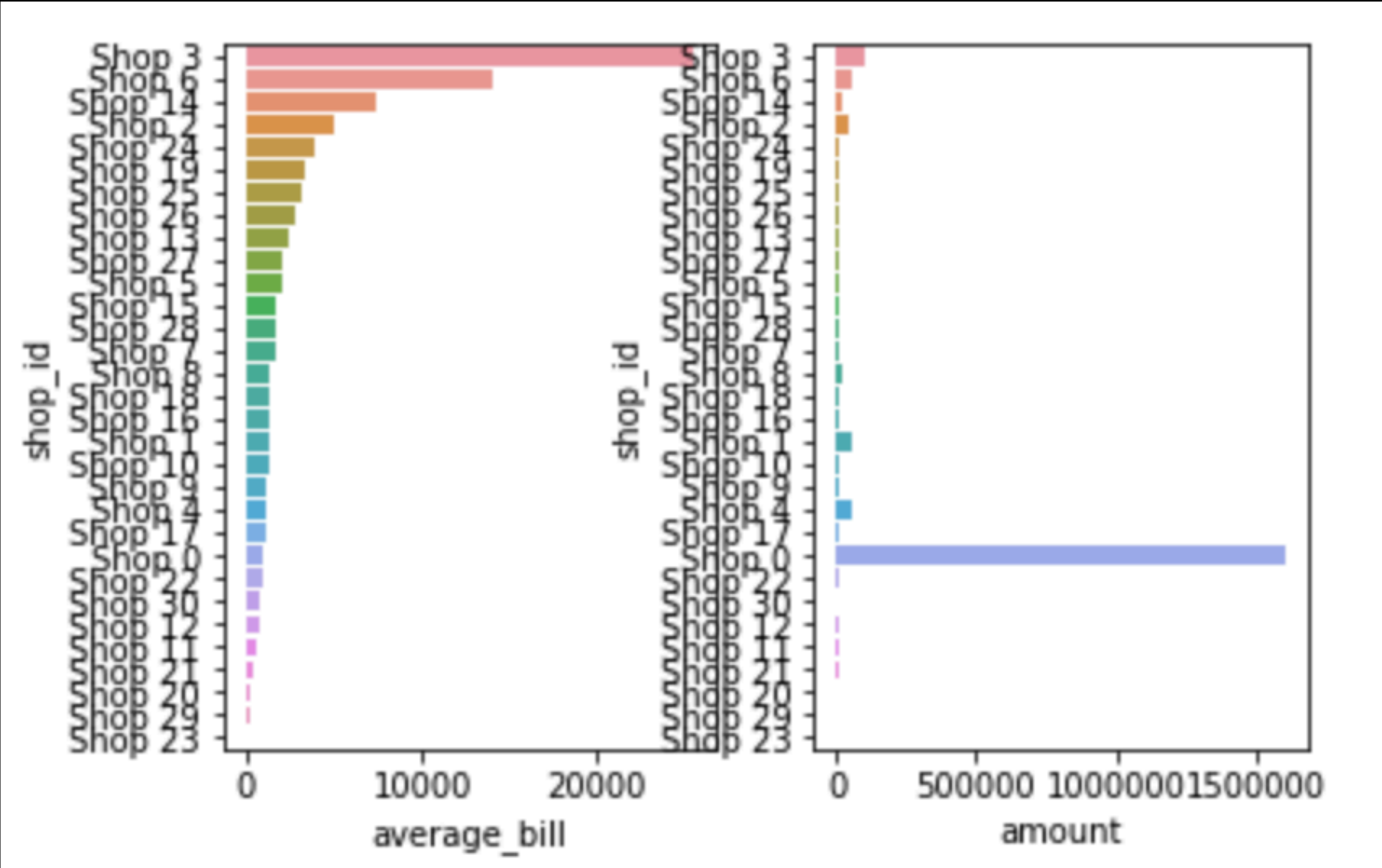
Количество уникальных покупателей, совершивших первую покупку, на каждой неделе



Перекрытия

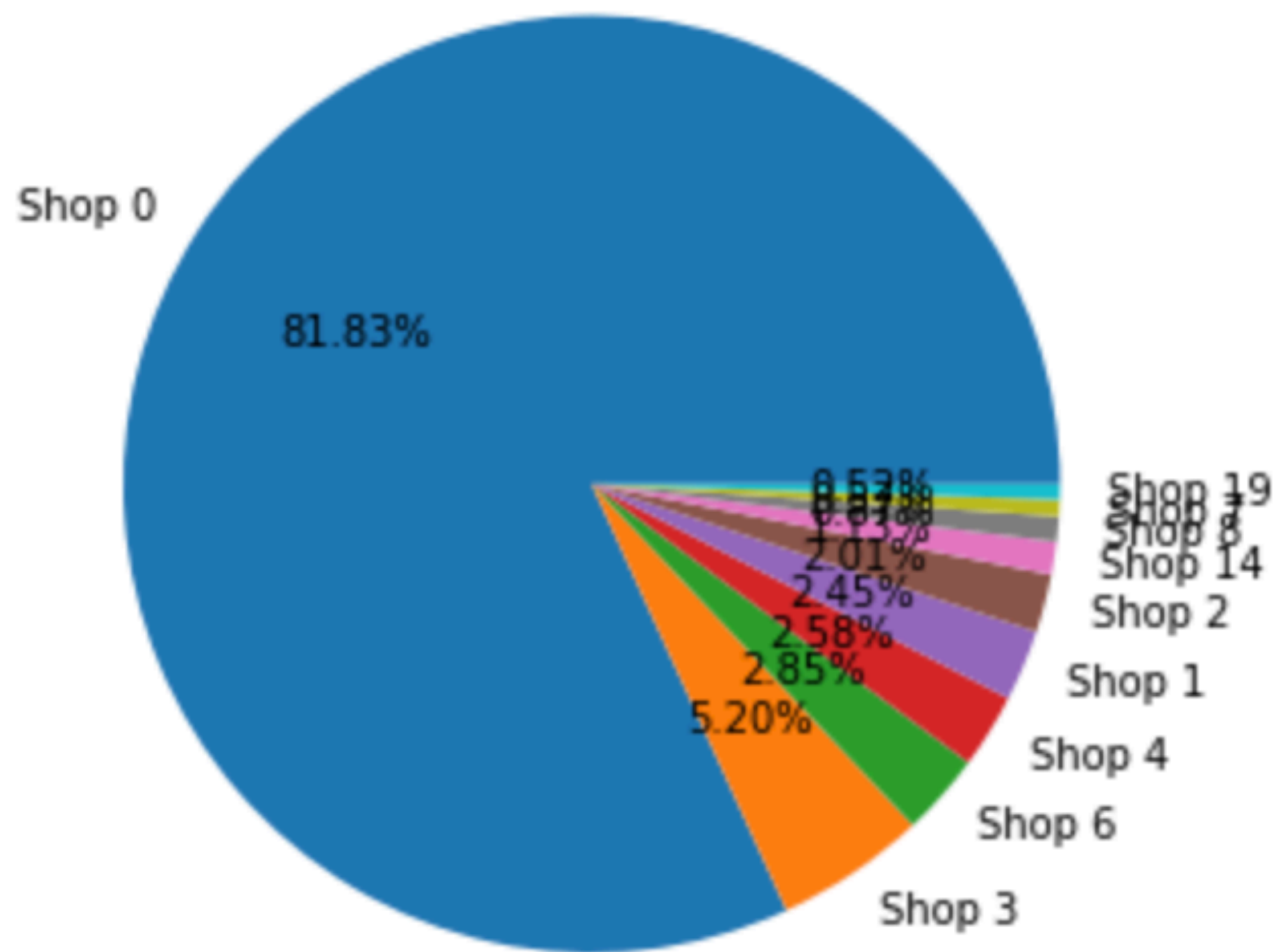


Слишком мелко



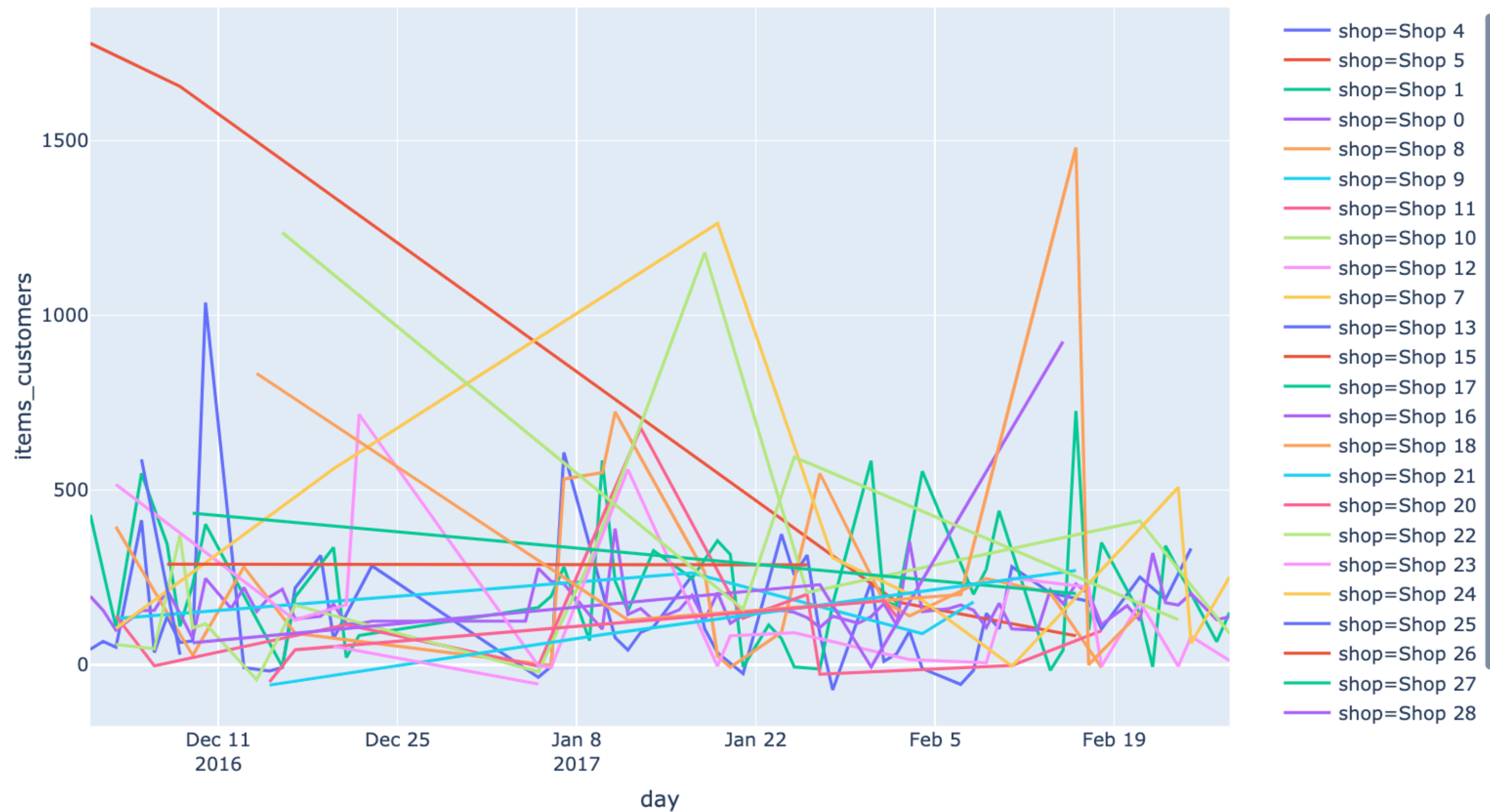
Слишком мелко

10 наиболее доходных магазинов, в % от общего дохода сети



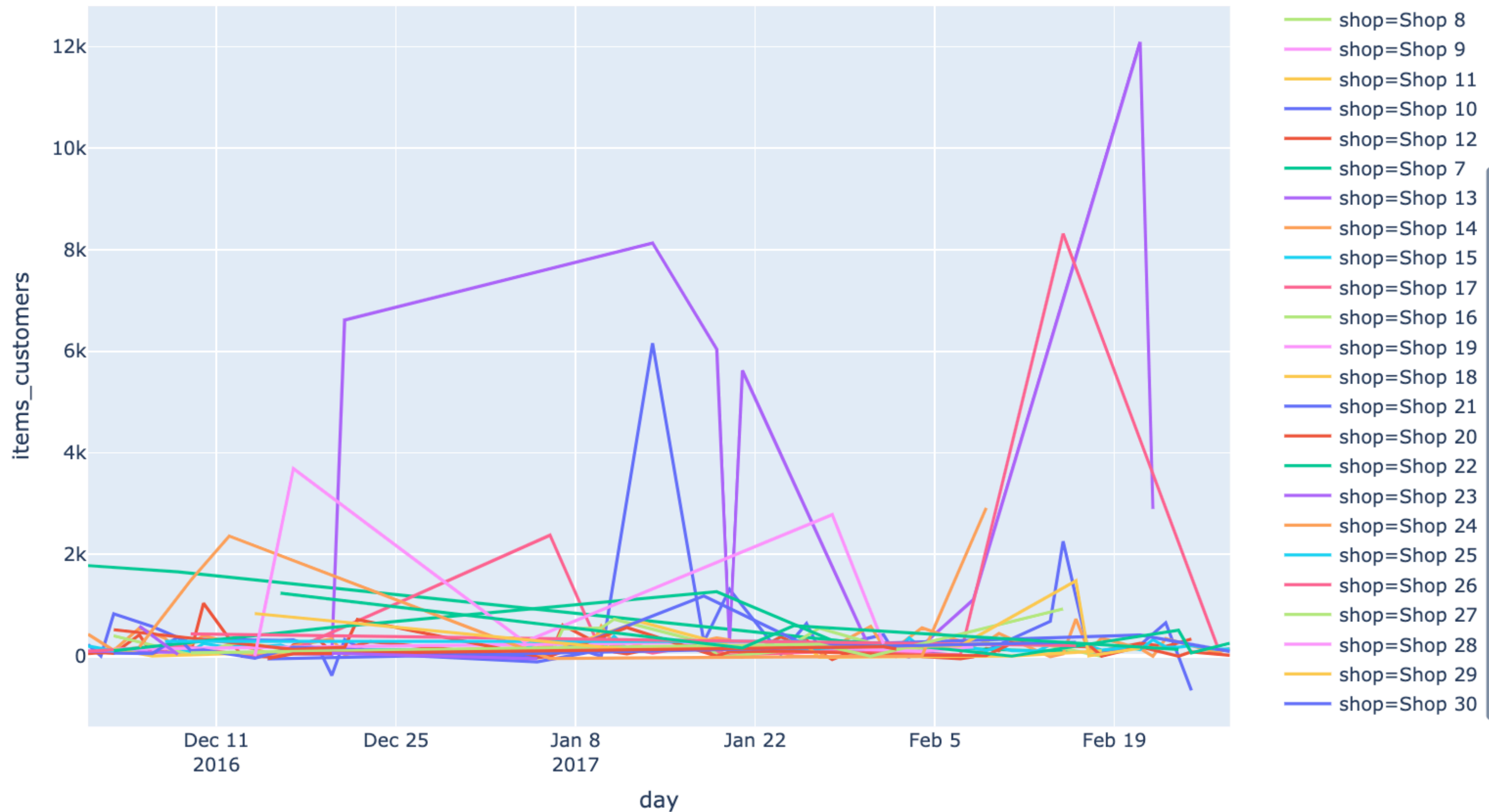
Современное искусство

Метрика: количество товаров на количество покупателей в день: аутсайдер точки



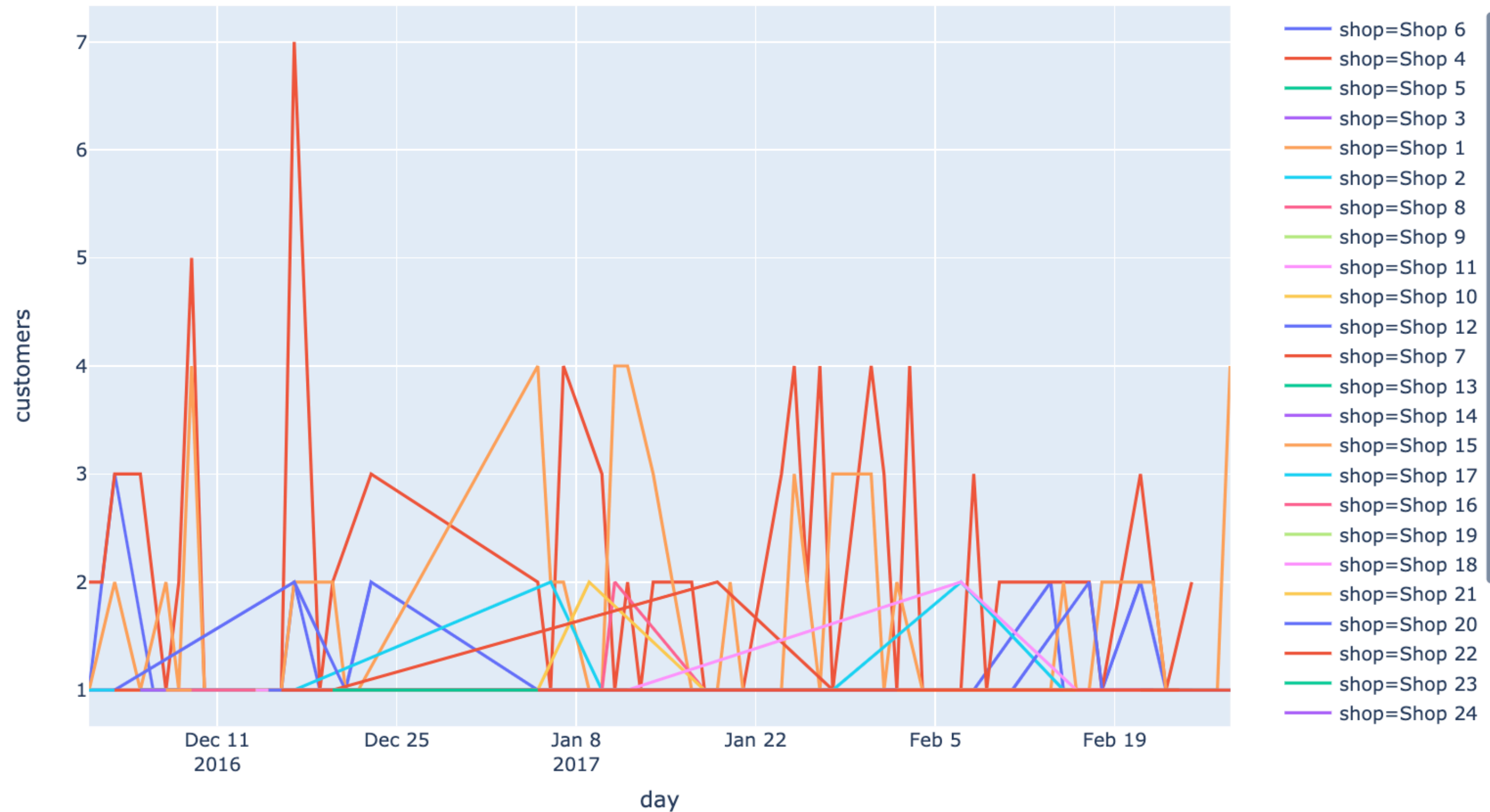
Современное искусство

Метрика: количество товаров на количество покупателей в день все точки



Современное искусство

Метрика: количество покупателей в день все точки кроме топовой точки 0



Ускоряем Pandas и строим правильные графики

Александр Ольферук,
наставник

Яндекс Практикум