

# Поиск клиентов с неоптимальными тарифами

Цели проекта:

- Изучить текущее положение дел с предоставлением услуг связи по тарифам А, В, С.
- Найти клиентов с неоптимальными тарифами и оценить возможности по предоставлению оптимальных.

Этапы проекта:

- Подготовка данных.** Данные будут *загружены, проверены* и при необходимости исправлены.
- Исследование данных.** На основе исходных данных *будет* выполнен расчёт объёма и стоимости услуг связи для каждого абонента за каждый месяц (*биллинг*), *далее будут* составлены портреты среднестатистического *клиента, тарифов* и их *пользователей*.
- Выбор оптимальных тарифов.** На основе биллинга *будет* проанализировано, все ли клиенты подключены к оптимальным тарифам, кто из них может сэкономить, сменить тариф, насколько снизится *доход*, если все клиенты перейдут на оптимальные тарифы, и как снизить *риски* и *компенсировать* снижение дохода.
- Гипотезы.** Будут проверены статистическая значимость *снижения доходов* после выбора пользователями оптимальных тарифов, а также то, что характеристики пользователей тарифа *В превосходят* характеристики пользователей тарифа *С*.
- Итоги.** В заключение будут подведены итоги проекта: сделаны основные выводы о положении дел с предоставлением услуг связи, а также даны рекомендации по улучшению ситуации с тарифами.

Примечания:

- Более детально структуру проекта можно посмотреть в автоматическом оглавлении тетради.
- По проекту подготовлена *презентация-очёт* для генерального директора.
- Часть исследования представлена в виде интерактивного *дашборда*.

## Подготовка данных

Подключим все необходимые для дальнейшей работы программные модули:

```
In [1]: import math
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from scipy import stats as st
```

## Загрузка данных

Загрузим данные из файла в память:

```
In [2]: clients = pd.read_csv('data/telecom_clients.csv')
calls = pd.read_csv('data/telecom_dataset.csv')
```

Проверим загрузку:

```
In [3]: clients.head(3)

Out[3]:   user_id  tariff_plan  date_start
0    166713      A    2019-08-15
1    166901      A    2019-08-23
2    168527      A    2019-10-29
```

```
In [4]: calls.head(3)

Out[4]:   user_id      date  direction  internal  operator_id  is_missed_call  calls_count  call_duration  total_call_duration
0    166377  2019-08-04 00:00:00-03:00      in      False      NaN              True           2           0           4
1    166377  2019-08-05 00:00:00-03:00      in      False    880022.0          True           3           0           5
2    166377  2019-08-05 00:00:00-03:00      out       True    880020.0          True           1           0           1
```

Структура данных соответствует описанию - загрузка произведена корректно.

## Предобработка данных

В каждой строке каждого набора проверим данные на наличие пропусков, некорректных значений, дублей и неверных типов. Исправим обнаруженные недостатки.

### Набор данных clients

```
user_id
Проверим тип значений:

In [5]: clients['user_id'].dtype

Out[5]: dtype('int64')
```

Тип значений соответствует хранимым в столбце данным.

Проверим наличие пропусков:

```
In [6]: clients['user_id'].isna().mean()

Out[6]: 0.0

Пропусков нет.
```

Проверим корректность значений - выисним их диапазон:

```
In [7]: clients['user_id'].min(), clients['user_id'].max()

Out[7]: (166373, 168606)
```

Идентификаторы клиентов находятся в подходящем диапазоне - некорректных значений нет.

Проверим значения на уникальность:

```
In [8]: clients['user_id'].duplicated().mean()

Out[8]: 0.0

Дублей нет - все значения уникальны.
```

tariff\_plan

Для удобства сократим название столбца:

```
In [9]: clients = clients.rename(columns={'tariff_plan': 'tariff'})

Проверим наличие пропущенных значений:
```

```
In [10]: clients['tariff'].isna().mean()

Out[10]: 0.0

Пропущенных значений нет.
```

Проверим корректность значений и их тип:

```
In [11]: clients['tariff'].unique()

Out[11]: array(['A', 'C', 'B'], dtype=object)

Значения и их тип корректны.
```

```
date_start
Проверим тип значений:

In [12]: clients['date_start'].dtype

Out[12]: dtype('O')
```

Тип неверен, изменим его:

```
In [13]: clients['date_start'] = pd.to_datetime(clients['date_start'], infer_datetime_format=True)

Проверим наличие пропущенных значений:
```

```
In [14]: clients['date_start'].isna().mean()

Out[14]: 0.0

Пропущенных значений нет.
```

Проверим корректность значений по диапазону:

```
In [15]: clients['date_start'].min(), clients['date_start'].max()

Out[15]: (Timestamp('2019-08-01 00:00:00'), Timestamp('2019-10-31 00:00:00'))
```

Диапазон не вызывает сомнений.

Проверим распределение:

```
In [16]: clients['date_start'].hist(figsize=(12,3), bins=30*3, color='purple')
plt.show()
```

Распределение вполне естественно: подключение к тарифам в основном по будням, примерно на одинаковом уровне. Заметных аномалий в распределении нет. Пик 01.08.2019 объясняется, скорее всего, либо ярким рекламным мероприятием в этот день (выставка/конференция на тему телекоммуникаций, обзором потенциальных клиентов), либо единовременной регистрацией в этот день ранее накопленных заявок на подключение.

### Набор данных calls

```
user_id
Проверим тип значений:

In [17]: calls['user_id'].dtype

Out[17]: dtype('int64')
```

Тип значений соответствует хранимым в столбце данным.

Проверим наличие пропусков:

```
In [18]: calls['user_id'].isna().mean()

Out[18]: 0.0

Пропусков нет.
```

Проверим корректность значений по диапазону:

```
In [19]: calls['user_id'].min(), calls['user_id'].max()

Out[19]: (166373, 168606)
```

Идентификаторы клиентов находятся в подходящем диапазоне - некорректных значений нет.

Проверим, все ли идентификаторы клентов из набора calls присутствуют в наборе clients:

```
In [20]: calls['user_id'].isin(clients['user_id']).mean()

Out[20]: 1.0
```

Да, все клиенты (100%) из набора calls, присутствуют в наборе clients.

Проверим обратное: для всех ли клиентов из набора clients, есть записи в наборе calls:

```
In [21]: clients['user_id'].isin(calls['user_id']).mean()

Out[21]: 0.41939896719382516
```

Для 58% клиентов из clients нет записей об использовании услуги связи в calls. Быть может эти клиенты только-только подключились и ещё не начали работать. Посмотрим моменты подключений:

```
In [22]: clients[clients['user_id'].isin(calls['user_id'])]['date_start'].hist(
figsize=(12,3), bins=100)
plt.show()
```

Многие из неактивных клиентов подключились уже давно (несколько месяцев назад), но всё ещё не начали пользоваться связью. Это странно: нерациональное поведение. И его демонстрирует 58% клиентов. Из чего следует, что пассивные клиенты включены в этот день по ошибке. Исключим их из набора clients:

```
In [23]: clients = clients[clients['user_id'].isin(calls['user_id'])]

date
Проверим тип значений:

In [24]: calls['date'].dtype

Out[24]: dtype('O')
```

Тип неверен, изменим его:

```
In [25]: calls['date'] = pd.to_datetime(calls['date'], infer_datetime_format=True)

Проверим наличие пропущенных значений:
```

```
In [26]: calls['date'].isna().mean()

Out[26]: 0.0

Пропущенных значений нет.
```

Проверим корректность значений по диапазону (сравним с диапазоном дат из calls):

```
In [27]: clients['date_start'].min(), clients['date_start'].max()

Out[27]: (Timestamp('2019-08-01 00:00:00'), Timestamp('2019-10-31 00:00:00'))

In [28]: calls['date'].min(), calls['date'].max()

Out[28]: (Timestamp('2019-08-02 00:00:00+0300', tz='pytz.FixedOffset(180)'),
Timestamp('2019-11-28 00:00:00+0300', tz='pytz.FixedOffset(180)'))
```

Время, за которое есть данные в calls и clients, почти совпадает. Отличие одно: в calls есть дополнительно данные за ноябрь.

Время, за которое есть данные в calls и clients, почти совпадает. Отличие одно: в calls есть дополнительно данные за ноябрь. А клиентов, зарегистрированных в ноябре, нам знать не надо, т.к. по ним нет полного месяца использования услуг связи. Законч. образом, диапазоны дат в calls и clients правильно синхронизированы и верны.

Проверим распределение:

```
In [29]: calls['date'].hist(figsize=(12,3), bins=30*4)
plt.show()
```

Распределение вполне естественно: по мере получения новых и новых клиентов объем вызовов растёт. Кроме того в данных четко просматривается недельный цикл: повышенное число звонков в будние дни и сниженное в выходные. Аномалий нет.

Данные верны.

```
direction
Проверим наличие пропущенных значений:

In [30]: calls['direction'].isna().mean()

Out[30]: 0.0

Пропущенных значений нет.
```

Проверим корректность значений и их тип:

```
In [31]: calls['direction'].unique()

Out[31]: array(['in', 'out'], dtype=object)

Значения и их тип корректны.
```

```
internal
Проверим, тип и значения в столбце:
```

```
In [32]: calls['internal'].unique()

Out[32]: array([False, True, nan], dtype=object)

В столбце есть пропуски. Тип столбца не соответствует типу хранимых значений.
```

Сначала оценим число записей с пропусками и их долю в общем числе записей:

```
In [33]: x = calls['internal'].isna()
x.sum(), x.mean()

Out[33]: (117, 0.062170605914473123)
```

Пропусков незначительное число. Восстановить пропущенные значения исходя из значений других столбцов не представляется возможным. Вероятная причина пропусков - технический сбой. Пропуски придётся удалить:

```
In [34]: calls = calls[~x]

Осталось исправить тип столбца:
```

```
In [35]: calls['internal'] = calls['internal'].astype('bool')

operator_id
Узнаем число пропущенных значений (в абсолютном и относительном выражении):
```

```
In [36]: x = calls['operator_id'].isna()
x.sum(), x.mean()

Out[36]: (8115, 0.15087849772241332)
```

Пропущено значительное число значений.

Разберёмся подробнее, в каких случаях оператор не указан. Рассмотрим направление вызовов:

```
In [37]: calls[x]['direction'].value_counts(normalize=True)

In [37]: In      0.975681
out      0.024399
Name: direction, dtype: float64
```

В подавляющем большинстве случаев оператор не указан для входящих вызовов. Посмотрим, состоялось ли соединение (разговор):

```
In [38]: calls[x]['is_missed_call'].mean()

Out[38]: 0.9649661121388016
```

Почти все вызовы, для которых оператор не указан (независимо от направления) закончились неудачей (разговор не состоялся).

Следовательно основной ситуацией, когда оператор не указывается, является поступление вызова в момент, когда все операторы заняты и не отвечают (звонящий положил трубку прежде, чем кто-либо из операторов освободился). Данные о таких вызовах можно удалить, т.к. несоответствие вызовов без указания оператора в данном исследовании нас не интересует:

```
In [39]: calls = calls[~(x & (calls['direction']=='in') & calls['is_missed_call'])]

В остальных ситуациях (исходящий вызов или состоявшийся разговор) оператор должен быть указан, а если это не так, то единственное объяснение этому - технический сбой. Выясним число таких записей:
```

```
In [40]: x = calls['operator_id'].isna()
x.sum(), x.mean()

Out[40]: (271, 0.0658988702901547636)
```

Число сбойных записей незначительно - удалим их:

```
In [41]: calls = calls[~x]

Пропуски исключены.
```

Проверим тип столбца:

```
In [42]: calls['operator_id'].dtype

Out[42]: dtype('float64')
```

Проверим, действительно ли значения представлены числами с плавающей запятой:

```
In [43]: (calls['operator_id'] != calls['operator_id'].round()).sum()

Out[43]: 0
```

Ни у одного идентификатора не обнаружено дробной части, значит идентификаторы - целые числа. Соответственно исправим тип:

```
In [44]: calls['operator_id'] = calls['operator_id'].astype('int')

is_missed_call
Проверим пропуск значений:
```

```
In [45]: calls['is_missed_call'].isna().sum()

Out[45]: 0
```

Пропущенных значений нет.

О том, состоялся разговор или нет, можно судить не только по столбцу is\_missed\_call, но и по столбцу call\_duration. Проверим, есть ли несогласованность значений столбцов:

```
In [46]: x = ~(calls['is_missed_call'] & (calls['call_duration'] == 0)) \
| (calls['is_missed_call'] & calls['call_duration'] > 0)
x.sum(), x.mean()

Out[46]: (344, 0.007532296912634114)
```

Рассинхронизация значений в столбцах is\_missing\_call и call\_duration встречается в незначительном числе строк. Причиной этого может быть только технический сбой. Противоречивым данным доверять нельзя - придётся удалить строки с ними:

```
In [47]: calls = calls[~x]

Теперь столбец is_missed_call больше не нужен, его можно удалить:
```

```
In [48]: calls = calls.drop(columns='is_missed_call')

calls_count
Для удобства дадим столбцу более короткое имя:
```

```
In [49]: calls = calls.rename(columns={'calls_count': 'count'})

Проверим тип столбца:
```

```
In [50]: calls['count'].dtype

Out[50]: dtype('int64')
```

Тип столбца соответствует характеру хранимой информации.

Проверим, нет ли пропущенных значений:

```
In [51]: calls['count'].isna().mean()

Out[51]: 0.0

Пропущенных значений нет.
```

Проверим корректность значений, исследовав их диапазон:

```
In [52]: calls['count'].min(), calls['count'].max()

Out[52]: (1, 4817)
```

Верхняя граница диапазона кажется слишком большой. Посмотрим диаграмму размаха:

```
In [53]: calls.boxplot(column='count', figsize=(12,1), vert=False)
plt.xticks([])
plt.show()
```

Согласно диаграмме значений, близких к верхней границе диапазона нет. Убедимся в этом:

```
In [54]: (calls['count'] > 4000).sum()

Out[54]: 1
```

Предположим, что больше значения - это выбросы. Узнаем их количество по стандартной методике полнотерного межквартильного размаха:

```
In [55]: q1 = calls['count'].quantile([0.25, 0.75])
q1, q3 = q1[0.25], q1[0.75]
limit = q3 + 1.5 * (q3 - q1)
x = calls[calls['count'] > limit].sum()
x, x / calls.shape[0]

Out[55]: (5289, 0.11668799364603996)
```

Методика определена как выбросы слишком большое число значений. Скорее всего, большое число вызовов характерно для операторов, специально обслуживающих большие группы людей в ходе программ информирования или в рекламных целях. Для этого могут использоваться роботы или поменно работающие сотрудники call-центров. Определим предельные возможности таких операторов:

```
In [56]: limit = 24 * 60 * 2
limit

Out[56]: 2880
```

Таким образом, если оператор совершает в среднем 2 звонка в минуту, то он успеет сделать 2880 вызовов за сутки. В среднем 30 сек. на вызов. Если большинство абонентов не желают слушать сообщение и почти сразу бросают трубку, то на вызов может тратиться меньше времени, например, введое, тогда число вызовов (в среднем) вырастет 2880 \* 2 = 5760. Самые большие числа вызовов в наборе данных не превосходят это число (4817 < 5760), значит их нельзя считать выбросами. Данные корректны.

call\_duration

Переименуем столбец для удобства:

```
In [57]: calls = calls.rename(columns={'call_duration': 'duration'})

Проверим наличие пропусков:
```

```
In [58]: calls['duration'].isna().mean()

Out[58]: 0.0

Пропусков нет.
```

Проверим тип столбца:

```
In [59]: calls['duration'].dtype

Out[59]: dtype('int64')
```

Тип соответствует хранимым данным (длительность соединений в секундах).

Проверим значения столбца по максимальной длительности соединений. При этом учтём, что длительность, равная нулю, установлена для вызовов, не закончившихся соединением (разговором), а если длительность больше нуля, то это суммарная длительность того числа звонков, которые агрегированы в этой строке данных:

```
In [60]: x = calls['duration'] > 0
y = calls.loc[x, 'duration'] / calls.loc[x, 'count']
y.min(), y.max()

Out[60]: (0.5, 3550.0)
```

Разговоры длятся от полсекунды до часа (60 \* 60 = 3600). Само по себе значение верхней границы не вызывает сомнений, однако, возможно, если просуммировать записи из расчёта на оператора, то сумма соединений превысит пределы человеческих возможностей. Проверим, нет ли таких операторов, сумма разговоров которых за сутки превышает 10 часов (рабочий день с запасом):

```
In [61]: wrong_operators = {
calls.drop_duplicates()
.groupby(['operator_id', 'date'])
.agg(['duration', 'sum'])
.reset_index()
.query('duration > (60*60*10)')['operator_id']
.unique()
}
calls['operator_id'].isin(wrong_operators).sum() / calls.shape[0]

Out[61]: 0.01922340378596565
```

Такие операторы существуют - их 2% от общего числа операторов. Их наличие можно объяснить только ошибкой в сборе данных.

Удалим данные по таким операторам:

```
In [62]: calls = calls[~calls['operator_id'].isin(wrong_operators)]

Данные по операторам со сверхкоэффициентами удалены!
```

Теперь проверим нижнюю границу длительности разговора:

```
In [63]: (y < 5).mean()

Out[63]: 0.013995621073241961
```

Соединения меньше полсекунды редки. Скорее всего, это случаи, 1) когда абонент случайно нажал подряд кнопки приёма и отбоя звонка подбод, 2) когда вызываемый уже поднял трубку, но вызывающий положил её почти в тот же момент, думая, что не дождался ответа, 3) когда связь обрывалась по техническим причинам: слабый сигнал, разряжен аккумулятор и т.п.

Проверим размах и распределение вызовов по длительности соединения:

```
In [64]: y.to_frame().boxplot(figsize=(12,1), vert=False)
plt.xticks([])
plt.show()
```

Распределение вызовов по длительности соединения (диапазон 0-480 сек, медиана - 80, среднее - 117)

Распределение для длительности звонков вполне естественно, аномалий не наблюдается. Данные корректны.

total\_call\_duration

Данные о длительности звонков с учётом ожидания не будут использоваться в настоящем проекте: для биллинга эти данные бесполезны, а технические аспекты (время использования линий связи) выходят за рамки исследования. Удаляем столбец:

```
In [66]: calls = calls.drop(columns='total_call_duration')

Дублирование строк
```

Проверим, нет ли дублирования строк в наборе данных:

```
In [67]: calls.duplicated(keep='first').sum()

Out[67]: 4859
```

Дубли имеются (возможно появились после удаления столбцов или были изначально). Удаляем дубли:

```
In [68]: calls = calls.drop_duplicates()

Синхронизация clients-calls
```

В ходе переработки набора calls возможно из набора исчезли некоторые клиенты, которые присутствуют в наборе clients. Проверим это:

```
In [69]: x = ~clients['user_id'].isin(calls['user_id'])
x.sum(), x.mean()

Out[69]: (18, 0.0565319218249423)
```

Да, небольшое количество клиентов исчезло из calls из-за некорректных данных по ним. Удалим этих клиентов также из набора clients:

```
In [70]: clients = clients[~x]
```

## Презентация

В целях создания презентации создадим функцию сохранения диаграмм в файлах на диске:

```
In [71]: savefig_file_number = 0
def savefig():
    global savefig_file_number
    plt.savefig(
        # '/home/alex/Documents/ya-praktikum/analyst/projects'
        # '/final/ides/(-.0f).png'.format(savefig_file_number),
        dpi=200, transparent=True)
    savefig_file_number += 1
```

## Дашборд

Выгрузим данные для интерактивного дашборда:

```
In [72]: dash = calls.drop(columns=['user_id', 'operator_id', 'duration'])
dash['date'] = dash['date'].dt.date.astype('datetime64[D]')
dash = dash.groupby(['date', 'direction', 'internal'])['count'].sum()
dash = dash.reset_index()
dash['direction'] = dash['direction'].map({'in': 'входящий', 'out': 'исходящий'})
dash['internal'] = dash['internal'].map({True: 'внутренний', False: 'внешний'})
dash.to_csv(
    # '/home/alex/Documents/ya-praktikum/analyst/projects/final/data/dash-final.csv',
    index=False)
```

## Итоги

Данные загружены и проверены. Обнаруженные дефекты исправлены. Данные готовы к использованию.

## Исследовательский анализ

Выполним биллинг - рассчитаем для каждого клиента за каждый месяц:

- сколько работало операторов,
- сколько среднее время, закончившихся соединением, и какова их суммарная продолжительность:
- исходящих внешних
- исходящих внутренних
- входящих внешних
- сколько стоили услуги (по тарифу клиента).

Данные биллинга будут основой дальнейшего анализа.

Скомпируем данные, необходимые для расчёта стоимости потребленных услуг (каждым клиентом в каждом месяце):

```
In [73]: # добавляем метку месяца для последующих группировок
calls['month'] = calls['date'].dt.date.astype('datetime64[M]')
# создаём базу биллинга
bills = clients[['user_id', 'tariff']]
# bills - клиенты
bills = bills.merge(
    calls.groupby(['user_id', 'month'])['operator_id'].nunique().reset_index(),
    how='left', on='user_id', left_index=True)
bills = bills.rename(columns={'operator_id': 'operators'})
# добавляем в биллинг счётчики внешних (вход и исход) вызовов
calls_out = (
    calls[~calls['internal'] & (calls['duration'] > 0)]
    .pivot_table(values='duration', count=1, index='user_id', month=1,
    columns='direction', aggfunc='sum')
    .reset_index()
)
calls_out.columns = ['user_id', 'month', 'in_count', 'out_count',
'in_duration', 'out_duration']
bills = bills.merge(calls_out, how='left', on='user_id', month=1)
# добавляем в биллинг счётчики внутренних исходящих вызовов
calls_internal = (
    calls[calls['internal'] & (calls['direction']=='out') & (calls['duration'] > 0)]
    .groupby(['user_id', 'month'], as_index=False)
    .agg(['count', 'sum', 'duration'])
)
calls_internal = calls_internal.rename(columns={'count': 'int_count',
'duration': 'int_duration'})
bills = bills.merge(calls_internal, how='left', on='user_id', month=1)
# заполним пропущенные значения нулями
bills = bills.fillna(0)
for column in bills.columns[4:]:
    bills[column] = bills[column].astype('int')
bills.head()
```

```
Out[73]:   user_id  tariff      month  operators  in_count  out_count  in_duration  out_duration  int_count  int_duration
0    166713      A    2019-08-01           2           6           0    1201              0           0           0
1    166713      A    2019-09-01           2           10          0    1171              0           1        141
2    166713      A    2019-10-01           2           10           1     536              70           0           0
3    166713      A    2019-11-01           2           17          1     6388             152           0           0
4    166901      A    2019-08-01           7           26          13     709             987          21        1201
```

Проверим, нет ли выровжденных строк (строк где все счётчики равны 0):

```
In [74]: (bills.sum(axis=1, numeric_only=True) == 0).sum()

Out[74]: 0
```

Выровженных строк нет.

Проверим, нет ли среди клиентов таких, для которых за всё время нет ни одного вызова, закончившегося соединением:



	in_count	out_count	in_duration	out_duration	in_count	out_count	in_duration	out_duration
167122	0	0	0	0	0	0	0	0
167364	0	0	0	0	0	0	0	0
168291	0	0	0	0	0	0	0	0

Несколько таких пассивных клиентов существуют. Узнаем, каким числом операторов они располагают:

user_id	tariff	month	operators	in_count	out_count	in_duration	out_duration	in_count	out_count	in_duration	out_duration
84	167364	C	2019-11-01	1	0	0	0	0	0	0	0
469	167122	B	2019-09-01	1	0	0	0	0	0	0	0
449	168291	B	2019-11-01	1	0	0	0	0	0	0	0

Операторов крайне мало. Данных клиентов можно удалить из статистики:

bills = bills[bills['user_id'].isin(bills_bill_total[x].index)]
---

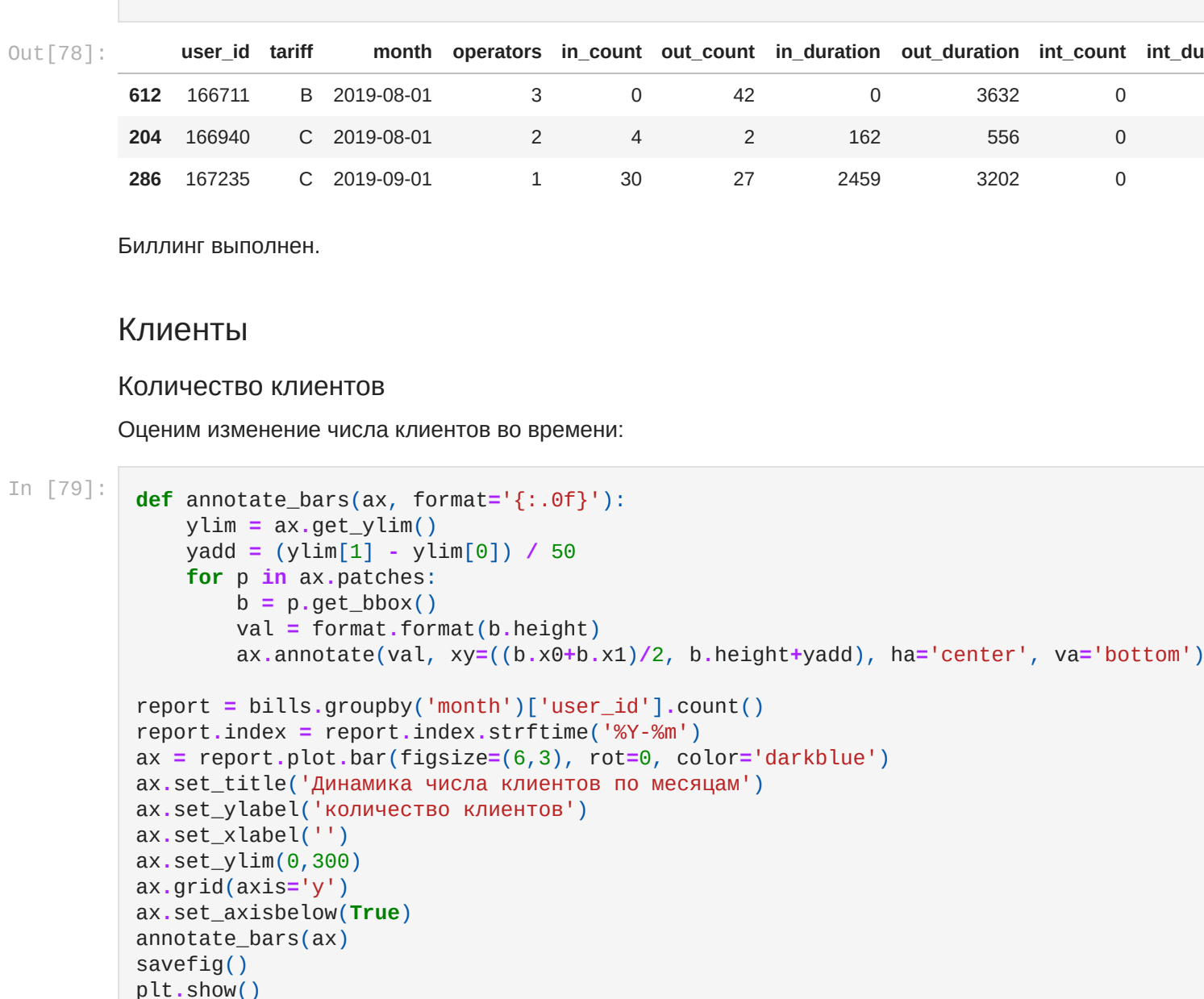
user_id	tariff	month	operators	in_count	out_count	in_duration	out_duration	in_count	out_count	in_duration	out_duration
612	166711	B	2019-08-01	3	4	42	0	3632	0	0	2330.5
204	166940	C	2019-08-01	2	4	2	162	556	0	0	1207.0
286	167325	C	2019-09-01	1	30	27	2459	3202	0	0	1137.8

Биллинг выполнен.

## КЛИЕНТЫ

### Количество клиентов

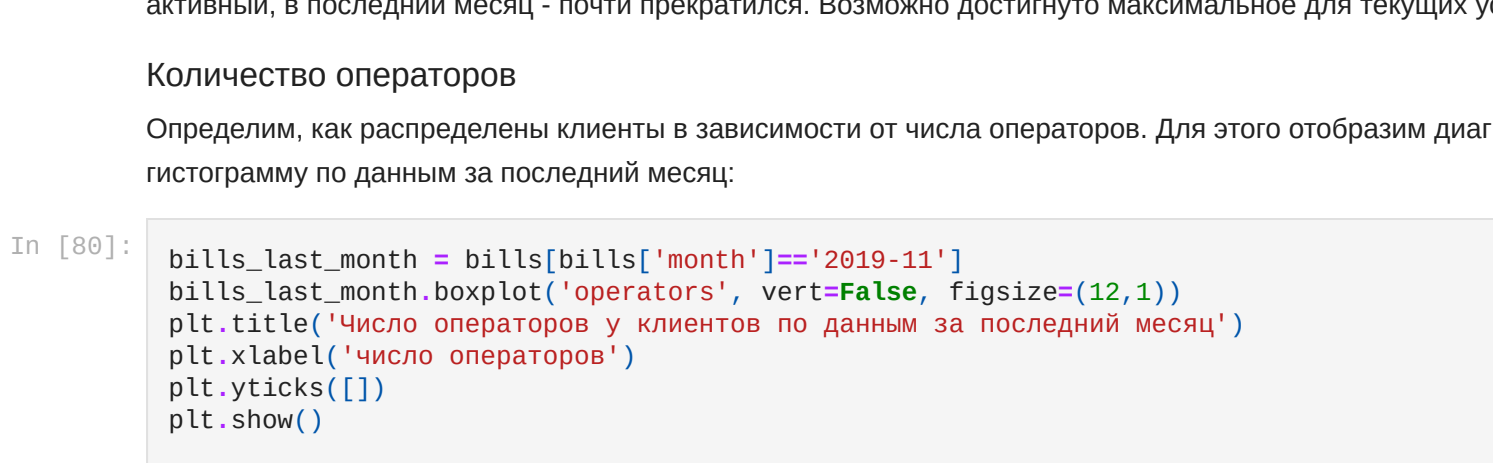
Оценим изменение числа клиентов во времени:



За несколько месяцев число клиентов практически выросло в несколько раз. В первые три месяца - рост был активным, в последний месяц - почти прекратился. Возможно достигнуто для текущих условий число клиентов.

### Количество операторов

Оценим, как распределены клиенты в зависимости от числа операторов. Для этого отобразим диаграмму размаха и гистограмму по данным за последний месяц:



Диаграммы показывают, что распределение далеко от нормального. Оно сильно смещено в область малых значений и имеет большое число выбросов. Для большинства клиентов характерно малое число операторов (1-4). В то же время отдельные клиенты располагают штатом от 20 до 50 операторов.

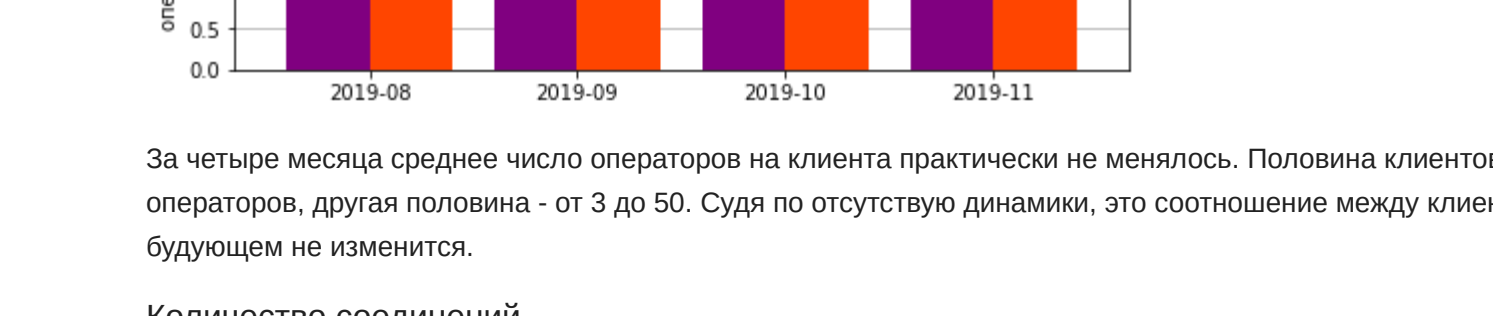
Оценим, как усреднительное значение числа операторов на клиента меняется со временем. Ввиду характера распределения основной мерой среднего выберем медиану, но для справки отобразим также и среднее арифметическое (на нём сильно сказывается появление клиентов с экстремально высокими значениями - исчезающее мало).



За четыре месяца среднее число операторов на клиента практически не менялось. Половина клиентов имеет одного-двух операторов, другая половина - от 3 до 5. Судя по отсутствию динамики, это соотношение между клиентами в ближайшем будущем не изменится.

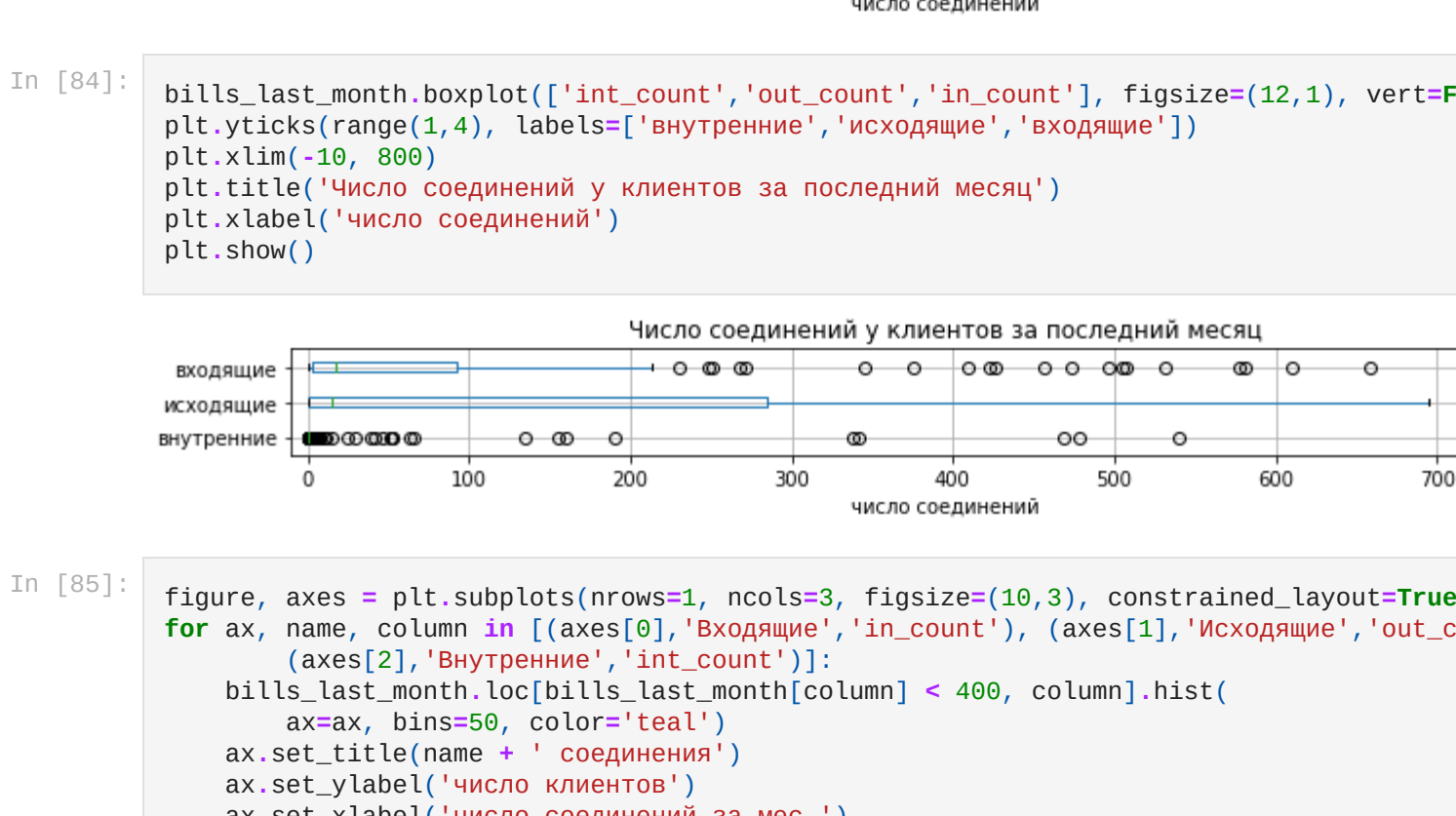
### Количество соединений

Узнаем распределение клиентов в зависимости от числа вызовов, закончившихся соединением (разговором). Для этого отобразим диаграмму размаха и гистограмму по данным за последний месяц:



Диаграммы показывают, что распределение далеко от нормального. Оно сильно смещено в область малых значений и имеет большое число выбросов. Для большинства клиентов совершено небольшое количество соединений. Среди клиентов есть и такие, у которых число соединений идет на тысячи в месяц (1000/30-33 в сутки) и даже есть один клиент примерно с 30 тыс. соединений в месяц (30000/30-1000 в сутки).

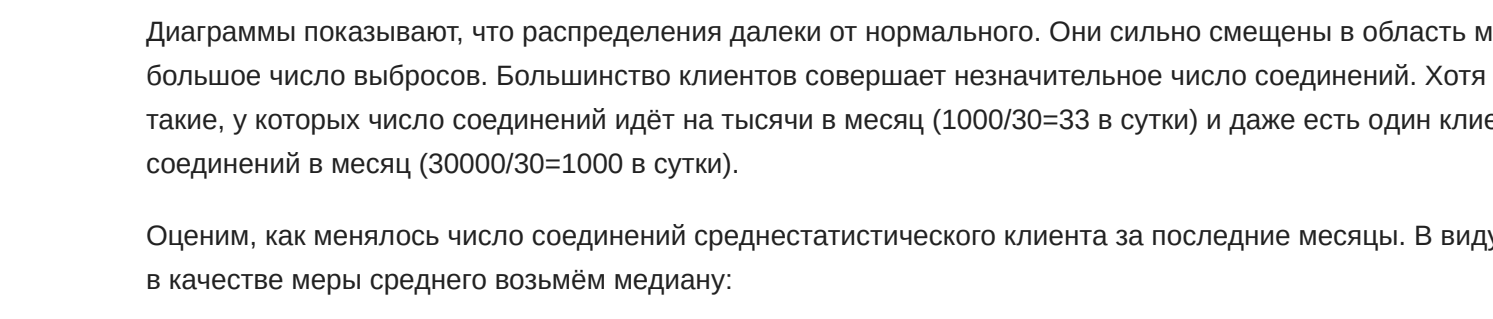
Оценим, как менялось число соединений усреднительного клиента за последние месяцы. Ввиду характера распределения в качестве меры среднего возьмем медиану, но для справки отобразим также и среднее арифметическое (на нём сильно сказывается появление клиентов с экстремально высокими значениями - исчезающее мало).



В первые месяцы после подключения к тарифам новых клиентов среднее число вызовов на клиента росло. В последний месяц практически прекратилось. Вместе со стабилизацией клиентской базы, стабилизировалось и среднее число вызовов на клиента: 50% входящих соединений, 50% исходящих, внутренних соединений - исчезающее мало.

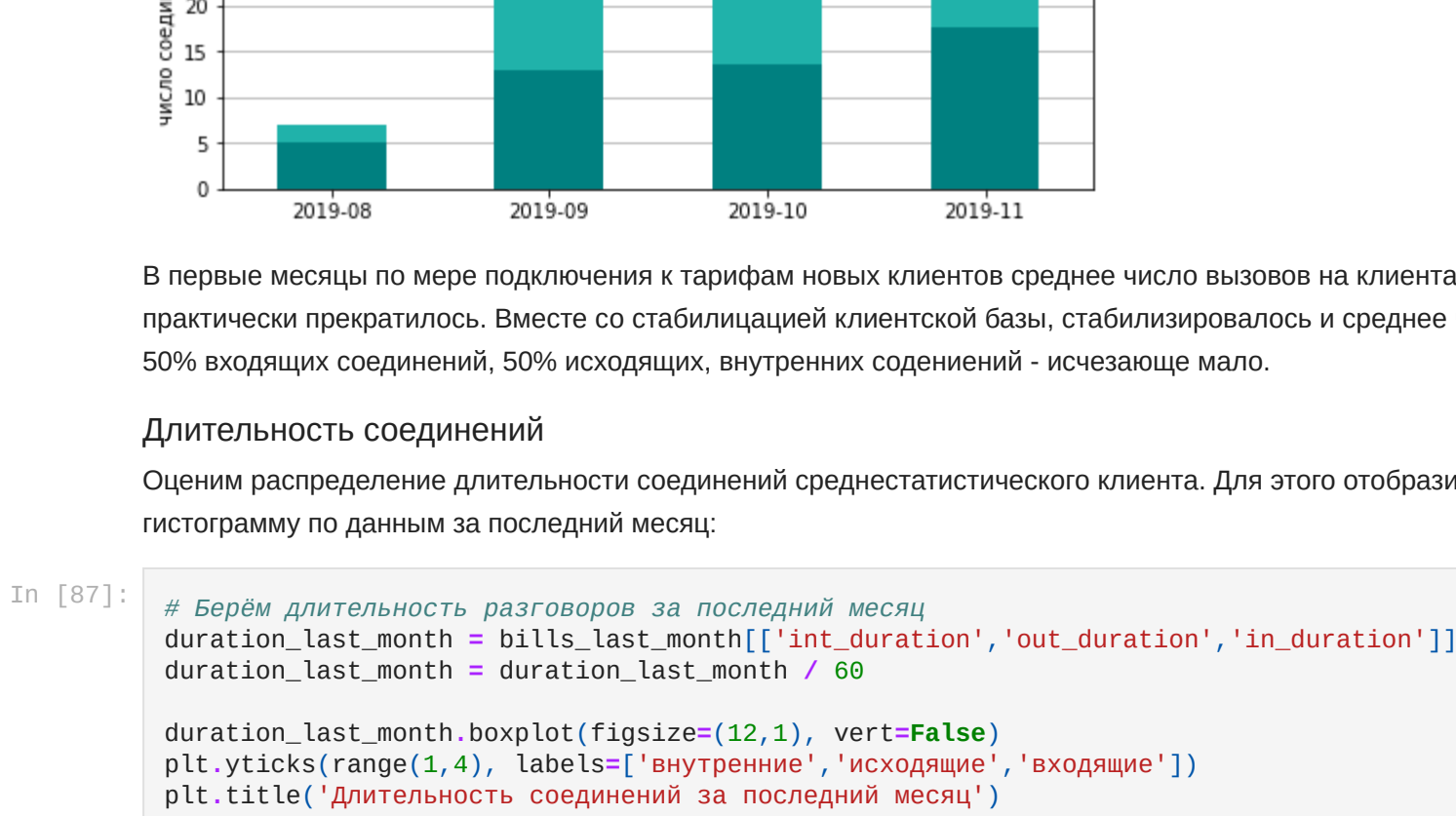
### Длительность соединений

Оценим распределение длительности соединений усреднительного клиента. Для этого отобразим диаграмму размаха и гистограмму по данным за последний месяц:



Диаграммы показывают, что распределение далеко от нормального. Оно сильно смещено в область малых значений и имеет большое число выбросов. Большинство клиентов совершает незначительное число соединений. Среди клиентов есть и такие, у которых число соединений идет на тысячи в месяц (1000/30-33 в сутки) и даже есть один клиент примерно с 30 тыс. соединений в месяц (30000/30-1000 в сутки).

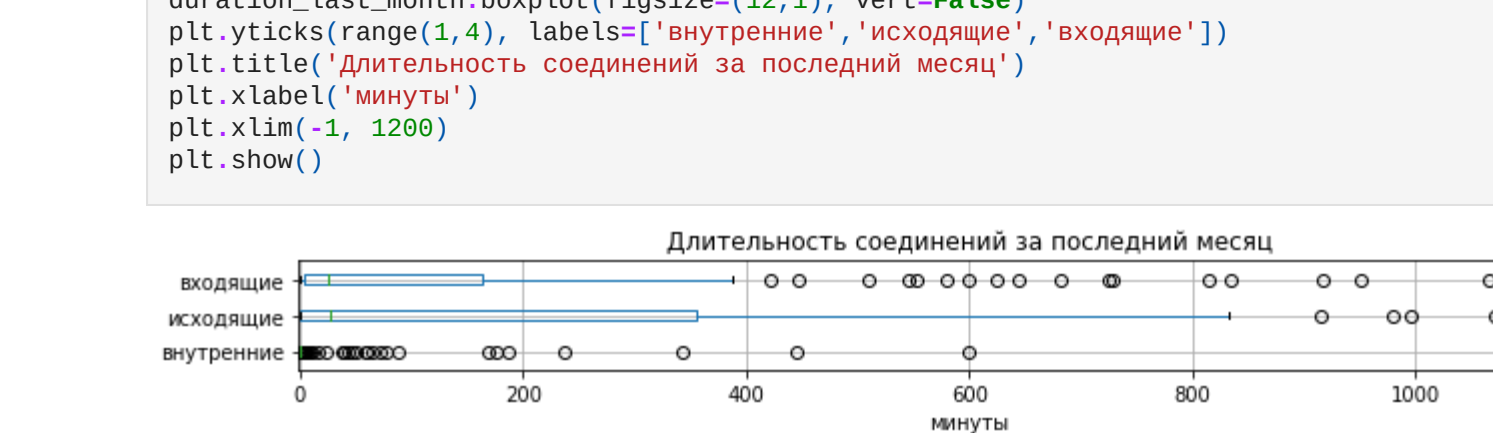
Оценим, как менялось число соединений усреднительного клиента за последние месяцы. Ввиду характера распределения в качестве меры среднего возьмем медиану, но для справки отобразим также и среднее арифметическое (на нём сильно сказывается появление клиентов с экстремально высокими значениями - исчезающее мало).



В первые два месяца среднее число соединений клиентов активно росло. В последние два месяца заметного изменения нет. Видно по мере роста числа клиентов показатель стабилизировался. На протяжении всего времени соотношение в длительности разговоров между разными видами соединений не изменилось: 50% длительности приходится на входящие соединения, 50% - на исходящие, доля внутренних соединений исчезающе мала.

### Суммарный доход

Оценим динамику суммарного дохода от клиентов:



Доходы росли по мере подключения к тарифам новых клиентов. В последний месяц рост почти прекратился.

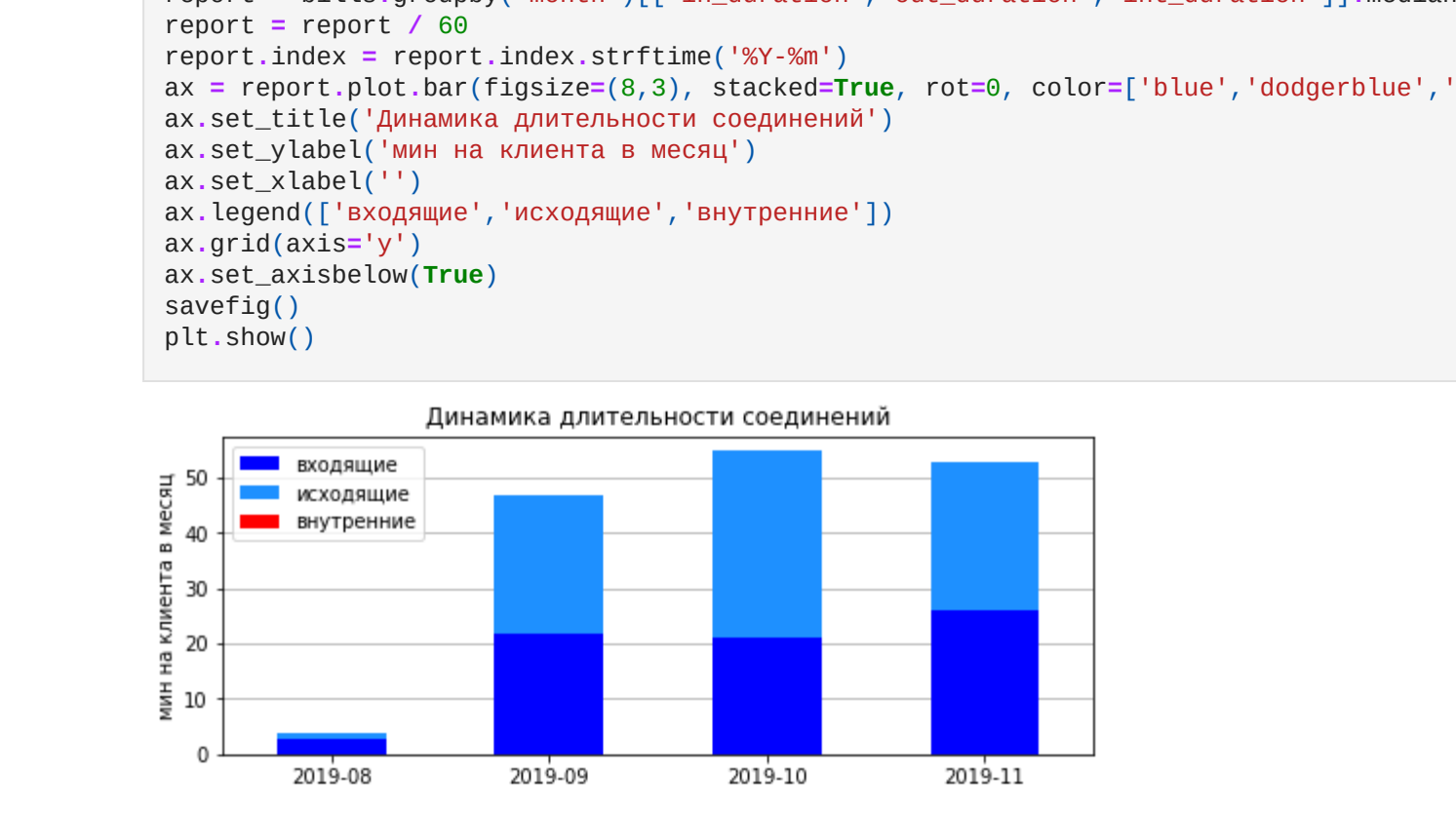
### Средняя доходность

Оценим, как распределены клиенты по сумме дохода, который они приносят. Для этого построим диаграмму размаха и гистограмму по данным за последний месяц:



Доходность от клиентов не имеет нормального распределения. Большинство клиентов приносит доход близкий размеру абонентской платы за тариф (1000, 2000 и 3000 руб. в мес.). В то же время существуют отдельные клиенты, которые приносят доход в размере от 10 до 30 тыс. руб. в месяц.

Оценим динамику средней доходности клиентов во времени. Ввиду характера распределения в качестве основной меры среднего возьмем медиану, но параллельно отобразим и среднее арифметическое (на нём сильно сказывается появление клиентов с экстремально высокими значениями - исчезающее мало).



Средняя доходность клиента практически не меняется с течением времени. Половина клиентов приносит доход от 1000 до 2100 руб. в месяц. Вторая половина клиентов имеет значительно больший разброс - от 2100 до 30 тыс. руб.

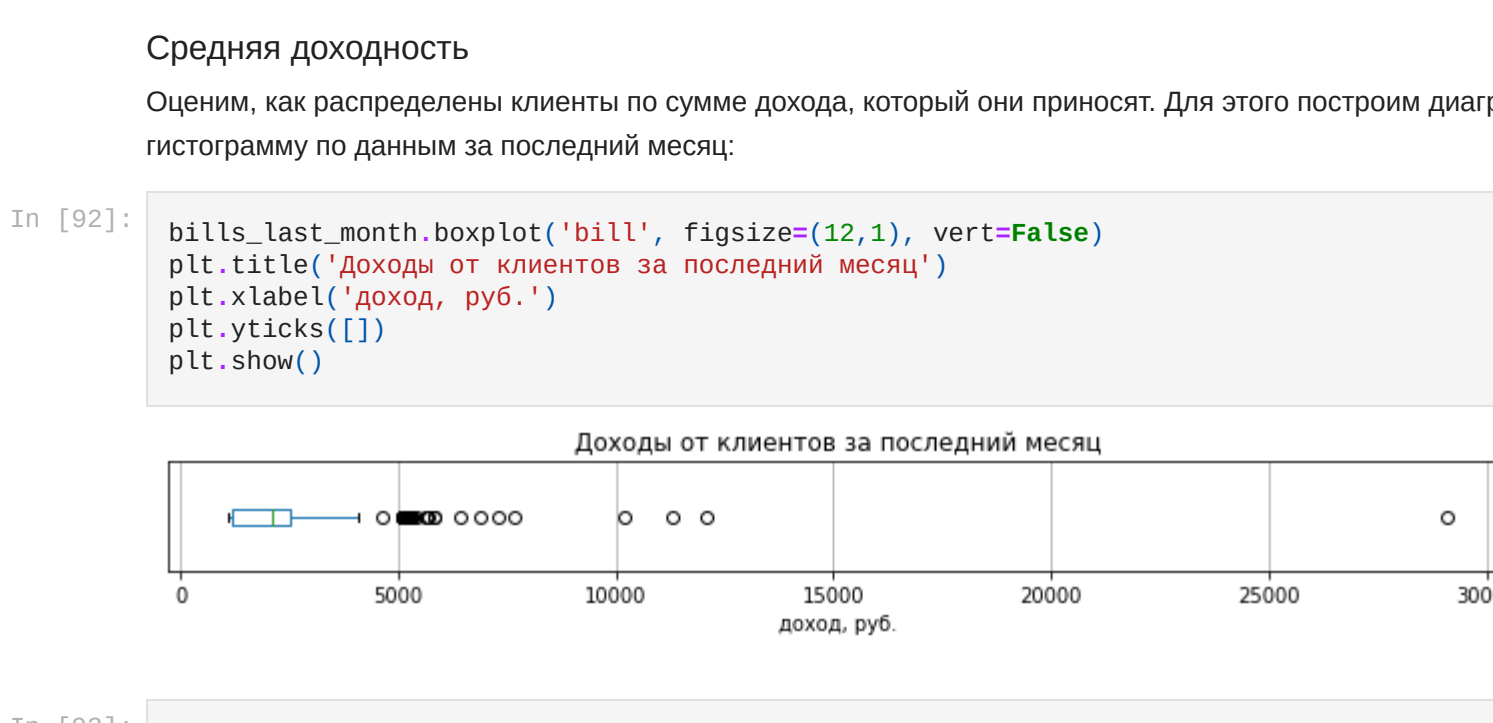
### Итого

По результатам исследования динамики и состава клиентов можно сделать следующие выводы:

- В первые месяцы после активное подключение к тарифам новых клиентов. В последнем месяце рост прекратился. По мере стабилизации числа клиентов стабилизируется и ряд их параметров: среднее число соединений и средняя длительность разговоров. Ряд параметров клиентов был стабилен независимо от роста клиентской базы: число операторов и средний доход (приносимый клиентами).
- В среднем клиенты имеют очень низкие значения параметров: 2 оператора на клиента, 35 звонков в месяц, длительность соединений - 50 минут в месяц. Средний доход от клиента равен 2100 руб. в мес. Большинство клиентов приносит доход, близкий к абонентской плате (1000, 2000 и 3000 руб. в мес.).
- Да половина клиентов, чьи характеристики выше среднего, имеют очень большой размах характеристик, на порядок превышающий средние значения: от 2 до 50 операторов, от 1 до 30 тыс. соединений в месяц, от 50 до 55 тыс. минут разговоров в месяц, от 2100 до 30 тыс. руб. дохода в месяц.
- Для клиентов в среднем характерно следующее соотношение соединений (как по количеству, так и по длительности): 50% - входящие соединения, 50% - исходящие соединения, число и длительность внутренних разговоров исчезающе мала.

### Тарифы

Узнаем, как менялась доходность и число пользователей тарифов во времени:



Из графиков видно, что доходы менялись по мере роста пользователей тарифов: в первые месяцы росли, в последний - перестали. Соотношение показателей разных тарифов на протяжении времени существенно не меняется: доходы от разных тарифов примерно равны, число пользователей тарифа А - наименьшее, и тарифа В - среднее, и тарифа С - наибольшее.

Узнаем, как менялись тарифицируемые параметры у пользователей разных тарифов во времени:



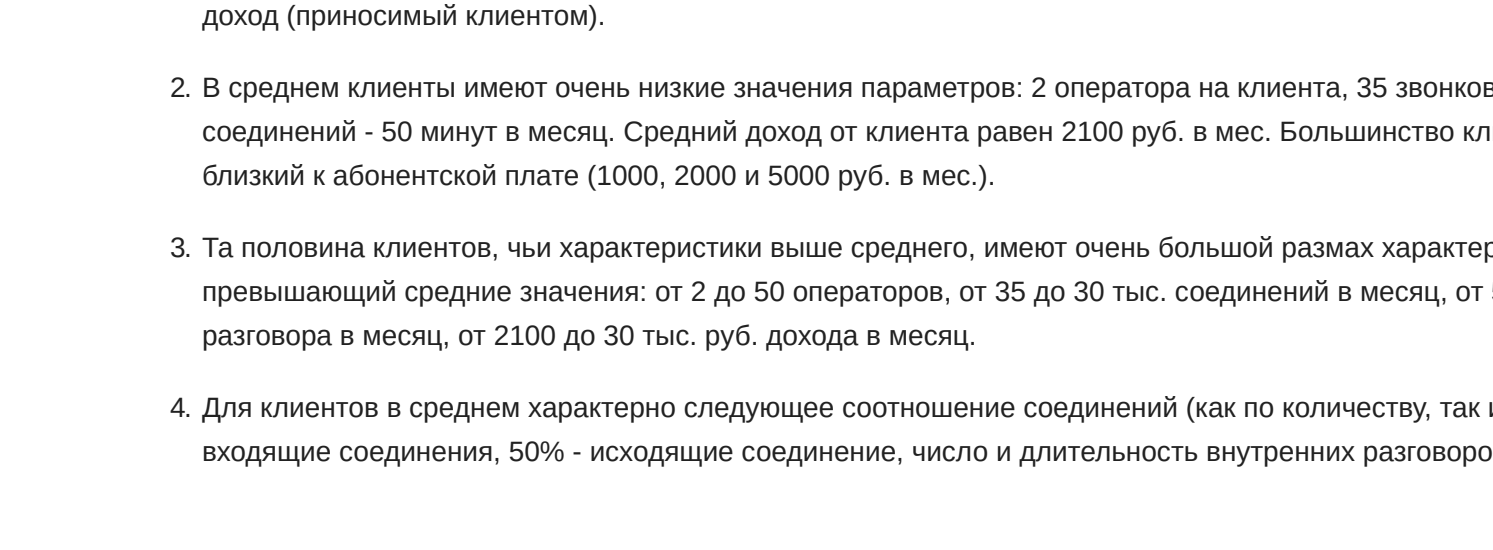
Среднее число операторов у пользователей разных тарифов и мало отличается и во времени практически не меняется. Средняя длительность соединений многократно выше у пользователей тарифа А и растёт, у пользователей тарифов В и С колеблется на низком уровне (у пользователей тарифа С - близка к нулю).

### Доходность и популярность

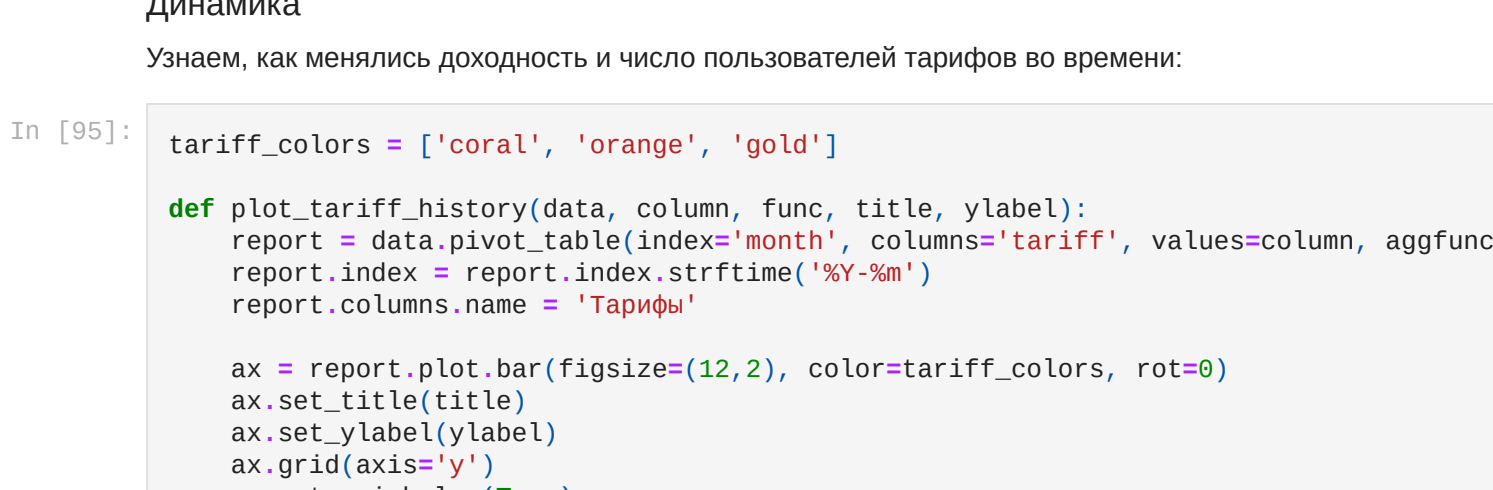
Узнаем как распределён доход и пользователи между тарифами по данным за последний месяц:



Расчитаем лучшие тарифы для клиентов по данным последнего месяца:



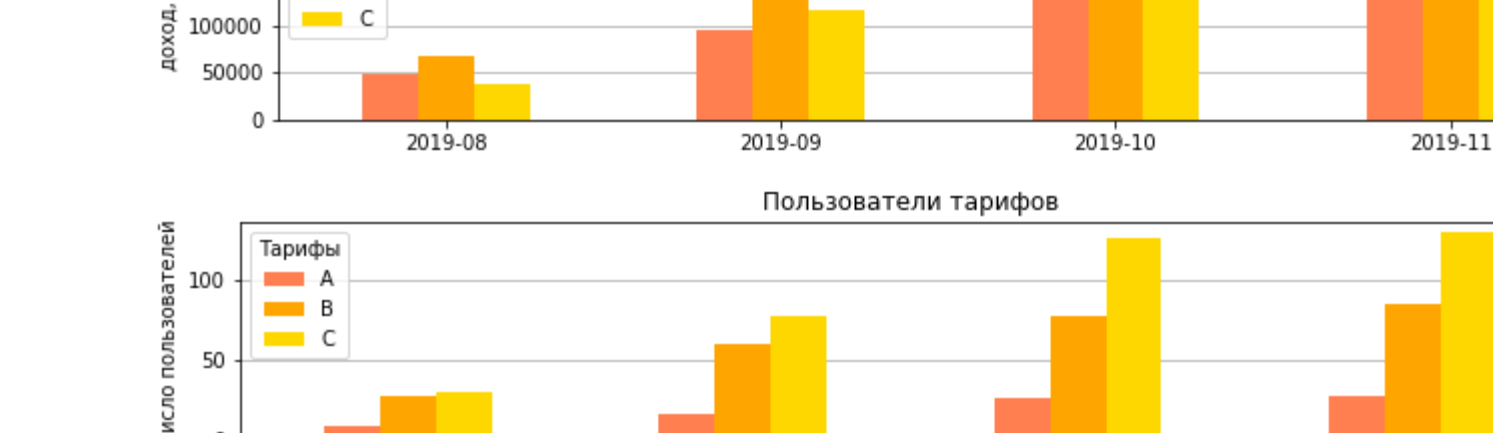
Выясним, сколько клиентов могут выбрать тариф, по которому будут платить меньше, чем по текущему:



Клиенты тарифов А и В либо переключаются своей потребностью в услугах связи, либо упрощают развитие, которое пока не состоялось. Так или иначе, но они переплативают, и алгоритм советует им - всем, кроме одного-двух - перейти на более дешёвые тарифы.

### Упущенный доход

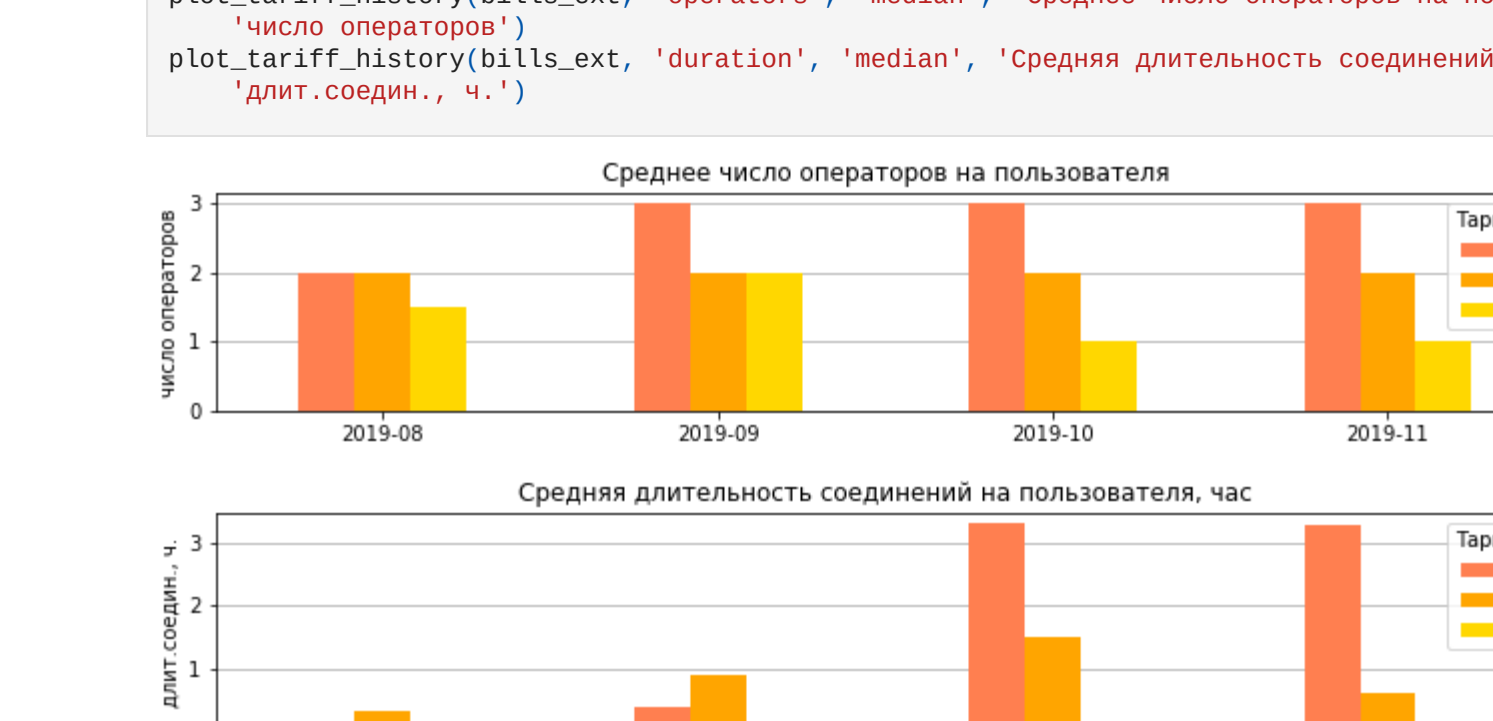
Расчитаем, насколько снижается доход, если все пользователи перейдут на оптимальные тарифы:



Доход сократится почти на треть, если все пользователи, которые могут, перейдут на оптимальные тарифы.

### Компенсация упущенного дохода

Постараемся компенсировать снижение доходов из-за возможного перехода всех клиентов на оптимальные тарифы. Запустим расчёт коэффициентов изменения. Значения коэффициентов покажу, в какой степени мы должны изменить стоимость тех или иных услуг в тарифных планах, чтобы компенсировать снижение дохода. Рассчитаем коэффициенты изменения абонентской платы (коэффициент fix), платы за оператора (oper) и стоимости минуты звонка (call):



fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff

fix oper call incomes\_diff



Компенсировать выплатающие доходы (с точностью до процента от текущего размера месячного дохода) поможет либо дуратурное увеличение на всех тарифах цен на подключение оператора и минут соединения одновременно, либо 70-процентное увеличение абонентской платы на всех тарифах, либо промежуточные варианты из таблицы выше.

## Снижение риска

Снизить риск быстрой потери доходов от перехода абонентов на оптимальные тарифы можно, если растянуть процесс и предлагать сменить тариф не всем абонентам, а только некоторым из них, например, только абонентам тарифов В и С, но только не абонентам тарифа А. Абоненты тарифа А (как было показано выше) приносят в средний доход многократно превышающий доход от абонентов тарифов В и С. Рассчитаем, насколько упадут доходы, если на более выгодные тарифы перейдут абоненты всех тарифов, кроме А:

```
In [104]: bills_a = bills_last_month[bills_last_month['tariff']=='A']
bills_b = bills_last_month[bills_last_month['tariff']=='B']
bills_c = count_bills_min(bills_b, tariffs)
x_min, x = bills_bc_min[bills_min]
sum() + bills_a[bills_min].sum(), bills_last_month[bills_min].sum()
round((x_min - x) / x_min, 2)
```

Out[104]: 0.15

Если тарифы сменяют только абоненты тарифов В и С, то доход упадёт не на 29%, а на 15%. Убыток сократится почти в два раза! При этом число перешедших на более выгодные тарифы клиентов сократится лишь на четверть: 27/(27+83)=0.25.

## Статистические гипотезы

### Переход на оптимальные тарифы снижает доход

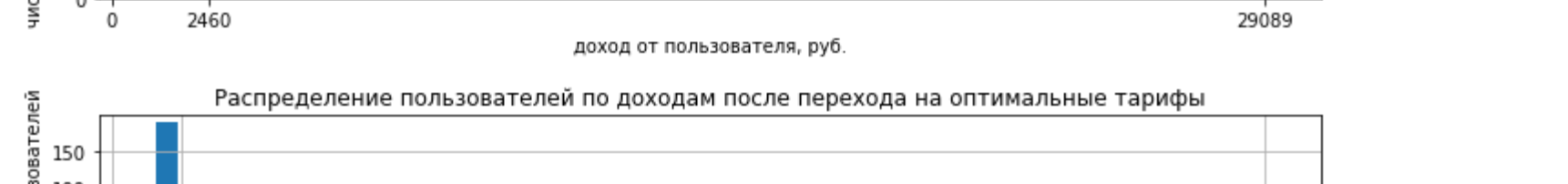
Проверим, верно ли, что переход пользователей на оптимальные тарифы снижает средний доход от пользователей? Как показано выше, суммарный доход снизился, но может быть это совпадение случайных факторов и на самом деле разница в средних доходах в расчёте на пользователя статистически незначима?

Сначала сохраним доход до смены тарифов и рассчитаем доход после смены тарифов:

```
In [105]: tariffs = pd.DataFrame(columns=['tariff', 'income'])
tariffs.loc[0, 'tariff', 'income'] = ['A', 5000, 100, 2000, 0.10, 0.40]
tariffs.loc[1, 'tariff', 'income'] = ['B', 2000, 100, 2000, 0.15, 0.50]
tariffs.loc[2, 'tariff', 'income'] = ['C', 1000, 100, 2000, 0.30, 0.70]
bills_last_month = bills[bills['month']=='bills_last_month'].max()
bills_min = count_bills_min(bills_last_month, tariffs)
```

Перед тем как проверить гипотезы изучим распределение сравниваемых величин и их средние значения:

```
In [106]: def show_hist(data, bins=50, title='', xlabel='', ylabel=''):
data.hist(figsize=(12,15), bins=bins)
plt.xticks(0, data.mean(), data.max())
plt.title(title)
plt.xlabel(xlabel)
plt.ylabel(ylabel)
plt.show()
show_hist(bills_last_month[bills_last_month['tariff']=='B'], bins=50,
title='Распределение пользователей по доходам до перехода на оптимальные тарифы',
xlabel='доход от пользователя, руб.',
ylabel='число пользователей')
show_hist(bills_min[bills_min['bill_min']], bins=50,
title='Распределение пользователей по доходам после перехода на оптимальные тарифы',
xlabel='доход от пользователя, руб.',
ylabel='число пользователей')
```



Так как распределение не является нормальным для проверки гипотез будем использовать критерий Манна-Уитни.

**Нулевая гипотеза H0:** средний доход на пользователя для выборки до и после перехода на оптимальные тарифы не различается.

**Альтернативная гипотеза H1:** средний доход на пользователя снизился после перехода на оптимальные тарифы.

Проверка:

```
In [107]: result = st.mannwhitneyu(bills_last_month[bills_min], bills_min[bills_min],
p_value = result[1])
```

Out[107]: 3.25943289976692e-10

Уровень статистической значимости значительно ниже 0.05 (и даже 0.01). Данные противоречат нулевой гипотезе, она должна быть отвергнута и вместо неё принята альтернативная гипотеза. Таким образом разница в выборках статистически значима и смена тарифов на оптимальные действительно снижает средний доход от пользователей тарифа.

### Пользователи тарифа В превосходят пользователей тарифа С

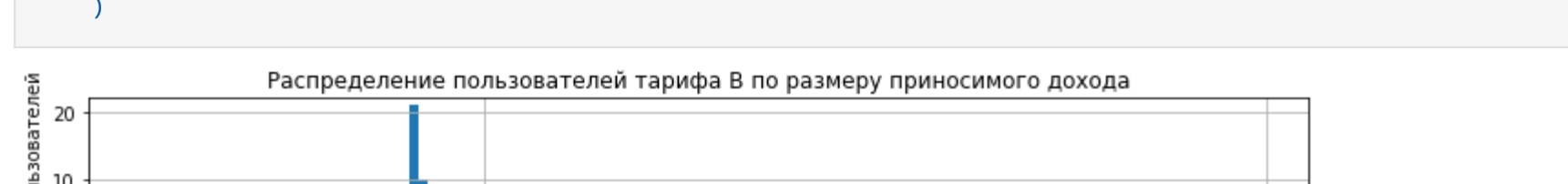
В разделах выше было показано, что пользователи тарифов В и С схожи в параметрах по сравнению с пользователями тарифа А. Пользователи тарифа А имеют в среднем в разы большие значения, чем у пользователей тарифов В и С. Последние же различаются между собой не более чем в полтора-два раза. Интересно узнать, значима ли такая разница или это случайное расхождение? Вопрос тем более актуален, что алгоритм выбора оптимальных тарифов посоветовал почти всем пользователям тарифа В перейти на тариф С.

#### Доход

Выясним, верно ли, что средний доход от пользователей тарифа В больше, чем средний доход от пользователей тарифа С.

Перед тем как проверить гипотезы, изучим распределение сравниваемых величин и их средние значения:

```
In [108]: show_hist(bills_last_month[bills_last_month['tariff']=='B']['bill'], bins=100,
title='Распределение пользователей тарифа В по размеру приносимого дохода',
xlabel='доход от пользователя в мес.',
ylabel='число пользователей')
show_hist(bills_last_month[bills_last_month['tariff']=='C']['bill'], bins=100,
title='Распределение пользователей тарифа С по размеру приносимого дохода',
xlabel='доход от пользователя в мес.',
ylabel='число пользователей')
```



Так как распределение не является нормальным для проверки гипотез будем использовать критерий Манна-Уитни.

**Нулевая гипотеза H0:** средний доход от пользователей тарифов В и С одинаков.

**Альтернативная гипотеза H1:** средний доход от пользователей тарифа В больше, чем от пользователей тарифа С.

Проверка:

```
In [109]: result = st.mannwhitneyu(
bills_last_month[bills_last_month['tariff']=='B']['bill'],
bills_last_month[bills_last_month['tariff']=='C']['bill'],
True, 'greater')
p_value = result[1]
```

Out[109]: 3.22904562882801e-27

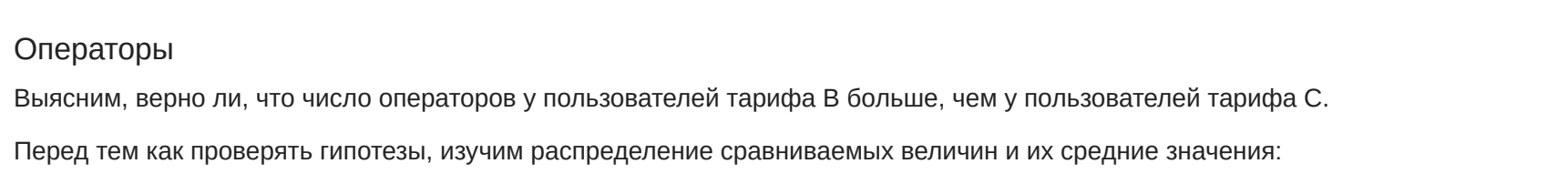
Уровень статистической значимости значительно ниже 0.05 (и даже 0.01). Данные противоречат нулевой гипотезе, она должна быть отвергнута и вместо неё принята альтернативная гипотеза. Таким образом разница в выборках статистически значима и средний доход от пользователей тарифов В больше среднего дохода от пользователей тарифа С.

#### Операторы

Выясним, верно ли, что число операторов у пользователей тарифа В больше, чем у пользователей тарифа С.

Перед тем как проверить гипотезы, изучим распределение сравниваемых величин и их средние значения:

```
In [110]: show_hist(bills_last_month[bills_last_month['tariff']=='B']['operators'], bins=16,
title='Распределение пользователей тарифа В по числу операторов',
xlabel='число операторов',
ylabel='число пользователей')
show_hist(bills_last_month[bills_last_month['tariff']=='C']['operators'], bins=13,
title='Распределение пользователей тарифа С по числу операторов',
xlabel='число операторов',
ylabel='число пользователей')
```



Так как распределение не является нормальным для проверки гипотез будем использовать критерий Манна-Уитни.

**Нулевая гипотеза H0:** среднее число операторов пользователей тарифов В и С одинаково.

**Альтернативная гипотеза H1:** среднее число операторов пользователей тарифов В больше, нежели у пользователей тарифа С.

Проверка:

```
In [111]: result = st.mannwhitneyu(
bills_last_month[bills_last_month['tariff']=='B']['operators'],
bills_last_month[bills_last_month['tariff']=='C']['operators'],
True, 'greater')
p_value = result[1]
```

Out[111]: 0.00921423913997827668

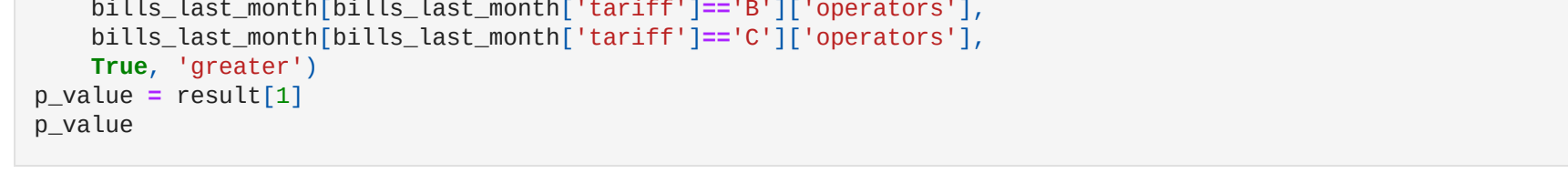
Уровень статистической значимости значительно ниже 0.05 (и даже 0.01). Данные противоречат нулевой гипотезе, она должна быть отвергнута и вместо неё принята альтернативная гипотеза. Таким образом разница в выборках статистически значима и среднее число операторов у пользователей тарифов В больше, чем у пользователей тарифа С.

#### Соединения

Выясним, совпадает средняя длительность разговоров у пользователей тарифов В и С.

Перед тем как проверить гипотезы, рассчитаем сравниваемые величины и изучим их распределение, а также средние значения:

```
In [112]: x = bills_last_month.copy()
x['tariff_minutes'] = (x['int_duration'] + x['out_duration']) / 60
show_hist(x[x['tariff']=='B']['tariff_minutes'], bins=100,
title='Распределение пользователей тарифа В по длительности соединений',
xlabel='длительность соединений, мин./мес.',
ylabel='число пользователей')
show_hist(x[x['tariff']=='C']['tariff_minutes'], bins=100,
title='Распределение пользователей тарифа С по длительности соединений',
xlabel='длительность соединений, мин./мес.',
ylabel='число пользователей')
```



Так как распределение не является нормальным для проверки гипотез будем использовать критерий Манна-Уитни.

**Нулевая гипотеза H0:** средняя длительность соединений для пользователей В и С совпадают.

**Альтернативная гипотеза H1:** средняя длительность соединений у пользователей тарифа В больше, чем у пользователей тарифа С.

Проверка:

```
In [113]: result = st.mannwhitneyu(
x[x['tariff']=='B']['tariff_minutes'],
x[x['tariff']=='C']['tariff_minutes'],
True, 'greater')
p_value = result[1]
```

Out[113]: 0.00426322457487433

Уровень статистической значимости значительно ниже 0.05 (и даже 0.01). Данные противоречат нулевой гипотезе, она должна быть отвергнута и вместо неё должна быть принята альтернативная гипотеза. Таким образом разница в выборках статистически значима и средняя длительность соединений у пользователей тарифа В больше, чем у пользователей тарифа С.

#### Итог

Проверка подтвердила статистическую значимость различий между пользователями тарифов В и С: параметры пользователей тарифа В - доход, число операторов, длительность соединений - все превосходят параметры пользователей тарифа С.

## Итоги

По итогам исследований следует сделать следующие выводы.

- Бурный рост подключений к тарифам А, В, С завершился.** Количественные показатели (число клиентов и суммарный доход), а также качественные показатели (средний доход, звонки и динамика тарифов) в течение четырех месяцев выросли и стабилизировались.
- Клиенты очень неоднородны.** Показатели потребления услуг и генерируемого дохода у разных клиентов могут различаться на многие порядки. Доход от клиента колеблется в диапазоне от 1 000 до 30 000 руб. в месяц. Количество соединений - от 1 до 30 000 соединений в месяц. Длительность разговоров - от 1 до 10 000 минут в месяц.
- Средний клиент потребляет мало услуг и приносит малый доход.** Средний доход, генерируемый половиной клиентов, не превышает 2 100 руб. в мес., среднее число соединений половины клиентов - не более 33 соединений в мес., средняя длительность соединений - не более 50 мин. в мес. Заменяет число клиентов даёт доход, близкий или почти совпадающий с размером абонентской платы по тарифу (1000, 2000 и 5000 руб. в мес.).
- Тарифы приносят одинаковый доход, но пользователи у них разные.** Каждый тариф приносит примерно треть (200 тыс. руб.) от суммарной выручки за месяц (600 тыс. руб.). При этом число пользователей существенно различается: тариф А - 12% пользователей, тариф В - 35%, тариф С - 53%. Соответственно существенно различается доход, приносимый пользователями, и объём потребляемых услуг: пользователи тарифа А приносят доход (и потребляют услуг) в несколько раз больше, чем пользователи тарифов В и С. Пользователи тарифа В опережают пользователей тарифа С в полтора-два раза. Статистическая проверка показывает, что эти различия не являются следствием случайных влияний текущих факторов.
- Пользователи тарифов переплываются.** Алгоритм выбора наиболее выгодного тарифа показывает, что подавляющему большинству пользователей тарифов А и В (более 90%) следует сменить свой тариф на тариф С. Это значит, что пользователи перенесли свои потребности в услугах связи и требуют отдельного изучения.
- Переход на оптимальные тарифы приведёт к снижению дохода.** Переход пользователей на более выгодные тарифы сократит суммарный доход на 29%. Статистическая проверка показывает, что это падение не может быть объяснено текущими случайными значениями различных факторов.

Рекомендации:

- Привлечение новых клиентов и стимулирование имеющихся.** Следует привлекать крупных клиентов, схожих по профилю с текущими пользователями тарифа А. Они дают в несколько раз больший доход, их потребности в связи растут. Имеющихся клиентов следует стимулировать. Способы стимулирования выйдут за рамки данного проекта и требуют отдельного изучения.
- Риск перехода пользователей на более дешёвые тарифы может быть снижен,** если предлагать сменить тариф не всем пользователям, которые могли бы это сделать, а только пользователям тарифа В. Если на более выгодные тарифы перейдут только они, то при максимальном числе довольных возможностью сэкономить пользователей общий доход упадёт не на 29%, а лишь на 15%.
- Компенсировать снижение дохода можно изменением тарифов.** Калькулятор изменения тарифов позволяет рассчитать разные комбинации изменения параметров текущих тарифов для компенсации упущенного дохода. Например, расчёты показывают, что увеличение абонентской платы (на всех тарифах) на 70% полностью компенсирует переход всех пользователей на более экономные тарифы. Другой вариант: абонентская плата остаётся неизменной, но плату за оператора и минуту разговоров придётся увеличить на 100%. Возможны и другие варианты.