

## Шаг 1. Откройте файл с данными и изучите общую информацию

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

data_calls = pd.read_csv("../datasets/calls.csv")
data_internet = pd.read_csv("../datasets/internet.csv")
data_messages = pd.read_csv("../datasets/messages.csv")
data_tariffs = pd.read_csv("../datasets/tariffs.csv")
data_users = pd.read_csv("../datasets/users.csv")

data_calls.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 202607 entries, 0 to 202606
Data columns (total 4 columns):
 id                202607 non-null object
 call_date         202607 non-null object
 call_duration     202607 non-null float64
 user_id           202607 non-null int64
 dtypes: float64(1), int64(1), object(2)
memory usage: 6.2+ MB

In [3]: data_internet.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 149396 entries, 0 to 149395
Data columns (total 8 columns):
 id                149396 non-null int64
 name              149396 non-null object
 mb_used           149396 non-null float64
 session_date      149396 non-null object
 user_id           149396 non-null int64
 reg_date          149396 non-null int64
 tariff            149396 non-null object
 dtypes: int64(1), int64(2), object(2)
memory usage: 5.7+ MB

In [4]: data_messages.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 123036 entries, 0 to 123035
Data columns (total 3 columns):
 id                123036 non-null object
 message_date      123036 non-null object
 user_id           123036 non-null int64
 dtypes: int64(1), object(2)
memory usage: 2.8+ MB

In [5]: data_tariffs.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2 entries, 0 to 1
Data columns (total 8 columns):
 messages_included 2 non-null int64
 mb_per_month_included 2 non-null int64
 minutes_included  2 non-null int64
 rub_monthly_fee   2 non-null int64
 rub_per_gb        2 non-null int64
 rub_per_message   2 non-null int64
 rub_per_minute    2 non-null int64
 tariff            2 non-null object
 dtypes: int64(7), object(1)
memory usage: 256.0 bytes

In [6]: data_users.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 8 columns):
 user_id           500 non-null int64
 age              500 non-null float64
 churn_date        38 non-null object
 city              500 non-null object
 first_name        500 non-null object
 last_name         500 non-null object
 reg_date          500 non-null object
 tariff            500 non-null object
 dtypes: int64(2), object(6)
memory usage: 51.4+ KB
```

## Вывод

Данные распределены по пяти датасетам: `data_calls`, `data_internet`, `data_messages`, `data_tariffs`, `data_users`. Названия столбцов не требуют изменений.

Датасет `data_calls` состоит из 4 столбцов и 202607 строк. Столбец `call_date` необходимо привести к формату даты. Датасет `data_internet` состоит из 5 столбцов и 149396 строк. Столбец `session_date` необходимо привести к формату даты.

Датасет `data_messages` состоит из 3 столбцов и 123036 строк. Столбец `message_date` необходимо привести к формату даты. Датасет `data_users` состоит из 8 столбцов и 500 строк. Столбцы `churn_date` и `reg_date` необходимо привести к формату даты. Столбец `churn_date` содержит 462 нулевых значения. Это означает, что пользователи ещё пользуются тарифом.

Датасеты `data_calls`, `data_internet`, `data_messages` и `data_users` связаны по `user_id`. Датасеты `data_tariffs` состоит из 8 столбцов и 2 строк, содержит информацию о тарифах.

## Шаг 2. Подготовьте данные

Приведем строки в столбце `call_date` таблицы `data_calls` к формату даты и добавим столбец месяца `call_month`.

```
In [7]: data_calls['call_date'] = pd.to_datetime(data_calls['call_date'], format='%Y-%m-%d')
data_calls['call_date'] = data_calls['call_date'].dt.month
```

Округлим время значения минут за исключением тех, чья длительность — 0.0. Для этого выберем из таблицы все значения кроме нулевых, отбросим дробную часть минутами `np.floor` и прибавим единицу. Сохраним результат в столбце `duration_minute`.

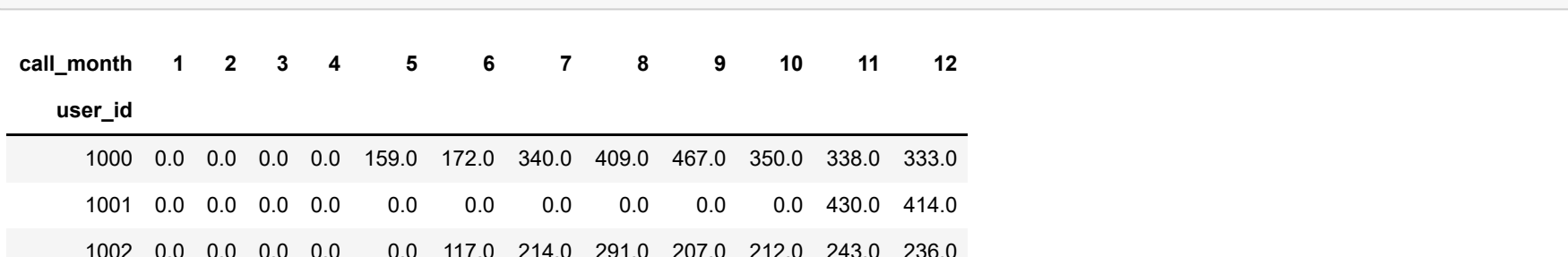
```
In [8]: data_calls['duration_minute'] = data_calls['duration']
data_calls.loc[data_calls['duration_minute'] != 0, 'duration_minute'] = data_calls.loc[data_calls['duration_minute'] != 0, 'duration_minute'].apply(lambda x: math.floor(x) + 1)
data_calls.head(10)
```

Out[8]:

	id	call_date	duration	user_id	call_month	duration_minute
0	1000_0	2018-07-25	0.00	1000	7	0.0
1	1000_1	2018-06-17	0.00	1000	8	0.0
2	1000_2	2018-06-11	2.85	1000	6	3.0
3	1000_3	2018-09-21	13.80	1000	9	14.0
4	1000_4	2018-12-15	5.18	1000	12	6.0
5	1000_5	2018-11-02	0.00	1000	11	0.0
6	1000_6	2018-10-18	0.00	1000	10	0.0
7	1000_7	2018-06-22	18.31	1000	8	18.0
8	1000_8	2018-09-15	18.44	1000	9	18.0
9	1000_9	2018-08-15	0.00	1000	8	0.0

Построим гистограмму округленной длительности звонка в минутах.

```
In [9]: data_calls['duration_minute'].hist(bins=len(data_calls['duration_minute'].value_counts()), figsize=(8, 5))
plt.title('Округленная длительность звонка в минутах');
```



Поскольку нулевые значения не влияют на выручку пользователей, оставим их в датасете. Максимальная длительность звонка не превышает 40 минут, что не является аномалией.

Посчитаем количество сделанных звонков по месяцам. Для этого создадим сводную таблицу `data_calls_pivot_count`, сгруппировав количество звонков по пользователям и месяцам и посчитав суммарное количество звонков.

```
In [10]: data_calls_pivot_count = data_calls.pivot_table(index='user_id', columns='call_month', \
values='id', aggfunc='count').fillna(0)
data_calls_pivot_count.head()
```

Поискать количество отправленных сообщений по месяцам. Для этого создадим сводную таблицу `data_messages_pivot_sum`, сгруппировав количество сообщений по пользователям и месяцам и посчитав количество сообщений.

```
data_messages_pivot = data_messages.pivot_table(index='user_id', columns='message_month', \
                                                values='count', aggfunc='count').fillna(0)
```

```
data_messages_pivot.head()
```

message_month	1	2	3	4	5	6	7	8	9	10	11	12	
user_id	1000	0.0	0.0	0.0	0.0	22.0	60.0	75.0	81.0	57.0	73.0	58.0	70.0

Посчитаем количество израсходованных минут разговора по месяцам. Для этого создадим сводную таблицу `data_calls_pivot`, сгруппировав количество звонков по пользователям и месяцам и посчитав суммарное количество минут разговора.

```
In [11]: data_calls_pivot = data_calls.pivot_table(index='user_id', columns='call_month', \
values='duration_minute', aggfunc='sum').fillna(0)
data_calls_pivot.head()
```

```
data_internet['session_date'] = pd.to_datetime(data_internet['session_date'], format='%Y-%m-%d')
data_internet['session_month'] = data_internet['session_date'].dt.month
```

Очистив нулев значения мегабит за исключением тех, чей айпи — 0.0. Для этого выберем из таблицы все значения кроме нулевых, отбросим дублирую часть из столбца np.floor и прибавим единицу. Сохраним результат в столбце mb\_used\_int.

```
data_internet['mb_used_int'] = data_internet['mb_used_int'] + 1
data_internet.loc[data_internet['mb_used_int'] != 0, 'mb_used_int'] = data_internet.loc[data_internet['mb_used_int'] != 0, 'mb_used_int'].np.floor + 1
```

Приведем строки в столбце `message_date` таблицы `data_messages` к формату даты и добавим столбец месяца `message_month`.

```
In [12]: data_messages['message_date'] = pd.to_datetime(data_messages['message_date'], format='%Y-%m-%d')
data_messages['message_date'] = data_messages['message_date'].dt.month
```

Посчитаем количество отправленных сообщений по месяцам. Для этого создадим сводную таблицу `data_messages_pivot_sum`, сгруппировав количество сообщений по пользователям и месяцам и посчитав суммарное количество сообщений.

```
In [13]: data_messages_pivot_sum = data_messages.pivot_table(index='user_id', columns='message_month', \
values='id', aggfunc='count').fillna(0)
data_messages_pivot_sum.head()
```

149392	149392	1499_15	490.13	2018-12-14	1499	12	491.0
149393	149393	1499_154	0.00	2018-10-27	1499	10	0.0
149394	149394	1499_155	1246.32	2018-11-26	1499	11	1247.0
149395	149395	1499_156	544.37	2018-10-26	1499	10	545.0

149396 rows x 7 columns

Построим гистограмму окружённого объёма интернет-графа (в мегабайтах).

Приведем строки в столбце `session_date` таблицы `data_internet` к формату даты и добавим столбец месяца `session_month`.

```
In [14]: data_internet['session_date'] = pd.to_datetime(data_internet['session_date'], format='%Y-%m-%d')
data_internet['session_date'] = data_internet['session_date'].dt.month
```

Округлим время значения минут за исключением тех, чей объем — 0.0. Для этого выберем из таблицы все значения кроме нулевых, отбросим дробную часть минутами `np.floor` и прибавим единицу. Сохраним результат в столбце `mb_used_int`.

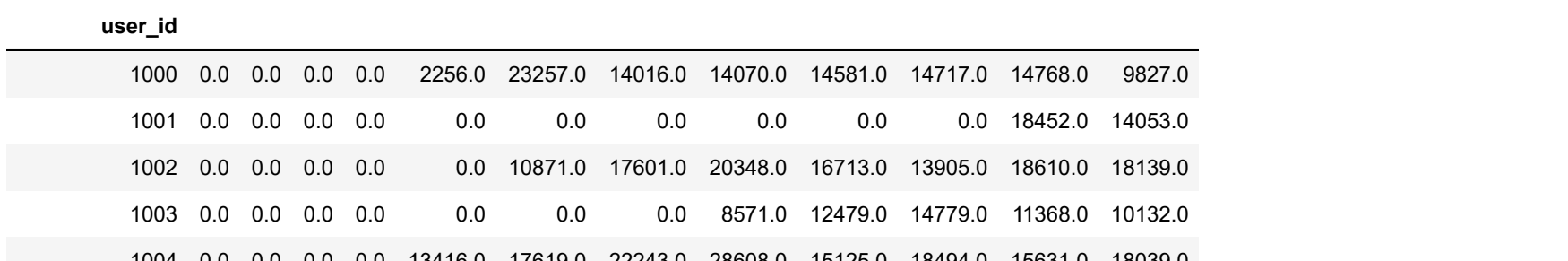
```
In [15]: data_internet['mb_used_int'] = data_internet['mb_used']
data_internet.loc[data_internet['mb_used_int'] != 0, 'mb_used_int'] = data_internet.loc[data_internet['mb_used_int'] != 0, 'mb_used_int'].apply(lambda x: math.floor(x) + 1)
data_internet
```

Out[15]:

Unnamed: 0	id	mb_used	session_date	user_id	session_month	mb_used_int	
0	0	1000_0	112.95	2018-11-25	1000	11	113.0
1	1	1000_1	1062.81	2018-09-07	1000	9	1063.0
2	2	1000_2	1197.26	2018-06-25	1000	6	1198.0
3	3	1000_3	550.27	2018-08-22	1000	8	551.0
4	4	1000_4	302.56	2018-09-24	1000	9	303.0
...	...	...	...	...	...	...	...
149391	149391	1499_152	318.90	2018-10-03	1499	10	319.0
149392	149392	1499_153	490.13	2018-12-14	1499	12	491.0
149393	149393	1499_154	0.00	2018-10-27	1499	10	0.0
149394	149394	1499_155	1248.32	2018-11-26	1499	11	1247.0
149395	149395	1499_156	544.37	2018-10-26	1499	10	545.0

Построим гистограмму округленного объема интернет-трафика (в мегабайтах).

```
In [16]: data_internet['mb_used_int'].hist(bins=50, figsize=(8, 5))
plt.title('Округленный объем интернет-трафика (в мегабайтах)');
```



Поскольку нулевые значения не влияют на выручку пользователей, оставим их в датасете. Максимальный объем потраченного за сессию интернет-трафика не превышает двух гигабайт, что не является уникальным.

Посчитаем объем израсходованного интернет-трафика по месяцам. Для этого создадим сводную таблицу `data_internet_pivot_sum`, сгруппировав объем интернет-трафика по пользователям по месяцам и посчитав суммарный объем интернет-трафика.

```
In [17]: data_internet_pivot_sum = data_internet.pivot_table(index='user_id', columns='session_month', \
values='mb_used_int', aggfunc='sum').fillna(0)
data_internet_pivot_sum.head()
```

```
data_users['churn_date_month'] = data_users['churn_date_month'].fillna(13)
data_users
```

user_id	age	churn_date	city	first_name	last_name	reg_date_month	tariff	reg_date_month	churn_date_month
0	1000	51	NaT	Краснодар	Воробей	2018-05-25	ultra	5	13.0
1	1001	41	NaT	Москва	Влад	2018-01-01	smart	11	13.0
2	1002	59	NaT	Свердловск	Евгений	2018-11-07	ultra	8	13.0
3	1003	23	NaT	Москва	Генна	2018-08-17	ultra	8	13.0
4	1004	68	NaT	Новосибирск	Татьяна	2018-05-14	ultra	5	13.0

Приведем строки в столбцах `reg_date` и `churn_date` таблицы `data_users` к формату даты и добавим столбцы месяца `reg_date_month` и `churn_date_month`. Заменим значения NaN в столбце `churn_date_month` на 13 (значит тариф еще действует).

```
In [18]: data_users['reg_date'] = pd.to_datetime(data_users['reg_date'], format='%Y-%m-%d')
data_users['reg_date_month'] = data_users['reg_date'].dt.month
data_users['churn_date'] = pd.to_datetime(data_users['churn_date'], format='%Y-%m-%d')
data_users['churn_date_month'] = data_users['churn_date'].dt.month
data_users['churn_date_month'] = data_users['churn_date_month'].fillna(13)
data_users
```

Out[18]:

user_id	age	churn_date	city	first_name	last_name	reg_date	tariff	reg_date_month	churn_date_month
0	1000	52	Краснодар	Радван	Верещанин	2018-05-25	ultra	5	13.0
1	1001	41	Москва	Иван	Евков	2018-11-01	smart	11	13.0
2	1002	59	НаТ	Стерлитамак	Евгений	2018-06-17	smart	6	13.0
3	1003	23	НаТ	Москва	Белла	2018-08-17	ultra	8	13.0
4	1004	68	НаТ	Новокузнецк	Татьяна	2018-05-14	ultra	5	13.0
...	...	...	...	...	...	...	...	...	...
495	1495	65	НаТ	Иркутск	Алексей	2018-08-28	ultra	8	13.0
496	1496	36	НаТ	Вологда	Трифон	2018-01-27	smart	1	13.0
497	1497	32	НаТ	Челябинск	Каролина	2018-10-09	smart	10	13.0
498	1498	68	2018-10-25	Владикавказ	Василий	2018-07-19	smart	7	10.0
499	1499	35	НаТ	Пермь	Гектор	2018-09-27	smart	9	13.0

500 rows × 10 columns

Создадим сводную таблицу `data_users_pivot`, где единицами будут отмечены месяцы подключения тарифа.

```
In [19]: data_users_pivot = data_users.pivot_table(index='user_id', columns='reg_date_month', \
values='tariff', aggfunc='count').fillna(0)
data_users_pivot.head()
```

Out[19]:

reg_date_month	1	2	3	4	5	6	7	8	9	10	11	12
user_id	1000	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
1001	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0
1002	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1003	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
1004	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Создадим сводную таблицу `data_users_zero`, заполненную нулями, она пригодится нам в дальнейшем.

```
In [20]: data_users_zero = data_users_pivot.replace(1, 0)
```

В цикле по каждому пользователю сравним каждый месяц с месяцем подключения тарифа и месяцем прекращения пользования тарифом. Если он внутри диапазона, меняем значение в сводной таблице на единицу.

```
In [21]: for i in range(data_users.shape[0]):
for j in range(1,13):
if (j >= data_users.loc[i, 'reg_date_month']) and (j <= data_users.loc[i, 'churn_date_month']):
data_users_pivot.loc[1000+i, j] = 1
data_users_pivot.tail()
```

Out[21]:

reg_date_month	1	2	3	4	5	6	7	8	9	10	11	12
user_id	1495	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0	1.0	1.0
1496	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
1497	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0
1498	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0	0.0	0.0
1499	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0	1.0	1.0

Чтобы найти помесичную выручку с каждого пользователя, преобразуем найденные сводные таблицы и просуммируем. Однако можно заметить, что они не совпадают по размеру. Некоторые пользователи не совершают звонки, другие не отправляют сообщения, третьи не пользуются интернетом. Приведем их к единому образцу, сложив с нулевой таблицей `data_users_zero` и обнулив значения NaN.

```
In [22]: data_calls_pivot = (data_users_zero + data_calls_pivot).fillna(0)
data_messages_pivot = (data_users_zero + data_messages_pivot).fillna(0)
data_internet_pivot = (data_users_zero + data_internet_pivot).fillna(0)
```

Посчитаем помесичную выручку с каждого пользователя, как если бы он пользовался обоими тарифами. Сначала найдем стоимость интернет-трафика сверх тарифного пакета по тарифу «Смарт». Для этого вычтем из таблиц `data_internet_pivot` бесплатный лимит в 15 Гб, обнулим отрицательные значения методом `mask()`, переведем в гигабайты, поделив на 1024, выберем все значения кроме нулевых, отбросим дробную часть методом `np.floor`, прибавим единицу и умножим на стоимости 1 Гб интернет-трафика сверх тарифа. Обнулим значения NaN и сохраним результат в таблице `pivot_data_smart`.

```
In [23]: pivot_data_smart = data_internet_pivot - 15360).mask((data_internet_pivot - 15360) < 0, 0) / 1024
pivot_data_smart = (pivot_data_smart[pivot_data_smart != 0]).apply(np.floor) + 1) * 200
pivot_data_smart = pivot_data_smart.fillna(0)
pivot_data_smart.head()
```

Out[23]:

reg_date_month	1	2	3	4	5	6	7	8	9	10	11	12
----------------	---	---	---	---	---	---	---	---	---	----	----	----







```
In [52]: data_messages_pivot_ultra_var = np.var(data_messages_pivot_ultra.loc[:, range(1, 13)])\
data_messages_pivot_ultra_var
Out [52]: 1 655.359184
2 668.808394
3 1753.600000
4 1372.147272
5 1529.620264
6 1906.782357
7 2160.720961
8 1999.978177
9 1886.492078
10 2403.047224
11 2229.390555
12 2342.475141
dtype: float64
```

Найдём среднее количество сообщений по месяцам для пользователей тарифа «Смарт», сохраним в переменную data\_messages\_pivot\_smart\_mean и выведем на экран.

```
In [53]: data_messages_pivot_smart_mean = data_messages_pivot_smart.loc[:, range(1, 13)]\
data_messages_pivot_smart_mean
Out [53]: 1 20.000000
2 27.940000
3 35.550725
4 35.864583
5 39.032787
6 39.577483
7 38.988506
8 38.750000
9 39.219731
10 39.082677
11 39.881481
12 42.264804
dtype: float64
```

Аналогично для тарифа «Ультра».

```
In [54]: data_messages_pivot_ultra_mean = data_messages_pivot_ultra.loc[:, range(1, 13)]\
data_messages_pivot_ultra_mean
Out [54]: 1 43.428571
2 37.937500
3 42.800000
4 44.515821
5 48.769231
6 52.313433
7 62.136986
8 58.897273
9 63.294737
10 69.776699
11 66.159292
12 72.008547
dtype: float64
```

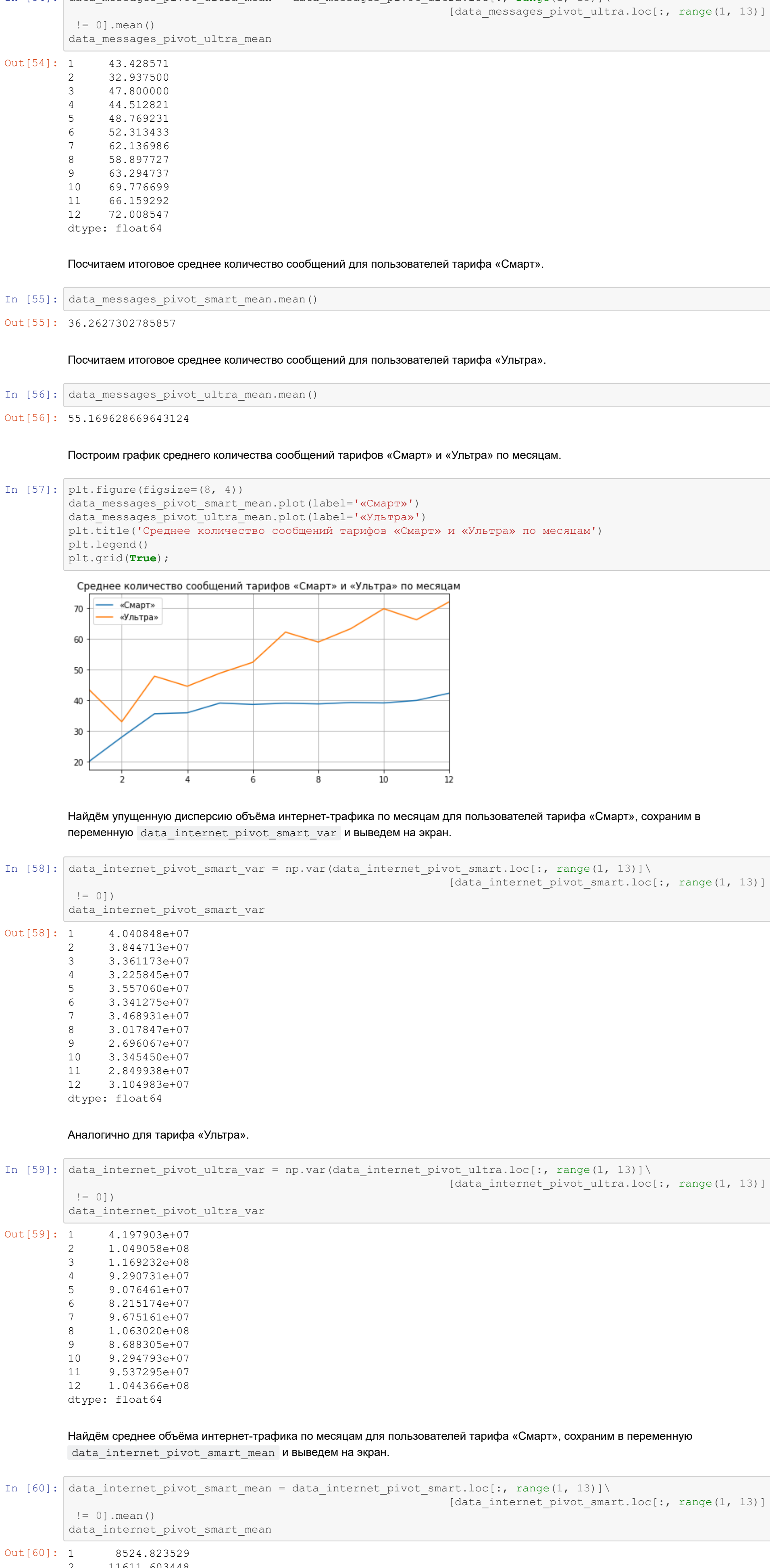
Посчитаем итоговое среднее количество сообщений для пользователей тарифа «Смарт».

```
In [55]: data_messages_pivot_smart_mean.mean()
Out [55]: 36.2627302785857
```

Посчитаем итоговое среднее количество сообщений для пользователей тарифа «Ультра».

```
In [56]: data_messages_pivot_ultra_mean.mean()
Out [56]: 55.169628669643124
```

Построим график среднего количества сообщений тарифов «Смарт» и «Ультра» по месяцам.



```
In [58]: data_internet_pivot_smart_var = np.var(data_internet_pivot_smart.loc[:, range(1, 13)]\
data_internet_pivot_smart_var
Out [58]: 1 4.040848e+07
2 3.844713e+07
3 3.681173e+07
4 3.229845e+07
5 3.557060e+07
6 3.341275e+07
7 3.468931e+07
8 3.017847e+07
9 2.696067e+07
10 3.345450e+07
11 2.849938e+07
12 3.104983e+07
dtype: float64
```

Аналогично для тарифа «Ультра».

```
In [59]: data_internet_pivot_ultra_var = np.var(data_internet_pivot_ultra.loc[:, range(1, 13)]\
data_internet_pivot_ultra_var
Out [59]: 1 4.197963e+07
2 1.849058e+08
3 1.169232e+08
4 9.290731e+07
5 9.076461e+07
6 8.215174e+07
7 9.679161e+07
8 1.050320e+08
9 8.688305e+07
10 9.294793e+07
11 9.537295e+07
12 1.044366e+08
dtype: float64
```

Найдём среднее объёма интернет-трафика по месяцам для пользователей тарифа «Смарт», сохраним в переменную data\_internet\_pivot\_smart\_mean и выведем на экран.

```
In [60]: data_internet_pivot_smart_mean = data_internet_pivot_smart.loc[:, range(1, 13)]\
data_internet_pivot_smart_mean
Out [60]: 1 8524.823529
2 11611.603448
3 15124.090909
4 13479.857143
5 15825.617021
6 15815.164634
7 15763.904948
8 16703.726496
9 16325.050000
10 16836.286689
11 16917.987342
12 18137.296736
dtype: float64
```

Аналогично для тарифа «Ультра».

```
In [61]: data_internet_pivot_ultra_mean = data_internet_pivot_ultra.loc[:, range(1, 13)]\
data_internet_pivot_ultra_mean
Out [61]: 1 13153.555556
2 12858.954545
3 17552.405405
4 14844.235364
5 19177.969231
6 19170.807229
7 20729.911111
8 19779.224512
9 19391.280729
10 20228.672131
11 19954.586466
12 21972.072464
dtype: float64
```

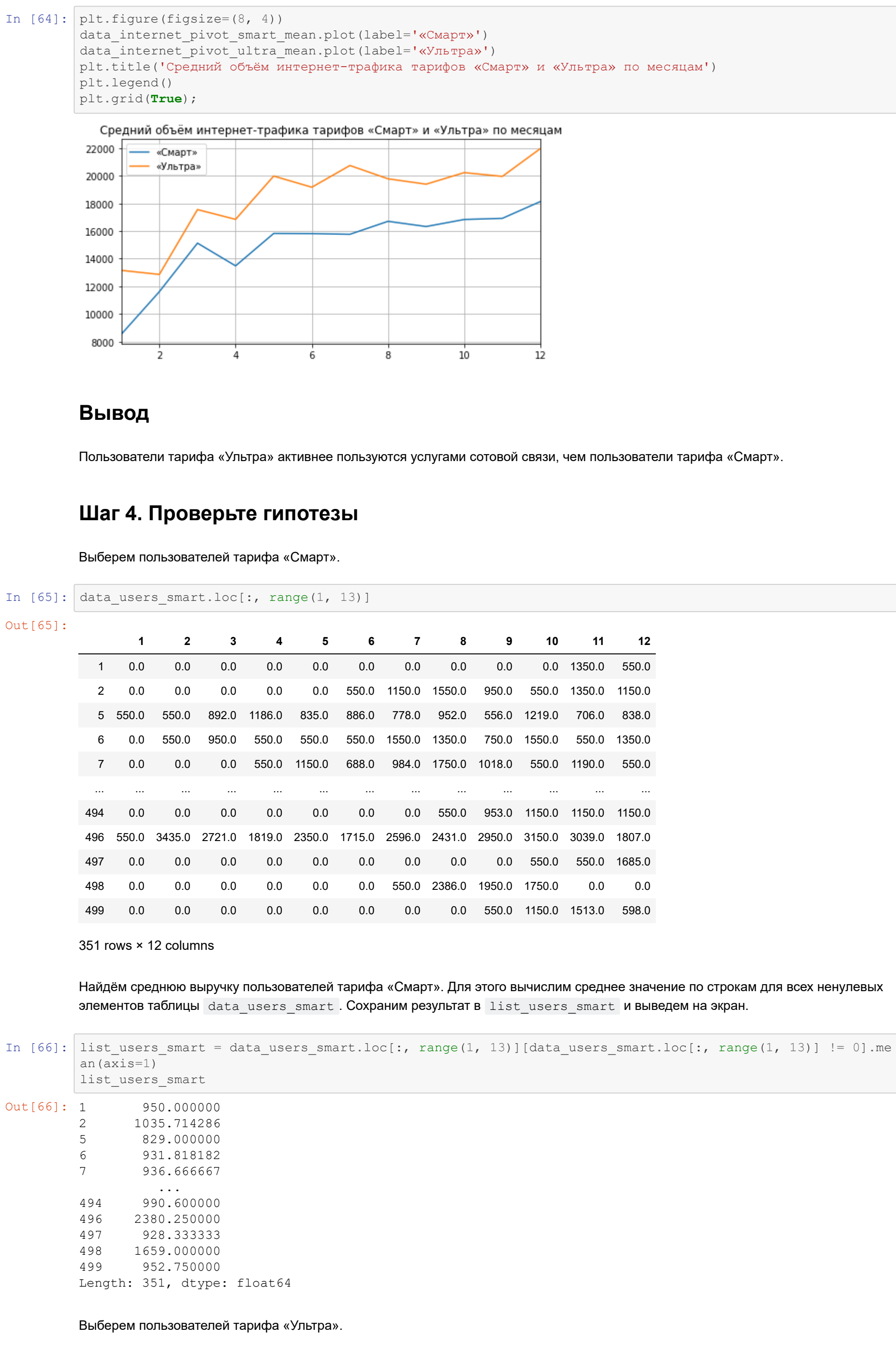
Посчитаем итоговый средний объём интернет-трафика для пользователей тарифа «Смарт».

```
In [62]: data_internet_pivot_smart_mean.mean()
Out [62]: 15088.78370378288
```

Посчитаем итоговый средний объём интернет-трафика для пользователей тарифа «Ультра».

```
In [63]: data_internet_pivot_ultra_mean.mean()
Out [63]: 18468.889801564805
```

Построим график среднего объёма интернет-трафика тарифов «Смарт» и «Ультра» по месяцам.



## Вывод

Пользователи тарифа «Ультра» активнее пользуются услугами сотовой связи, чем пользователи тарифа «Смарт».

## Шаг 4. Проверьте гипотезы

Выберем пользователей тарифа «Смарт».

```
In [65]: data_users_smart.loc[:, range(1, 13)]
Out [65]: 1 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1350.0 550.0
2 0.0 0.0 0.0 0.0 0.0 0.0 550.0 1150.0 1550.0 950.0 550.0 1150.0
3 550.0 550.0 892.0 1186.0 835.0 886.0 778.0 952.0 556.0 1219.0 706.0 838.0
6 0.0 550.0 950.0 550.0 550.0 550.0 1550.0 1350.0 750.0 1550.0 550.0 1350.0
7 0.0 0.0 0.0 0.0 550.0 1150.0 688.0 984.0 1750.0 1018.0 550.0 1190.0 550.0
... ..
494 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 953.0 1150.0 1150.0 1150.0
495 550.0 3435.0 2721.0 1819.0 2350.0 1715.0 2596.0 2431.0 2950.0 3150.0 3039.0 1807.0
497 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 550.0 1685.0
498 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 2386.0 1950.0 1750.0 0.0 0.0
499 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 1150.0 1513.0 598.0
```

351 rows x 12 columns

Найдём среднюю выручку пользователей тарифа «Смарт». Для этого вычислим среднее значение по строкам для всех ненулевых элементов таблицы data\_users\_smart. Сохраним результат в list\_users\_smart и выведем на экран.

```
In [66]: list_users_smart = data_users_smart.loc[:, range(1, 13)][data_users_smart.loc[:, range(1, 13)] != 0].mean(axis=1)
list_users_smart
Out [66]: 1 950.000000
2 1035.714286
3 829.000000
4 931.818182
5 936.666667
...
494 990.600000
495 2380.250000
497 928.333333
498 1659.000000
499 952.750000
Length: 351, dtype: float64
```

Выберем пользователей тарифа «Ультра».

```
In [67]: data_users_ultra.loc[:, range(1, 13)]
Out [67]: 1 0.0 0.0 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
2 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 1950.0 1950.0 1950.0 1950.0
3 0.0 0.0 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
4 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 1950.0
13 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 1950.0
16 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 1950.0
... ..
476 0.0 0.0 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
477 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0
485 0.0 0.0 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
491 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 2750.0
493 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 1950.0
495 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 2400.0 1950.0 3300.0 1950.0
```

149 rows x 12 columns

Найдём среднюю выручку пользователей тарифа «Ультра». Для этого вычислим среднее значение по строкам для всех ненулевых элементов таблицы data\_users\_ultra. Сохраним результат в list\_users\_ultra и выведем на экран.

```
In [68]: list_users_ultra = data_users_ultra.loc[:, range(1, 13)][data_users_ultra.loc[:, range(1, 13)] != 0].mean(axis=1)
list_users_ultra
Out [68]: 0 1950.0
3 1950.0
4 1950.0
13 1950.0
16 1950.0
...
476 1950.0
485 1950.0
491 1950.0
493 2850.0
495 2310.0
Length: 149, dtype: float64
```

Сравним среднюю выручку пользователей тарифов «Смарт» и «Ультра». Для этого проверим гипотезу о равенстве среднего двух генеральных совокупностей методом scipy.stats.ttest\_ind().

```
In [70]: alpha = 0.05 # пороговое значение
results = st.ttest_ind(
    list_users_smart,
    list_users_ultra,
    equal_var = False
)
print('p-значение:', results.pvalue)
if (results.pvalue < alpha):
    print("Отвергаем нулевую гипотезу")
else:
    print("Не получилось отвергнуть нулевую гипотезу")
p-значение: 7.689671550143974e-59
Отвергаем нулевую гипотезу
```

Средняя выручка пользователей тарифов «Ультра» и «Смарт» различается.

```
In [71]: np.var(list_users_smart)
Out [71]: 368543.83725846589
```

```
In [72]: np.var(list_users_ultra)
Out [72]: 89420.18554420327
```

Выберем пользователей тарифа «Смарт», проживающих в Москве, сохраним их в data\_users\_smart\_moscow и выведем на экран.

```
In [73]: data_users_smart_moscow = data_users_smart.query('city == "Москва"').loc[:, range(1, 13)]
data_users_smart_moscow
Out [73]: 1 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1350.0 550.0
7 0.0 0.0 0.0 0.0 550.0 1150.0 688.0 984.0 1750.0 1018.0 550.0 1190.0 550.0
31 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1048.0 738.0
33 0.0 0.0 0.0 0.0 0.0 0.0 550.0 2031.0 2831.0 1640.0 2618.0 1816.0
38 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 950.0 550.0 950.0
... ..
447 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 1606.0 1827.0 931.0
449 0.0 0.0 0.0 0.0 0.0 0.0 550.0 1023.0 1192.0 1005.0 1926.0 314.0
450 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 1634.0 891.0
481 0.0 550.0 1150.0 1550.0 1150.0 550.0 1350.0 1750.0 950.0 1150.0 1550.0 950.0
490 0.0 0.0 0.0 0.0 550.0 1750.0 1550.0 586.0 950.0 550.0 1750.0 1350.0 750.0
```

67 rows x 12 columns

Выберем пользователей тарифа «Смарт», проживающих в других регионах, сохраним их в data\_users\_smart\_not\_moscow, и выведем на экран.

```
In [74]: data_users_smart_not_moscow = data_users_smart.query('city != "Москва"').loc[:, range(1, 13)]
data_users_smart_not_moscow
Out [74]: 1 0.0 0.0 0.0 0.0 0.0 550.0 1150.0 1550.0 950.0 550.0 1150.0 550.0
2 0.0 0.0 0.0 0.0 0.0 550.0 1150.0 1550.0 1550.0 1550.0 1550.0 1550.0
5 550.0 550.0 892.0 1186.0 835.0 886.0 778.0 952.0 556.0 1219.0 706.0 838.0
6 0.0 550.0 950.0 550.0 550.0 550.0 1550.0 1350.0 750.0 1550.0 550.0 1350.0
8 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 568.0
9 0.0 0.0 0.0 550.0 550.0 550.0 750.0 2550.0 2350.0 1950.0 1950.0 950.0 1350.0
... ..
494 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 953.0 1150.0 1150.0 1150.0
496 550.0 3435.0 2721.0 1819.0 2350.0 1715.0 2596.0 2431.0 2950.0 3150.0 3039.0 1807.0
497 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 550.0 1685.0
498 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 2386.0 1950.0 1750.0 0.0 0.0
499 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 550.0 1150.0 1513.0 598.0
```

284 rows x 12 columns

Найдём среднюю выручку пользователей тарифа «Смарт», проживающих в Москве. Для этого вычислим среднее значение по строкам для всех ненулевых элементов таблицы data\_users\_smart\_moscow. Сохраним результат в list\_users\_smart\_moscow и выведем на экран.

```
In [75]: list_users_smart_moscow = data_users_smart_moscow[data_users_smart_moscow != 0].mean(axis=1)
list_users_smart_moscow
Out [75]: 1 950.000000
7 936.666667
31 892.000000
33 1914.000000
38 816.666667
...
447 1832.750000
449 1104.500000
450 1025.700000
457 928.333333
490 1087.333333
Length: 67, dtype: float64
```

Найдём среднюю выручку пользователей тарифа «Смарт», проживающих в других регионах. Для этого вычислим среднее значение по строкам для всех ненулевых элементов таблицы data\_users\_smart\_not\_moscow. Сохраним результат в list\_users\_smart\_not\_moscow и выведем на экран.

```
In [76]: list_users_smart_not_moscow = data_users_smart_not_moscow[data_users_smart_not_moscow != 0].mean(axis=1)
list_users_smart_not_moscow
Out [76]: 2 1035.714286
5 829.000000
6 931.818182
8 568.000000
9 1310.000000
...
494 990.600000
496 2380.250000
497 928.333333
498 1659.000000
499 952.750000
Length: 284, dtype: float64
```

Сравним среднюю выручку пользователей тарифа «Смарт» из Москвы и выручку пользователей из других регионов. Для этого проверим гипотезу о равенстве среднего двух генеральных совокупностей методом scipy.stats.ttest\_ind().

```
In [77]: results = st.ttest_ind(
    list_users_smart_moscow,
    list_users_smart_not_moscow,
    equal_var = False
)
print('p-значение:', results.pvalue)
if (results.pvalue < alpha):
    print("Отвергаем нулевую гипотезу")
else:
    print("Не получилось отвергнуть нулевую гипотезу")
p-значение: 0.14993780434254916
Не получилось отвергнуть нулевую гипотезу
```

Средняя выручка пользователей тарифа «Смарт» из Москвы не отличается от выручки пользователей тарифа «Смарт» из других регионов.

Выберем пользователей тарифа «Ультра», проживающих в Москве, сохраним их в data\_users\_ultra\_moscow и выведем на экран.

```
In [82]: data_users_ultra_moscow = data_users_ultra.query('city == "Москва"').loc[:, range(1, 13)]
data_users_ultra_moscow
Out [82]: 1 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 1950.0
49 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
57 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
80 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
81 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
```

Выберем пользователей тарифа «Ультра», проживающих в других регионах, сохраним их в data\_users\_ultra\_not\_moscow, и выведем на экран.

```
In [83]: data_users_ultra_not_moscow = data_users_ultra.query('city != "Москва"').loc[:, range(1, 13)]
data_users_ultra_not_moscow
Out [83]: 0 0.0 0.0 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
4 0.0 0.0 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
13 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 1950.0
16 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 1950.0
18 0.0 0.0 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
... ..
474 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 3750.0
485 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0
491 0.0 0.0 0.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0 1950.0
493 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 3750.0
495 0.0 0.0 0.0 0.0 0.0 0.0 0.0 1950.0 2400.0 1950.0 3300.0 1950.0
```

117 rows x 12 columns

Найдём среднюю выручку пользователей тарифа «Ультра», проживающих в Москве. Для этого вычислим среднее значение по строкам для всех ненулевых элементов таблицы data\_users\_ultra\_moscow. Сохраним результат в list\_users\_ultra\_moscow и выведем на экран.

```
In [84]: list_users_ultra_moscow = data_users_ultra_moscow[data_users_ultra_moscow != 0].mean(axis=1)
list_users_ultra_moscow.head()
Out [84]: 3 1950.0
49 1950.0
80 1950.0
81 1950.0
dtype: float64
```

Найдём среднюю выручку пользователей тарифа «Ультра», проживающих в других регионах. Для этого вычислим среднее значение по строкам для всех ненулевых элементов таблицы data\_users\_ultra\_not\_moscow. Сохраним результат в list\_users\_ultra\_not\_moscow и выведем на экран.

```
In [85]: list_users_ultra_not_moscow = data_users_ultra_not_moscow[data_users_ultra_not_moscow != 0].mean(axis=1)
list_users_ultra_not_moscow
Out [85]: 0 1950.0
3 1950.0
4 1950.0
16 1950.0
18 1950.0
...
474 2850.0
485 1950.0
491 1950.0
493 2810.0
495 2310.0
Length: 117, dtype: float64
```

Сравним среднюю выручку пользователей тарифа «Ультра» из Москвы и выручку пользователей из других регионов. Для этого проверим гипотезу о равенстве среднего двух генеральных совокупностей методом scipy.stats.ttest\_ind().

```
In [86]: results = st.ttest_ind(
    list_users_ultra_moscow,
    list_users_ultra_not_moscow,
    equal_var = False
)
print('p-значение:', results.pvalue)
if (results.pvalue < alpha):
    print("Отвергаем нулевую гипотезу")
else:
    print("Не получилось отвергнуть нулевую гипотезу")
p-значение: 0.48164162894509566
Не получилось отвергнуть нулевую гипотезу
```

Средняя выручка пользователей тарифа «Ультра» из Москвы не отличается от выручки пользователей тарифа «Ультра» из других регионов.

Выберем пользователей тарифа «Ультра», проживающих в Москве, сохраним их в data\_users\_ultra\_moscow и выведем на экран.

```
In [87]: np.var(list_users_ultra_moscow)
Out [87]: 40193.8668492251
```

```
In [88]: np.var(list_users_ultra_not_moscow)
Out [88]: 102649.57007101324
```

Средняя выручка пользователей тарифа «Ультра» из Москвы не отличается от выручки пользователей тарифа «Ультра» из других регионов.

## Вывод

Средняя выручка пользователей тарифов «Ультра» и «С