

Procrastinate Pro+

Расследование

Почему компания терпит убытки

Описание проекта

Вы — маркетинговый аналитик развлекательного приложения Procrastinate Pro+.

Несмотря на огромные вложения в рекламу, последние несколько месяцев **компания терпит убытки**.

0.1 Ваша задача — разобраться в причинах и помочь компании выйти в плюс.

Есть данные о пользователях, привлечённых с 1 мая по 27 октября 2019 года:

- лог сервера с данными об их посещениях,
- выгрузка их покупок за этот период,
- рекламные расходы.

0.2 Вам предстоит изучить:

- откуда приходят пользователи и какими устройствами они пользуются,
- сколько стоит привлечение пользователей из различных рекламных каналов;
- сколько денег приносит каждый клиент,
- когда расходы на привлечение клиента окупаются,
- какие факторы мешают привлечению клиентов.

1 План работы

1. Загрузите данные и подготовьте их к анализу

- Положите данные о визитах, заказах и рекламных тратах в переменные.
- Подготовьте данные к анализу. Убедитесь, что тип данных во всех колонках соответствует значениям. Проверьте отсутствие дубликатов.

2. Задайте функции для расчета и анализа LTV, ROI, удержания и конверсии

3. Проведите исследовательский анализ данных

- Создайте пользовательские профили. Определите минимальную и максимальную даты привлечения пользователей.
- Выясните, из каких стран пользователи приходят в приложение и на какую страну приходится больше всего платящих пользователей. Постройте таблицу, отражающую количество пользователей и долю платящих из каждой страны.
- Узнайте, какими устройствами пользуются клиенты и какие устройства предпочитают платящие пользователи. Постройте таблицу, отражающую количество пользователей и долю платящих для каждого устройства.

- Изучите рекламные источники привлечения и определите каналы, из которых пришло больше всего платящих пользователей. Постройте таблицу, отражающую количество пользователей и долю платящих для каждого канала привлечения.

4. Маркетинг

- Посчитайте общую сумму расходов на маркетинг. Выясните, как траты распределены по источникам. Визуализируйте изменения метрик во времени.
- Узнайте, сколько в среднем стоило привлечение одного пользователя из каждого источника. Рассчитайте средний CAC на одного пользователя для всего проекта и для каждого источника трафика. Используйте профили пользователей.

5. Оцените окупаемость рекламы

Используя графики LTV, ROI и CAC, проанализируйте окупаемость рекламы. Считайте, что на календаре 1 ноября 2019 года, а в бизнес-плане заложено, что пользователи должны окупаться не позднее чем через две недели после привлечения. Необходимость включения в анализ органических пользователей определите самостоятельно.

- Проанализируйте общую окупаемость рекламы. Постройте графики LTV и ROI, а также графики динамики LTV, CAC и ROI.
- Проанализируйте окупаемость рекламы с разбивкой по рекламным каналам. Постройте графики LTV и ROI, а также графики динамики LTV, CAC и ROI.
- Проанализируйте окупаемость рекламы с разбивкой по странам. Постройте графики LTV и ROI, а также графики динамики LTV, CAC и ROI.
- Постройте и изучите графики конверсии и удержания с разбивкой по устройствам, странам, рекламным каналам.

Ответьте на такие вопросы:

- Окупается ли реклама в целом?
- Какие устройства, страны и каналы могут снижать окупаемость рекламы?
- Чем могут быть вызваны проблемы окупаемости?
 - Опишите возможные причины обнаруженных проблем и промежуточные рекомендации для рекламного отдела.

6. Напишите выводы

- Выделите причины неэффективности привлечения пользователей.
- Сформулируйте рекомендации для отдела маркетинга.

In [2]:



```
1 # импорт библиотек
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 # importing a library for operating with a data type "NaN"
6 import numpy as np
7 import datetime
8 from datetime import datetime, timedelta
```

In [3]:



```
1 # отключаем некритичные уведомления
2 import warnings
3 warnings.filterwarnings('ignore')
4 # показывать до 40ка колонок
5 pd.set_option('display.max.columns', 40)
6 # установка формата вывода на дисплей численных значений
7 pd.options.display.float_format = '{:,.2f}'.format
8 # Библиотека для отображения картинок
9 from IPython.display import Image
10 # библиотека для укоругления в БОльшую сторону
11 import math
```

In [4]:



```
1 # считываем данные
2 try:
3     # Yandex path
4     visits = pd.read_csv('/datas...ort.csv')
5     orders = pd.read_csv('/datas...ort.csv')
6     costs = pd.read_csv('/datas...rt.csv')
7
8
9 except:
10     # Local patch
11     visits = pd.read_csv(r"C:\Users\eddyd...ort.csv")
12     orders = pd.read_csv(r"C:\Users\eddyd...ort.csv")
13     costs = pd.read_csv(r"C:\Users\eddyd...ort.csv")
```

In [5]:



```
1 # сведения о данных
2 def full_info(dataframe):
3     print(dataframe.columns)
4     display(dataframe.head(3))
5     display(dataframe.info())
6     print("describe")
7     display(dataframe.describe())
8     print("Доли отсутствующих значений")
9     print(round(dataframe.isna().sum() * 100 / len(dataframe), 2))
10     return
```

2 Предобработка данных

2.1 visits

```
1 full_info(visits)
```

```
Index(['User Id', 'Region', 'Device', 'Channel', 'Session Start',  
      'Session End'],  
      dtype='object')
```

	User Id	Region	Device	Channel	Session Start	Session End
0	981449118918	United States	iPhone	organic	2019-05-01 02:36:01	2019-05-01 02:45:01
1	278965908054	United States	iPhone	organic	2019-05-01 04:46:31	2019-05-01 04:47:35
2	590706206550	United States	Mac	organic	2019-05-01 14:09:25	2019-05-01 15:32:08

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 309901 entries, 0 to 309900  
Data columns (total 6 columns):  
#   Column          Non-Null Count  Dtype  
---  -  
0   User Id         309901 non-null int64  
1   Region          309901 non-null object  
2   Device          309901 non-null object  
3   Channel         309901 non-null object  
4   Session Start   309901 non-null object  
5   Session End     309901 non-null object  
dtypes: int64(1), object(5)  
memory usage: 14.2+ MB
```

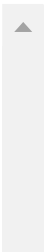
None

describe

	User Id
count	309,901.00
mean	499,766,449,382.69
std	288,789,916,887.83
min	599,326.00
25%	249,369,122,776.00
50%	498,990,589,687.00
75%	749,521,111,616.00
max	999,999,563,947.00

Доли отсутствующих значений

User Id	0.00
Region	0.00
Device	0.00
Channel	0.00
Session Start	0.00



Session End 0.00
dtype: float64



2.1.1 пропусков данных нет

2.1.2 Названия колонок приведём к змеинному формату

In [7]:



```
1 # всё - нижним регистром
2 visits.columns = visits.columns.str.lower()
3 visits.columns
```

Out[7]:

```
Index(['user id', 'region', 'device', 'channel', 'session start',
      'session end'],
      dtype='object')
```

In [8]:



```
1 # заменим пробелы символами подчеркивания
2 visits.columns = [_name.replace(' ', '_') for _name in visits.columns]
3 visits.columns
```

Out[8]:

```
Index(['user_id', 'region', 'device', 'channel', 'session_start',
      'session_end'],
      dtype='object')
```

2.1.3 user_id

посмотрим, сколько уникальных значений

In [9]:



```
1 visits['user_id'].value_counts()
```

Out[9]:

```
33606575057    46
943775408561    36
901180916748    35
870784569391    34
764086596354    33
..
214203066007     1
369265191867     1
346271445800     1
133742530598     1
279181973476     1
Name: user_id, Length: 150008, dtype: int64
```

В колонке **user_id** есть записи с одной строкой на значение user_id

и есть с несколькими строками на значение user_id
150000 уникальных пользователей

2.1.4 region

посмотрим, сколько уникальных значений
нет ли неявных дубликатов

In [10]:



```
1 visits['region'].value_counts()
```

Out[10]:

```
United States    207327
UK               36419
France          35396
Germany         30759
Name: region, dtype: int64
```

В **region** встречается **4** уникальных значения.

Лидер упоминаемости - **United States**

2.1.5 device

посмотрим, сколько уникальных значений
нет ли неявных дубликатов

In [11]:



```
1 visits['device'].value_counts()
```

Out[11]:

```
iPhone    112603
Android   72590
PC        62686
Mac       62022
Name: device, dtype: int64
```

В **device** встречается **4** уникальных значения.

Лидер упоминаемости - **iPhone**

2.1.6 channel

посмотрим, сколько уникальных значений
нет ли неявных дубликатов

In [12]:



```
1 visits['channel'].value_counts()
```

Out[12]:

```
organic          107760
TipTop           54794
FaceBoom        49022
WahooNetBanner  20465
LeapBob         17013
OppleCreativeMedia 16794
RocketSuperAds  12724
YRabbit         9053
MediaTornado    8878
AdNonSense      6891
lambdaMediaAds  6507
Name: channel, dtype: int64
```

В **channel** встречается **11** уникальных значения.

Лидер упоминаемости - **organic**

А среди платных каналов - **TipTop** и **FaceBoom**

2.1.7 session start и session_end

переведём в формат даты

In [13]:



```
1 visits['session_start'] = pd.to_datetime(visits['session_start'])
2 visits['session_end'] = pd.to_datetime(visits['session_end'])
```

минимальные и максимальные значения (временной интервал)

In [14]:



```
1 print("Минимальное значение session_start", min(visits['session_start']))
2 print("Минимальное значение session_end", min(visits['session_end']))
3 print()
4 print("Максимальное значение session_start", max(visits['session_start']))
5 print("Максимальное значение session_end", max(visits['session_end']))
6
```

```
Минимальное значение session_start 2019-05-01 00:00:41
Минимальное значение session_end 2019-05-01 00:07:06
```

```
Максимальное значение session_start 2019-10-31 23:59:23
Максимальное значение session_end 2019-11-01 01:38:46
```

2.2 orders

In [15]:



```
1 full_info(orders)
```

Index(['User Id', 'Event Dt', 'Revenue'], dtype='object')

	User Id	Event Dt	Revenue
0	188246423999	2019-05-01 23:09:52	4.99
1	174361394180	2019-05-01 12:24:04	4.99
2	529610067795	2019-05-01 11:34:04	4.99

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 40212 entries, 0 to 40211
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   User Id     40212 non-null  int64
1   Event Dt    40212 non-null  object
2   Revenue     40212 non-null  float64
dtypes: float64(1), int64(1), object(1)
memory usage: 942.6+ KB
```

None

describe

	User Id	Revenue
count	40,212.00	40,212.00
mean	499,029,531,203.23	5.37
std	286,093,675,967.16	3.45
min	599,326.00	4.99
25%	251,132,440,436.75	4.99
50%	498,283,972,665.00	4.99
75%	743,332,711,780.00	4.99
max	999,895,427,370.00	49.99

Доли отсутствующих значений
User Id 0.00
Event Dt 0.00
Revenue 0.00
dtype: float64

2.2.1 Пропусков данных нет

2.2.2 Названия колонок приведём к змеиному формату

In [16]:

```
1 # всё - нижним регистром
2 orders.columns = orders.columns.str.lower()
3 orders.columns
```

Out[16]:

```
Index(['user id', 'event dt', 'revenue'], dtype='object')
```

In [17]:

```
1 # заменим пробелы символами подчеркивания
2 orders.columns = [_name.replace(' ', '_') for _name in orders.columns]
3 orders.columns
```

Out[17]:

```
Index(['user_id', 'event_dt', 'revenue'], dtype='object')
```

2.2.3 user_id

посмотрим, сколько уникальных значений

In [18]:

```
1 orders['user_id'].value_counts()
```

Out[18]:

```
901180916748    22
883098437811    20
75337957494     19
512471511263    19
295795879965    19
..
829252887757     1
499471996783     1
847348136580     1
390188868722     1
168548862926     1
Name: user_id, Length: 8881, dtype: int64
```

В колонке user_id есть записи с одной строкой на значение user_id
и есть с несколькими строками на значение user_id
8881 уникальных пользователей

интересно, это количество значительно меньше, чем в данных о визитах

2.2.4 event_dt

переведём в формат даты

посмотрим мин и макс значения

In [19]:



```
1 orders['event_dt'] = pd.to_datetime(orders['event_dt'])
2 print("Минимальное значение event_dt", min(orders['event_dt']))
3 print("Максимальное значение event_dt", max(orders['event_dt']))
```

Минимальное значение event_dt 2019-05-01 00:28:11
Максимальное значение event_dt 2019-10-31 23:56:56

Интервал времени **совпадает** с данными в табл о визитах

2.2.5 revenue

In [20]:



```
1 print("Минимальное значение revenue", min(orders['revenue']))
2 print("Среднее значение revenue", orders['revenue'].mean())
3 print("Медиана значение revenue", orders['revenue'].median())
4 print("Максимальное значение revenue", max(orders['revenue']))
```

Минимальное значение revenue 4.99
Среднее значение revenue 5.370607778770249
Медиана значение revenue 4.99
Максимальное значение revenue 49.99

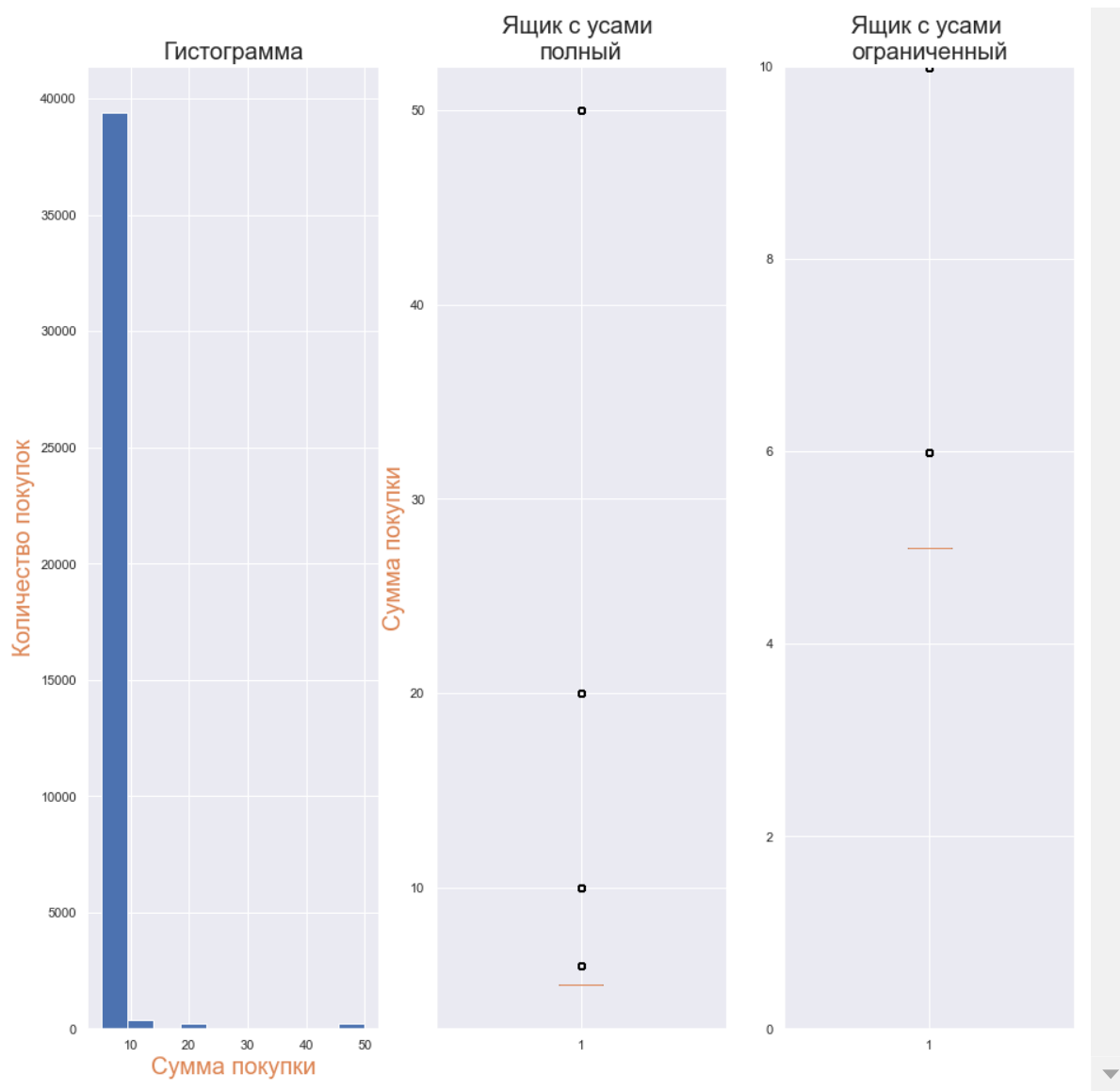
Распределение и ящик с усами

In [21]:



```
1 # setup size plot
2 sns.set(rc = {'figure.figsize':(15,15)})
3 fig=plt.figure()
4
5 ax=fig.add_subplot(1,3,1)
6 ax2=fig.add_subplot(1,3,2)
7 ax3=fig.add_subplot(1,3,3)
8 # plot line graph
9 ax.hist(orders['revenue'])
10 ax.set_ylabel("Количество покупок", color="C1", fontsize=20)
11 ax.set_xlabel("Сумма покупки", color="C1", fontsize=20)
12 ax.set_title("Гистограмма", fontsize=20)
13
14 # plot box
15 ax2.boxplot(orders['revenue'])
16 ax2.set_ylabel("Сумма покупки", color="C1", fontsize=20)
17 ax2.set_title("Ящик с усами \nполный", fontsize=20)
18
19 # plot box 2
20 ax3.boxplot(orders['revenue'])
21 ax3.set_ylabel("", color="C1", fontsize=20)
22 ax3.set_title("Ящик с усами \nограниченный", fontsize=20)
23 ax3.set_ylim(0, 10.0)
24
25 plt.show()
```





Подавляющее число продаж - минимальный пакет за базовую стоимость 4,99

2.2.5.1 Посмотрим, какие покупки по цене встречаются

In [22]:

```
1 orders['revenue'].unique()
2
```

Out[22]:

```
array([ 4.99,  5.99,  9.99, 49.99, 19.99])
```

Подписки: * 4,99 * 5,99 * 9,99 * 19,99 * 49,99

2.3 costs

In [23]:



```
1 full_info(costs)
```

Index(['dt', 'Channel', 'costs'], dtype='object')

	dt	Channel	costs
0	2019-05-01	FaceBoom	113.30
1	2019-05-02	FaceBoom	78.10
2	2019-05-03	FaceBoom	85.80

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1800 entries, 0 to 1799
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   dt          1800 non-null   object
1   Channel     1800 non-null   object
2   costs       1800 non-null   float64
dtypes: float64(1), object(2)
memory usage: 42.3+ KB
```

None

describe

	costs
count	1,800.00
mean	58.61
std	107.74
min	0.80
25%	6.50
50%	12.29
75%	33.60
max	630.00

Доли отсутствующих значений

dt	0.00
Channel	0.00
costs	0.00

dtype: float64

Пропусков данных нет

2.3.1 Названия колонок приведём к змеиному формату

In [24]:



```
1 # всё - нижним регистром
2 costs.columns = costs.columns.str.lower()
3 costs.columns
```

Out[24]:

```
Index(['dt', 'channel', 'costs'], dtype='object')
```

2.3.2 dt

переведём в формат даты

- название месяца
- в коротком формате

посмотрим мин и макс значения

In [25]:



```
1 # выделяем месяц
2 costs['month'] = pd.to_datetime(costs['dt']).dt.strftime('%b')
3 # выделяем дату
4 costs['dt'] = pd.to_datetime(costs['dt']).dt.date
5 print("Минимальное значение event_dt", min(costs['dt']))
6 print("Максимальное значение event_dt", max(costs['dt']))
```

Минимальное значение event_dt 2019-05-01

Максимальное значение event_dt 2019-10-27

Интервал времени совпадает с другими таблицами

Добавим колонку с названием месяца

2.3.3 channel

Посмотрим уникальные значения и количества записей с ними

In [26]:



```
1 costs['channel'].value_counts()
```

Out[26]:

```
FaceBoom          180
MediaTornado      180
RocketSuperAds    180
TipTop            180
YRabbit           180
AdNonSense        180
LeapBob           180
OpplCreativeMedia 180
WahooNetBanner    180
lambdaMediaAds    180
Name: channel, dtype: int64
```

10 значений - все каналы за исключением органического
Количество записей одинаковое - **180** для каждого
Значит, во все каналы вкладывались равномерно

2.3.4 costs

In [27]:



```
1 print("Минимальное значение costs", min(costs['costs']))
2 print("Среднее значение costs", costs['costs'].mean())
3 print("Медиана значение costs", costs['costs'].median())
4 print("Максимальное значение costs", max(costs['costs']))
```

```
Минимальное значение costs 0.8
Среднее значение costs 58.609611111111118
Медиана значение costs 12.285000000000002
Максимальное значение costs 630.0
```

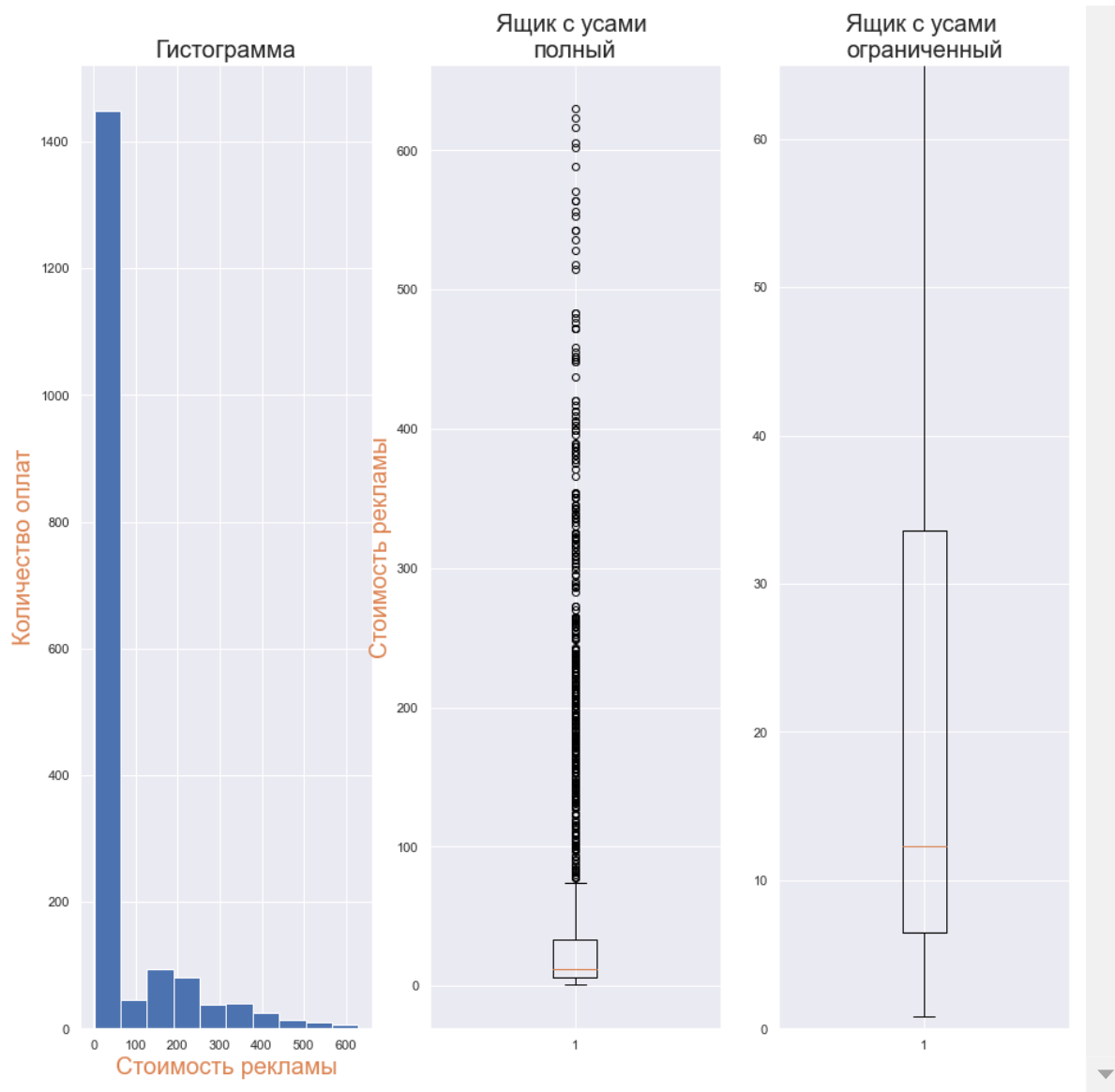
Распределение и ящик с усами

In [28]:



```
1 # setup size plot
2 sns.set(rc = {'figure.figsize':(15,15)})
3 fig=plt.figure()
4
5 ax=fig.add_subplot(1,3,1)
6 ax2=fig.add_subplot(1,3,2)
7 ax3=fig.add_subplot(1,3,3)
8 # plot line graph
9 ax.hist(costs['costs'])
10 ax.set_ylabel("Количество оплат", color="C1", fontsize=20)
11 ax.set_xlabel("Стоимость рекламы", color="C1", fontsize=20)
12 ax.set_title("Гистограмма", fontsize=20)
13
14 # plot box
15 ax2.boxplot(costs['costs'])
16 ax2.set_ylabel("Стоимость рекламы", color="C1", fontsize=20)
17 ax2.set_title("Ящик с усами \nполный", fontsize=20)
18
19 # plot box 2
20 ax3.boxplot(costs['costs'])
21 ax3.set_ylabel("", color="C1", fontsize=20)
22 ax3.set_title("Ящик с усами \nограниченный", fontsize=20)
23 ax3.set_ylim(0, 65.0)
24
25 plt.show()
```





Интересная картина.

Основная масса платежей за рекламные каналы в интервале от 7 до 35 единиц.

При этом наблюдается достаточно мощный горб гистограммы в области 150 - 250 единиц за канал. В дальнейшем можно будет внимательнее рассматривать эту историю, при необходимости.

2.4 Вывод по Предобработка данных

Данные поступили в хорошем состоянии.

Названия колонок приведены к змеиную формат.

Данные без пропусков и дубликатов.

Также отсутствуют неявные дубликаты.

Количества каналов и временные интервалы по всем таблицам коррелируют друг с другом (то есть совпадают).

Проблем не выявлено.

2.4.1 Структура данных

2.4.1.1 visits

- user_id - идентификатор пользователя, целочисленное
- region - регион, в котором находится пользователь, строковое
- device - тип устройства, строковое
- channel - канал привлечения, строковое
- session_start - тайм-код начала сессии, полная дата
- session_end - тайм-код окончания сессии, полная дата

2.4.1.2 orders

- user_id - идентификатор пользователя, целочисленное
- event_dt - тайм-код покупки, полная дата
- revenue - сумма покупки, численное

2.4.1.3 costs

- dt - дата оплаты рекламного канала, дата без времени
- channel - рекламный канал, строковое
- costs - сумма оплаты, численное

2.4.2 Основные особенности на предварительном анализе:

- * 150000 уникальных пользователей
 - * United States - лидер упоминаемости
 - * iPhone - лидер упоминаемости
 - * 10 платных каналов трафика
 - * **TipTop** и **FaceBoom** - лидеры по упоминаемости
 - * интервалы времени с мая 2019 по ноябрь 2019
 - * покупателей лишь 8881 из 150000 визитёров
 - * Подавляющее число продаж - минимальный пакет за базовую стоимость 4,99
 - * Основная масса платежей за рекламные каналы в интервале от 7 до 35 единиц.
- * * Приэтом наблюдается достаточно мощный горб гистограммы в области 150 - 250 единиц за канал.

3 Функции для расчётов

для расчета и анализа LTV, ROI, удержания и конверсии

3.1 LTV и ROI

get_ltv(profiles, purchases, observation_date, horizon_days, dimensions=[], ignore_horizon=False)

return (
 result_raw, # сырые данные
 result_grouped, # таблица LTV
 result_in_time, # таблица динамики LTV

```
roi_grouped, # таблица ROI
roi_in_time, # таблица динамики ROI
)
```

```

1  def get_ltv(
2      profiles,
3      purchases,
4      observation_date,
5      horizon_days,
6      dimensions=[],
7      ignore_horizon=False,
8  ):
9
10     # исключаем пользователей, не «доживших» до горизонта анализа
11     last_suitable_acquisition_date = observation_date
12     if not ignore_horizon:
13         last_suitable_acquisition_date = observation_date - timedelta(
14             days=horizon_days - 1
15         )
16     result_raw = profiles.query('dt <= @last_suitable_acquisition_date')
17     # добавляем данные о покупках в профили
18     result_raw = result_raw.merge(
19         purchases[['user_id', 'event_dt', 'revenue']], on='user_id', how='left'
20     )
21     # рассчитываем лайфтайм пользователя для каждой покупки
22     result_raw['lifetime'] = (
23         result_raw['event_dt'] - result_raw['first_ts']
24     ).dt.days
25     # группируем по cohort, если в dimensions ничего нет
26     if len(dimensions) == 0:
27         result_raw['cohort'] = 'All users'
28         dimensions = dimensions + ['cohort']
29
30     # функция группировки по желаемым признакам
31     def group_by_dimensions(df, dims, horizon_days):
32         # строим «треугольную» таблицу выручки
33         result = df.pivot_table(
34             index=dims, columns='lifetime', values='revenue', aggfunc='sum'
35         )
36         # находим сумму выручки с накоплением
37         result = result.fillna(0).cumsum(axis=1)
38         # вычисляем размеры когорт
39         cohort_sizes = (
40             df.groupby(dims)
41             .agg({'user_id': 'nunique'})
42             .rename(columns={'user_id': 'cohort_size'})
43         )
44         # объединяем размеры когорт и таблицу выручки
45         result = cohort_sizes.merge(result, on=dims, how='left').fillna(0)
46         # считаем LTV: делим каждую «ячейку» в строке на размер когорты
47         result = result.div(result['cohort_size'], axis=0)
48         # исключаем все лайфтаймы, превышающие горизонт анализа
49         result = result[['cohort_size'] + list(range(horizon_days))]
50         # восстанавливаем размеры когорт
51         result['cohort_size'] = cohort_sizes
52
53         # собираем датафрейм с данными пользователей и значениями CAC,
54         # добавляя параметры из dimensions
55         cac = df[['user_id', 'acquisition_cost'] + dims].drop_duplicates()
56
57         # считаем средний CAC по параметрам из dimensions
58         cac = (
59             cac.groupby(dims)

```

```

60         .agg({'acquisition_cost': 'mean'})
61         .rename(columns={'acquisition_cost': 'cac'})
62     )
63
64     # считаем ROI: делим LTV на CAC
65     roi = result.div(cac['cac'], axis=0)
66
67     # удаляем строки с бесконечным ROI
68     roi = roi[~roi['cohort_size'].isin([np.inf])]
69
70     # восстанавливаем размеры когорт в таблице ROI
71     roi['cohort_size'] = cohort_sizes
72
73     # добавляем CAC в таблицу ROI
74     roi['cac'] = cac['cac']
75
76     # в финальной таблице оставляем размеры когорт, CAC
77     # и ROI в лайфтаймы, не превышающие горизонт анализа
78     roi = roi[['cohort_size', 'cac'] + list(range(horizon_days))]
79
80     # возвращаем таблицы LTV и ROI
81     return result, roi
82
83     # получаем таблицы LTV и ROI
84     result_grouped, roi_grouped = group_by_dimensions(
85         result_raw, dimensions, horizon_days
86     )
87
88     # для таблиц динамики убираем 'cohort' из dimensions
89     if 'cohort' in dimensions:
90         dimensions = []
91
92     # получаем таблицы динамики LTV и ROI
93     result_in_time, roi_in_time = group_by_dimensions(
94         result_raw, dimensions + ['dt'], horizon_days
95     )
96
97     return (
98         result_raw, # сырые данные
99         result_grouped, # таблица LTV
100         result_in_time, # таблица динамики LTV
101         roi_grouped, # таблица ROI
102         roi_in_time, # таблица динамики ROI
103     )
104

```

3.2 удержания

get_retention(profiles, sessions, observation_date, horizon_days, dimensions=[], ignore_horizon=False)

return result_raw, result_grouped, result_in_time

```

1  def get_retention(
2      profiles,
3      sessions,
4      observation_date,
5      horizon_days,
6      dimensions=[],
7      ignore_horizon=False,
8  ):
9
10     # добавляем столбец payer в передаваемый dimensions список
11     dimensions = ['payer'] + dimensions
12
13     # исключаем пользователей, не «доживших» до горизонта анализа
14     last_suitable_acquisition_date = observation_date
15     if not ignore_horizon:
16         last_suitable_acquisition_date = observation_date - timedelta(
17             days=horizon_days - 1
18         )
19     result_raw = profiles.query('dt <= @last_suitable_acquisition_date')
20
21     # собираем «сырые» данные для расчёта удержания
22     result_raw = result_raw.merge(
23         sessions[['user_id', 'session_start']], on='user_id', how='left'
24     )
25     result_raw['lifetime'] = (
26         result_raw['session_start'] - result_raw['first_ts']
27     ).dt.days
28
29     # функция для группировки таблицы по желаемым признакам
30     def group_by_dimensions(df, dims, horizon_days):
31         result = df.pivot_table(
32             index=dims, columns='lifetime', values='user_id', aggfunc='nunique'
33         )
34         cohort_sizes = (
35             df.groupby(dims)
36             .agg({'user_id': 'nunique'})
37             .rename(columns={'user_id': 'cohort_size'})
38         )
39         result = cohort_sizes.merge(result, on=dims, how='left').fillna(0)
40         result = result.div(result['cohort_size'], axis=0)
41         result = result[['cohort_size'] + list(range(horizon_days))]
42         result['cohort_size'] = cohort_sizes
43         return result
44
45     # получаем таблицу удержания
46     result_grouped = group_by_dimensions(result_raw, dimensions, horizon_days)
47
48     # получаем таблицу динамики удержания
49     result_in_time = group_by_dimensions(
50         result_raw, dimensions + ['dt'], horizon_days
51     )
52
53     # возвращаем обе таблицы и сырые данные
54     return result_raw, result_grouped, result_in_time

```

3.3 конверсии

```
get_conversion( profiles, purchases, observation_date, horizon_days, dimensions=[], ignore_horizon=False )
```

```
return result_raw, result_grouped, result_in_time
```

```

1  def get_conversion(
2      profiles,
3      purchases,
4      observation_date,
5      horizon_days,
6      dimensions=[],
7      ignore_horizon=False,
8  ):
9
10     # исключаем пользователей, не «доживших» до горизонта анализа
11     last_suitable_acquisition_date = observation_date
12     if not ignore_horizon:
13         last_suitable_acquisition_date = observation_date - timedelta(
14             days=horizon_days - 1
15         )
16     result_raw = profiles.query('dt <= @last_suitable_acquisition_date')
17
18     # определяем дату и время первой покупки для каждого пользователя
19     first_purchases = (
20         purchases.sort_values(by=['user_id', 'event_dt'])
21         .groupby('user_id')
22         .agg({'event_dt': 'first'})
23         .reset_index()
24     )
25
26     # добавляем данные о покупках в профили
27     result_raw = result_raw.merge(
28         first_purchases[['user_id', 'event_dt']], on='user_id', how='left'
29     )
30
31     # рассчитываем лайфтайм для каждой покупки
32     result_raw['lifetime'] = (
33         result_raw['event_dt'] - result_raw['first_ts']
34     ).dt.days
35
36     # группируем по cohort, если в dimensions ничего нет
37     if len(dimensions) == 0:
38         result_raw['cohort'] = 'All users'
39         dimensions = dimensions + ['cohort']
40
41     # функция для группировки таблицы по желаемым признакам
42     def group_by_dimensions(df, dims, horizon_days):
43         result = df.pivot_table(
44             index=dims, columns='lifetime', values='user_id', aggfunc='nunique'
45         )
46         result = result.fillna(0).cumsum(axis = 1)
47         cohort_sizes = (
48             df.groupby(dims)
49             .agg({'user_id': 'nunique'})
50             .rename(columns={'user_id': 'cohort_size'})
51         )
52         result = cohort_sizes.merge(result, on=dims, how='left').fillna(0)
53         # делим каждую «ячейку» в строке на размер когорты
54         # и получаем conversion rate
55         result = result.div(result['cohort_size'], axis=0)
56         result = result[['cohort_size'] + list(range(horizon_days))]
57         result['cohort_size'] = cohort_sizes
58         return result
59

```



```
60     # получаем таблицу конверсии
61     result_grouped = group_by_dimensions(result_raw, dimensions, horizon_days)
62
63     # для таблицы динамики конверсии убираем 'cohort' из dimensions
64     if 'cohort' in dimensions:
65         dimensions = []
66
67     # получаем таблицу динамики конверсии
68     result_in_time = group_by_dimensions(
69         result_raw, dimensions + ['dt'], horizon_days
70     )
71
72     # возвращаем обе таблицы и сырые данные
73     return result_raw, result_grouped, result_in_time
```

3.4 Служебные функции

- `get_profiles(sessions, orders, events, ad_costs, event_names=[])`

return profiles

```

1 def get_profiles(sessions, orders, events, ad_costs, event_names=[]):
2
3     # находим параметры первых посещений
4     profiles = (
5         sessions.sort_values(by=['user_id', 'session_start'])
6         .groupby('user_id')
7         .agg(
8             {
9                 'session_start': 'first',
10                'channel': 'first',
11                'device': 'first',
12                'region': 'first',
13            }
14        )
15        .rename(columns={'session_start': 'first_ts'})
16        .reset_index()
17    )
18
19    # для когортного анализа определяем дату первого посещения
20    # и первый день месяца, в который это посещение произошло
21    profiles['dt'] = profiles['first_ts'].dt.date
22    profiles['month'] = profiles['first_ts'].astype('datetime64[M]')
23
24    # добавляем признак платящих пользователей
25    profiles['payer'] = profiles['user_id'].isin(orders['user_id'].unique())
26
27    # добавляем флаги для всех событий из event_names
28    for event in event_names:
29        if event in events['event_name'].unique():
30            profiles[event] = profiles['user_id'].isin(
31                events.query('event_name == @event')['user_id'].unique()
32            )
33
34    # считаем количество уникальных пользователей
35    # с одинаковыми источником и датой привлечения
36    new_users = (
37        profiles.groupby(['dt', 'channel'])
38        .agg({'user_id': 'nunique'})
39        .rename(columns={'user_id': 'unique_users'})
40        .reset_index()
41    )
42
43    # объединяем траты на рекламу и число привлечённых пользователей
44    ad_costs = ad_costs.merge(new_users, on=['dt', 'channel'], how='left')
45
46    # делим рекламные расходы на число привлечённых пользователей
47    ad_costs['acquisition_cost'] = ad_costs['costs'] / ad_costs['unique_users']
48
49    # добавляем стоимость привлечения в профили
50    profiles = profiles.merge(
51        ad_costs[['dt', 'channel', 'acquisition_cost']],
52        on=['dt', 'channel'],
53        how='left',
54    )
55
56    # стоимость привлечения органических пользователей равна нулю
57    profiles['acquisition_cost'] = profiles['acquisition_cost'].fillna(0)
58
59    return profiles

```

3.5 Графические функции

- `filter_data(df, window)`
- ▪ `return df`
- `plot_retention(retention, retention_history, horizon, window=7)`
- `plot_conversion(conversion, conversion_history, horizon, window=7)`
- `plot_ltv_roi(ltv, ltv_history, roi, roi_history, horizon, window=7)`

In [33]:



```
1 def filter_data(df, window):
2     # для каждого столбца применяем скользящее среднее
3     for column in df.columns.values:
4         df[column] = df[column].rolling(window).mean()
5     return df
6
```

```

1 def plot_retention(retention, retention_history, horizon, window=7):
2
3     # задаём размер сетки для графиков
4     plt.figure(figsize=(15, 10))
5
6     # исключаем размеры когорт и удержание первого дня
7     retention = retention.drop(columns=['cohort_size', 0])
8     # в таблице динамики оставляем только нужный лайфтайм
9     retention_history = retention_history.drop(columns=['cohort_size'])[
10         [horizon - 1]
11     ]
12
13     # если в индексах таблицы удержания только payer,
14     # добавляем второй признак – cohort
15     if retention.index.nlevels == 1:
16         retention['cohort'] = 'All users'
17         retention = retention.reset_index().set_index(['cohort', 'payer'])
18
19     # в таблице графиков – два столбца и две строки, четыре ячейки
20     # в первой строим кривые удержания платящих пользователей
21     ax1 = plt.subplot(2, 2, 1)
22     retention.query('payer == True').droplevel('payer').T.plot(
23         grid=True, ax=ax1
24     )
25     plt.legend()
26     plt.xlabel('Лайфтайм')
27     plt.title('Удержание платящих пользователей')
28
29     # во второй ячейке строим кривые удержания неплатящих
30     # вертикальная ось – от графика из первой ячейки
31     ax2 = plt.subplot(2, 2, 2, sharey=ax1)
32     retention.query('payer == False').droplevel('payer').T.plot(
33         grid=True, ax=ax2
34     )
35     plt.legend()
36     plt.xlabel('Лайфтайм')
37     plt.title('Удержание неплатящих пользователей')
38
39     # в третьей ячейке – динамика удержания платящих
40     ax3 = plt.subplot(2, 2, 3)
41     # получаем названия столбцов для сводной таблицы
42     columns = [
43         name
44         for name in retention_history.index.names
45         if name not in ['dt', 'payer']
46     ]
47     # фильтруем данные и строим график
48     filtered_data = retention_history.query('payer == True').pivot_table(
49         index='dt', columns=columns, values=horizon - 1, aggfunc='mean'
50     )
51     filter_data(filtered_data, window).plot(grid=True, ax=ax3)
52     plt.xlabel('Дата привлечения')
53     plt.title(
54         'Динамика удержания платящих пользователей на {}-й день'.format(
55             horizon
56         )
57     )
58
59     # в четвёртой ячейке – динамика удержания неплатящих

```

```

60 ax4 = plt.subplot(2, 2, 4, sharey=ax3)
61 # фильтруем данные и строим график
62 filtered_data = retention_history.query('payer == False').pivot_table(
63     index='dt', columns=columns, values=horizon - 1, aggfunc='mean'
64 )
65 filter_data(filtered_data, window).plot(grid=True, ax=ax4)
66 plt.xlabel('Дата привлечения')
67 plt.title(
68     'Динамика удержания неплатящих пользователей на {}-й день'.format(
69         horizon
70     )
71 )
72
73 plt.tight_layout()
74 plt.show()
75

```

In [35]:



```

1 def plot_conversion(conversion, conversion_history, horizon, window=7):
2
3     # задаём размер сетки для графиков
4     plt.figure(figsize=(15, 5))
5
6     # исключаем размеры когорт
7     conversion = conversion.drop(columns=['cohort_size'])
8     # в таблице динамики оставляем только нужный лайфтайм
9     conversion_history = conversion_history.drop(columns=['cohort_size'])[
10         [horizon - 1]
11     ]
12
13     # первый график – кривые конверсии
14     ax1 = plt.subplot(1, 2, 1)
15     conversion.T.plot(grid=True, ax=ax1)
16     plt.legend()
17     plt.xlabel('Лайфтайм')
18     plt.title('Конверсия пользователей')
19
20     # второй график – динамика конверсии
21     ax2 = plt.subplot(1, 2, 2, sharey=ax1)
22     columns = [
23         # столбцами сводной таблицы станут все столбцы индекса, кроме даты
24         name for name in conversion_history.index.names if name not in ['dt']
25     ]
26     filtered_data = conversion_history.pivot_table(
27         index='dt', columns=columns, values=horizon - 1, aggfunc='mean'
28     )
29     filter_data(filtered_data, window).plot(grid=True, ax=ax2)
30     plt.xlabel('Дата привлечения')
31     plt.title('Динамика конверсии пользователей на {}-й день'.format(horizon))
32
33     plt.tight_layout()
34     plt.show()
35

```

```

1 def plot_ltv_roi(ltv, ltv_history, roi, roi_history, horizon, window=7):
2
3     # задаём сетку отрисовки графиков
4     plt.figure(figsize=(20, 10))
5
6     # из таблицы ltv исключаем размеры когорт
7     ltv = ltv.drop(columns=['cohort_size'])
8     # в таблице динамики ltv оставляем только нужный лайфтайм
9     ltv_history = ltv_history.drop(columns=['cohort_size'])[[horizon - 1]]
10
11     # стоимость привлечения запишем в отдельный фрейм
12     cac_history = roi_history[['cac']]
13
14     # из таблицы roi исключаем размеры когорт и cac
15     roi = roi.drop(columns=['cohort_size', 'cac'])
16     # в таблице динамики roi оставляем только нужный лайфтайм
17     roi_history = roi_history.drop(columns=['cohort_size', 'cac'])[
18         [horizon - 1]
19     ]
20
21     # первый график – кривые ltv
22     ax1 = plt.subplot(2, 3, 1)
23     ltv.T.plot(grid=True, ax=ax1)
24     plt.legend()
25     plt.xlabel('Лайфтайм')
26     plt.title('LTV')
27
28     # второй график – динамика ltv
29     ax2 = plt.subplot(2, 3, 2, sharey=ax1)
30     # столбцами сводной таблицы станут все столбцы индекса, кроме даты
31     columns = [name for name in ltv_history.index.names if name not in ['dt']]
32     filtered_data = ltv_history.pivot_table(
33         index='dt', columns=columns, values=horizon - 1, aggfunc='mean'
34     )
35     filter_data(filtered_data, window).plot(grid=True, ax=ax2)
36     plt.xlabel('Дата привлечения')
37     plt.title('Динамика LTV пользователей на {}-й день'.format(horizon))
38
39     # третий график – динамика cac
40     ax3 = plt.subplot(2, 3, 3, sharey=ax1)
41     # столбцами сводной таблицы станут все столбцы индекса, кроме даты
42     columns = [name for name in cac_history.index.names if name not in ['dt']]
43     filtered_data = cac_history.pivot_table(
44         index='dt', columns=columns, values='cac', aggfunc='mean'
45     )
46     filter_data(filtered_data, window).plot(grid=True, ax=ax3)
47     plt.xlabel('Дата привлечения')
48     plt.title('Динамика стоимости привлечения пользователей')
49
50     # четвёртый график – кривые roi
51     ax4 = plt.subplot(2, 3, 4)
52     roi.T.plot(grid=True, ax=ax4)
53     plt.axhline(y=1, color='red', linestyle='--', label='Уровень окупаемости')
54     plt.legend()
55     plt.xlabel('Лайфтайм')
56     plt.title('ROI')
57
58     # пятый график – динамика roi
59     ax5 = plt.subplot(2, 3, 5, sharey=ax4)

```

```

60 # столбцами сводной таблицы станут все столбцы индекса, кроме даты
61 columns = [name for name in roi_history.index.names if name not in ['dt']]
62 filtered_data = roi_history.pivot_table(
63     index='dt', columns=columns, values=horizon - 1, aggfunc='mean'
64 )
65 filter_data(filtered_data, window).plot(grid=True, ax=ax5)
66 plt.axhline(y=1, color='red', linestyle='--', label='Уровень окупаемости')
67 plt.xlabel('Дата привлечения')
68 plt.title('Динамика ROI пользователей на {}-й день'.format(horizon))
69
70 plt.tight_layout()
71 plt.show()
72

```

4 Исследовательский анализ данных

4.1 Пользовательские профили.

Минимальную и максимальную даты привлечения пользователей

Вызовем функцию `get_profiles()`, передав ей данные о посещениях, покупках, и тратах на рекламу. Укажем, что данных о событиях у нас нет.

In [37]:

»

```

1 events = None
2 profiles = get_profiles(visits, orders, events, costs)
3 display(profiles.head(5))

```

	user_id	first_ts	channel	device	region	dt	month	payer	acquisition_cost
0	599326	2019-05-07 20:58:57	FaceBoom	Mac	United States	2019-05-07	2019-05-01	True	1.09
1	4919697	2019-07-09 12:46:07	FaceBoom	iPhone	United States	2019-07-09	2019-07-01	False	1.11
2	6085896	2019-10-01 09:58:33	organic	iPhone	France	2019-10-01	2019-10-01	False	0.00
3	22593348	2019-08-22 21:35:48	AdNonSense	PC	Germany	2019-08-22	2019-08-01	False	0.99
4	31989216	2019-10-02 00:07:44	YRabbit	iPhone	United States	2019-10-02	2019-10-01	False	0.23

4.1.1 Минимальные и максимальные даты

In [38]:

```
1 print(min(profiles['first_ts'].dt.date), "- минимальная дата первой сессии пользователей")
2 print(max(profiles['first_ts'].dt.date), "- максимальная дата первой сессии пользователей")
```

2019-05-01 - минимальная дата первой сессии пользователей
2019-10-27 - максимальная дата первой сессии пользователей

4.1.2 Выясним, из каких стран пользователи приходят в приложение

In [39]:

```
1 profiles[['region', 'user_id']].groupby(by='region').count().sort_values(by='user_id', ascending=False)
```

Out[39]:

user_id	
region	
United States	100002
UK	17575
France	17450
Germany	14981

4.1.3 На какую страну приходится больше всего платящих пользователей

In [76]:

```
1 profiles[['region', 'payer']][profiles['payer'] == True]\
2 .groupby(by='region').count().sort_values(by='payer', ascending=False)
```

Out[76]:

payer	
region	
United States	6902
UK	700
France	663
Germany	616

Рейтинги общего количества пользователей и количества платящих пользователей по странам совпадают

4.1.4 Доля платящих пользователей по странам

In [77]:



```
1 # all users
2 all_users = profiles[['region', 'user_id']].groupby(by='region').count().sort_values(by='region')
3 # payers
4 payers = profiles[['region', 'payer']][profiles['payer'] == True]\
5     .groupby(by='region').count().sort_values(by='payer', ascending=False)
6 # объединим две таблицы
7 percentage_of_buyers = all_users.merge(payers, on='region')
8 # добавим колонку долей
9 percentage_of_buyers['percentage'] = (percentage_of_buyers['payer'] / percentage_of_buyers['user_id'])
10 display(percentage_of_buyers)
```

	user_id	payer	percentage
region			
United States	100002	6902	6.90
UK	17575	700	3.98
France	17450	663	3.80
Germany	14981	616	4.11

United States - безусловный лидер.

И по количеству пользователей, и доле платящих из них.

Germany занимает последнее место по количеству пользователей, но на **втором** месте по доле платящих.

хорошая точка роста бизнеса

4.1.5 Какими устройствами пользуются клиенты

какие устройства предпочитают платящие пользователи

Построим таблицу, отражающую количество пользователей и долю платящих для каждого устройства

In [78]:



```
1 # all devices
2 all_users_d = profiles[['device', 'user_id']].groupby(by='device').count().sort_values(t
3 # payers
4 devices = profiles[['device', 'payer']][profiles['payer'] == True]\
5     .groupby(by='device').count().sort_values(by='payer', ascending=False)
6 # объединим две таблицы
7 percentage_of_devices = all_users_d.merge(devices, on='device')
8 # добавим колонку долей
9 percentage_of_devices['percentage'] = (percentage_of_devices['payer'] / percentage_of_c
10 display(percentage_of_devices)
```

	user_id	payer	percentage
device			
iPhone	54479	3382	6.21
Android	35032	2050	5.85
PC	30455	1537	5.05
Mac	30042	1912	6.36

Пользователи **iPhone** и **Mac** подтверждают устоявшееся мнение о них)))

4.1.6 Рекламные источники привлечения

- каналы, из которых пришло больше всего платящих пользователей
- таблица, отражающая количество пользователей и долю платящих для каждого канала привлечения

In [79]:



```
1 # channels
2 all_users_ch = profiles[['channel', 'user_id']].groupby(by='channel').count().sort_values
3 # payers
4 channels = profiles[['channel', 'payer']][profiles['payer'] == True]\
5     .groupby(by='channel').count().sort_values(by='payer', ascending=False)
6 # объединим две таблицы
7 percentage_of_channels = all_users_ch.merge(channels, on='channel')
8 # добавим колонку долей
9 percentage_of_channels['percentage'] = (percentage_of_channels['payer'] / percentage_of
10 display(percentage_of_channels.sort_values('payer', ascending=False)))
```

	user_id	payer	percentage
channel			
FaceBoom	29144	3557	12.20
TipTop	19561	1878	9.60
organic	56439	1160	2.06
WahooNetBanner	8553	453	5.30
AdNonSense	3880	440	11.34
RocketSuperAds	4448	352	7.91
LeapBob	8553	262	3.06
OppleCreativeMedia	8605	233	2.71
lambdaMediaAds	2149	225	10.47
YRabbit	4312	165	3.83
MediaTornado	4364	156	3.57

FaceBoom - рекорсмен конверсии (12,20%). Это позволило каналу FaceBoom опередить органический канал по количеству платящих пользователей.

Перспективные каналы для инвестиций - **AdNonSense** и **lambdaMediaAds**. Эти каналы характеризуются высоким коэффициентов конверсии.

4.1.7 Выводы по исследованию профилей пользователей

- 2019-05-01 - минимальная дата первой сессии пользователей
- 2019-10-27 - максимальная дата первой сессии пользователей
- United States - страна-лидер по количеству и качеству пользователей
- Germany занимает последнее место по количеству пользователей, но на втором месте по доле платящих.
- *хорошая точка роста бизнеса**
- Пользователи **iPhone** и **Mac** лидируют
- **FaceBoom** - рекорсмен конверсии (12,20%). Это позволило каналу FaceBoom опередить органический канал по количеству платящих пользователей.
- Перспективные каналы для инвестиций - **AdNonSense** и **lambdaMediaAds**. Эти каналы характеризуются высоким коэффициентов конверсии.

4.2 Маркетинг

4.2.1 Общая сумма расходов на маркетинг.

Выясним, как траты распределены по источникам.

Таблица **costs**

4.2.2 Всего расходов на рекламу

In [44]:



```
1 print(int(costs['costs'].sum()), "общая сумма расходов на рекламу")
```

105497 общая сумма расходов на рекламу

4.2.3 Расходы на рекламу по каналам

In [45]:



```
1 costs[['channel', 'costs']].groupby(by='channel').sum().sort_values(by='costs', ascending=True)
```

Out[45]:

costs	
channel	
TipTop	54,751.30
FaceBoom	32,445.60
WahooNetBanner	5,151.00
AdNonSense	3,911.25
OppleCreativeMedia	2,151.25
RocketSuperAds	1,833.00
LeapBob	1,797.60
lambdaMediaAds	1,557.60
MediaTornado	954.48
YRabbit	944.22

Рекламный бизнес в плену ожиготажных тенденций.

Самые большие вливания в новомодный канал **TipTop**

4.2.4 Изменение затрат на рекламу в течении времени

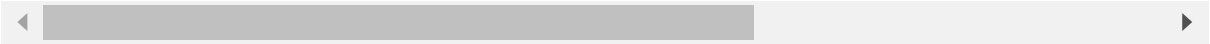
In [46]:

▶

```
1 costs.pivot_table(index = 'dt', columns = 'channel', values = 'costs', aggfunc = 'sum')
2
```

Out[46]:

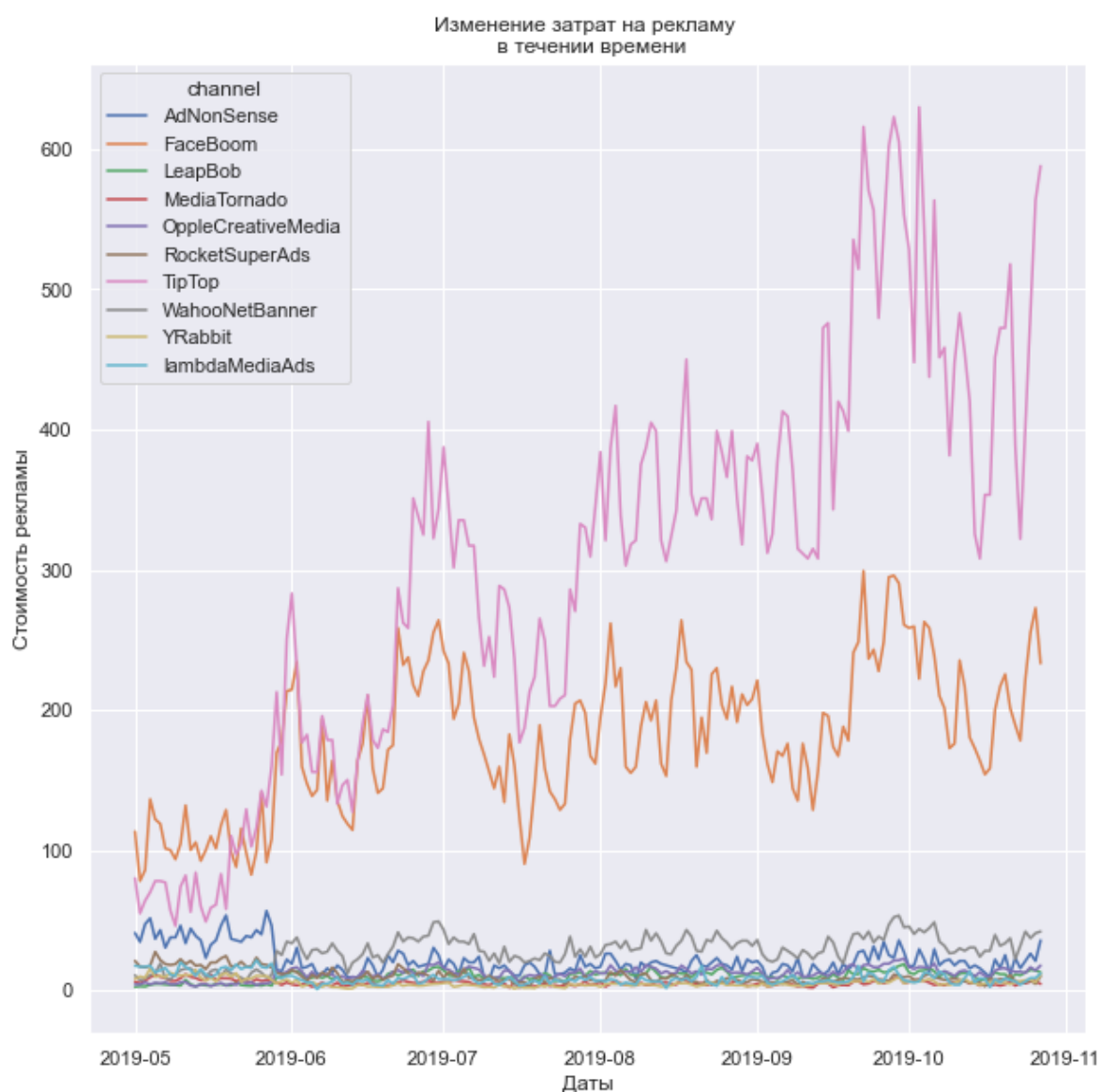
channel	AdNonSense	FaceBoom	LeapBob	MediaTornado	OppleCreativeMedia	RocketSuperAc
dt						
2019-05-01	40.95	113.30	2.52	6.24	4.25	21.0
2019-05-02	34.65	78.10	2.94	5.04	4.25	16.9
2019-05-03	47.25	85.80	2.73	6.96	5.75	16.3
2019-05-04	51.45	136.40	3.99	9.36	4.25	17.5
2019-05-05	36.75	122.10	4.62	11.04	5.50	27.5

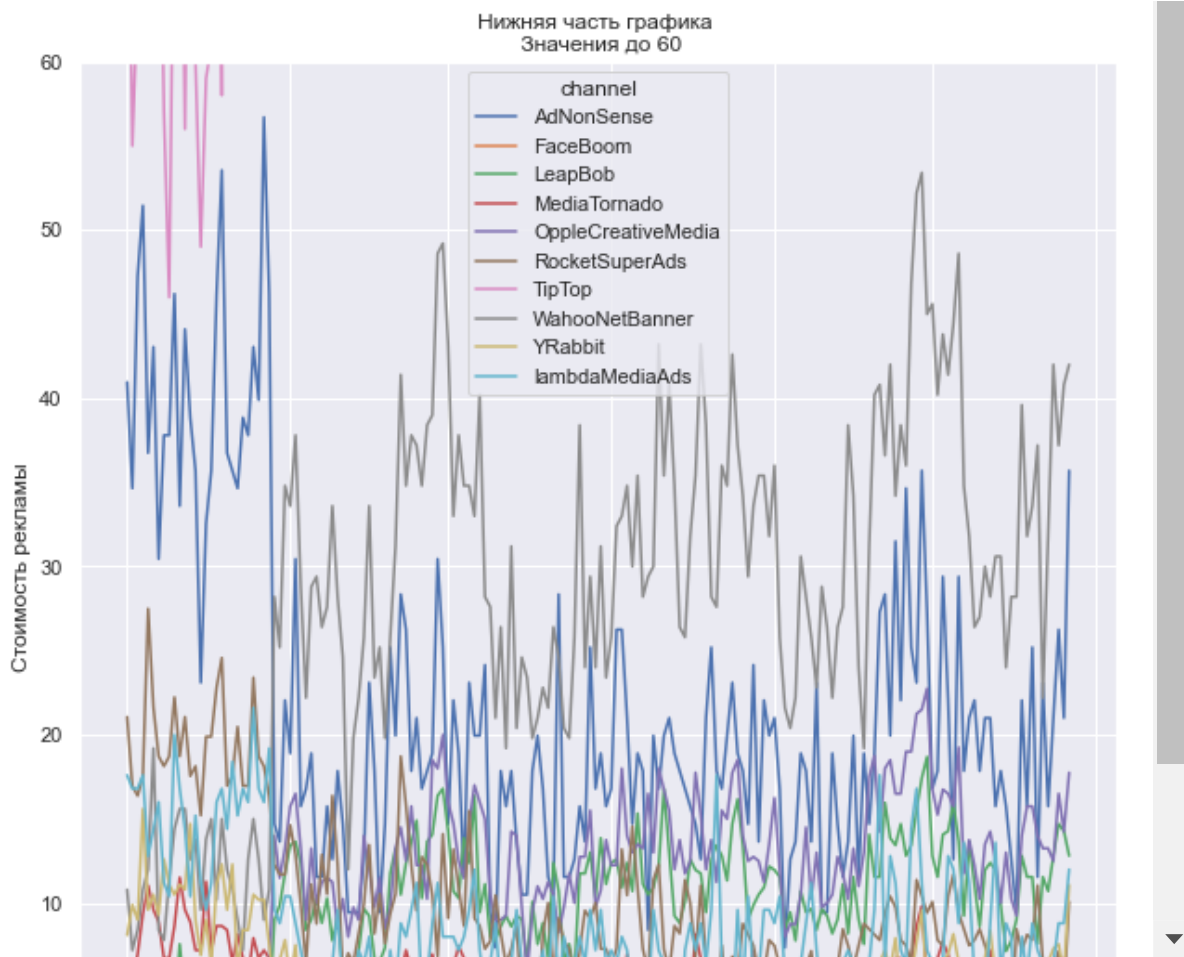


In [47]:



```
1 sns.set(rc = {'figure.figsize':(10,10)})
2 ax=fig.add_subplot(2,2,1)
3 ax2=fig.add_subplot(2,2,2)
4 # полноформатный график
5 ax = costs.pivot_table(index = 'dt', columns = 'channel', values = 'costs', aggfunc = 'sum')
6 ax.set_ylabel("Стоимость рекламы");
7 ax.set_xlabel("Даты");
8 ax.set_title("Изменение затрат на рекламу \n в течении времени");
9 # график усечённый по оси Y
10 ax2 = costs.pivot_table(index = 'dt', columns = 'channel', values = 'costs', aggfunc = 'sum')
11 ax2.set_ylim(0, 60)
12 ax2.set_ylabel("Стоимость рекламы");
13 ax2.set_xlabel("Даты");
14 ax2.set_title("Нижняя часть графика \n Значения до 60");
15
16 plt.show()
```





По графикам видим: Вливания в каналы - лидеры по оплатам **TipTop** и **FaceBoom** идут по **нарастающей** в течении года.

Канал **AdNonSense** начали достаточно активно, но потом значительно снизили. Опрданно ли это?...

4.2.5 Сколько в среднем стоило привлечение одного пользователя из каждого источника

- средний CAC на одного пользователя для всего проекта
- средний CAC на одного пользователя для каждого источника трафика

Группируем данные по источникам.

Суммируем суммы платежей.

Суммируем количество пользователей.

Делим платежи на пользователей.

In [48]:



```
1 # пользователи и каналы
2 users_and_channels = profiles[['channel', 'user_id']].groupby(by='channel').count()
3 users_and_channels
```

Out[48]:

	user_id
channel	
AdNonSense	3880
FaceBoom	29144
LeapBob	8553
MediaTornado	4364
OppleCreativeMedia	8605
RocketSuperAds	4448
TipTop	19561
WahooNetBanner	8553
YRabbit	4312
lambdaMediaAds	2149
organic	56439

In [49]:



```
1 # каналы и оплаты
2 channels_and_costs = costs[['channel', 'costs']].groupby(by='channel').sum()
3 channels_and_costs
```

Out[49]:

	costs
channel	
AdNonSense	3,911.25
FaceBoom	32,445.60
LeapBob	1,797.60
MediaTornado	954.48
OpplCreativeMedia	2,151.25
RocketSuperAds	1,833.00
TipTop	54,751.30
WahooNetBanner	5,151.00
YRabbit	944.22
lambdaMediaAds	1,557.60

In [50]:



```
1 # объединим результаты
2 # только по платным каналам
3 cac_by_channels = channels_and_costs.merge(users_and_channels, on='channel')
4 cac_by_channels
```

Out[50]:

	costs	user_id
channel		
AdNonSense	3,911.25	3880
FaceBoom	32,445.60	29144
LeapBob	1,797.60	8553
MediaTornado	954.48	4364
OpplCreativeMedia	2,151.25	8605
RocketSuperAds	1,833.00	4448
TipTop	54,751.30	19561
WahooNetBanner	5,151.00	8553
YRabbit	944.22	4312
lambdaMediaAds	1,557.60	2149

4.2.5.1 Средний САС в целом

In [51]:

```
1 print("Стоимость привлечения одного пользователя в целом по проекту:")
2 print()
3 print((sum(costs['costs']) / len(visits['user_id'].value_counts())) , "по всему проекту")
4 print()
5 print((sum(cac_by_channels['costs']) / sum(cac_by_channels['user_id']))) , "за вычетом г
```

Стоимость привлечения одного пользователя в целом по проекту:

0.7032778251826577 по всему проекту, включая органический трафик

1.127481323942759 за вычетом пользователей органического трафика

4.2.5.2 Средний CAC по каналам

In [52]:

```
1 # вычисляем CAC и сортируем по нему
2 cac_by_channels['cac'] = cac_by_channels['costs'] / cac_by_channels['user_id']
3 display(cac_by_channels.sort_values(by='cac', ascending=False))
```

	costs	user_id	cac
channel			
TipTop	54,751.30	19561	2.80
FaceBoom	32,445.60	29144	1.11
AdNonSense	3,911.25	3880	1.01
lambdaMediaAds	1,557.60	2149	0.72
WahooNetBanner	5,151.00	8553	0.60
RocketSuperAds	1,833.00	4448	0.41
OpplCreativeMedia	2,151.25	8605	0.25
YRabbit	944.22	4312	0.22
MediaTornado	954.48	4364	0.22
LeapBob	1,797.60	8553	0.21

Самые дорогие привлечённые пользователи из **TipTop**

Вдвое дешевле из **FaceBoom** и **AdNonSense**

В группе дешевых каналов **4** участника:

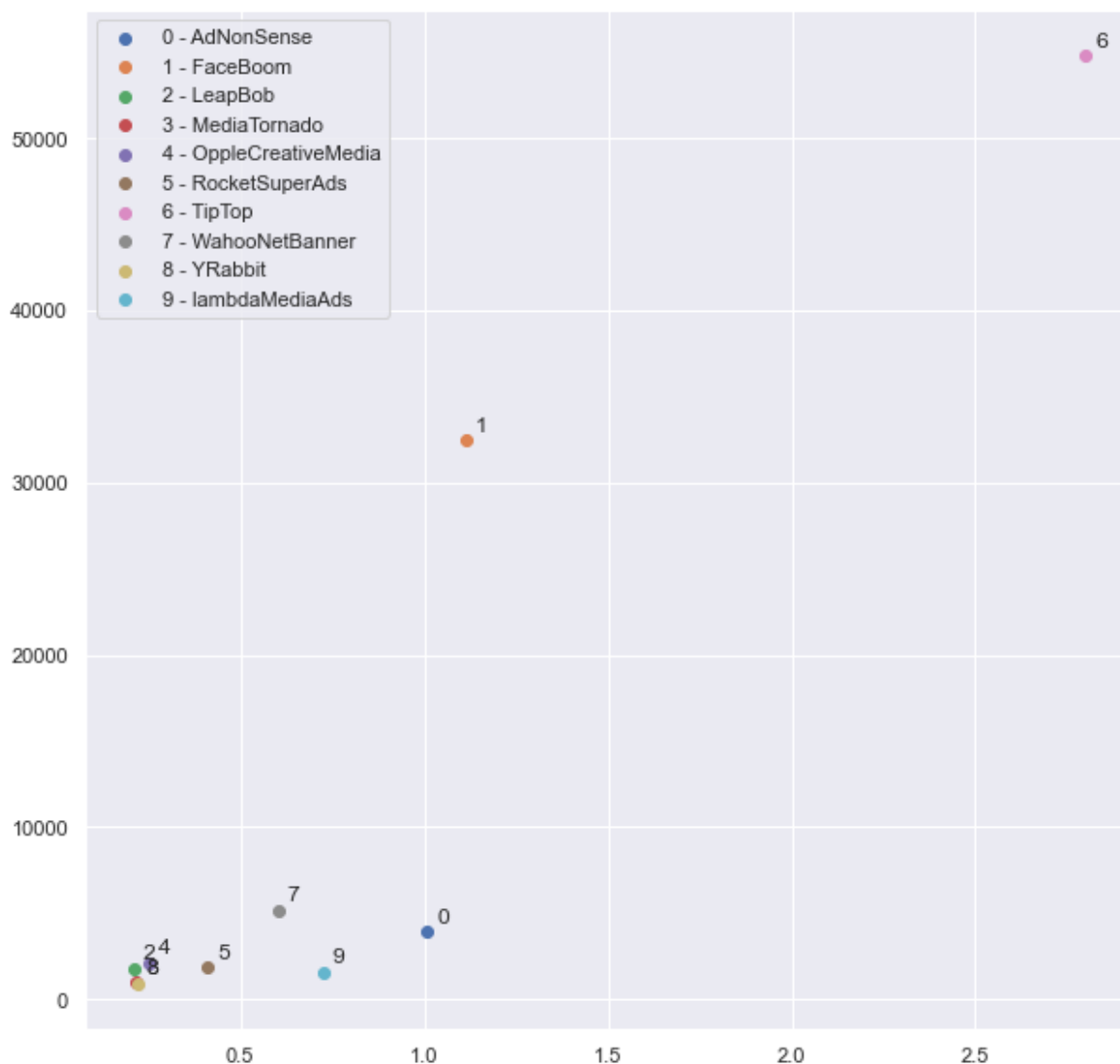
- OpplCreativeMedia
- YRabbit
- MediaTornado
- LeapBob

4.2.5.3 Графически соотношение CAC и вливаний в каналы

In [80]:



```
1 # Код ревьюера:
2 fig, ax = plt.subplots()
3
4 for i, txt in enumerate(cac_by_channels.index.values):
5     legend = (str(i) + ' - ' + str(cac_by_channels.reset_index()['channel'][i]))
6     ax.annotate(i, (cac_by_channels['cac'][i], cac_by_channels['costs'][i]),
7                xytext=(5,5), textcoords='offset points')
8     plt.scatter(cac_by_channels['cac'][i], cac_by_channels['costs'][i], label = legend)
9
10 plt.grid(True)
11 plt.legend(loc="upper left")
12
13 plt.show()
```



Здесь видно, что в каналы с **низкой** стоимостью привлечения, **мало** денежных вливаний.

А в каналы с **большой** стоимостью привлечения пользователя - самые большие денежные транши.

4.2.6 Может быть, надо сделать наоборот?

4.2.6.1 Динамика изменения САС для каждого источника во времени

Группируем профили пользователей по-месячно
(за вычетом органического трафика)

Группируем оплаты рекламы по-месячно

In [57]:

```
1 profiles.pivot_table(index='dt', columns='channel', values='acquisition_cost', aggfunc=
2 plt.ylabel('CAC')
3 plt.xlabel('Дата')
4 plt.title('Динамика CAC')
5 plt.show()
```



4.3 Окупаемость рекламы

Сегодня - 01.11.2019 г.

Срок окупаемости - ДВЕ недели.

Сделаем ДВА анализа:

- с органикой
- без органики

4.3.1 Установим момент и горизонт анализа данных.

In [60]:

```
1 observation_date = datetime(2019, 11, 1).date() # момент анализа
2 horizon_days = 14 # горизонт анализа
```

Используя графики LTV, ROI и CAC

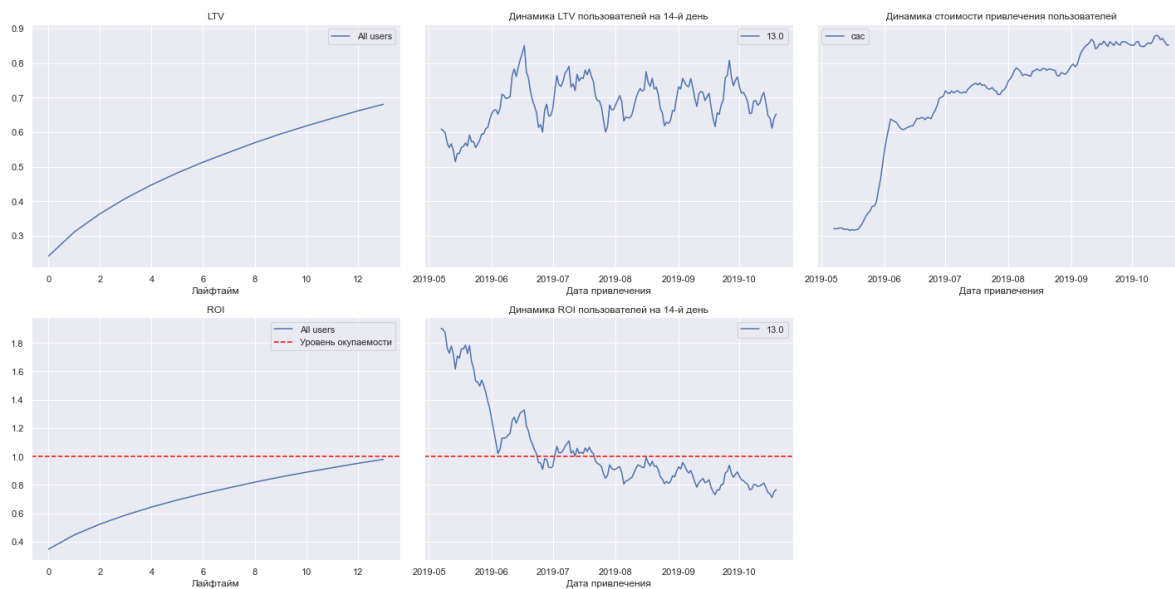
4.3.2 Считаем бизнес-показатели

Для начала оценим общую ситуацию — посмотрим на окупаемость рекламы. Рассчитаем и визуализируем LTV и ROI, вызвав функции `get_ltv()` и `plot_ltv_roi()`.

4.3.2.1 Включая органический трафик

In [61]:

```
1 # считаем LTV и ROI
2 ltv_raw, ltv_grouped, ltv_history, roi_grouped, roi_history = get_ltv(
3     profiles, orders, observation_date, horizon_days
4 )
5
6 # строим графики
7 plot_ltv_roi(ltv_grouped, ltv_history, roi_grouped, roi_history, horizon_days)
```

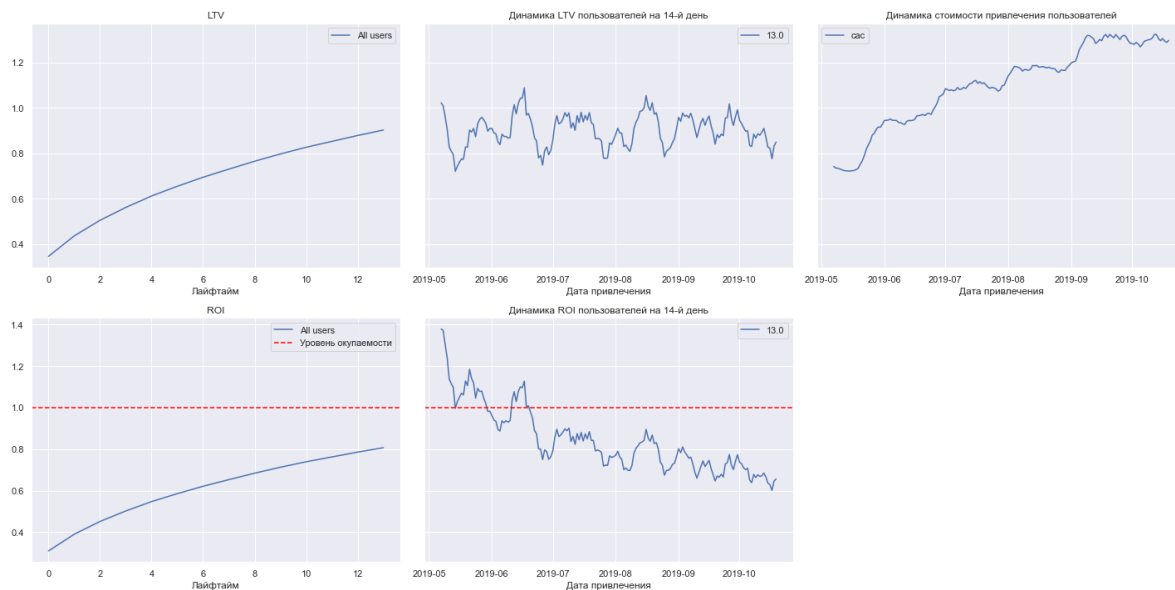


4.3.2.2 Исключая органический трафик

```

1 # profiles_not_organic
2 profiles_not_organic = profiles.query('channel != "organic"')
3
4 # считаем LTV и ROI
5 ltv_raw, ltv_grouped, ltv_history, roi_grouped, roi_history = get_ltv(
6     profiles_not_organic, orders, observation_date, horizon_days
7 )
8
9 # строим графики
10 plot_ltv_roi(ltv_grouped, ltv_history, roi_grouped, roi_history, horizon_days)

```



Органический трафик сильно **улучшает** показатели.

Поэтому для чистоты исследования **исключим** данные органического трафика в данных для анализа.

4.3.3 По графикам можно сделать такие выводы:

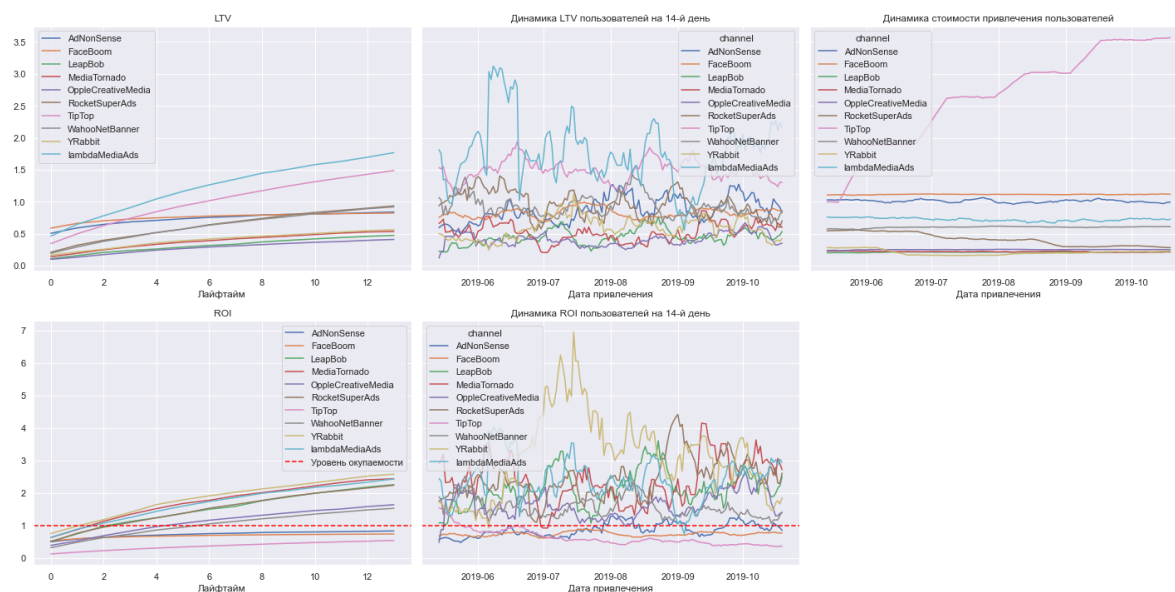
- Реклама не окупается. ROI в конце недели — чуть выше 80%.
- САС стабильно увеличивается. Значит, эффективность рекламных кампаний падает.
- LTV +/- стабилен. Значит, дело не в ухудшении качества пользователей.
- Чтобы разобраться в причинах, пройдем по всем доступным характеристикам пользователей — рекламным каналам, стране и устройству первого посещения.

4.3.4 Начнём с разбивки по рекламным каналам: передадим параметру `dimensions` столбец `channel`.

In [63]:



```
1 # смотрим окупаемость с разбивкой по источникам привлечения
2
3 dimensions = ['channel']
4
5 ltv_raw, ltv_grouped, ltv_history, roi_grouped, roi_history = get_ltv(
6     profiles_not_organic, orders, observation_date, horizon_days, dimensions=dimensions
7 )
8
9 plot_ltv_roi(
10     ltv_grouped, ltv_history, roi_grouped, roi_history, horizon_days, window=14
11 )
```



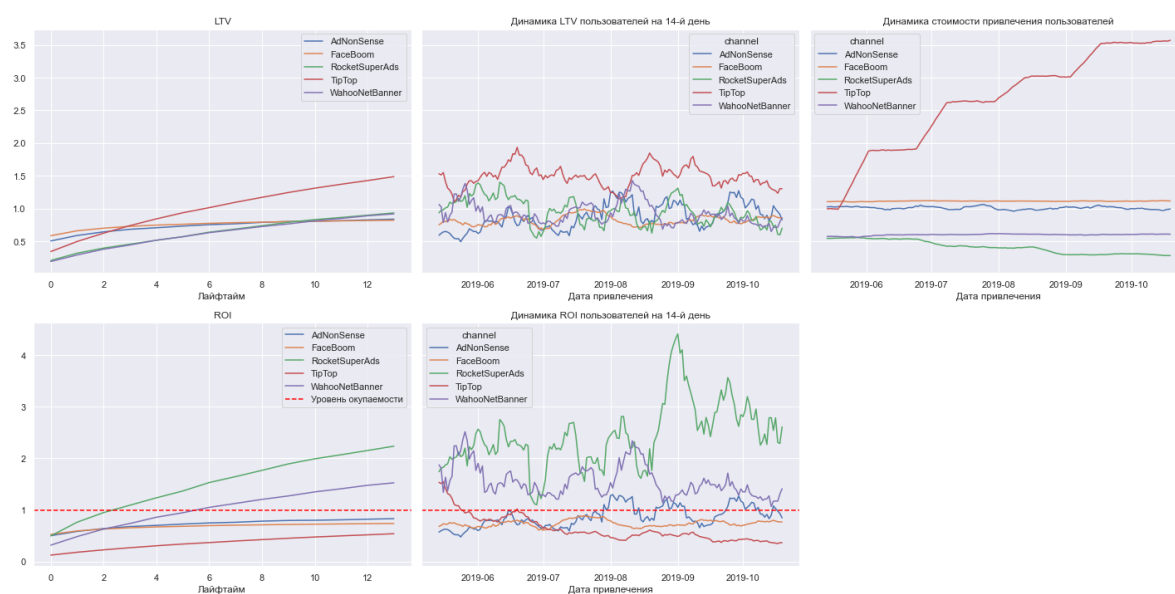
каша-малаша

Выведем данные двумя группами. По 5 штук.

In [64]:



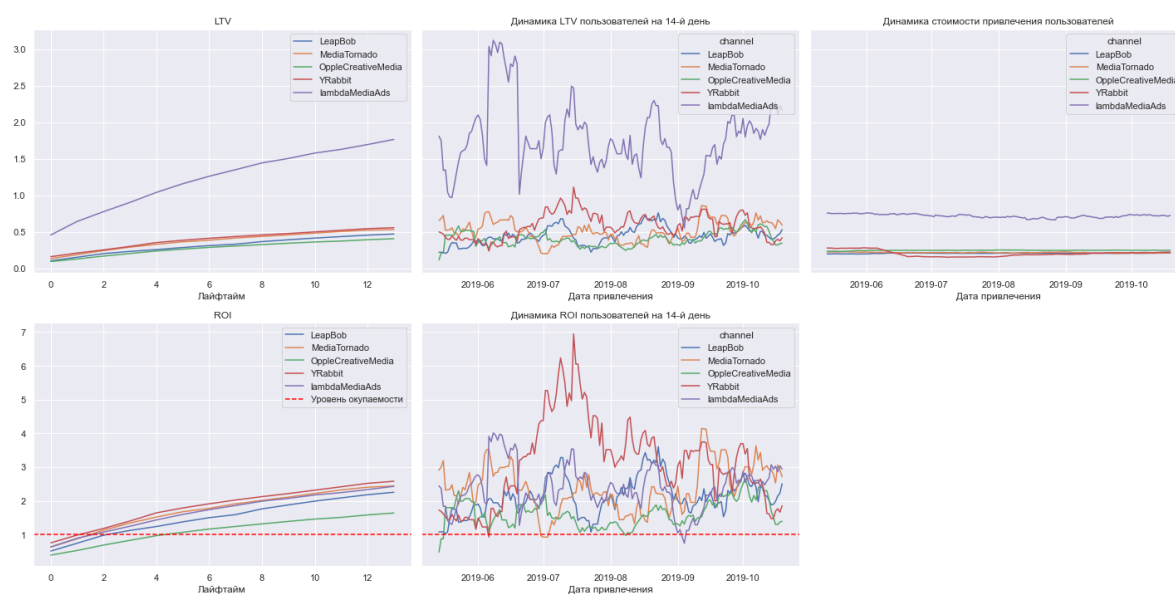
```
1 # смотрим окупаемость с разбивкой по источникам привлечения
2
3 dimensions = ['channel']
4
5 # каналы группа раз
6 channels_1 = ['FaceBoom', 'TipTop', 'WahooNetBanner', 'AdNonSense', 'RocketSuperAds']
7
8
9 ltv_raw, ltv_grouped, ltv_history, roi_grouped, roi_history = get_ltv(
10     profiles_not_organic.query('channel in @channels_1'), orders, observation_date, ho
11 )
12
13 plot_ltv_roi(
14     ltv_grouped, ltv_history, roi_grouped, roi_history, horizon_days, window=14
15 )
```



In [65]:



```
1 # смотрим окупаемость с разбивкой по источникам привлечения
2
3 dimensions = ['channel']
4
5 # каналы группа раз
6 channels_2 = ['LeapBob', 'OpplCreativeMedia', 'lambdaMediaAds', 'YRabbit', 'MediaTornado']
7
8
9 ltv_raw, ltv_grouped, ltv_history, roi_grouped, roi_history = get_ltv(
10     profiles_not_organic.query('channel in @channels_2'), orders, observation_date, horizon_days
11 )
12
13 plot_ltv_roi(
14     ltv_grouped, ltv_history, roi_grouped, roi_history, horizon_days, window=14
15 )
```

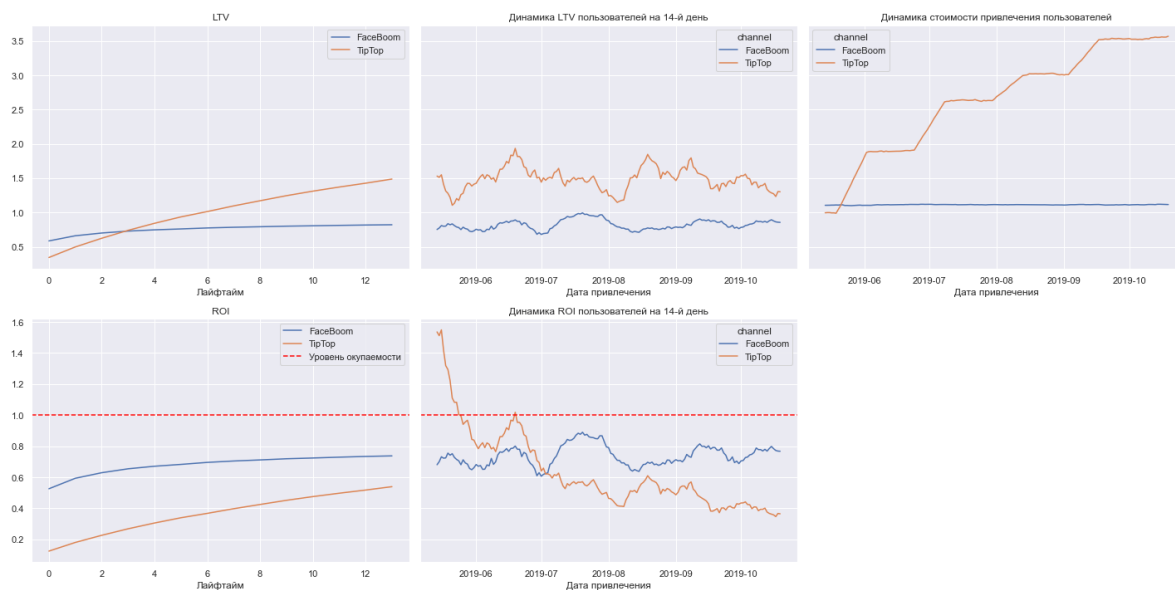


а теперь - два самых хитовых канала
FaceBoom и TipTop

In [66]:



```
1 # смотрим окупаемость с разбивкой по источникам привлечения
2
3 dimensions = ['channel']
4
5 # каналы группа раз
6 channels_3 = ['FaceBoom', 'TipTop']
7
8
9 ltv_raw, ltv_grouped, ltv_history, roi_grouped, roi_history = get_ltv(
10     profiles_not_organic.query('channel in @channels_3'), orders, observation_date, h
11 )
12
13 plot_ltv_roi(
14     ltv_grouped, ltv_history, roi_grouped, roi_history, horizon_days, window=14
15 )
```



4.3.4.1 Вывод промежуточный по рекламным каналам:

- TipTop
 - самый большой объём вливаний
 - хорошая конверсия в покупателей
 - окупаемость **низкая**, около 0,5
 - окупаемость **снижается** в течении отчётного периода
 - стоимость привлечения **растёт** с огромной скоростью
 - LTV при этом растёт и стабилен на 14й день
- FaceBoom
 - стабильные показатели жизни пользователей и стоимости их привлечения
 - стабильные показатели по времени
 - окупаемость **низкая** до 0,8
- WahooNetBanner
 - стабильная стоимость привлечения пользователей
 - окупаемость **более 1,5**
 - с течением времени окупаемость нестабильна, но всегда **положительная**
- AdNonSense
 - стабильная стоимость привлечения пользователей
 - окупаемость **отрицательная**

- с течением времени окупаемость нестабильна, то положительная, то отрицательная
- OpplCreativeMedia
 - стабильная стоимость привлечения пользователей
 - окупаемость **более 1,5**
 - с течением времени окупаемость нестабильна, но всегда **положительная**
- RocketSuperAds
 - стоимость привлечения пользователей ** снижается** с течением времени
 - окупаемость **более 2**
 - с течением времени окупаемость нестабильна, но всегда **положительная**
- LeapBob
 - стоимость привлечения пользователей ** снижается** с течением времени
 - окупаемость **более 2**
 - с течением времени окупаемость нестабильна, но всегда **положительная**
- lambdaMediaAds
 - стоимость привлечения пользователей ** снижается** с течением времени
 - окупаемость **более 2**
 - с течением времени окупаемость нестабильна, но всегда **положительная**
- MediaTornado
 - стоимость привлечения пользователей ** снижается** с течением времени
 - окупаемость **более 2**
 - с течением времени окупаемость нестабильна, но всегда **положительная**
- YRabbit
 - стоимость привлечения пользователей ** снижается** с течением времени
 - окупаемость **более 2**
 - с течением времени окупаемость нестабильна, но всегда **положительная**

4.3.5 По странам: передадим параметру dimensions столбец region.

In [67]:



```
1 # смотрим окупаемость с разбивкой по источникам привлечения
2
3 dimensions = ['region']
4
5 ltv_raw, ltv_grouped, ltv_history, roi_grouped, roi_history = get_ltv(
6     profiles_not_organic, orders, observation_date, horizon_days, dimensions=dimensions
7 )
8
9 plot_ltv_roi(
10     ltv_grouped, ltv_history, roi_grouped, roi_history, horizon_days, window=14
11 )
```



4.3.5.1 United States - портит всю картину

- рост расходов и стоимости привлечения со временем
- падение окупаемости со временем
- отрицательная окупаемость

4.3.6 Перейдём к устройствам.

In [68]:



```
1 # смотрим окупаемость с разбивкой по источникам привлечения
2
3 dimensions = ['device']
4
5 ltv_raw, ltv_grouped, ltv_history, roi_grouped, roi_history = get_ltv(
6     profiles_not_organic, orders, observation_date, horizon_days, dimensions=dimensions
7 )
8
9 plot_ltv_roi(
10     ltv_grouped, ltv_history, roi_grouped, roi_history, horizon_days, window=14
11 )
```



4.3.6.1 PC - лидер по показателям

Android идёт за ними, но не достигает окупаемости всё равно

- iPhone и Mac идут вровень, но **очень плохо**, что очень странно, ведь считается, что эта аудитория наиболее **платёжеспособна и лояльна к платным услугам и товарам**

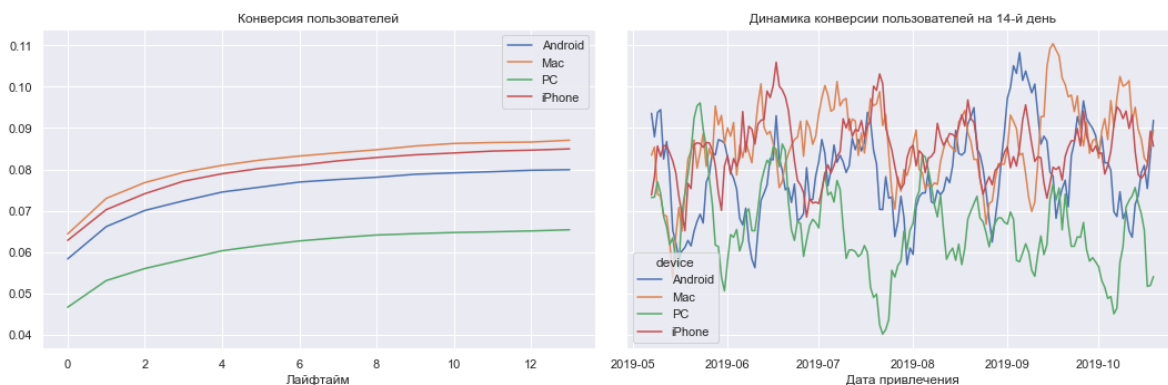
4.3.7 Графики конверсии и удержания с разбивкой по устройствам, странам, рекламным

каналам.

Посчитаем и визуализируем конверсию, вызвав функции `get_conversion()` и `plot_conversion()`.

In [69]:

```
1 # смотрим конверсию с разбивкой по устройствам
2
3 conversion_raw, conversion_grouped, conversion_history = get_conversion(
4     profiles_not_organic, orders, observation_date, horizon_days, dimensions=dimensions
5 )
6
7 plot_conversion(conversion_grouped, conversion_history, horizon_days)
```



4.3.7.1 Судя по графикам, пользователи iPhone и Mac конвертируются очень хорошо, причём постоянно.

Видимо, дело в удержании. .

4.3.8 Вывод по окупаемости рекламы

- **страны** - проблемная страна **United States**
- **каналы** - проблемные каналы
 - ◦ TipTop
 - ◦ FaceBoom

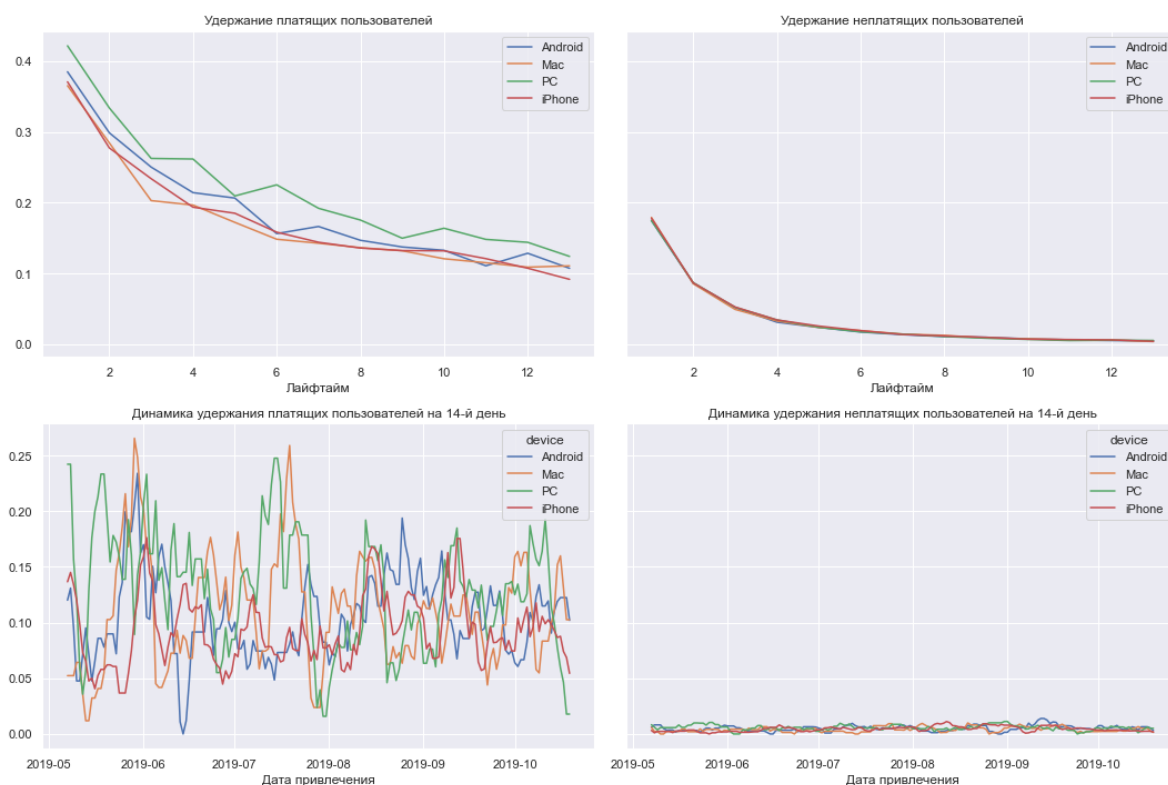
Вызовем функции `get_retention()` и `plot_retention()`, чтобы рассчитать и отразить на графиках этот показатель.

4.3.9 Удержание по устройствам

In [70]:



```
1 # смотрим удержание с разбивкой по устройствам
2
3 dimensions = ['device']
4
5 retention_raw, retention_grouped, retention_history = get_retention(
6     profiles_not_organic, visits, observation_date, horizon_days, dimensions=dimensions
7 )
8
9 plot_retention(retention_grouped, retention_history, horizon_days)
```



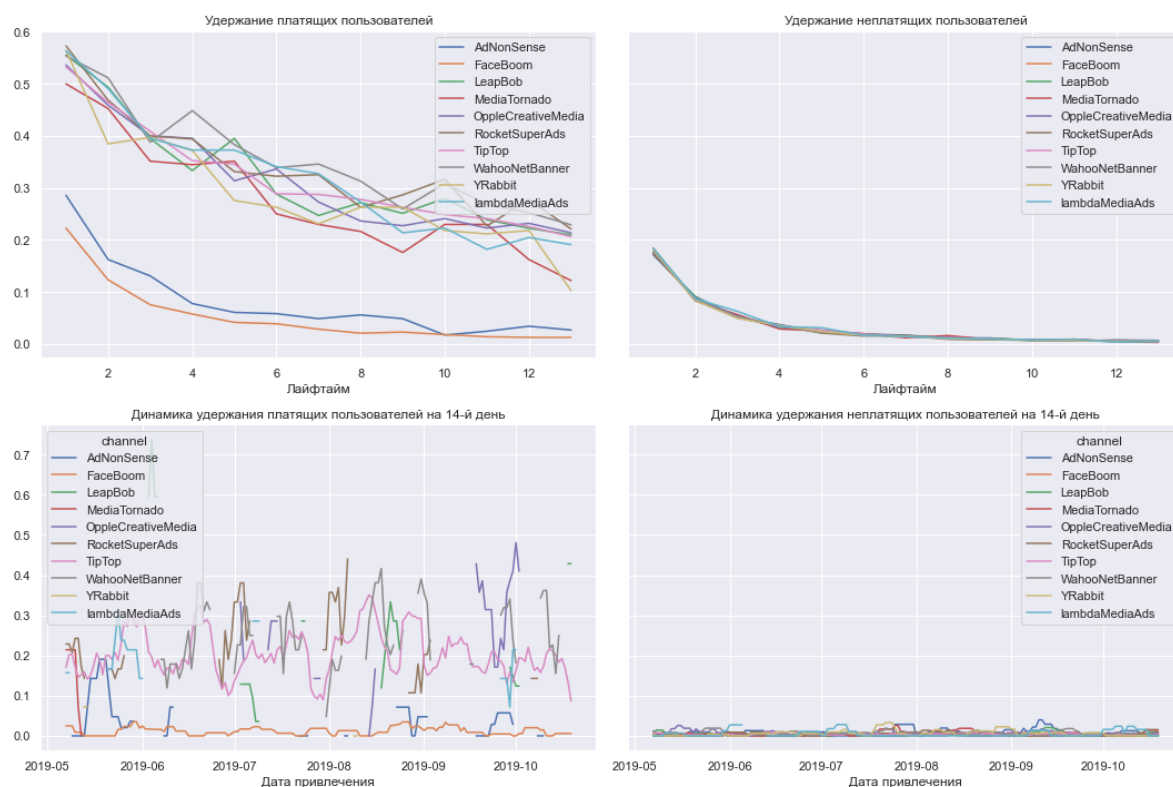
4.3.9.1 Удержание пользователей всех платформ примерно схоже

4.3.10 Удержание по рекламным каналам

```

1 # смотрим удержание с разбивкой по рекламным каналам
2
3 dimensions = ['channel']
4
5 retention_raw, retention_grouped, retention_history = get_retention(
6     profiles_not_organic, visits, observation_date, horizon_days, dimensions=dimensions
7 )
8
9 plot_retention(retention_grouped, retention_history, horizon_days)

```



4.3.10.1 Проблемное удержание пользователей

- AdNonSense
- FaceBoom

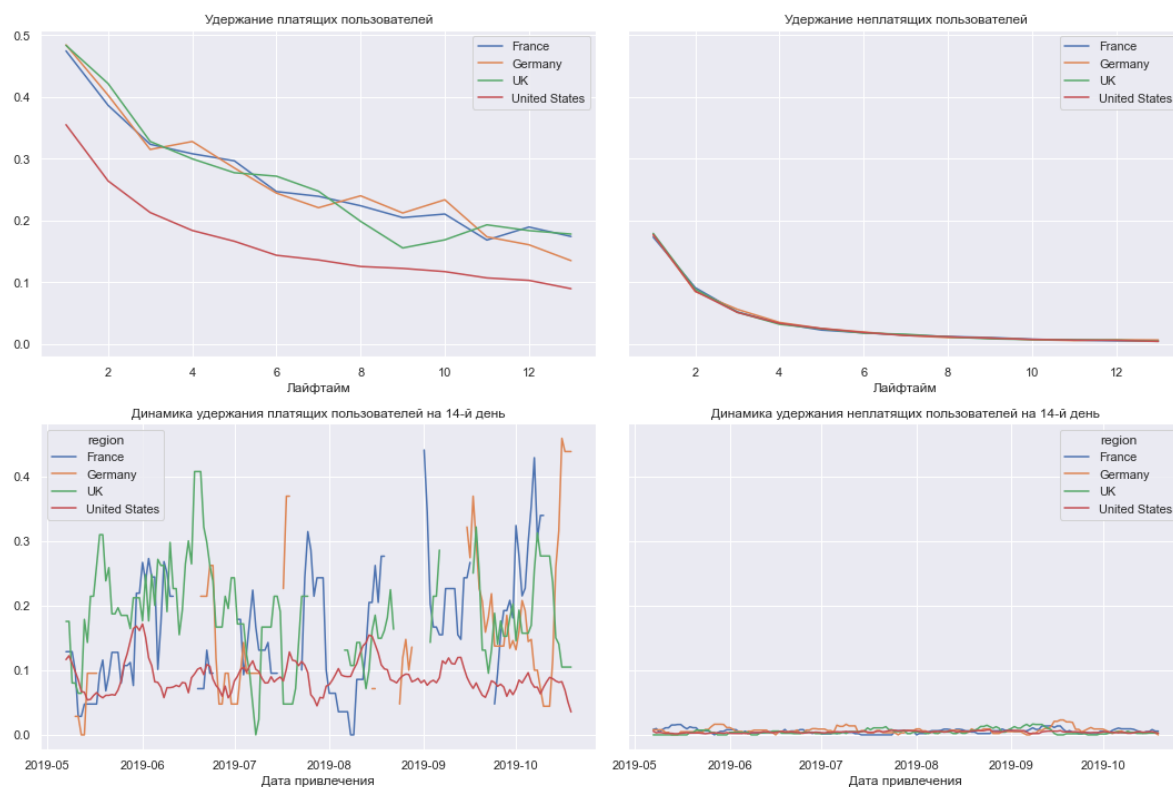
удержание канала TipTop хорошее

4.3.11 Удержание по странам

In [72]:



```
1 # смотрим удержание с разбивкой по странам
2
3 dimensions = ['region']
4
5 retention_raw, retention_grouped, retention_history = get_retention(
6     profiles_not_organic, visits, observation_date, horizon_days, dimensions=dimensions
7 )
8
9 plot_retention(retention_grouped, retention_history, horizon_days)
```



4.3.11.1 United States - отстают от других стран

4.3.12 итог

4.3.12.1 Реклама в целом не окупается

4.3.12.2 Провальные моменты

- Страна - аутсайдер - **United States**
- Провальные рекламные каналы - **FaceBoom** и **TipTop**

4.3.12.3 Возможные причины

- рекламные кампании построены **без учёта особенностей** пользователей страны **United States** и каналов **FaceBoom** и **TipTop**

4.3.12.4 Стратегия выхода на окупаемость

- **отключить** рекламу в **United States** и в каналах **FaceBoom** и **TipTop**
- изучить особенности пользователей в **United States** и в каналах **FaceBoom** и **TipTop**
- (пригласить к сотрудничеству специалистов по **United States** и каналам **FaceBoom** и **TipTop**)
- запустить в тестовом режиме обновлённые рекламные кампании в **United States** и в каналах **FaceBoom** и **TipTop**
- проверить результаты

5 Вывод

- Реклама в целом не окупается

5.1 Тянут вниз

- United States
- FaceBoom
- TipTop

5.2 Драйверы успеха

Рекламные каналы

- RocketSuperAds
- LeapBob
- lambdaMediaAds
- MediaTornado
- YRabbit

Платформы

- PC
- Android

5.3 Рекомендации по исправлению ситуации

- **отключить** рекламу в **United States** и в каналах **FaceBoom** и **TipTop**
- изучить особенности пользователей в **United States** и в каналах **FaceBoom** и **TipTop**
- (пригласить к сотрудничеству специалистов по **United States** и каналам **FaceBoom** и **TipTop**)
- **запустить в тестовом режиме** обновлённые рекламные кампании в **United States** и в каналах **FaceBoom** и **TipTop**
- проверить результаты

Анализ провёл Эдуард Дементьев

(<https://eddydewrussia.ru/category/%d0%b0%d0%bd%d0%b0%d0%bb%d0%b8%d0%b7-%d0%b4%d0%b0%d0%bd%d0%bd%d1%8b%d1%85/>)