



+ Код + Текст



ОЗУ



Редактирование



Подготовка



```
[39] 1 # imports
2 import numpy as np
3 import pandas as pd
4 import matplotlib
5 import plotly.express as px
6 from sqlalchemy import create_engine
7 from datetime import datetime, timedelta
8 from plotly.subplots import make_subplots
```

```
[40] 1 # connection
2 HOST = '37.139.42.145'
3 DBNAME = 'game-analytics'
4 USER = 'analytics'
5 PASSWORD = 'BRTtqYiJyr29WXN'
6 TABLE_SCHEMA = 'data_viz_1068.project_dataset'
7 engine = create_engine(f'postgresql://{{USER}}:{{PASSWORD}}@{{HOST}}/{{DBNAME}}')
```



```
[41] 1 # getting table from database
2 project_dataset = pd.read_sql(f"""
3     SELECT *
4     FROM {TABLE_SCHEMA}
5     """, con=engine)
6
7 # converting some types to datetime
8 project_dataset['event_date'] = pd.to_datetime(project_dataset['event_date'])
9 project_dataset['cohort_date'] = pd.to_datetime(project_dataset['cohort_date'])
10
11 # dropping full duplicates (if they exists)
12 project_dataset.drop_duplicates()
13
14 # printing
15 project_dataset.head(2)
```

	event_time	event_date	event_name	revenue_usd	region	country	device_type	platform	cohort_date	user_id	user_type	content_id	currency	revenue
0	2020-12-14 21:16:20	2020-12-14	LaunchApp	0.0	EU	RU	xiaomi-Redmi Note 7	android	2020-09-29	120512	non_organic	None	None	0
1	2020-12-14 21:16:41	2020-12-14	LaunchApp	0.0	AS	KG	samsung-SM-G928F	android	2020-12-12	212364	organic	None	None	0



4

▶

Дополнение основного датасета

```
[42] 1 project_dataset = project_dataset.dropna(subset=['event_time'])
```

```
[43] 1 # copy of main df
2 refined_df = project_dataset.copy()
```

```
[44] 1 # adding CIS region
2 refined_df['region'] = (np.where(
3     refined_df['country'].isin(['RU', 'UA', 'BE', 'AZ', 'AM', 'KZ', 'KG', 'MD', 'TJ', 'TM', 'UZ']),
4     'CIS',
5     refined_df['region']))
```

```
[45] 1 # adding install datetime
2 install_times = refined_df[refined_df['event_name']=='FirstLaunchApp']
3 install_times = install_times.drop_duplicates(subset='user_id')
4 install_times = install_times[['user_id', 'event_time']]
5 install_times = install_times.rename(columns={'event_time': 'install_time'})
6 install_times['install_date'] = install_times['install_time'].dt.date
7 refined_df = refined_df.merge(install_times, on='user_id')
8 refined_df.head(1)
```

	event_time	event_date	event_name	revenue_usd	region	country	device_type	platform	cohort_date	user_id	user_type	content_id	currency	revenue
0	2020-12-14 21:16:41	2020-12-14	LaunchApp	0.0	CIS	KG	samsung-SM-G928F	android	2020-12-12	212364	organic	None	None	0.0

```
[46] 1 # Removing events which happened before install times.
2 refined_df = refined_df[refined_df['event_time'] >= refined_df['install_time']].copy()
3 # adding living hours
4 refined_df['liffehour'] = ((refined_df['event_time'] - refined_df['install_time']).dt.total_seconds() / 3600).astype('int')
5 # adding living days
6 refined_df['lifeday'] = (refined_df['event_time'] - refined_df['install_time']).dt.days
7 refined_df.head(2)
```

	event_time	event_date	event_name	revenue_usd	region	country	device_type	platform	cohort_date	user_id	user_type	content_id	currency	revenue
0	2020-12-14 21:16:41	2020-12-14	LaunchApp	0.0	CIS	KG	samsung-SM-G928F	android	2020-12-12	212364	organic	None	None	None
1	2020-12-14 23:56:25	2020-12-14	achieve_level_6	0.0	CIS	KG	samsung-SM-G928F	android	2020-12-12	212364	organic	None	None	None

```
[47] 1 refined_df = refined_df[refined_df['install_date'] > pd.to_datetime('2020-11-10')]
2 refined_df = refined_df[refined_df['install_date'] < pd.to_datetime('2021-02-10')]
3 refined_df = refined_df[refined_df['lifeday'] < 90]
```

/usr/local/lib/python3.7/dist-packages/pandas/core/ops/array_ops.py:73: FutureWarning:

Comparison of Timestamp with datetime.date is deprecated in order to match the standard library behavior. In a future version these will be considered

Функции

AVERAGE CARPU

```
[48] 1 def GetAverageCarpu(df, cohort_name):
2     install_users = df.groupby(['install_date']).agg(install_users=('user_id', 'nunique'))
3     df = df.merge(install_users, on='install_date')
4     df = df.groupby(['install_date', 'lifeday']).apply(lambda g: pd.Series({
5         'total_revenue' : g['revenue_usd'].sum(),
6         'total_users' : g['install_users'].iloc[0],
7     }).reset_index())
8     df = df.groupby(['install_date']).apply(lambda g: g.assign(
9         cum_sum=g['total_revenue'].cumsum(),
10        CARPU=g['total_revenue'].cumsum()/g['total_users'].iloc[0]
11    )).droplevel(0)
12     df.sort_values(by=['install_date', 'lifeday'])
13     df = df.groupby('lifeday')[['CARPU']].mean().reset_index()
14     df['cohort'] = cohort_name
15     return df
16
```

```
[49] 1 def ShowAverageCarpuForCohort(cohort, name):
2     result = GetAverageCarpu(cohort[cohort['region']=='CIS'], name + '_cis')
3     result = result.append(GetAverageCarpu(cohort[cohort['region']=='NA'], name + '_na'))
4
5     display(px.line(result, x='lifeday', y='CARPU', color='cohort').show())
```

```
[50] 1 def GetCarpuPerDay(df):
2     install_users = df.groupby(['install_date']).agg(install_users=('user_id', 'nunique'))
3     df = df.merge(install_users, on='install_date')
4     df = df.groupby(['install_date', 'lifeday']).apply(lambda g: pd.Series({
5         'total_revenue' : g['revenue_usd'].sum(),
6         'total_users' : g['install_users'].iloc[0],
7     }).reset_index())
8     df = df.groupby(['install_date']).apply(lambda g: g.assign(
9         cum_sum=g['total_revenue'].cumsum(),
10        CARPU=g['total_revenue'].cumsum()/g['total_users'].iloc[0]
11    )).droplevel(0)
12     df.sort_values(by=['install_date', 'lifeday'])
13     #pivot = df.pivot(index='install_date', columns='lifeday', values='CARPU')
14     return df
```

```
[51] 1 df = GetCarpuPerDay(refined_df)
2 pivot = df.pivot(index='install_date', columns='lifeday', values='CARPU')
```

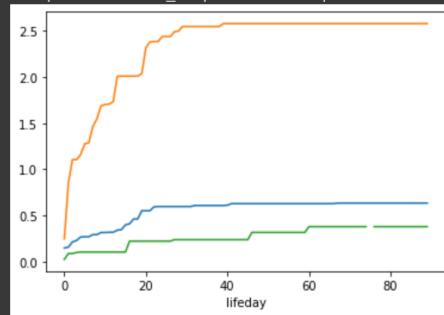
Ради интереса построим CARPU для нескольких случайных дней инсталла.

На графике видно что в конце жизни CARPU выходит на плато и begins to plateau

На графиках видно что в конце концов CACR выходит на плато и больше значимо не растет.

[52] 1 display(pivot.iloc[11].plot())
2 display(pivot.iloc[31].plot())
3 display(pivot.iloc[60].plot())

↳ <matplotlib.axes._subplots.AxesSubplot at 0x7fc26bda1dd0>
<matplotlib.axes._subplots.AxesSubplot at 0x7fc26bda1dd0>
<matplotlib.axes._subplots.AxesSubplot at 0x7fc26bda1dd0>



▼ ROLLING

[53] 1 def get_arpdau(df):
2 result = df.groupby('event_date')['revenue_usd'].sum() / df.groupby('event_date')['user_id'].nunique()
3 result = result.mean()
4 return result

[54] 1 def get_rolling_retention(data_frame):
2 copy_df = data_frame
3 copy_df = copy_df.groupby('user_id', as_index = False).agg(max_lifeday=('lifeday', 'max'))
4
5 df_1_to_24 = pd.DataFrame({'lifeday' : range(0, 90)})
6 df_1_to_24['for_merge'] = 1
7 copy_df['for_merge'] = 1
8
9 df_rolling = pd.merge(df_1_to_24, copy_df, on = 'for_merge')
10 df_rolling = df_rolling[df_rolling['max_lifeday'] >= df_rolling['lifeday']]
11 df_rolling = (df_rolling.groupby('lifeday', as_index = False).agg(unique_users=('user_id', 'nunique'))
12 .sort_values(by='lifeday'))
13 df_rolling['rolling_retention'] = round(df_rolling['unique_users'] / df_rolling['unique_users'][0], 4) * 100
14
15 return df_rolling

[55] 1 def show_ltv_via_rolling(df, name):
2 retention = get_rolling_retention(df)
3 lifetime = (retention['rolling_retention'] / 100).sum()
4 arpdau = get_arpdau(df)
5 result = lifetime * arpdau
6 print(name + " " + str(result))
7 return result

[56] 1 def show_rolling_ltv(cohort, name):
2 show_ltv_via_rolling(cohort[cohort['region']=='CIS'], name + '_cis')
3 show_ltv_via_rolling(cohort[cohort['region']=='NA'], name + '_na')
4 print('----')

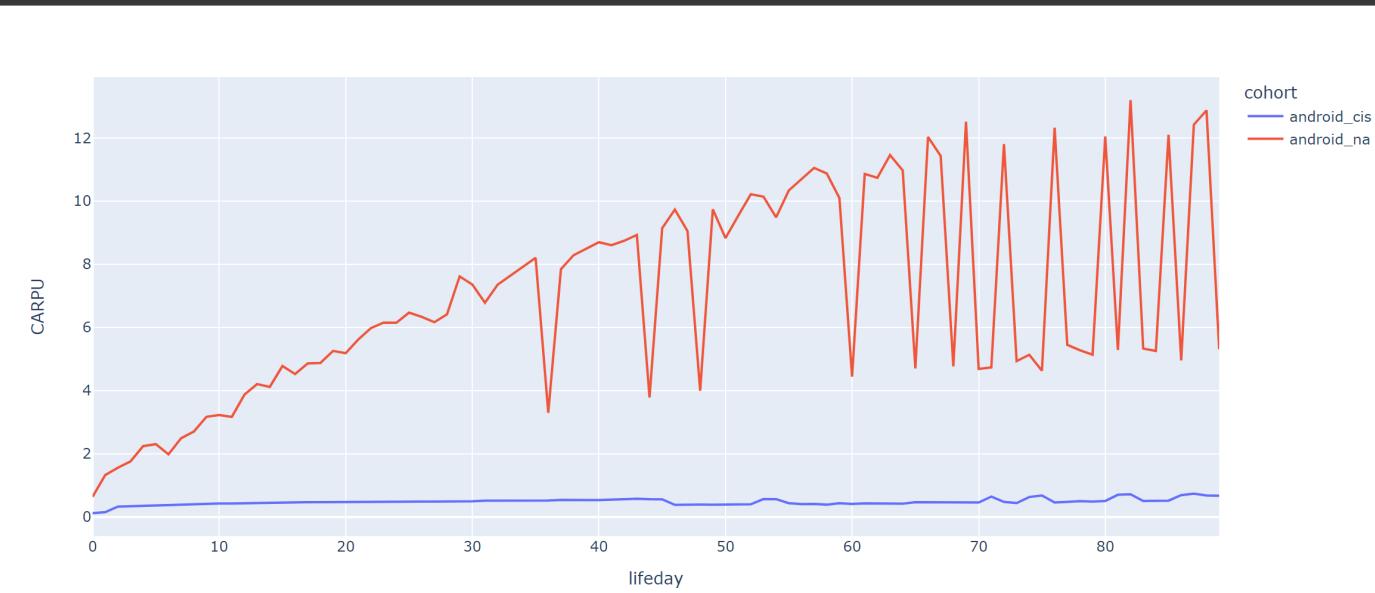
▼ Retention

[57] 1 def get_retention_rate(data_frame):
2 copy_df = data_frame
3
4 copy_df = copy_df.groupby('lifeday', as_index = False).agg(unique_users=('user_id', 'nunique')).sort_values(by='lifeday')
5 copy_df['retention'] = round(copy_df['unique_users'] / copy_df['unique_users'][0], 4) * 100
6
7 return copy_df
8
9
10 def show_ltv_via_retention(df, name):
11 retention = get_retention_rate(df)
12 lifetime = (retention['retention'] / 100).sum()
13 arpdau = get_arpdau(df)
14 result = lifetime * arpdau
15 print(name + " " + str(result))
16 return result
17
18
19 def show_retention_ltv(cohort, name):

```
20     show_ltv_via_retention(cohort[cohort['region']=='CIS'], name + '_cis')
21     show_ltv_via_retention(cohort[cohort['region']=='NA'], name + '_na')
22     print('----')
```

▼ 1. Android

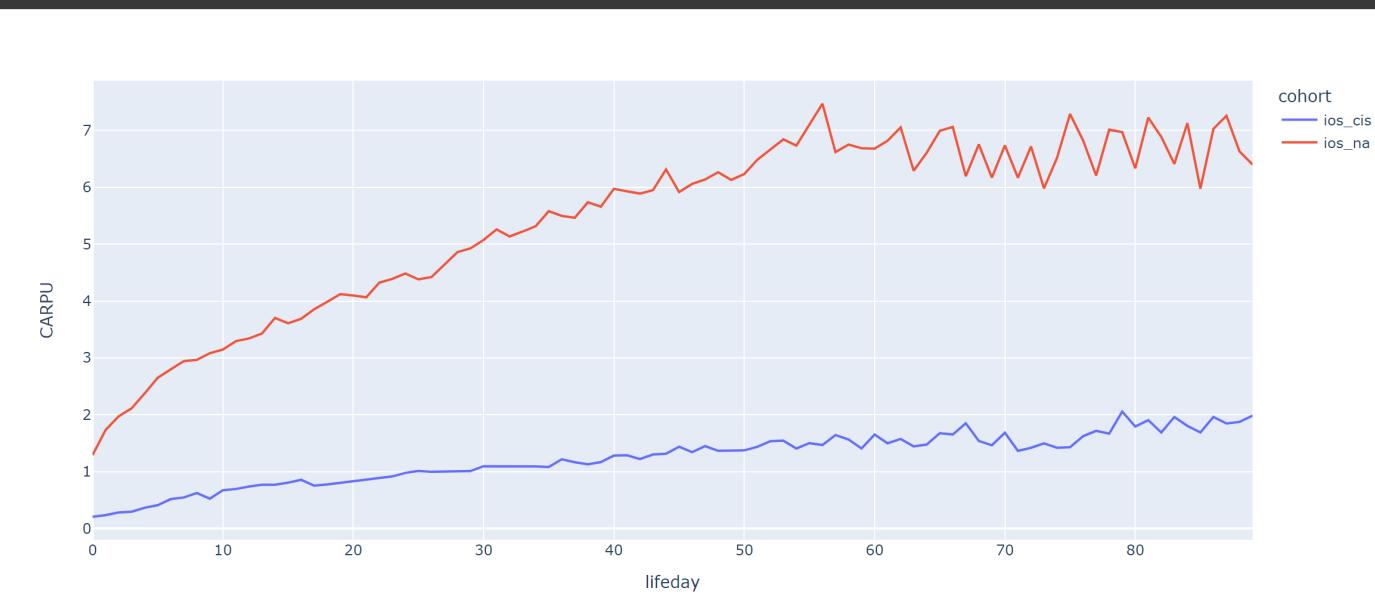
```
[58] 1 ShowAverageCarpuForCohort(refined_df[refined_df['platform']=='android'], 'android')
2 show_rolling_ltv(refined_df[refined_df['platform']=='android'], 'android')
3 show_retention_ltv(refined_df[refined_df['platform']=='android'], 'android')
```



```
None
android_cis 0.7012674812181283
android_na 8.71079309853721
----
android_cis 0.34849130010341867
android_na 4.209621212843139
----
```

▼ 1. iOS

```
[59] 1 ShowAverageCarpuForCohort(refined_df[refined_df['platform']=='ios'], 'ios')
2 show_rolling_ltv(refined_df[refined_df['platform']=='ios'], 'ios')
3 show_retention_ltv(refined_df[refined_df['platform']=='ios'], 'ios')
```



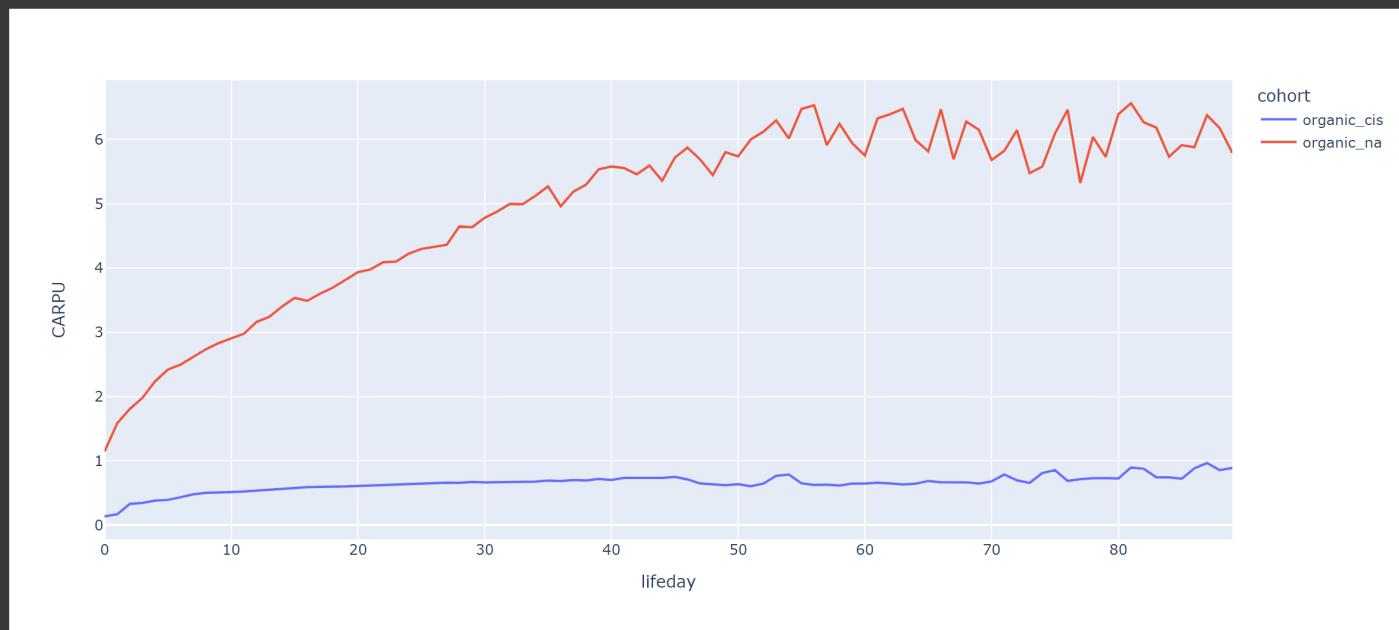
```
None
ios_cis 1.7815401796828456
ios_na 10.79870833512914
----
ios_cis 0.8579563440828234
ios_na 5.19355050300654
```

```
ios_na 5.10865059799854
```

```
----
```

▼ 1. ORGANIC

```
✓ [60] 1 ShowAverageCarpuForCohort(refined_df[refined_df['user_type']=='organic'], 'organic')
2 show_rolling_ltv(refined_df[refined_df['user_type']=='organic'], 'organic')
3 show_retention_ltv(refined_df[refined_df['user_type']=='organic'], 'organic')
```

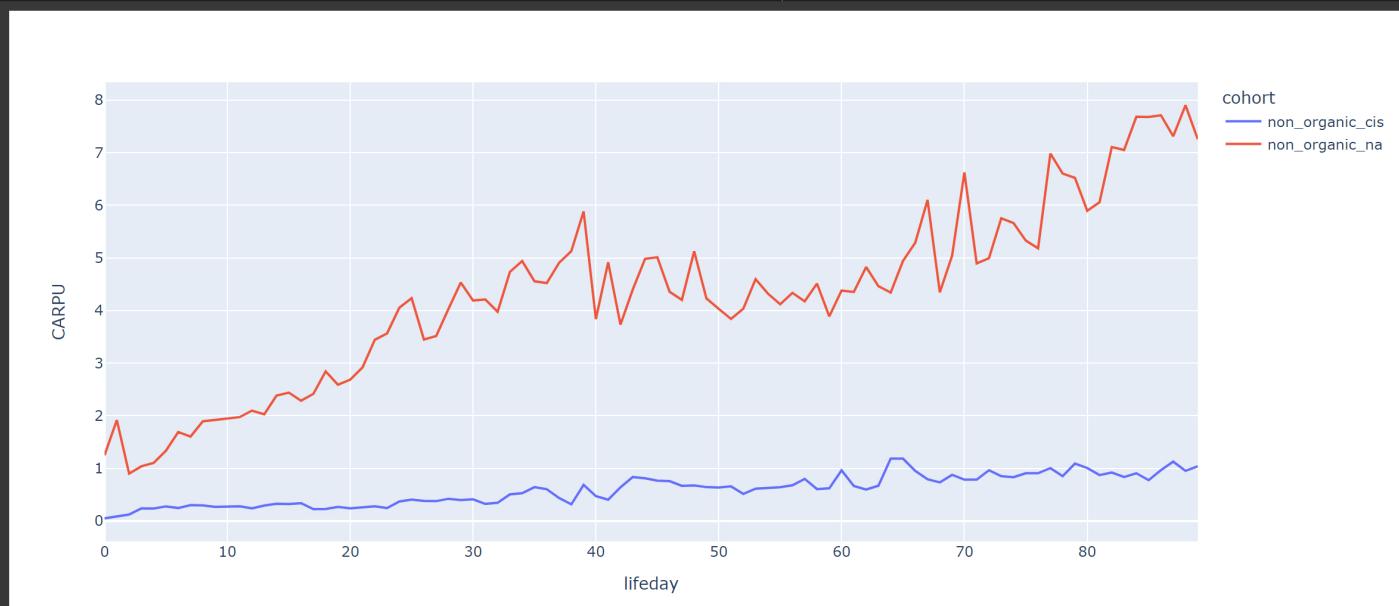


```
None
organic_cis 0.9665162635374484
organic_na 10.188857490966747
```

```
-----
organic_cis 0.4771642295991386
organic_na 4.835725610876053
-----
```

▼ 1. NON ORGANIC

```
✓ [61] 1 ShowAverageCarpuForCohort(refined_df[refined_df['user_type']=='non_organic'], 'non_organic')
2 show_rolling_ltv(refined_df[refined_df['user_type']=='non_organic'], 'non_organic')
3 show_retention_ltv(refined_df[refined_df['user_type']=='non_organic'], 'non_organic')
```



```
None
non_organic_cis 0.647981809089102
non_organic_na 9.531639423966949
```

```
-----
non_organic_cis 0.3095068071540245
non_organic_na 4.575676407937269
-----
```

▼ Вывод

Все методы показывают разные значения LTV.

Разница весьма существенна. LTV от retention rate самый маленький он примерно в 2 раза меньше чем LTV от rolling retention.

Среднее CARPU N-ного где то по середине между ними (если считать максимально стабильное значение = LTV).

Все методы показывают что среди всех когорт NA - безусловный лидер по LTV.

Если смотреть графики то среднее CARPU N-ного дня стабилизируется примерно на 60 день.

Какой метод самый практичный?

Если любому из этих методов можно доверять. И мы не знаем какой метод самый точный. То самым практическим будет метод который проще всего выполнить. Для меня это LTV через Retention Rate.