

Был запущен сплит-тест (A/B-тест), направленный на улучшение платёжной активности пользователей. Вам дан датасет с транзакциями пользователей до и во время теста в контрольной и тестовых группах

- Какие вы можете сделать выводы? Какая группа показала лучшие результаты?
- Можем ли мы как-то оценить из этих данных равномерность и валидность распределения юзеров по разным группам?
- Если не ограничиваться теми данными, которые приведены в рамках этого задания, что еще вы бы посчитали для оценки результатов групп?

Описание данных:

В таблице `users` приведена информация о том, какой юзер в какой момент времени попал в A/B тест:

- `tag` - лейбл группы (`control` - контрольная, остальные - тестовые)
- `ts` - время, когда впервые был выдан tag. То есть, все события до наступления времени ts происходили с юзером до попадания в A/B тест
- `user_uid` - внутренний id юзера (для мажннга со второй таблицей)
- `registration_time` - время регистрации пользователя в сервисе
- `conv_ts` - время совершения первой покупки пользователем в сервисе

В таблице `purchases` приведена информация о транзакциях пользователей из таблицы `users` до и во время A/B теста:

- `user_uid` - внутренний id юзера (для мажннга со второй таблицей)
- `time` - время совершения транзакции
- `consumption_mode` - вид потребления контента (`dto` - единица контента куплена навсегда, `rent` - единица контента взята в аренду, `subscription` - оформлена подписка)
- `element_uid` - уникальный id единицы контента или подписки
- `price` - цена (преобразованная)

Значения в полях `price` и всех полей, указывающих на время - преобразованы. Это значит, что значение в таблице не настоящее, но является линейным преобразованием реального значения, где ко всем значениям одного поля применено одно и то же преобразование - между ними сохранено отношение порядка. Ко всем полям, обозначающим время, применено одно и то же преобразование.

Выводы

1. Какие вы можете сделать выводы? Какая группа показала лучшие результаты?

Тестовая группа №3 показала лучшие результаты по ARPU

69.59 рубля против 68.34 рубля в контрольной

2. Можем ли мы как-то оценить из этих данных равномерность и валидность распределения юзеров по разным группам?

Всего 155104 пользователя совершили покупки после деления их на группы. 22% пользователей сервиса от общего числа прил. Пользователи распределились по группам не совсем равномерно

In [1]: *# Необходимые библиотеки для исследования*

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
from scipy import stats as st
```

In [2]:

```
sns.set(
    font_scale=2,
    style='whitegrid',
    rc={'figure.figsize':(16,6)}
)
```

In [3]: *#Информация о том, какой юзер в какой момент времени попал в A/B тест*
`users = pd.read_csv("Users/gaidovsky/Downloads/Apxim/users_...csv")`
`users.head()`

Out[3]:

	tag	ts	user_uid	registration_time	conv_ts
0	test4	30162.843868	6018855047f1168205c02a33369a	27410.154590	29405.523691
1	control	30133.146235	6e7b0c9f19f913a421b6767c703a4e1	27410.171795	27632.424734
2	test1	30140.447675	32096c27863c95a0407a9d1a7b1a	27410.217636	27414.028899
3	control	30136.790712	9502a3294c0a07d7f9a5db7c0ee12c60	27410.227367	27673.751236
4	test2	30150.336264	64112261a0214e999f154d28ba4622	27410.230786	29545.830319

In [4]: `users.isna().sum()`

Out[4]:

	tag	ts	user_uid	registration_time	conv_ts
tag	0	0	0	0	0
ts	0	0	0	0	0
user_uid	0	0	0	0	0
registration_time	0	0	0	0	0
conv_ts	0	0	0	0	0
dtype:	int64	int64	int64	int64	int64

In [5]: `users.dtypes`

Out[5]:

	tag	ts	user_uid	registration_time	conv_ts
tag	object	float64	object	float64	float64
ts	float64	float64	object	float64	float64
user_uid	object	float64	object	float64	float64
registration_time	float64	float64	object	float64	float64
conv_ts	float64	float64	object	float64	float64
dtype:	object	float64	object	float64	float64

In [6]: `users.shape`
`print(f"users.shape[0] количество строк")`
696982 количество строк

In [7]: `print(f"({users.user_uid.nunique()}) уникальных пользователей")`

694819 уникальных пользователей

Комментарий

Заметим, что количество строк в `users` больше, чем уникальных пользователей. Посмотрим на уникальные id. Не было ли такого, что один пользователь попал в разные группы

In [8]: `dupl = users.groupby('user_uid', as_index=False).agg({'registration_time':'count'}) \`
`.sort_values(by='registration_time', ascending=False).rename(columns={"registration_time": "number_of_tags"})`
`dupl[more_than_one] = dupl.number_of_tags > 2`
`dupl.head()`

Out[8]:

	user_uid	number_of_tags	more_than_one
293127	6c2289cfa50d5b66b159fe12ca4ed	4	True
496483	b6c9ba2b34242e1863b395d9b1323d	4	True
104642	26890a48724800518c97c310a1a005	4	True
132086	31c496417e66203117967a11c225e25	3	True
345378	7f98317238a090a0162656532a03	3	True

Комментарий

Видим, что некоторые пользователям tag был присвоен несколько раз. Посмотрим, что именно происходило с первым пользователем по количеству присвоенных тэгов

In [9]: `users.query("user_uid == '6c2289cfa50d5b66b159fe12ca4ed")`.sort_values(by='ts', ascending=False)

Out[9]:

	tag	ts	user_uid	registration_time	conv_ts
97224	test4	30140.554481	6c2289cfa50d5b66b159fe12ca4ed	28992.364445	28906.943711
90820	test3	30135.593332	6c2289cfa50d5b66b159fe12ca4ed	28992.364445	28906.943711
103896	control	30136.492787	6c2289cfa50d5b66b159fe12ca4ed	28992.364445	28906.943711
93037	test3	30133.872325	6c2289cfa50d5b66b159fe12ca4ed	28992.364445	28906.943711

In [10]: `print(f"({dupl.more_than_one.sum()}) пользователям был присвоен тэг больше 1 раза")`

2684 пользователям был присвоен тэг больше 1 раза

In [13]: `dupl.query("more_than_one == True").number_of_tags.sum()`

Out[13]: 4234

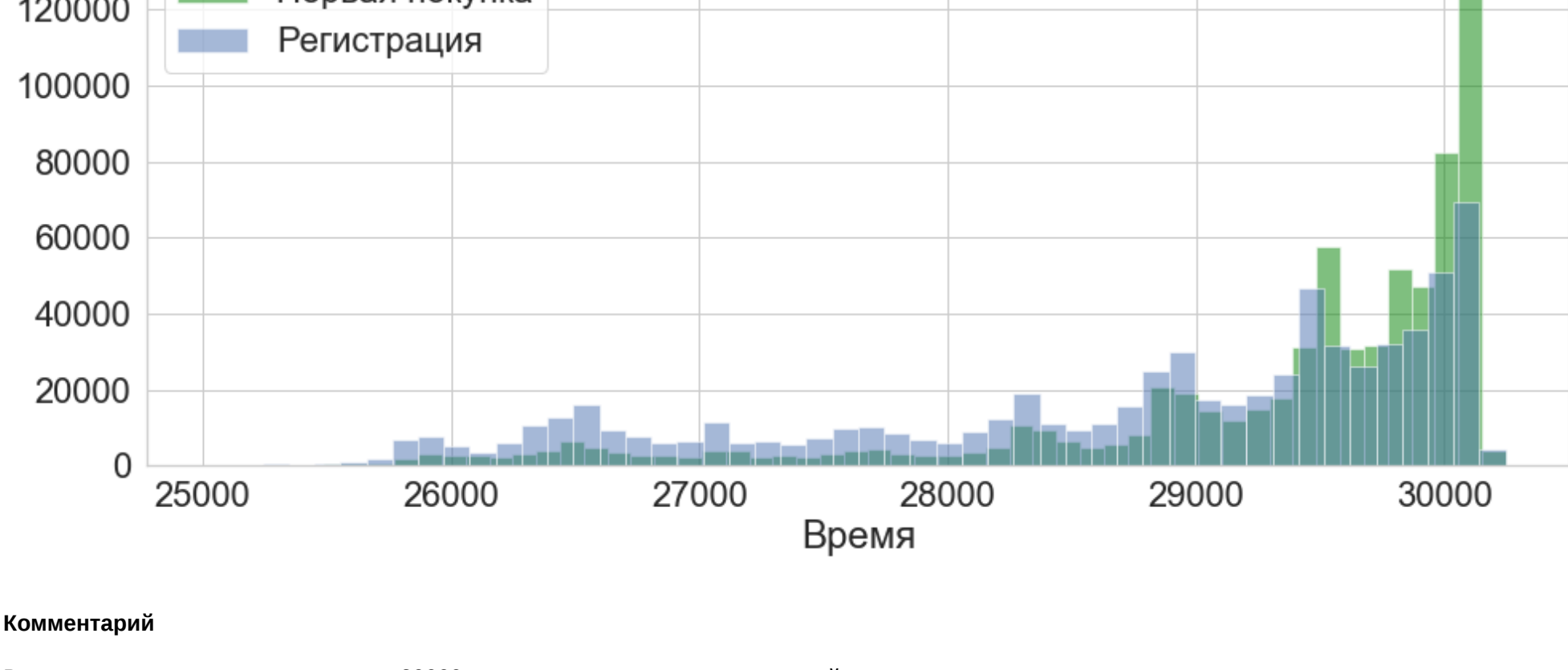
In [15]: `users.shape[0] - dupl.query("more_than_one == True").number_of_tags.sum() = dupl.more_than_one.sum()`

Out[15]: 694832

In []:

In [16]: `users.conv_ts.hist(alpha=0.5, color='green', bins=50, label='Первая покупка')`
`users.registration_time.hist(alpha=0.5, bins=50, label='Регистрация')`

#Пл. X.Ticks(rotation=90)
`plt.title("Распределение регистраций и первых покупок по времени")`
`plt.xlabel("Время")`
`plt.legend()`



Комментарий

Видим, что после условного времени 29000 произошел рост числа регистраций

In [17]: *#данные о транзакциях пользователей из таблицы users_до и во время A/B теста:*
`purchases = pd.read_csv("Users/gaidovsky/Downloads/Apxim/purchases_...csv", sep=';',)`
`purchases['price'] = purchases.price.round(1)`
`purchases.head()`

Out[17]:

	user_uid	time	consumption_mode	element_uid	price
0	d00a708c7b7e99146fe40e8f6535862b	30156.645112	dto	2ba66ac9785731d67b2b6155efac5c	44.5
1	0906074e1a1a229d5e74989b0646962	30156.645015	dto	e5642227569c0a672b1db55efac5c	38.6
2	efb6eca3325d573739c5d48ce330	30156.644990	dto	5447f0d33b061558dc073baedf0a7a	50.4
3	89a56a37245d591c4e49b44839712f	30156.644789	dto	8256d0c7b25382aac0c006c9c0e3	44.5
4	cc05e437f1899826205c70e0c1760	30156.644200	dto	a579f32dc1166240c661d843b66d5e5	44.5

In [18]: `purchases.shape`

Out[18]: (663849, 5)

In [19]: `purchases.isna().sum()`

Out[19]:

	user_uid	time	consumption_mode	element_uid	price
user_uid	0	0	0	0	0
time	0	0	0	0	0
consumption_mode	0	0	0	0	0
element_uid	0	0	0	0	0
price	0	0	0	0	0
dtype:	int64	int64	int64	int64	int64

In [20]: `purchases.dtypes`

Out[20]:

	user_uid	time	consumption_mode	element_uid	price
user_uid	object	float64	object	float64	float64
time	float64	float64	object	float64	float64
consumption_mode	object	float64	object	float64	float64
element_uid	object	float64	object	float64	float64
price	float64	float64	object	float64	float64
dtype:	object	float64	object	float64	float64

In [21]: `print(f"({purchases.user_uid.nunique()}) уникальных пользователей")`

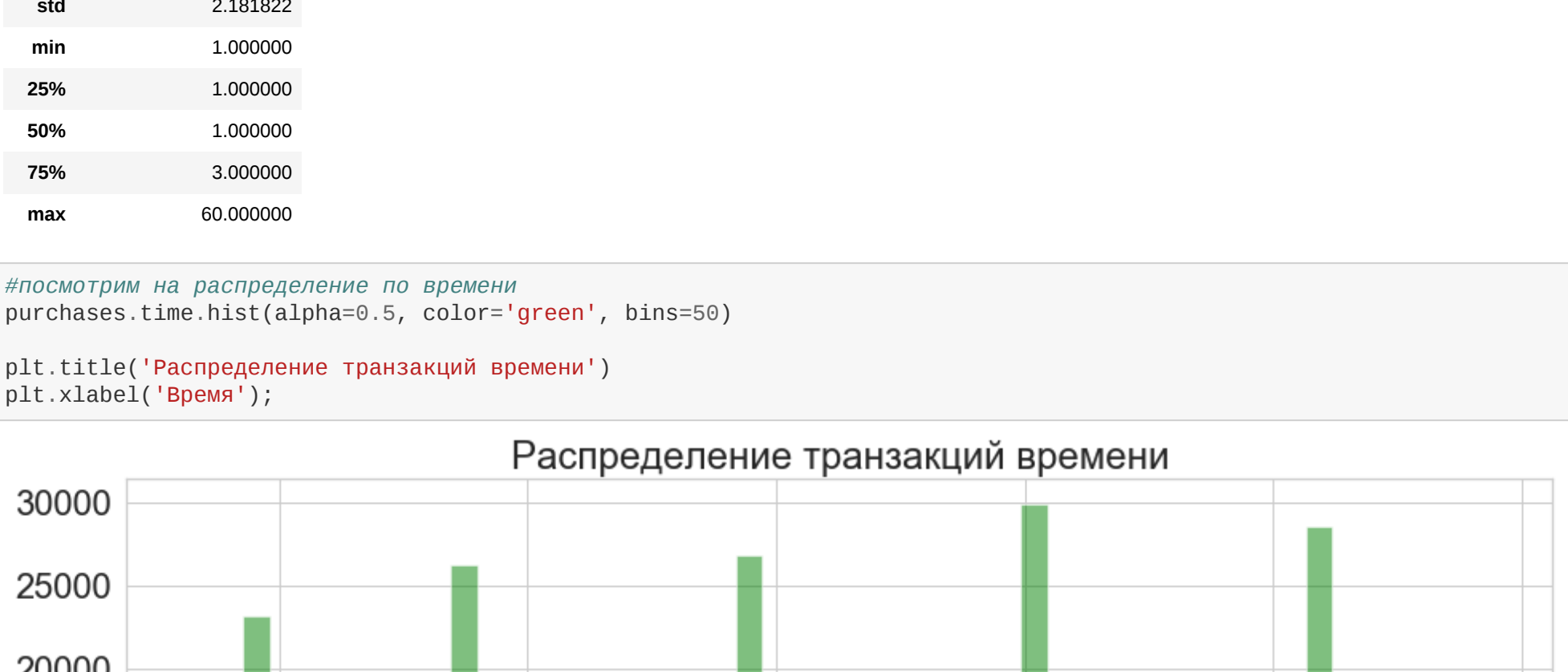
295290 уникальных пользователей

In [22]: *#зададим отдельную таблицу по количеству транзакций*
`number_of_purchases = purchases.groupby('user_uid', as_index=False).agg({'time':'count'}) \`
`.rename(columns={"time": "number_of_purchases"})`
`.sort_values(by="number_of_purchases", ascending=False)`
#посмотрим на описательные статистики покупок пользователей
`number_of_purchases.describe()`

Out[22]:

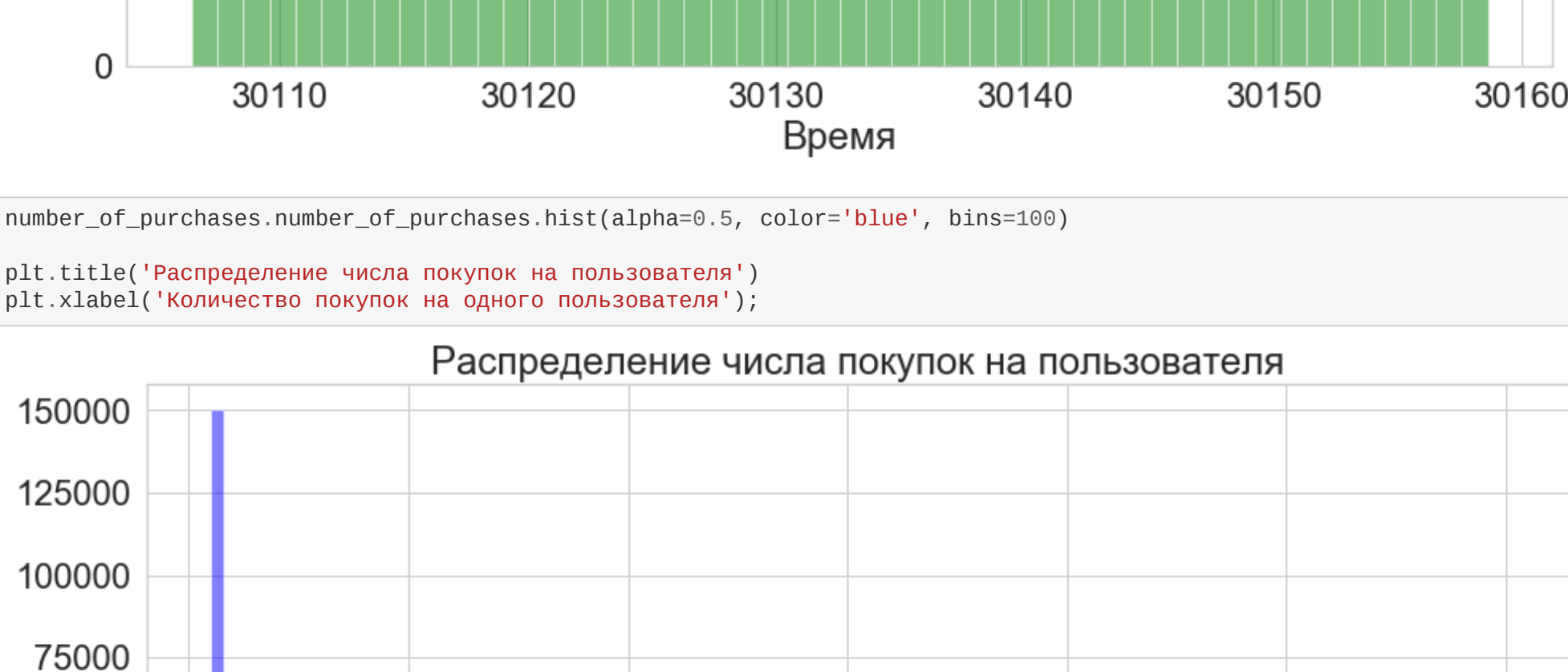
	number_of_purchases
count	295200.000000
mean	2.244811
std	2.181822
min	1.000000
25%	1.000000
50%	1.000000
75%	3.000000
max	60.000000

In [23]: *#посмотрим на распределение по времени*
`purchases.time.hist(alpha=0.5, color='green', bins=50)`
`plt.title("Распределение транзакций времени")`
`plt.xlabel("Время")`



In [24]: `number_of_purchases.number_of_purchases.hist(alpha=0.5, color='blue', bins=100)`

`plt.title("Распределение числа покупок на пользователя")`
`plt.xlabel("Количество покупок на одного пользователя")`



Комментарий

- В среднем пользователи совершили 2.2 транзакции за рассмотренный условный период
- 25% пользователей совершили 3 и более транзакции
- Максимальное количество покупок на одного юзера - 60
- Большинство пользователей совершили 1 покупку
- Распределение покупок по времени имеет явную цикличность (возможно, по дня недели)

In [25]: *#количество приобретений контента по виду потребления контента*
`popular_consumption_mode = purchases.pivot_table(index = 'consumption_mode', values = 'user_uid',`
`.agg({'time':'count', 'price':'mean'}) \`
`.rename(columns={"time": "number_of_transactions", 'price':'mean_price'})`
`consumption_mode['revenue'] = consumption_mode.number_of_transactions * consumption_mode.mean_price`
`consumption_mode = consumption_mode.sort_values(by='revenue', ascending=False)`
`consumption_mode['share_of_revenue'] = consumption_mode.revenue / consumption_mode.revenue.sum()`
`consumption_mode['share_of_transactions'] = consumption_mode.number_of_transactions / consumption_mode.number_of_transactions.sum()`
`round(consumption_mode,2)`

Out[25]:

	consumption_mode	number_of_transactions	mean_price	revenue	share_of_revenue	share_of_transactions
0	dto	245276	45.81	10989759.4	0.45	0.36
2	subscription	239897	34.93	8367984.9	0.35	0.37
1	rent	178676	27.81	4969602.0	0.20	0.27

- Подписка принесла 45% выручку за рассмотренный период
- Так же подписка - самый дорогой вид потребления. Средняя цена - 45.8
- Покупка контента и подписка почти одинаково популярны 37% и 36% соответственно
- Меньше всего клиенты сервиса пользовались 'арендой'. Так же аренда контента внесла наименьший вклад в выручку сервиса 20%

In [26]: `purchases.pivot_table(index = 'consumption_mode', values = 'price',`
`.aggfunc = 'mean', fill_value=0).reset_index() \`
`.rename(columns={"user_uid": "number_of_transactions"}) \`
`.sort_values(by="number_of_transactions", ascending=False)`

Out[26]:

	consumption_mode	price
0	dto	34.830315
1	rent	27.813484
2	subscription	45.810324

In [27]: *#количество транзакций, средняя цена и выручка по каждому виду потребления*
`consumption_mode = purchases.groupby('consumption_mode', as_index=False) \`
`.agg({'time':'count', 'price':'mean'}) \`
`.rename(columns={"time": "number_of_transactions", 'price':'mean_price'})`
`consumption_mode['revenue'] = consumption_mode.number_of_transactions * consumption_mode.mean_price`
`consumption_mode = consumption_mode.sort_values(by='revenue', ascending=False)`
`consumption_mode['share_of_revenue'] = consumption_mode.revenue / consumption_mode.revenue.sum()`
`consumption_mode['share_of_transactions'] = consumption_mode.number_of_transactions / consumption_mode.number_of_transactions.sum()`
`round(consumption_mode,2)`

Out[27]:

	consumption_mode	number_of_transactions	mean_price	revenue	share_of_revenue	share_of_transactions
2	subscription	239897	34.93	10989759.4	0.45	0.36
0	dto	245276	45.81	8367984.9	0.35	0.37
1	rent	178676	27.81	4969602.0	0.20	0.27

- Подписка принесла 45% выручку за рассмотренный период
- Так же подписка - самый дорогой вид потребления. Средняя цена - 45.8
- Покупка контента и подписка почти одинаково популярны 37% и 36% соответственно
- Меньше всего клиенты сервиса пользовались 'арендой'. Так же аренда контента внесла наименьший вклад в выручку сервиса 20%

In [28]: *#объединяем таблицы*
`us_pur = purchases.merge(users, on='user_uid')`

In [29]: *#ставим транзакции, которые были совершены после определения тэга для юзера*
`us_pur_ab = us_pur.query("time > ts")`
`us_pur_ab.head()`

Out[29]:

	user_uid	time	consumption_mode	element_uid	price	tag	ts	registration_time
0	d00a708c7b7e99146fe40e8f6535862b	30156.645112	dto	2ba66ac9785731d67b2b6155efac5c	44.5	control	30145.756945	29592.948896
1	d00a708c7b7e99146fe40e8f6535862b	30157.124337	dto	5aa724c3930870c3a18c6250310a45	44.5	control	30145.756945	29592.948896
2	d00a708c7b7e99146fe40e8f6535862b	30153.692015	dto	d1a6e6c3f9d6b15e12b0109898e4e6	38.6	control	30145.756945	29592.948896
3	d00a708c7b7e99146fe40e8f6535862b	30153.584597	dto	6e03393238a09a05a2608a0a76775	38.6	control	30145.756945	29592.948896
8	0906074e1a1a229d5e74989b0646962	30156.645015	dto	e5642227569c0a672b1db55efac5c	38.6	control	30137.415922	28079.964466

In [30]: `us_pur_ab.shape`

Out[30]: (279550, 9)

In []:

In []:

In [31]: `us_pur_user = us_pur_ab.groupby(['user_uid', 'tag'], as_index = False) \`
`.agg({'price': 'sum'})`
`sort_values(by='price')`

Out[31]:

	user_uid	tag	price
5170	540ab0d47380a531ba9f13520249597	control	68.0
80601	83d21e131c303950fbc20632acdb91	control	103.6
135975	e01c49cbb45c460ef5423a3834c13e	control	68.0
38034	3eb1a6e84d1ad0c79392929c7995dc	control	38.6
111887	b7ef16048d477036d38246a201561	control	68.0
...
60164	62223a0b78c3205ca3918460010c86a	test4	29.8
60163	63223a0b78c3205ca3918460010c86a	test4	29.8
60149	6314d651e6b90909506a4962d5e4e27	test4	38.6
60177	6320e6630c70c1b615e0158c03a601	test4	197.4
155104	fff09368a050a054eae7821b1	test4	68.0

In [43]: *#посчитаем ARPU и распределение пользователей по группам*
`groups = sum_per_user.groupby(['tag'], as_index = False).agg({'price':'mean', 'user_uid':'count'})`
`.rename(columns={"price": "ARPU", "user_uid":"number_of_users"})`
`groups['test_user_share'] = groups.number_of_users / groups.number_of_users.sum()`
`groups = round(groups, 2)`
`groups.head()`

Out[43]:

	tag	ARPU	number_of_users	test_users_share
0	control	68.34	30117	0.19
1	test1	69.15	31326	0.20
2	test2	69.40	31104	0.20
3	test3	69.60	31162	0.20
4	test4	69.20	31195	0.20

Комментарий

- на первый взгляд пользователи по тестовым группам распределились относительно равномерно, но в контрольной группе значения сильно меньше
- проверим равномерность распределения
- наибольшее различие по ARPU у контрольной группы с test3
- сравним ARPU

In [39]: `from scipy.stats import chisquare`

In [34]: `chi_pvalue = chisquare(groups.number_of_users)[1]`
`if 0.05 == chi_pvalue:`
`print('отклонем H0, пользователи распределены равномерно по группам')`
`else:`
`print('не отклоняем H0, пользователи распределены неравномерно по группам')`

Отклоняем H0, пользователи распределены не равномерно по группам</