

Банки — Сегментация пользователей по потреблению

Задача:

Анализ и сегментация клиентов регионального банка по количеству потребляемых продуктов.

- Исследовательский анализ данных;
- Сегментация пользователей на основе данных о количестве потребляемых продуктов;
- Формулировка и проверка статистических гипотез различия дохода между теми клиентами, которые пользуются двумя продуктами банка, и теми, которые пользуются одним.

Описание данных

Датасет содержит данные о клиентах банка «Метантром». Банк располагается в Ярославле и областных городах: Ростов Великий и Рыбинск.

- `userid` — идентификатор пользователя,
- `score` — баллы кредитного скоринга,
- `city` — город,
- `Age` — пол,
- `Gender` — возраст,
- `Objects` — количество объектов в собственности,
- `Balance` — баланс на счёте,
- `Products` — количество продуктов, которыми пользуется клиент,
- `CreditCard` — есть ли кредитная карта,
- `Loyalty` — активный клиент,
- `estimated_salary` — заработная плата клиента,
- `Churn` — ушёл или нет.

Оглавление

- [Шаг 1. Предобработка и исследовательский анализ данных](#)
- [Шаг 2. Сегментация пользователей на основе данных о количестве потребляемых продуктов](#)
- [Шаг 3. Формулировка и проверка статистических гипотез различия дохода](#)
- [Выводы](#)

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.pyplot as px
import seaborn as sns
import math
import datetime
from IPython.display import display
from plotly import graph_objects as go
from scipy import stats as st
import warnings
warnings.filterwarnings('ignore')
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import Lasso, Ridge
from sklearn.ensemble import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from sklearn.metrics import LogisticRegression
from sklearn.metrics import accuracy_score, precision_score, recall_score
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
from scipy.cluster.hierarchy import dendrogram, linkage
from sklearn.cluster import KMeans
```

Шаг 1. Предобработка и исследовательский анализ данных

```
In [2]: bank_dataset = pd.read_csv('datasets/bank_dataset.csv')

In [3]: bank_dataset.info()
```

```
Out[3]: <class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 12 columns):
userid      10000 non-null int64
score       10000 non-null int64
city        10000 non-null object
Age         10000 non-null object
Gender      10000 non-null int64
Age         10000 non-null int64
Objects     10000 non-null int64
Balance     6383 non-null float64
Products    10000 non-null int64
CreditCard  10000 non-null int64
Loyalty     10000 non-null int64
estimated_salary 10000 non-null float64
Churn       10000 non-null int64
dtypes: float64(2), int64(8), object(2)
memory usage: 937.6+ KB
```

```
In [4]: bank_dataset.duplicated().sum()
Out[4]: 0
```

Столбец `Balance` содержит 3617 пропусков. Предположим, что пропущенное значение — это нулевой баланс на счёте. Заменим пропущенные значения:

```
bank_dataset['Balance'] = bank_dataset['Balance'].fillna(0)
```

Вывод: Датасет состоит из 12 столбцов и 10000 строк. Названия столбцов: одноимённые, но станут заменами. Дубликаты отсутствуют.

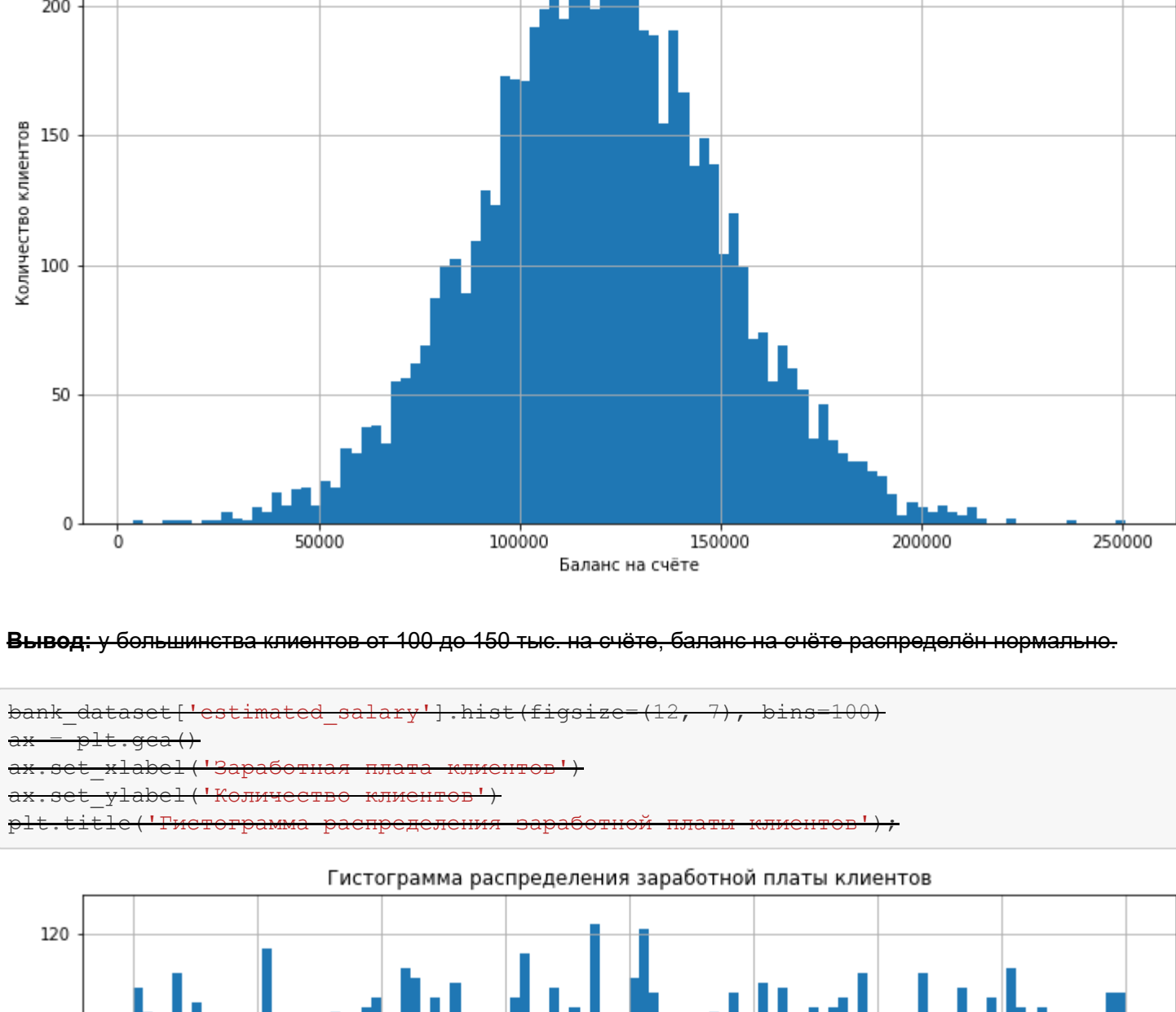
Проанализируем данные методом `describe()` и построим гистограммы распределения основных параметров.

```
bank_dataset.describe()
```

	count	mean	std	min	25%	50%	75%	max
userid	10000.0	4.673838e+02	7.100848e+01	14688422.00	14688422.00	14688422.00	14688422.00	14688422.00
score	10000.0	4.673838e+02	7.100848e+01	14688422.00	14688422.00	14688422.00	14688422.00	14688422.00
Age	10000.0	2.802180e+01	4.487896	18.00	3.200000e+00	3.700000e+01	4.400000e+01	92.00
Objects	6383.0	6.014280e+00	2.802174	0.00	0.000000e+00	6.000000e+00	7.000000e+00	10.00
Balance	6383.0	4.108275e+05	8.000545e+02	3789.69	4.001402e+05	4.108307e+05	4.205423e+05	2.68089e+06
Products	10000.0	4.630200e+00	0.681654	4.00	4.000000e+00	4.000000e+00	4.000000e+00	4.00
CreditCard	10000.0	7.554000e-01	0.455049	0.00	0.000000e+00	1.000000e+00	1.000000e+00	1.00
Loyalty	10000.0	5.545000e-01	0.499732	0.00	0.000000e+00	1.000000e+00	1.000000e+00	1.00
estimated_salary	10000.0	4.000000e+05	6.751040e+04	44.58	6.160214e+04	4.001030e+05	4.403882e+05	4.00000e+05
Churn	10000.0	2.802180e-01	0.402759	0.00	0.000000e+00	0.000000e+00	0.000000e+00	1.00

```
bank_dataset['score'].hist(figsize=(10,7), bins=100)
plt.show()
```

Гистограмма распределения баллов кредитного скоринга

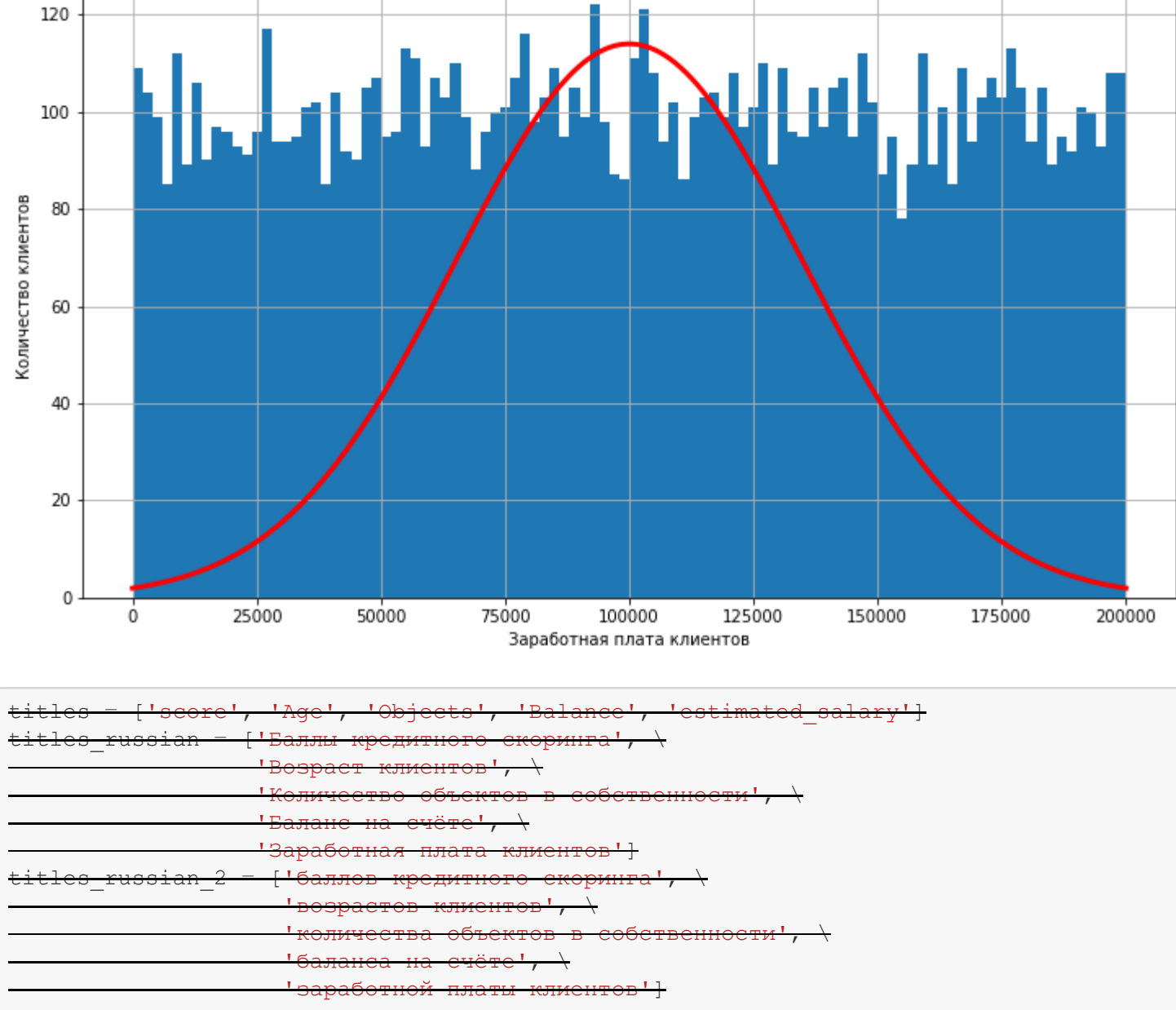


```
bank_dataset['score'].value_counts()
```

Вывод: Баллы кредитного скоринга от 400 до 800. Баллы распределены нормально, однако присутствуют аномально высокие значения (253) с максимальным значением баллов (850). Пока оставим выбор без изменений.

```
bank_dataset['Age'].hist(figsize=(10,7), bins=100)
plt.show()
```

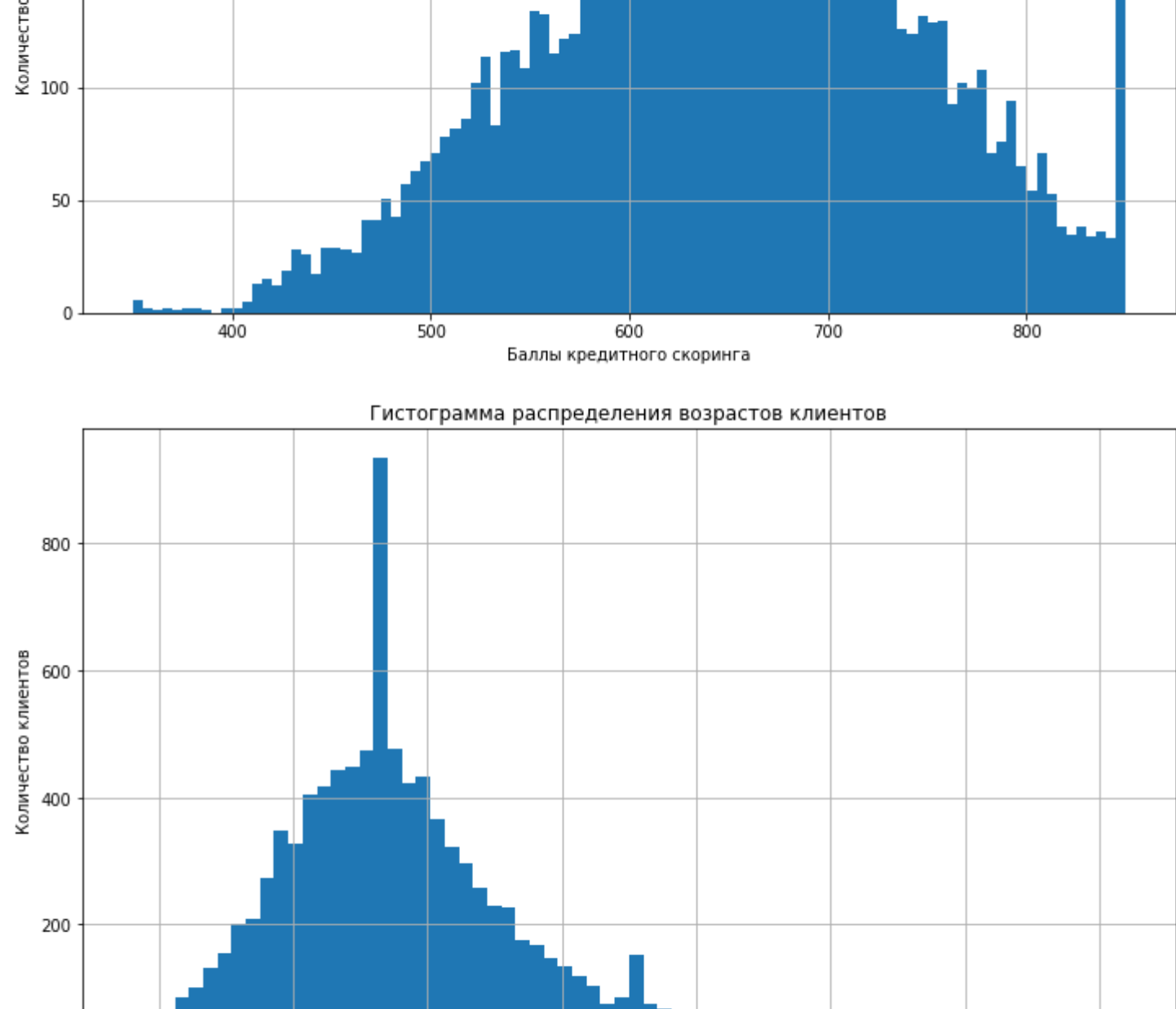
Гистограмма распределения возрастов клиентов



Вывод: Большинство клиентов старше от 30 до 40 лет. Данные распределены нормально.

```
bank_dataset['Age'].hist(figsize=(10,7), bins=100)
plt.show()
```

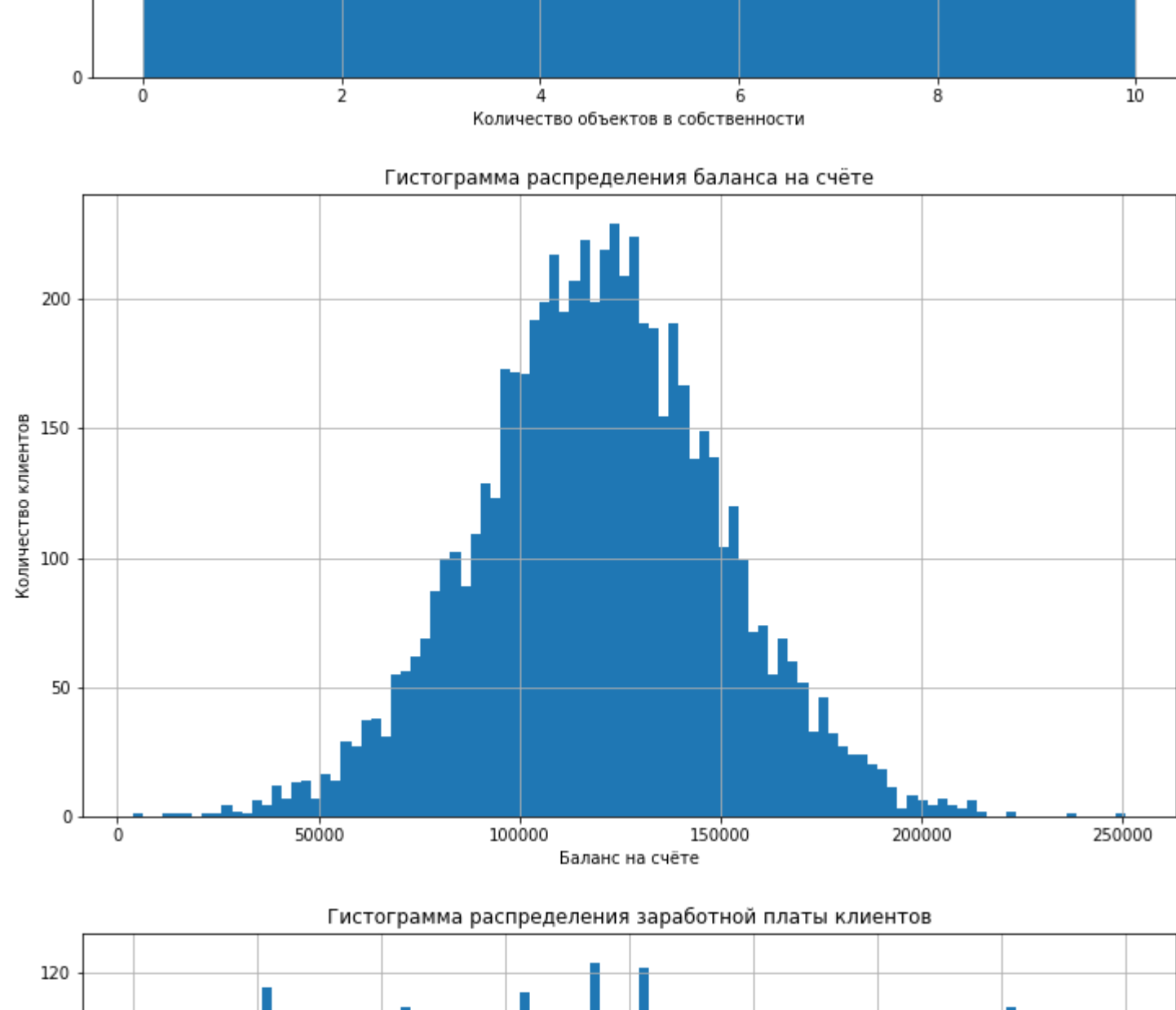
Гистограмма распределения количества объектов в собственности



Вывод: Количество объектов в собственности распределено нормально.

```
bank_dataset['Balance'].hist(figsize=(10,7), bins=100)
plt.show()
```

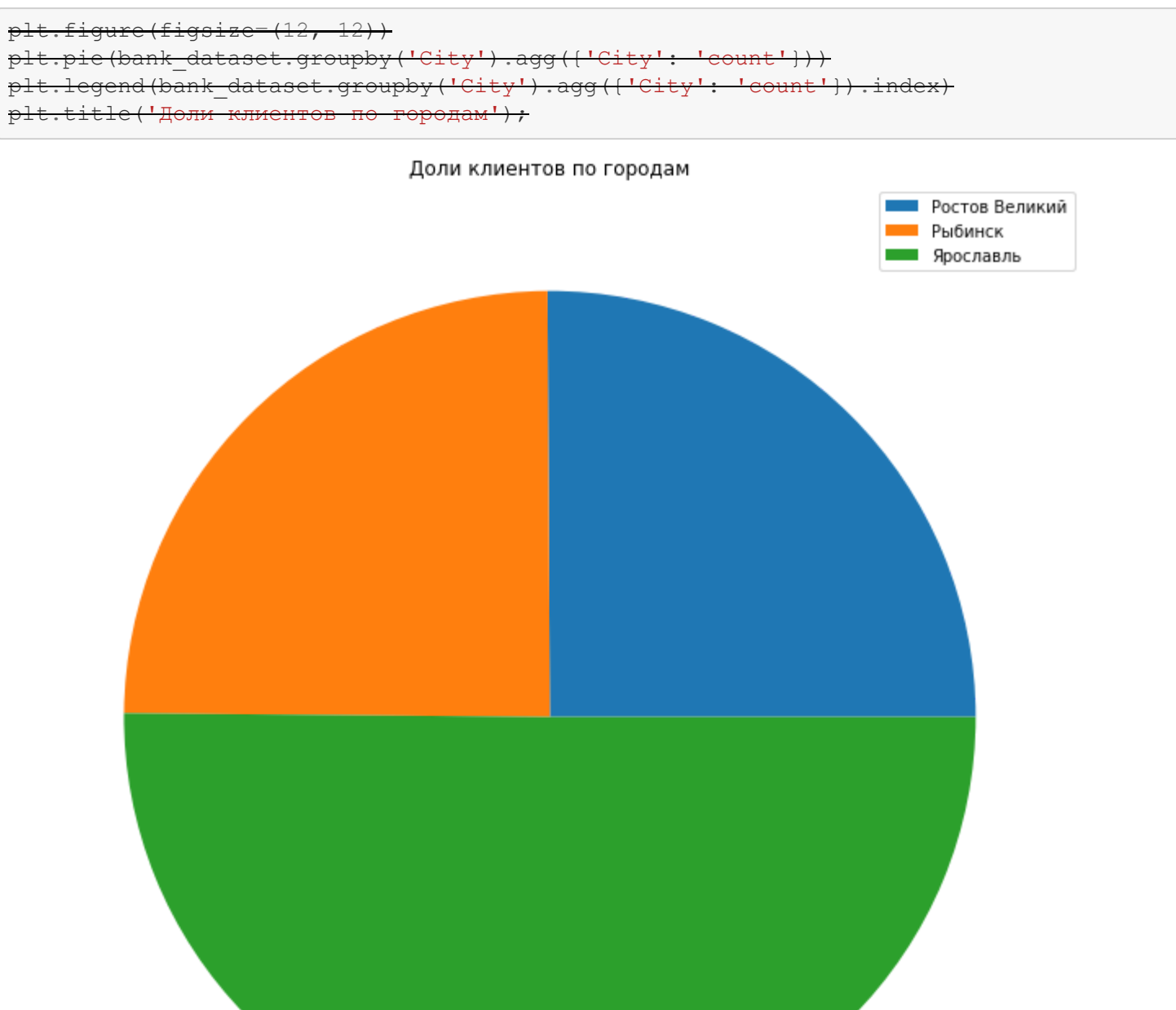
Гистограмма распределения баланса на счёте



Вывод: у большинства клиентов от 100 до 150 тыс. на счёте. Баланс на счёте распределён нормально.

```
bank_dataset['Balance'].hist(figsize=(10,7), bins=100)
plt.show()
```

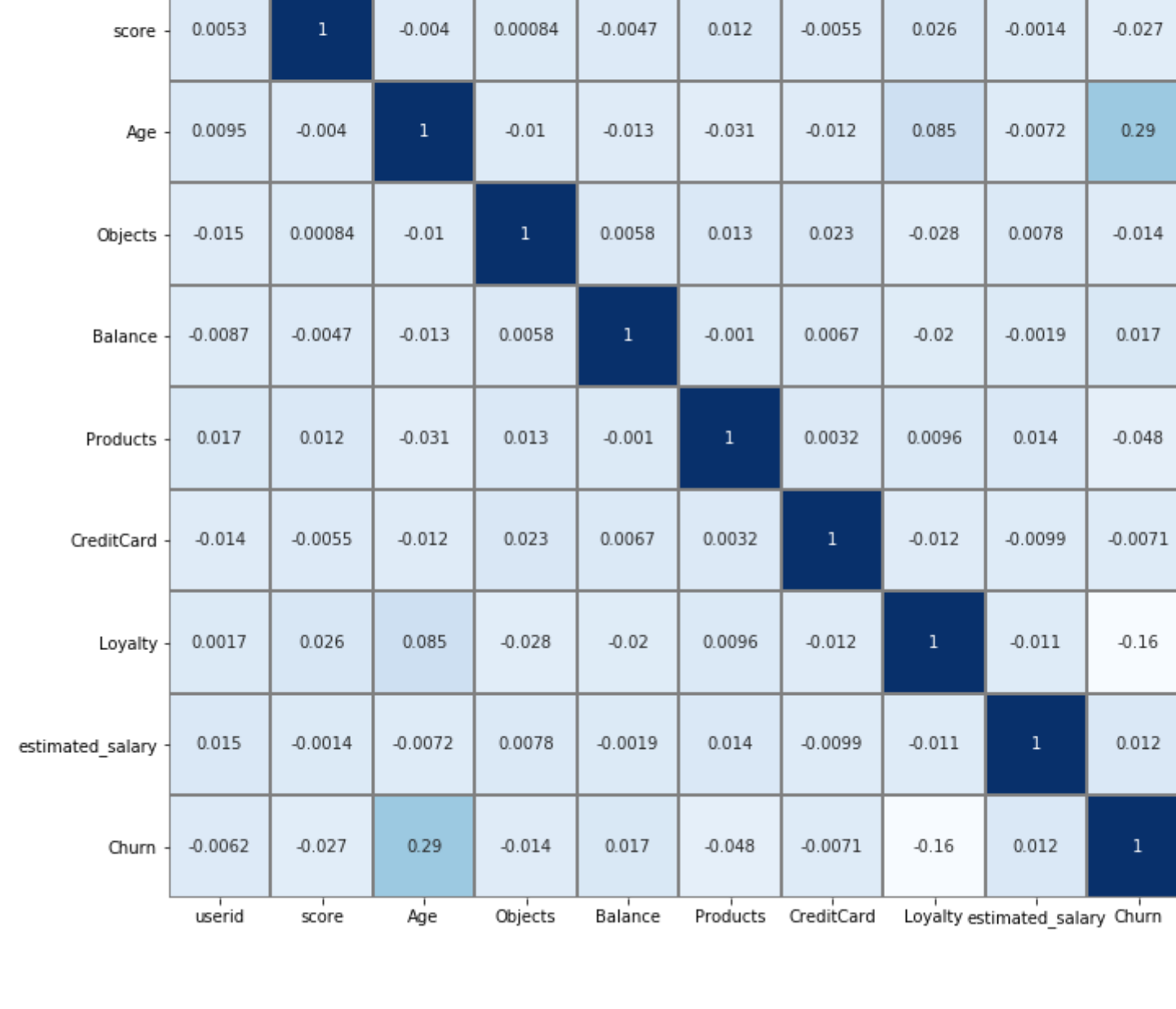
Гистограмма распределения заработной платы клиентов



Вывод: средняя заработная плата клиента — 10000. Распределение заработной плат равномерное. Есть ли аномально низкие (<5000), так и аномально высокие (>150000)? (Почему? Предположим, что данные о зарплате предоставлены сотрудниками доходов и примерно равном количестве)

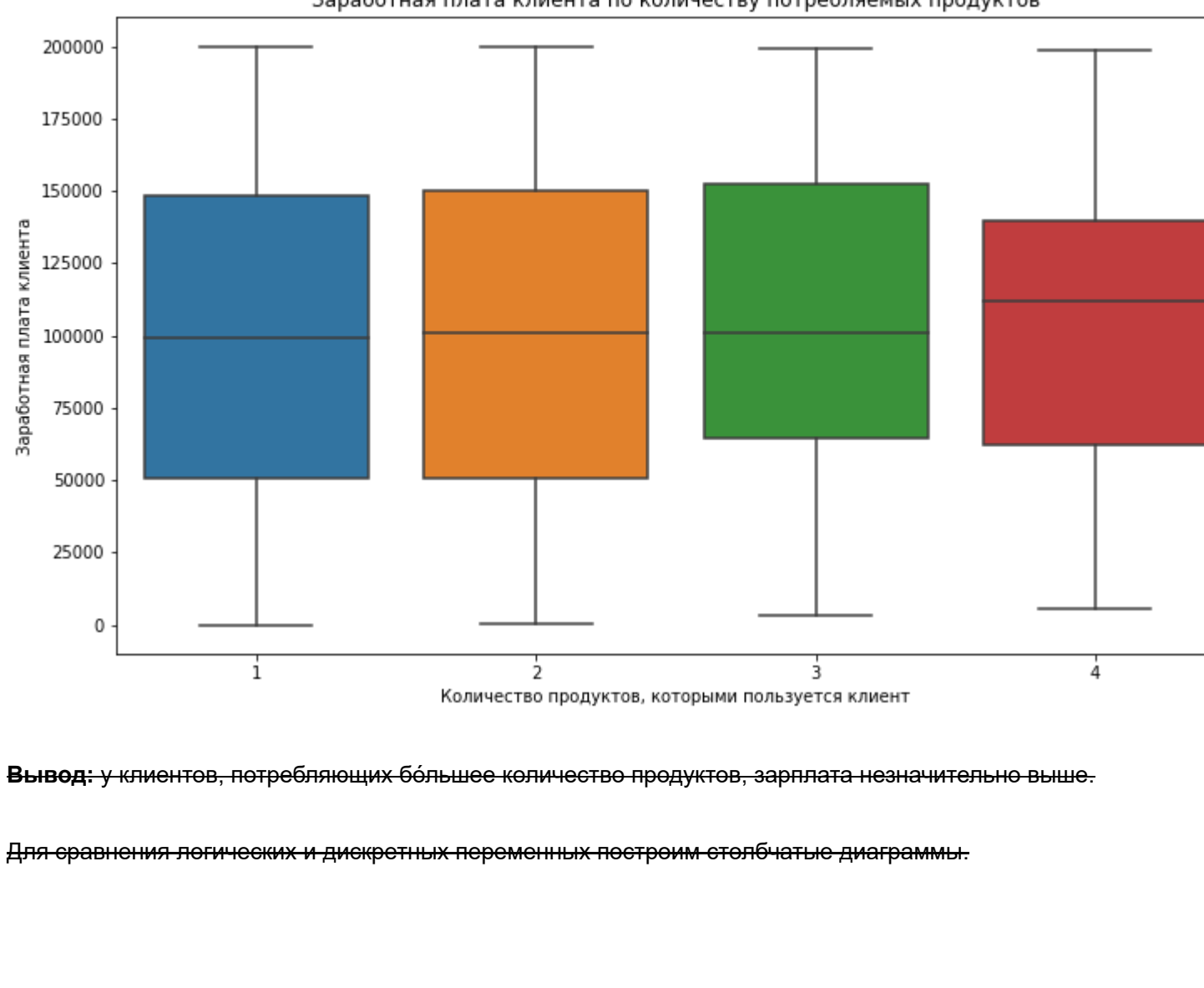
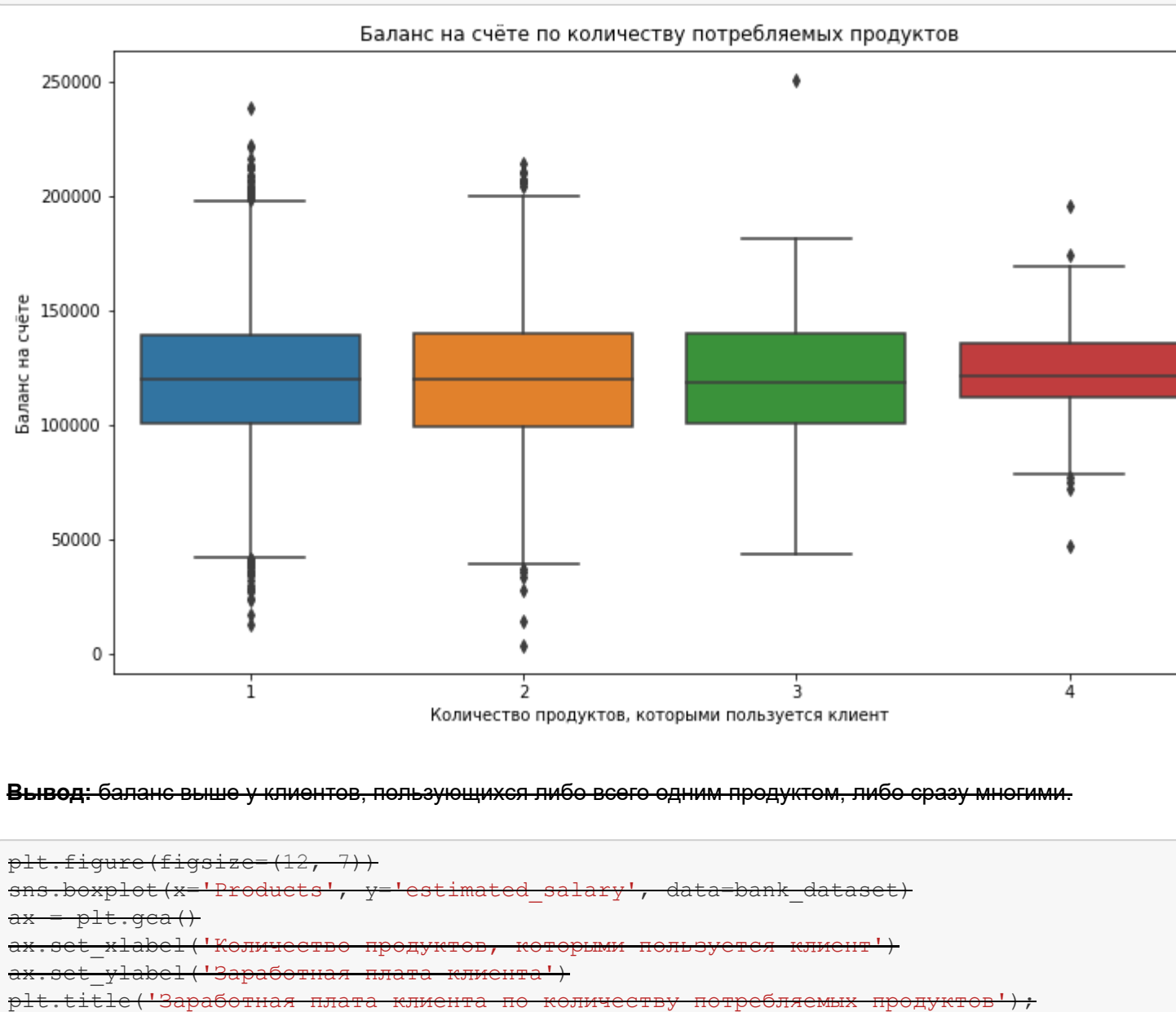
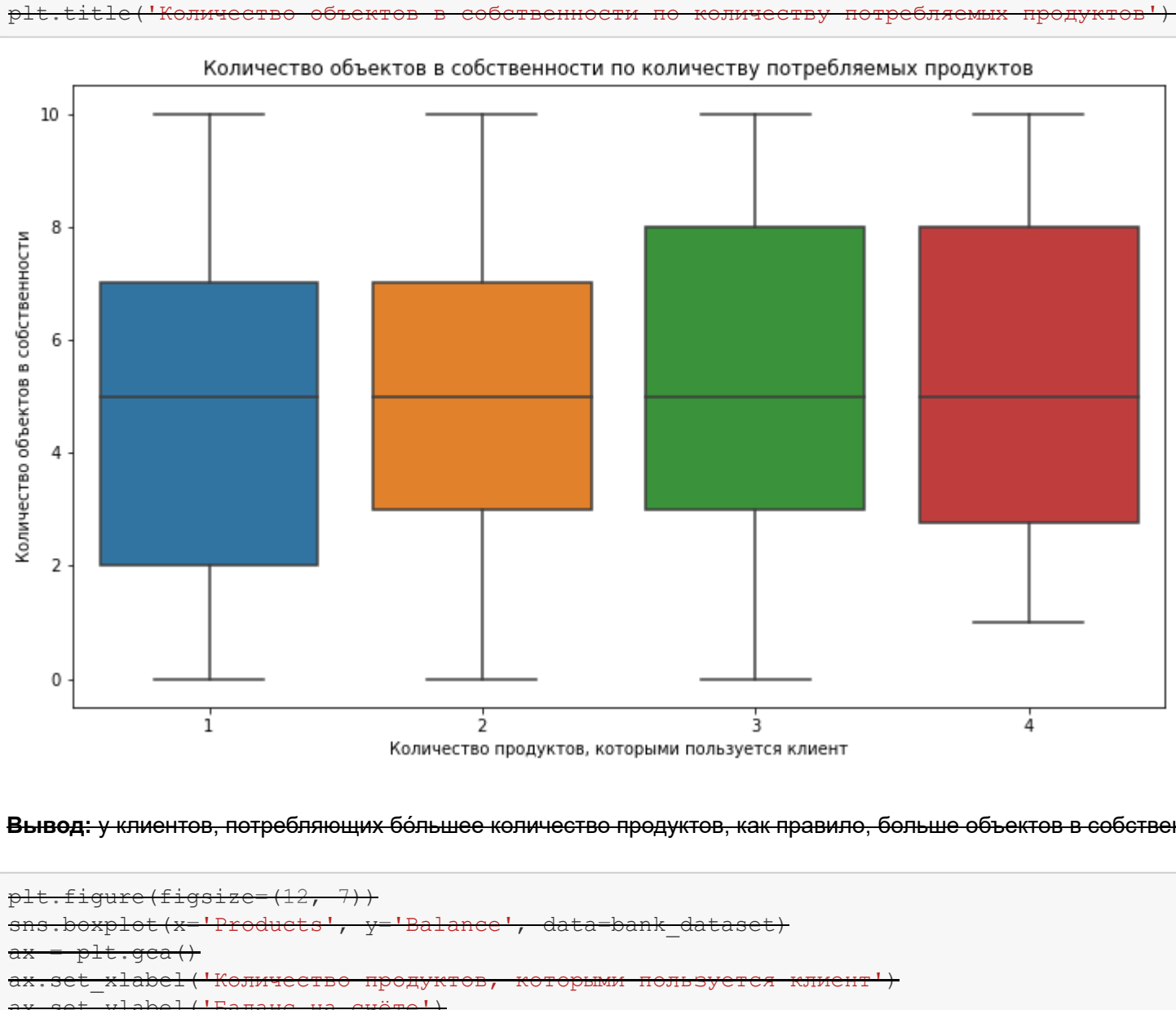
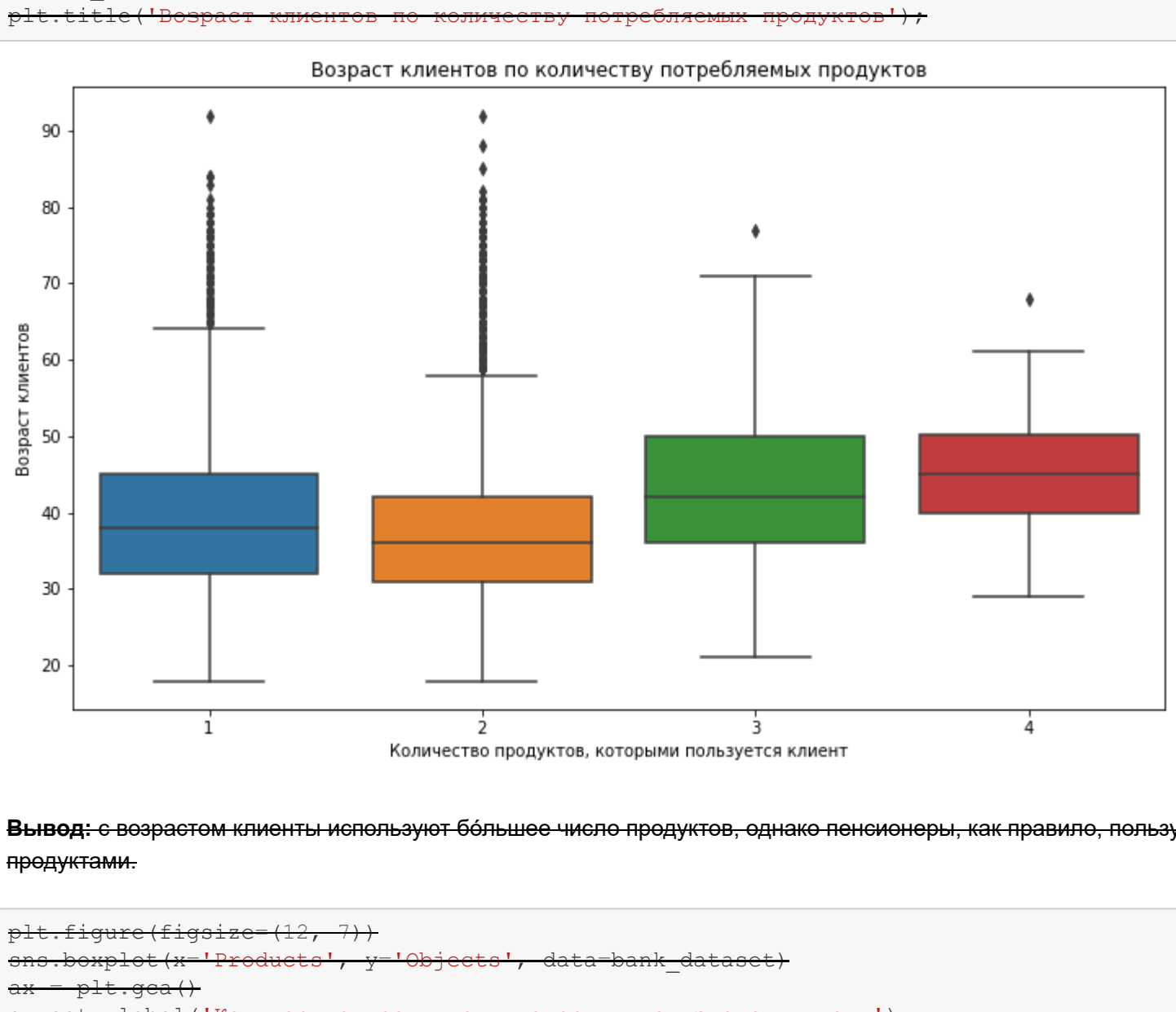
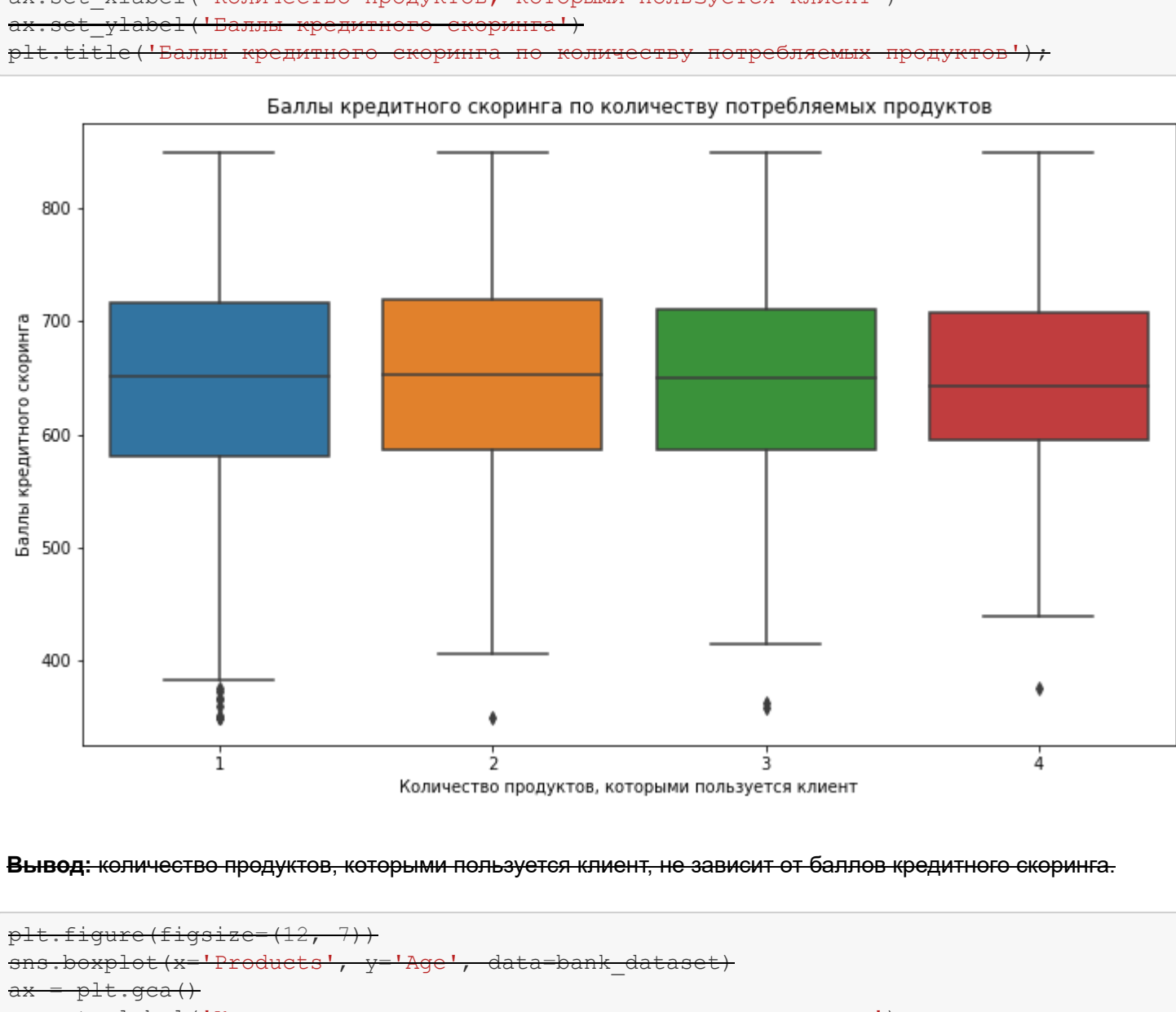
```
from scipy.stats import norm
plt.figure(figsize=(10,7))
plt.hist(bank_dataset['estimated_salary'], bins=100)
plt.show()
```

Гистограмма распределения заработной платы клиентов



```
bank_dataset['score'].hist(figsize=(10,7), bins=100)
plt.show()
```

Гистограмма распределения баллов кредитного скоринга



Оценим, как распределены клиенты по городам:

```
bank_dataset['city'].value_counts()
```

Город: Ростов Великий, Рыбинск, Ярославль

Доли клиентов по городам

Вывод: клиентов в Ярославле в два раза больше, чем в Ростове Великом и Рыбинске.

Построим матрицу корреляций и отобразим её при помощи тепловой карты:

```
bank_dataset[['score', 'Age', 'Objects', 'Balance', 'Products', 'CreditCard', 'Loyalty', 'estimated_salary', 'Churn']].corr()
```

Тепловая карта матрицы корреляций признаков

Вывод: сильно коррелирующие признаки отсутствуют. Самые сильные связи наблюдаются между оттоком и возрастом, количеством используемых продуктов и балансом на счёте.

Оглавление

Шаг 2. Сегментация пользователей на основе данных о количестве потребляемых продуктов

Сгруппируем пользователей по количеству потребляемых продуктов и сравним параметры каждой группы методом `boxplot()`.

```
bank_dataset.groupby('Products').boxplot()
```

Баллы кредитного скоринга по количеству потребляемых продуктов

Вывод: количество продуктов, которыми пользуется клиент, не зависит от баллов кредитного скоринга.

```
bank_dataset.groupby('Age').boxplot()
```

Возраст клиентов по количеству потребляемых продуктов

Вывод: с возрастом клиенты используют большее число продуктов, однако пенсионеры, как правило, пользуются одним-двумя продуктами.

```
bank_dataset.groupby('Objects').boxplot()
```

Количество объектов в собственности по количеству потребляемых продуктов

Вывод: у клиентов, потребляющих большее количество продуктов, как правило, больше объектов в собственности.

```
bank_dataset.groupby('Balance').boxplot()
```

Баланс на счёте по количеству потребляемых продуктов

Вывод: баланс выше у клиентов, пользующихся либо одним продуктом, либо сразу многими.

```
bank_dataset.groupby('estimated_salary').boxplot()
```

Зарплата клиента по количеству потребляемых продуктов

Вывод: у клиентов, потребляющих большее количество продуктов, зарплата незначительно выше.

Для сравнения признаков и дисперсии переменных построим scatter-диаграммы.


```

[3]: pd.concat([bank_dataset[bank_dataset['City'] == 'Иркутск'].groupby('Product').agg({'Product': 'count', 'City': 'count'}),
            bank_dataset[bank_dataset['City'] == 'Иркутск'].groupby('Product').agg({'Product': 'count', 'City': 'count'})
        ].reset_index()

[4]: bank_dataset[bank_dataset['City'] == 'Иркутск'].groupby('Product').agg({'Product': 'count', 'City': 'count'})

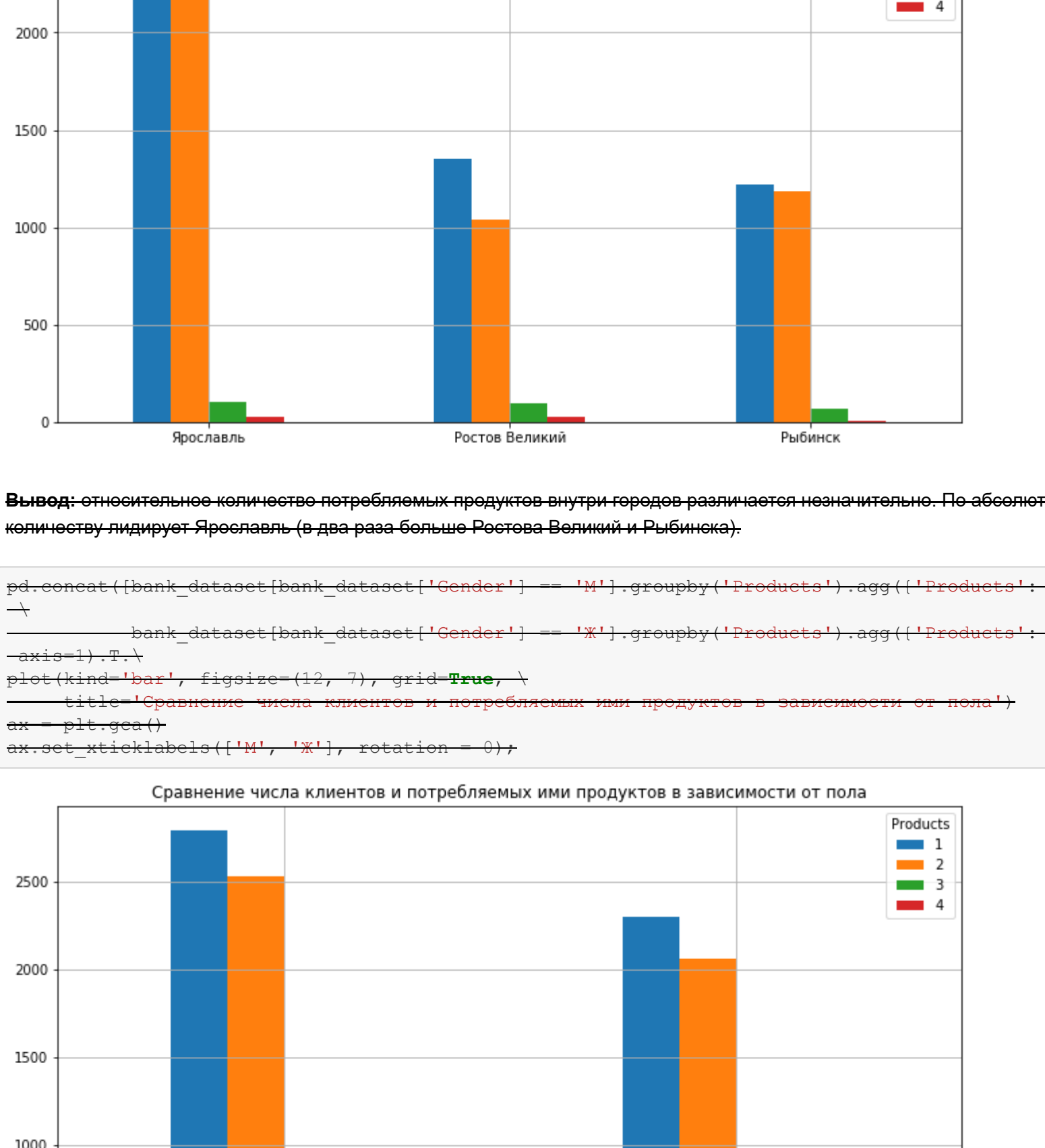
[5]: pd.concat([bank_dataset[bank_dataset['City'] == 'Иркутск'].groupby('Product').agg({'Product': 'count', 'City': 'count'}),
            bank_dataset[bank_dataset['City'] == 'Иркутск'].groupby('Product').agg({'Product': 'count', 'City': 'count'})
        ].reset_index()

[6]: plot(kind='bar', figsize=(10, 7), grid=True, title='Сравнение числа клиентов и потребляемых ими продуктов по городам')

[7]: ax.set_xticklabels(['Иркутск', 'Иркутск', 'Иркутск', 'Иркутск', 'Иркутск', 'Иркутск', 'Иркутск', 'Иркутск', 'Иркутск', 'Иркутск'])

```

Category	Product 1 (blue)	Product 2 (orange)
Иркутск	2200	1800
Иркутск	2000	1600
Иркутск	1800	1400
Иркутск	1600	1200
Иркутск	1400	1000
Иркутск	1200	800
Иркутск	1000	600
Иркутск	800	400
Иркутск	600	200
Иркутск	400	100



Category	Male	Female
No	~550	~550
Yes	~550	~550
Don't know	~100	~100

Вывод: количество потребляемых продуктов не зависит от пола. Общее число клиентов-мужчин незначительно превышает число клиентов-женщин.

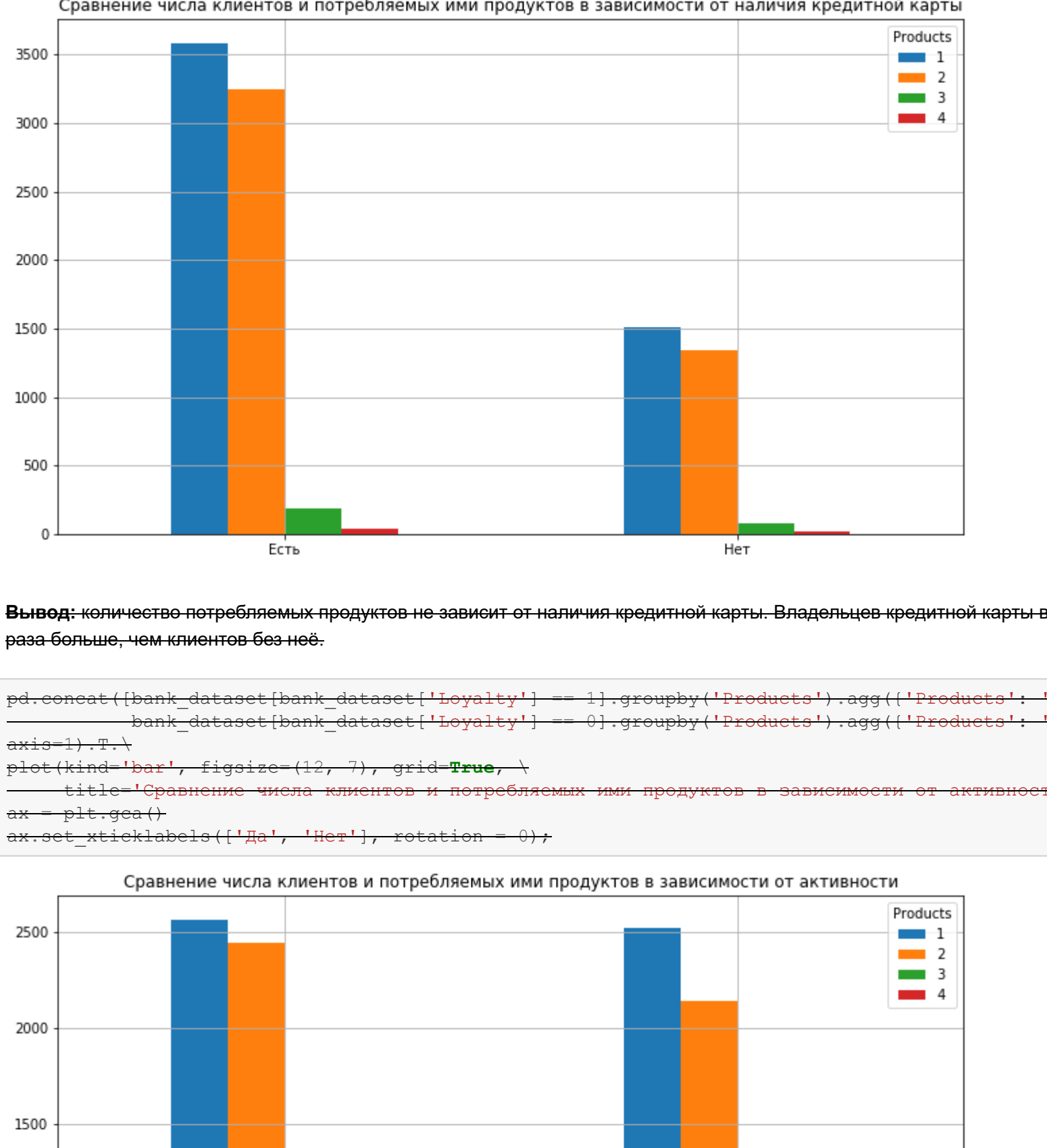
```

from itertools import groupby
from collections import Counter

bank_data = bank_data[bank_data['gender'] != 'other']

for gender, group in groupby(bank_data, lambda x: x['gender']):
    print(f'Gender: {gender}')
    count = Counter()
    for product in group['products']:
        count[product] += 1
    print(count)

```



Активность клиента	Категория 1 (Blue)	Категория 2 (Orange)	Категория 3 (Green)
Да	10	8	3
Нет	10	8	3

Вывод: количество потребляемых продуктов не зависит от активности клиента

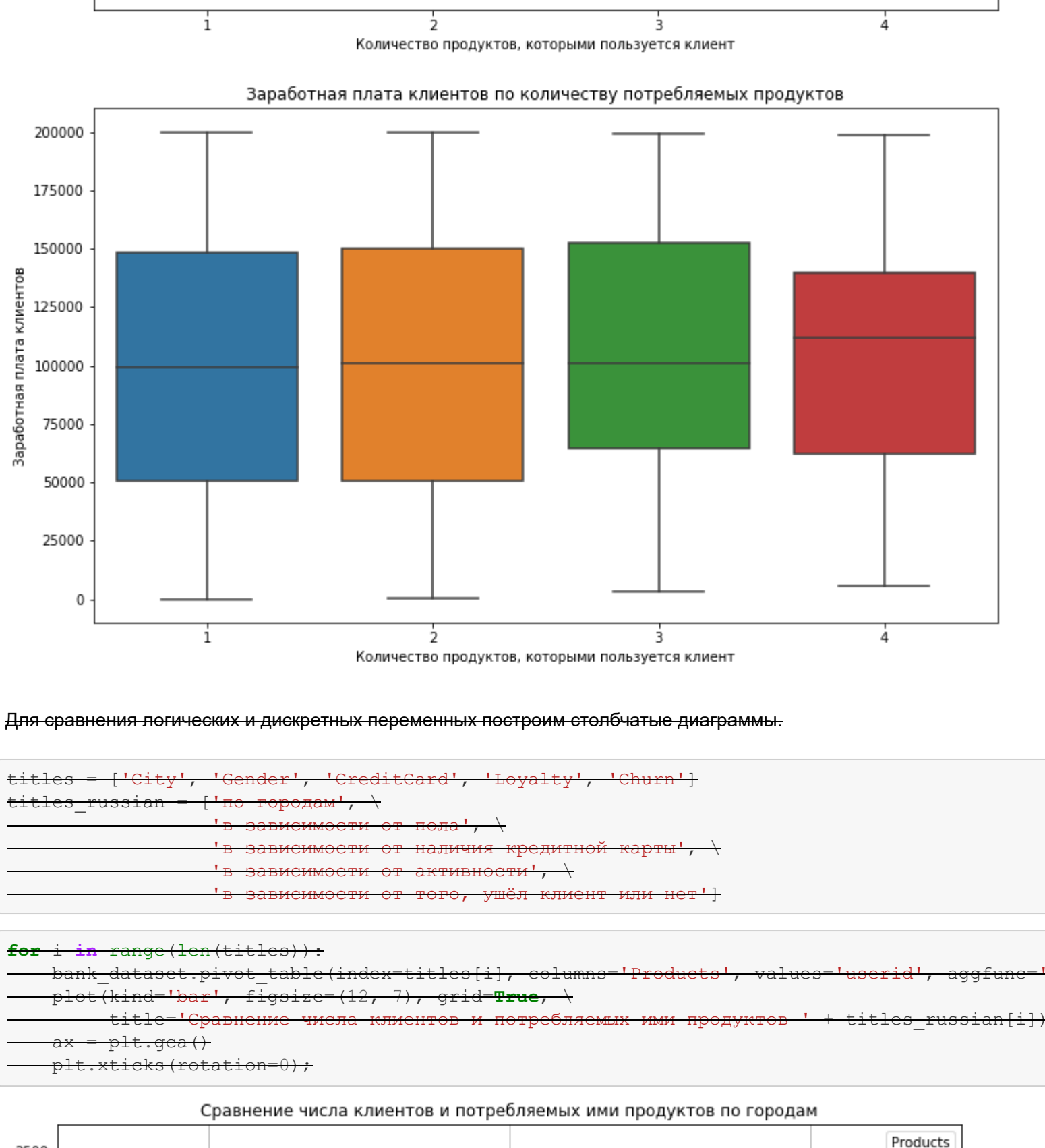
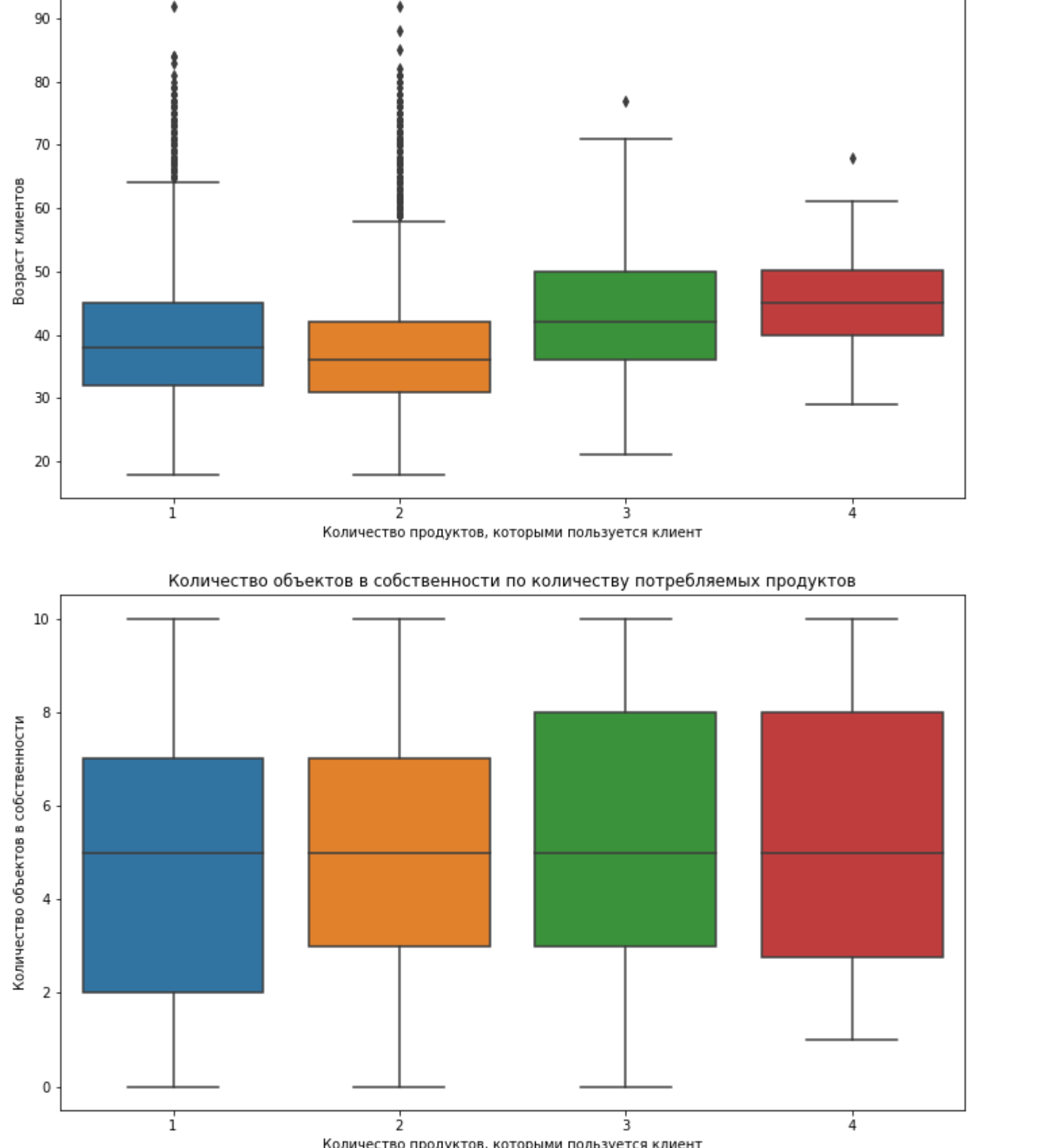
```

In [122]: pd.concat([bank_data, credit_data], axis=1).groupby(['ProductType', 'Activity']).count().unstack().fillna(0)
           bank_data credit_data
ProductType Activity
1             10         8
2             10         8
3              3         3

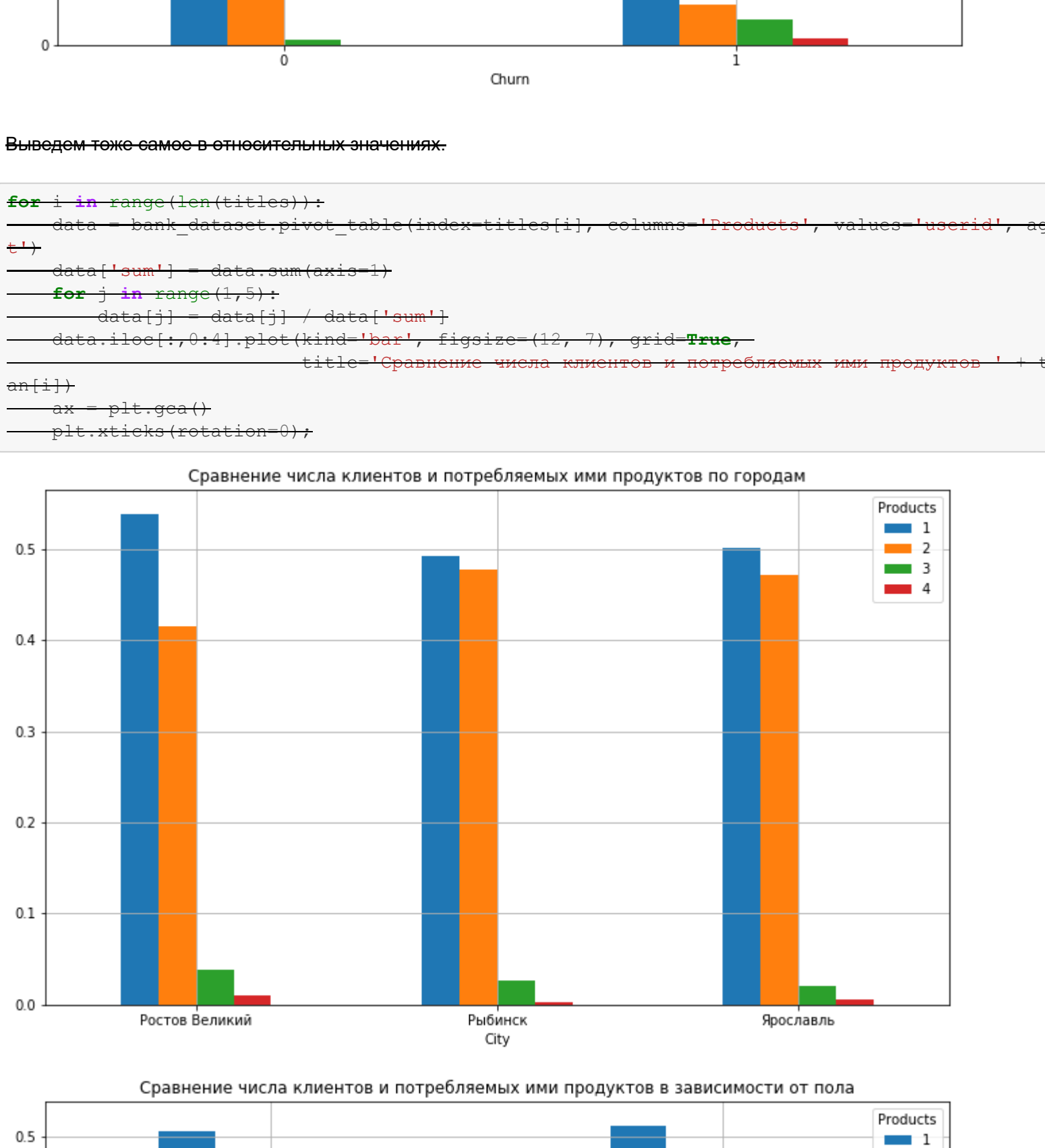
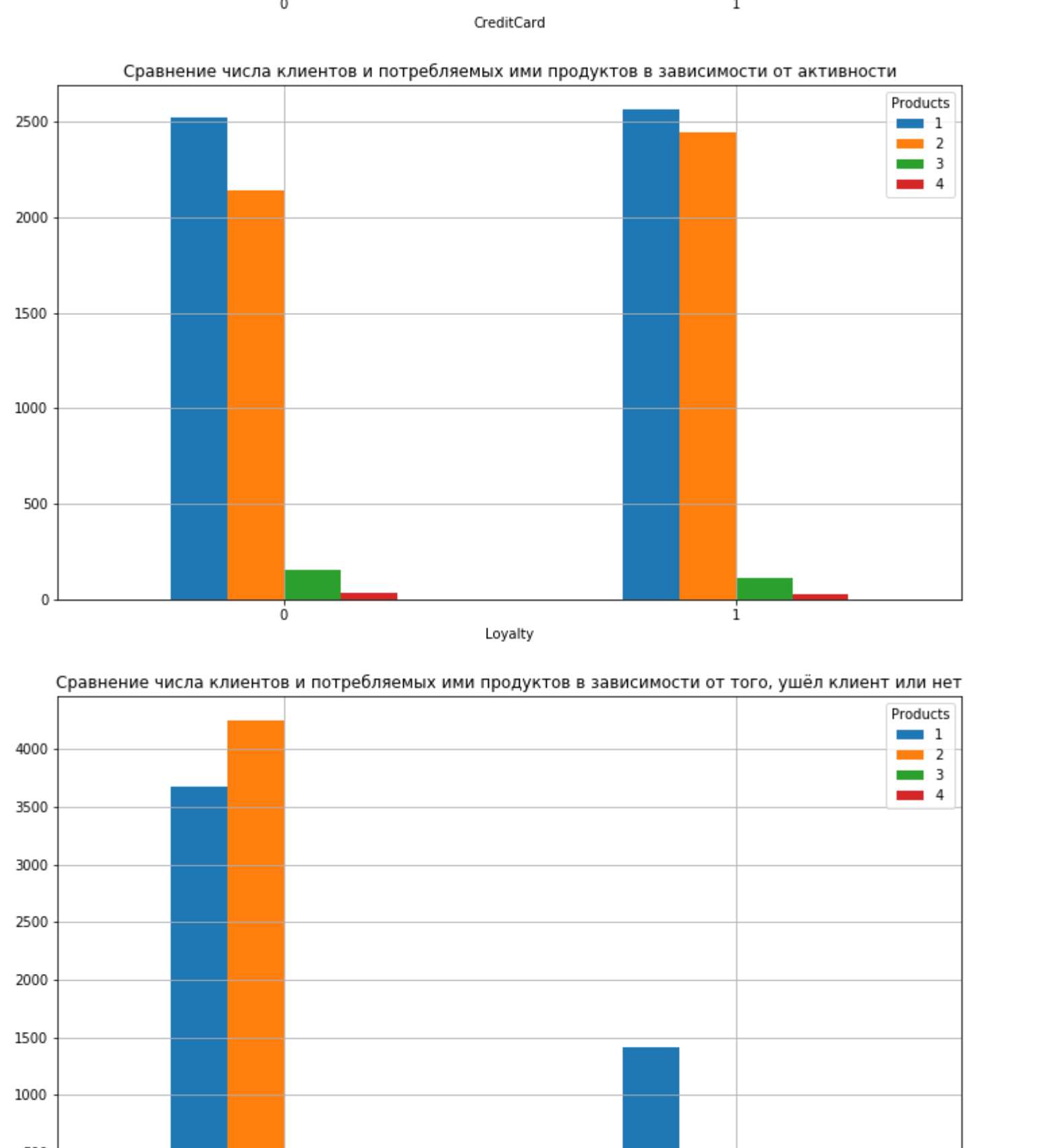
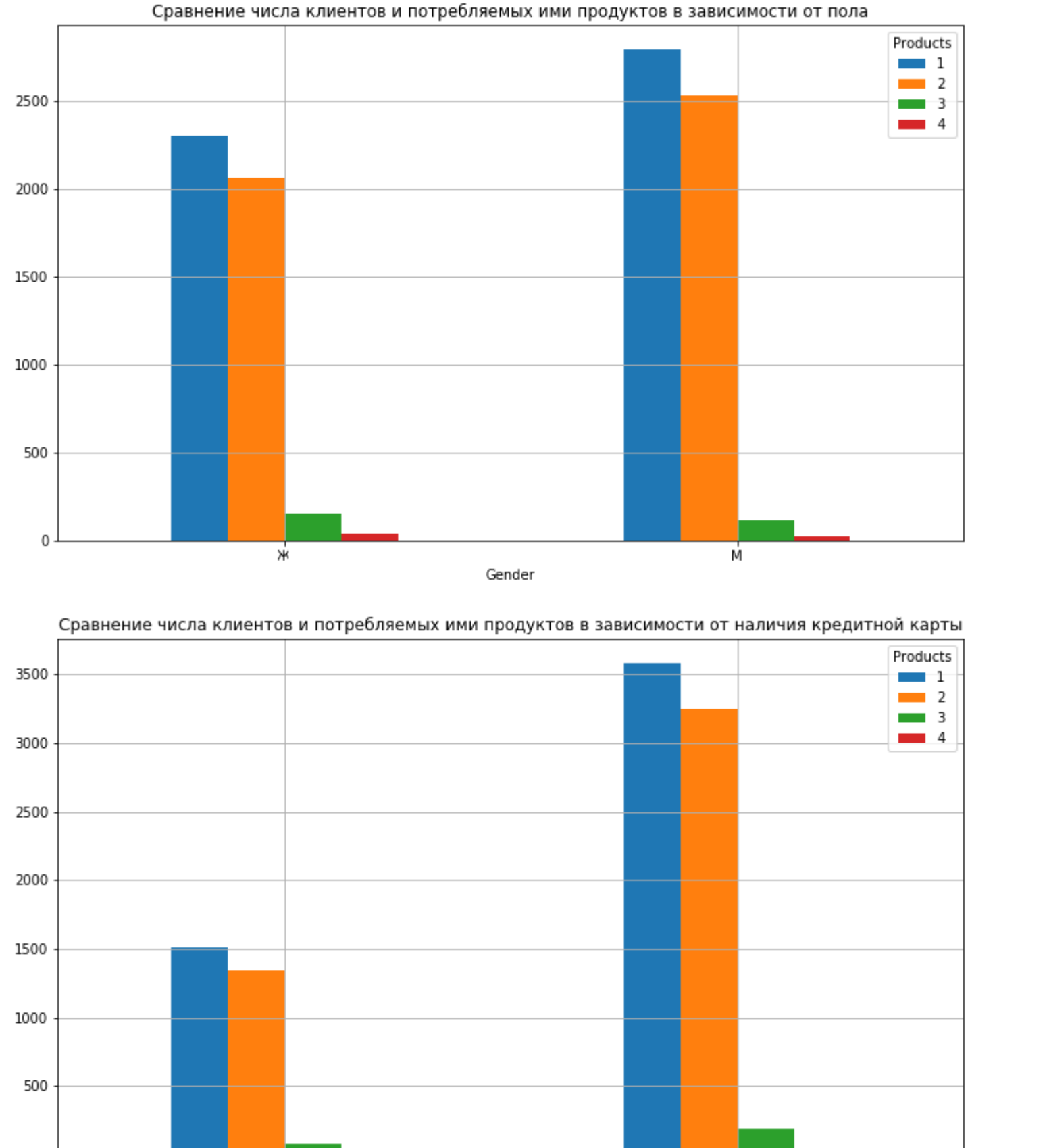
```



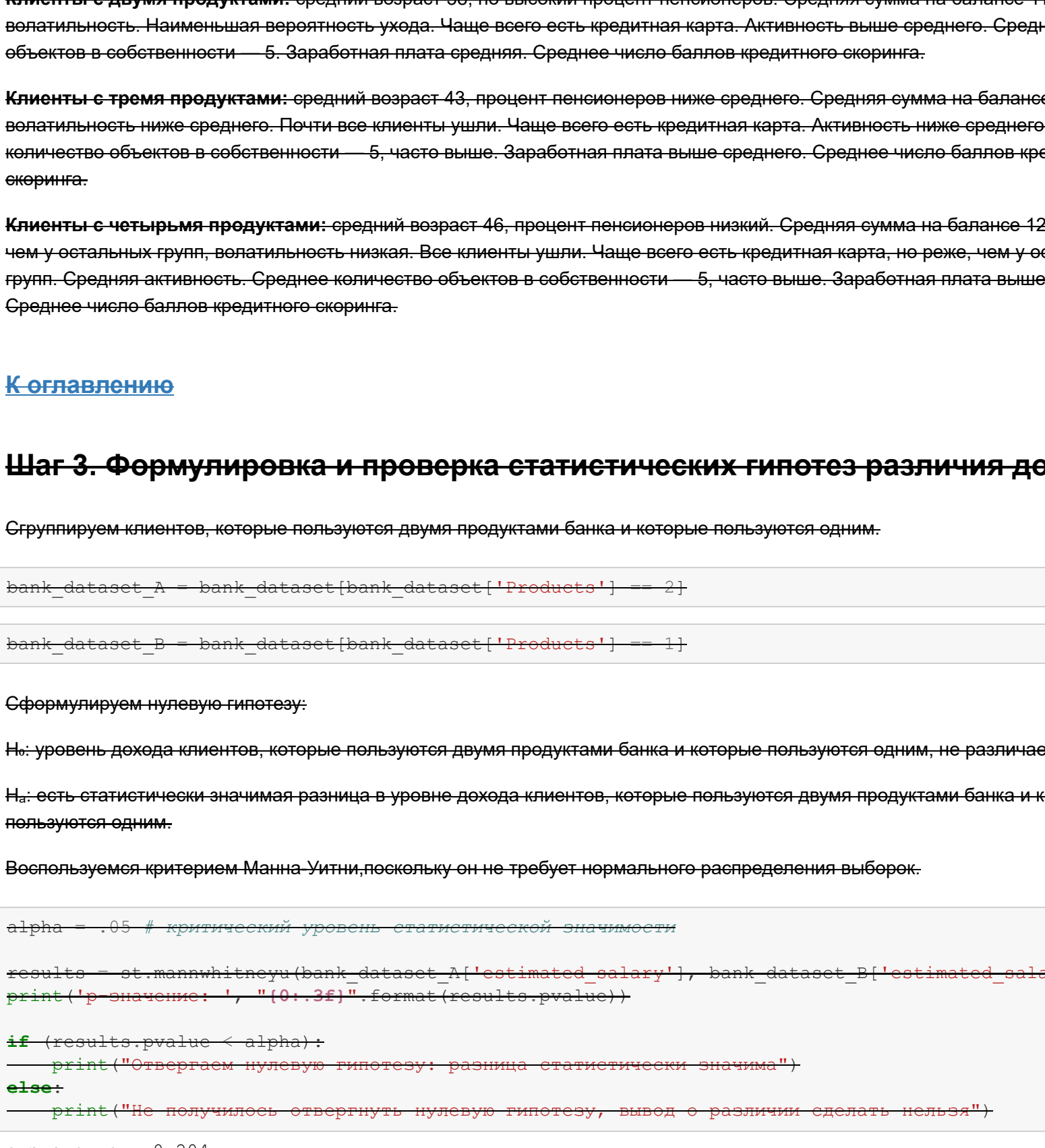
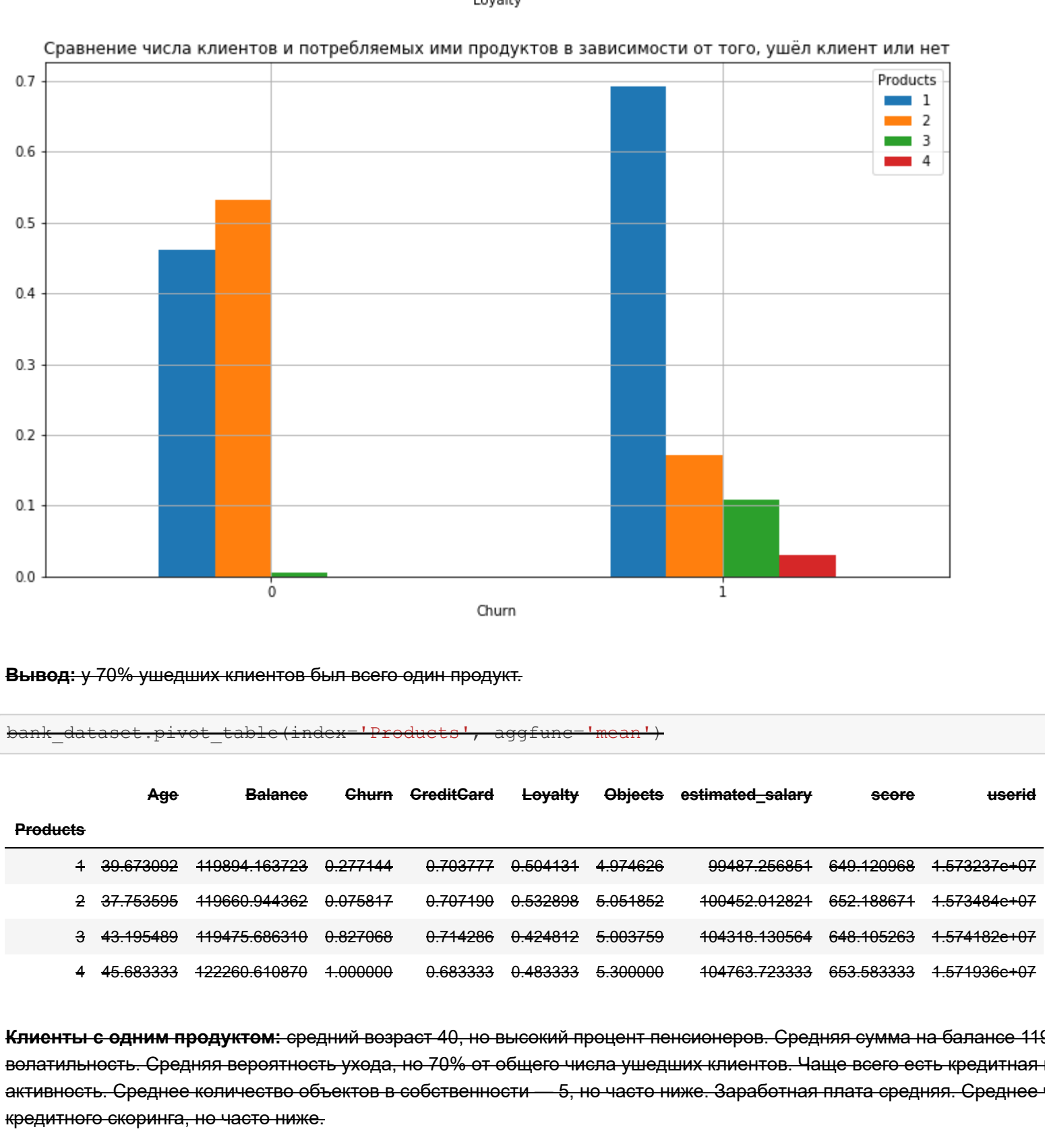
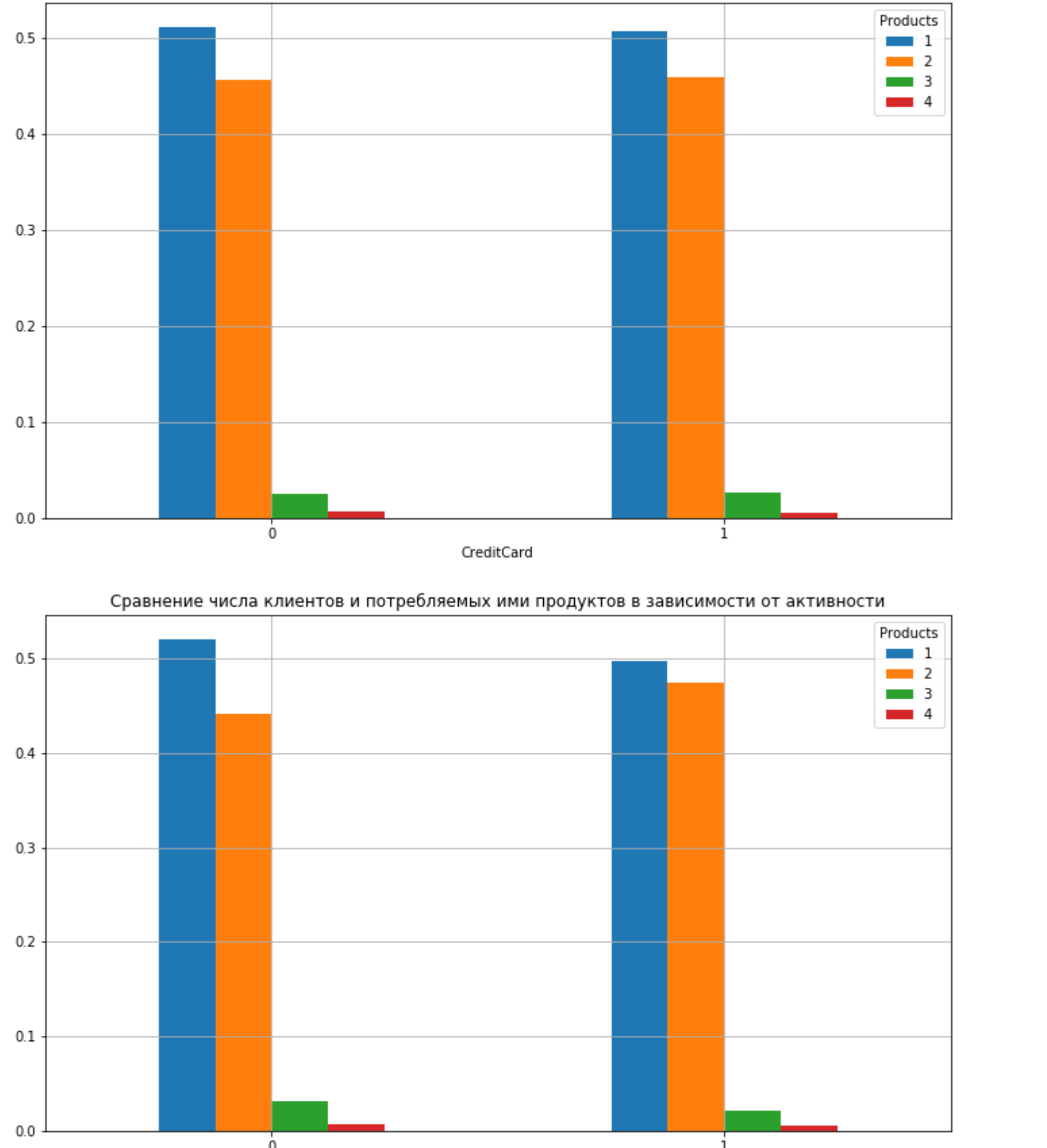
Возраст клиентов по количеству потребляемых продуктов



Category	Blue Bar (Approx. Count)	Orange Bar (Approx. Count)	Green Bar (Approx. Count)
Protein Structure	1350	1050	100
Protein City	1200	1150	100
Biochemicals	1650	1600	100



Category	No (Blue)	Yes (Orange)
Before	0.65	0.35
After	0.65	0.35



Вывод: уровень дохода клиентов, которые пользуются д

```
bank_dataset %>% 'estimated salary'.mean()
```

```

--> 300450-01020135073
-->
dim<-dim(bank_dataset)[!is.na(bank_dataset$y)]-1
-->
Oval<-Oval(300450-01020135073)
-->

```

Проверим, различаются ли балансы клиентов, которые используют одну продукцию банка и которые используют одним:
 Сформируем новую переменную:

Н₁-баланс на счетах клиентов, которые используют двумя продуктами банка и которые используют одним; не различается.
 Н₂-есть статистически значимая разница в балансах счетов клиентов, которые используют двумя продуктами банка и которые используют одним:

[illegible]

```
bank_data = df[['balance']], mean()
+-----+-----+
+0894-16372337679
```

Средний баланс на счет клиента, который пользуется двумя продуктами банка, в полтора раза выше среднего баланса на счетах клиентов, которые пользуются одним продуктом.

Количество продуктов, которыми пользуется клиент, не зависит от баллов кредитного скоринга.

С возрастом клиенты используют большее число продуктов, однако пенсионеры, как правило, пользуются одним-двумя продуктами.

У клиентов, потребляющих большее количество продуктов, как правило, больше сбережений в собственности.

Баланс выше у клиентов, пользующихся либо одним продуктом, либо сразу многими.

У клиентов, потребляющих большее количество продуктов, зарплата noticeably выше.

Относительное количество потребляемых продуктов скоринг скоринг различает незначительно. По абсолютному количеству

Количество потребляемых продуктов не зависит от пола. Общее число клиентов-мужчин ненамного превышает число клиентов-женщин.

Количество потребляемых продуктов не зависит от наличия кредитной карты. Владетель кредитной карты в два с лишним раза больше, чем клиент без нее.

Количество потребляемых продуктов не зависит от активности клиента.

У 70% владельцев клиентов был всего один продукт.

[К оглавлению](#)