

Подключаем необходимые для обработки данных библиотеки

```
In [1]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.patches as mpatches
import numpy as np
```

Читаем и смотрим датасет

```
In [2]: path = '_data.csv'
pd.set_option('display.max_columns', 30)
#pd.set_option('display.max_info_columns', 30)
#pd.set_option('max_info_columns', 30)
r1st = pd.read_csv(path, index_col=0)
with pd.option_context('display.max_colwidth', 25):
    display(r1st)
```

	ID	Количество объявления	Количество комнат	Тип	Метро	Адрес	Площадь, м2	Дом	Парков
0	271271157	4	Квартира	м. Смоленский (9 мин ...	Москва, улица Новый А...	200.0/20.0	5/16,	подземн	
1	271634126	4	Квартира	м. Смоленский (8 мин ...	Москва, улица Новый А...	198.0/95.0/18.0	5/16,	подземн	
2	271173086	4, 0ба варианта	Квартира	м. Смоленская (7 мин ...	Москва, улица Новый А...	200.0/116.0/4.0	5/16	подзем	
3	272197456	4, 0ба варианта	Квартира	м. Смоленская (3 мин ...	Москва, перуелок Плот...	170.0/95.0/17.0	5/6	подзем	
4	273614615	2	Квартира	м. Арбатская (7 мин п...	Москва, улица Новый А...	58.0/38.0/5.0	12/26, Панельный		N
...	
23363	215565511	NaN	Квартира	м.Говорово (8 мин пе...	Москва, Боровское шос...	35.0/16.4/8.0	10/14		N
23364	274654844	1	Квартира	м. Солнцево (7 мин пе...	Москва, Производствен...	38.7/16.5/11.0	5/18, Монолитный		N
23365	268679909	2, 0ба варианта	Квартира	м. Солнцево (6 мин пе...	Москва, Боровский про...	43.1	5/5, Кирпичный		N
23366	274807526	2	Квартира	м. Солнцево (11 мин п...	Москва, улица Богдано...	52.5/10.0	8/23, Монолитный	назем	
23367	274672423	2, 0ба варианта	Квартира	м.Говорово (10 мин п...	Москва, улица Богдано...	90.0/48.2/15.5	6/23, Панельный		N

23368 rows x 24 columns

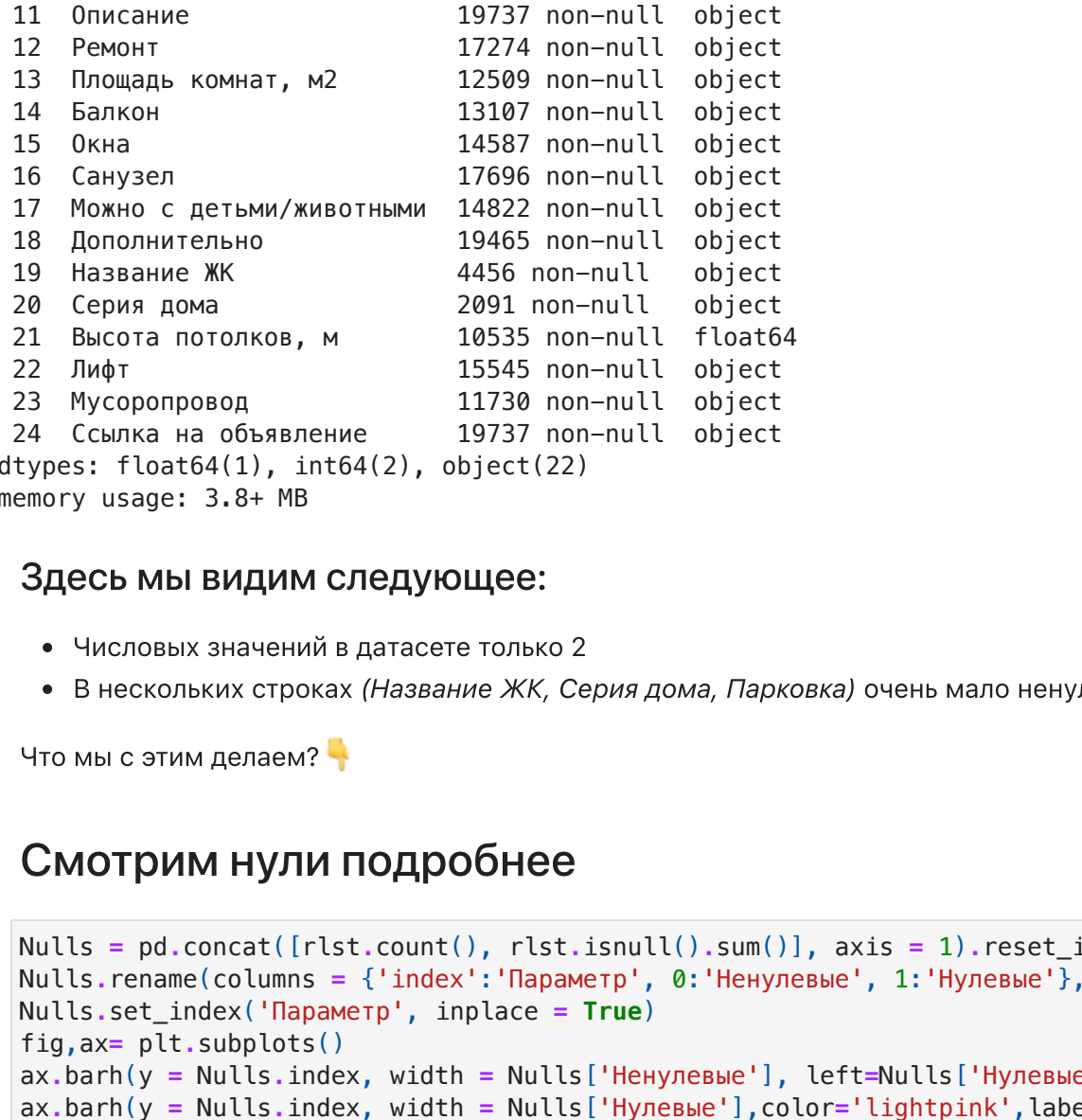
В датареиме мы видим **24 столбца** с данными. Можно сразу обратить внимание на следующее:

- В некоторых столбцах присутствуют **NaN (not a number)** значения
- Часть данных, которую хотелось бы видеть в числовом формате, представлена в виде строк, не состоящих полностью из цифр (**Количество комнат, Расстояние до метро, Цена, Площадь, Площадь комнат**)
- В части строк представлены данные, которые стоило бы разделить на отдельные столбцы для удобства обработки (**Количество комнат, Метро, Цена, Площадь, Площадь комнат, Лифт, Санузел**)
- В таблице присутствуют данные, которые очевидно не могут влиять на стоимость квартиры (**ID, Ссылка**)

Москва?

В качестве базового региона выбрана **Москва**, поэтому посмотрим, есть ли у нас квартиры **не в Москве**

```
In [3]: config InlineBackend.figure_format='retina'
bars=sns.countplot(x=r1st['Адрес'].apply(lambda x: x.split(sep=',')[0])!= 'Москва').index, inplace=True)
r1st.reset_index(inplace=True)
print("Размер датареиме только с Москвой стан (r1st.shape)")
for i in bars.containers[0]:
    bars_bar_label(i)
```



Видно, что в датасете **3631** квартира не в Москве. Что делать? 🤔

Удаляем всё кроме Москвы

```
In [4]: r1st.drop(r1st.loc[r1st['Адрес'].apply(lambda x: x.split(sep=',')[0])!= 'Москва'].index, inplace=True)
r1st.reset_index(inplace=True)
print("Размер датареиме только с Москвой стан (r1st.shape)")
```

Размер датареиме только с Москвой стан (19737, 25)

Инфо

Посмотрим общую информацию обо всех столбцах нашего набора данных, чтобы узнать их **тип** и **количество ненулевых значений**

```
In [5]: formatted_r1st = r1st.copy()
r1st.info(show_counts=True)

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 19737 entries, 0 to 19736
Data columns (total 25 columns):
# Column Non-Null Count Dtype
---
0 index 19737 non-null int64
1 ID объявления 19737 non-null int64
2 Количество комнат 19202 non-null object
3 Тип 19737 non-null object
4 Метро 19391 non-null object
5 Адрес 19737 non-null object
6 Площадь, м2 19737 non-null object
7 Дом 19737 non-null object
8 Парковка 8563 non-null object
9 Цена 19737 non-null object
10 Телефоны 19737 non-null object
11 Описание 19737 non-null object
12 Ремонт 17274 non-null object
13 Площадь комнат, м2 12509 non-null object
14 Балкон 13187 non-null object
15 Окна 14587 non-null object
16 Санузел 17696 non-null object
17 Можно с детьми/животными 14822 non-null object
18 Дополнительно 19465 non-null object
19 Название ЖК 4456 non-null object
20 Серия дома 2091 non-null object
21 Высота потолков, м 10535 non-null float64
22 Лифт 15545 non-null object
23 Мусоропровод 11730 non-null object
24 Ссылка на объявление 19737 non-null object
dtypes: float64(1), int64(2), object(22)
memory usage: 136.0 MB
```

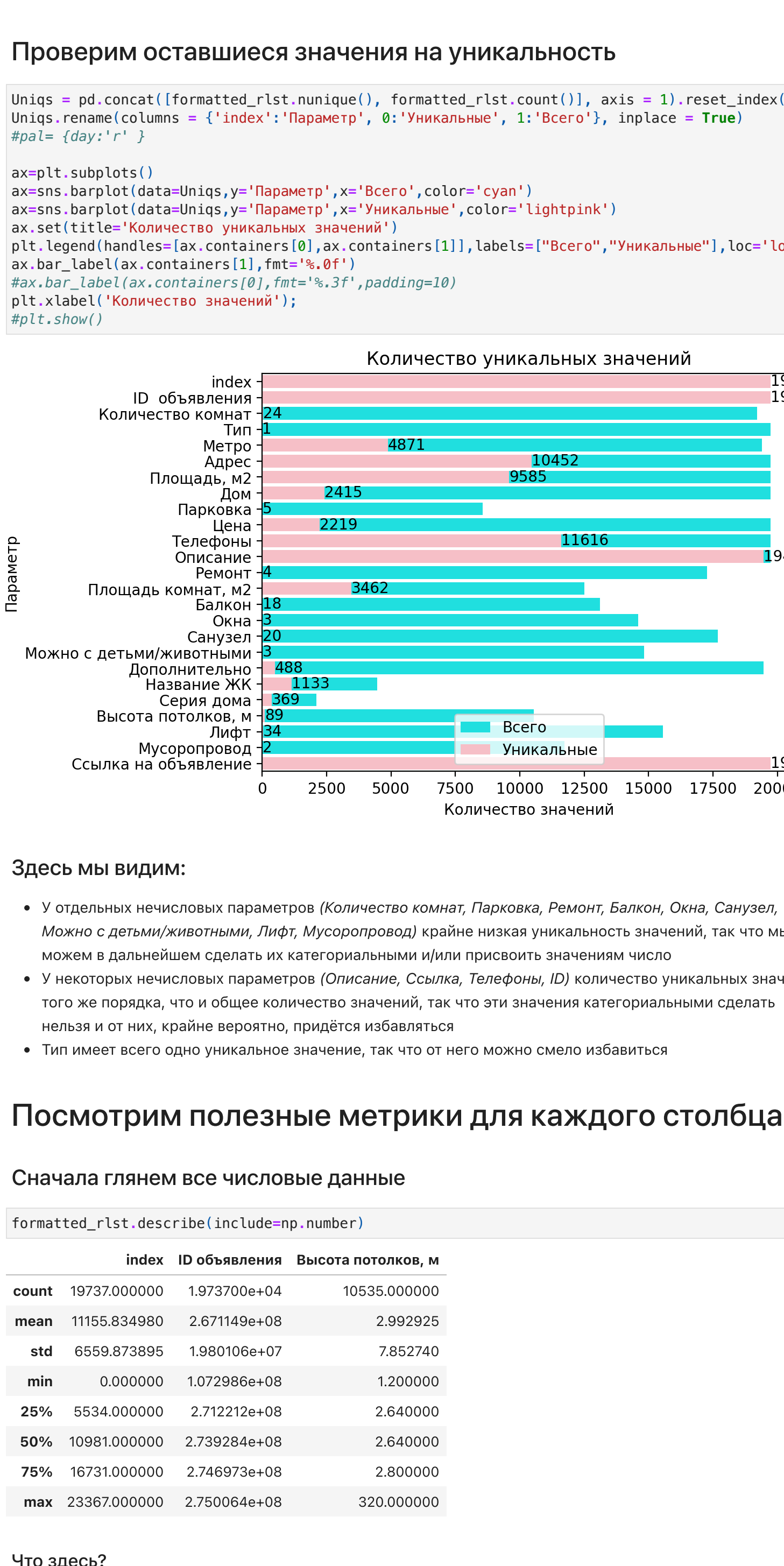
Здесь мы видим следующее:

- Числовых значений в датасете только 2
- В нескольких строках (**Название ЖК, Серия дома, Парковка**) очень мало ненулевых значений

Что мы с этим делаем? 🤔

Смотрим нули подробнее

```
In [6]: Nulls = pd.concat([r1st.count(), r1st.isnull().sum()], axis = 1).reset_index()
Nulls.rename(columns = {'index': 'Параметр', 0: 'Ненулевые', 1: 'Нулевые', inplace = True)
fig=plt.subplot(1,1,figsize=(10,5))
ax=Nulls.groupby('Параметр').sum(numeric_only=True)
ax.barh(y = Nulls.index, width = Nulls['Ненулевые'], left=Nulls['Нулевые'],color='cyan',label='Нуле')
ax.barh(y = Nulls.index, width = Nulls['Нулевые'],color='lightpink',label='Нулевые')
ax.set_title('Соотношение нулевых и ненулевых значений',label='Кол-во значений')
plt.legend()
```

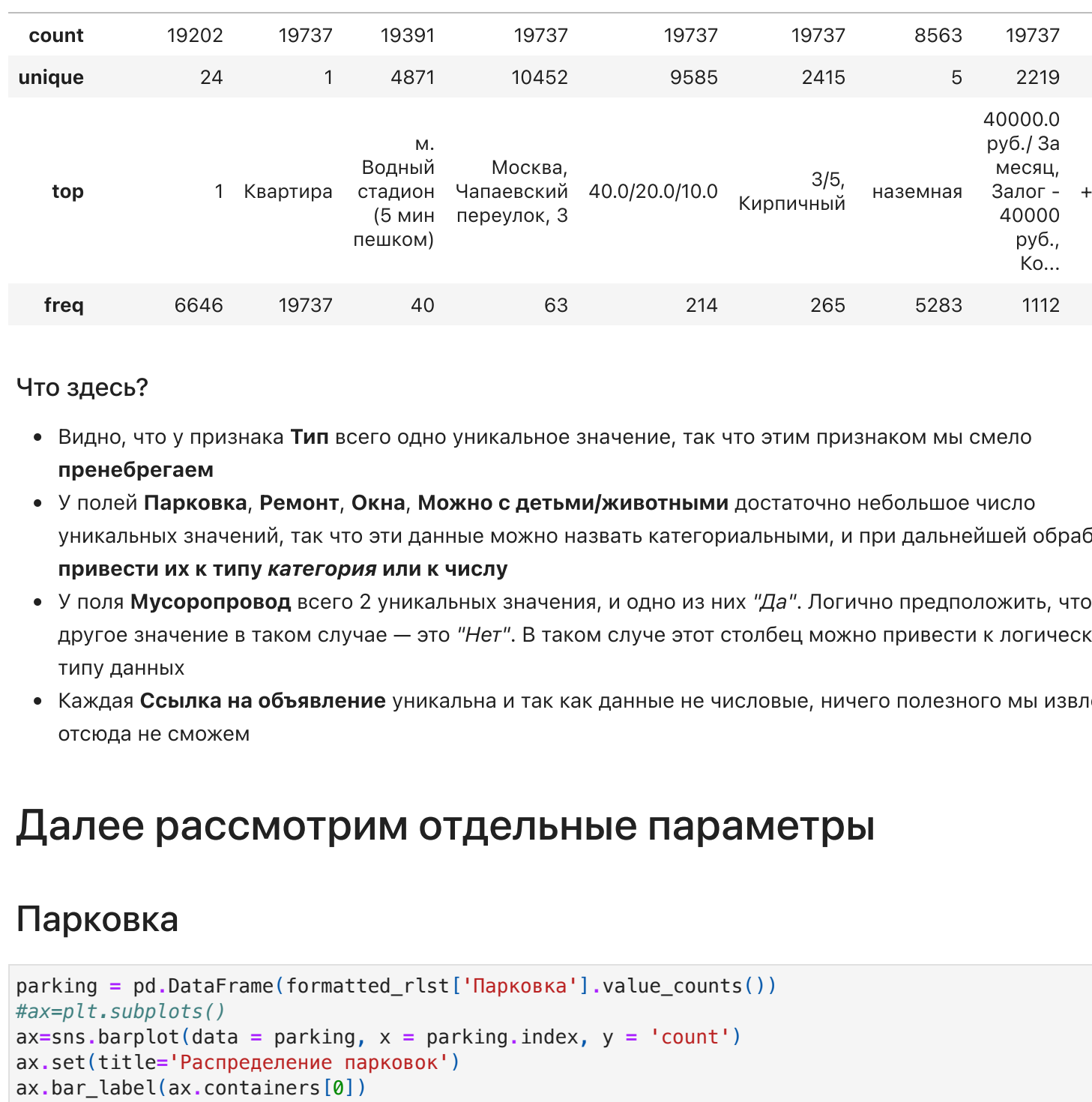


Что мы видим?

Видим, что у **названий ЖК и серии дома** больше значений нулевых, так что эти столбцы в дальнейшем можем смело выкинуть, к тому же мы полагаем их маловероятна корреляция с ценой

Проверим оставшиеся значения на уникальность

```
In [7]: Uniqs = pd.concat([formatted_r1st.nunique(), formatted_r1st.count()], axis = 1).reset_index()
Uniqs.rename(columns = {'index': 'Параметр', 0: 'Уникальные', 1: 'Всего'}, inplace = True)
#pal = (day::r')
ax=plt.subplot(1,1,figsize=(10,5))
ax=Uniqs.groupby('Параметр').sum(numeric_only=True)
ax=Uniqs.groupby('Параметр').sum(numeric_only=True)
ax.set_title('Количество уникальных значений')
plt.legend(handles=[ax.containers[0],ax.containers[1]],labels=['Всего','Уникальные'],loc='lower center')
ax.set_xlabel('Количество значений')
plt.xlabel('Количество значений')
plt.show()
```



Здесь мы видим:

- У отдельных чиселых параметров (**Количество комнат, Парковка, Ремонт, Балкон, Окна, Санузел, Можно с детьми/животными, Лифт, Мусоропровод**) крайне низкая уникальность значений, так что мы можем в дальнейшем сдлать их категориальными и/или присвоить значениям число
- У некоторых числовых параметров (**Описание, Ссылка, Телефоны, ID**) количество уникальных значений того же порядка, что и общее количество значений, так что эти значения категориальными сделать нельзя и от них, крайне вероятно, придётся избавиться
- Тип имеет всего одно уникальное значение, так что от него можно смело избавляться

Посмотрим полезные метрики для каждого столбца

Сначала глянем все числовые данные

```
In [8]: formatted_r1st.describe(include=np.number)

Out[8]:
```

	index	ID объявления	Высота потолков, м
count	19737.000000	1.973700e+04	1.053500e+00
mean	11155.834980	2.671149e+08	2.992925
std	6569.873895	1.980106e+07	7.852740
min	0.000000	1.078288e+08	1.200000
25%	5534.000000	2.712212e+08	2.640000
50%	10981.000000	2.739284e+08	2.640000
75%	16731.000000	2.746973e+08	2.800000
max	23367.000000	2.750064e+08	320.000000

Что здесь?

- Индекс и ID объявления для нас интереса не представляют
- А вот высота потолков нас интересует. Тут мы можем заметить, что максимальная и минимальная высоты неадекватны, поэтому глянем, насколько много таких выбросов

Потолок

```
In [9]: print("Top 5 самых низких потолков:\n",formatted_r1st['Высота потолков, м'].sort_values().head(),'\n')
print("Top 10 самых высоких потолков:\n",formatted_r1st['Высота потолков, м'].sort_values(ascending=False).head(10),'\n')

Top 5 самых низких потолков:
index Высота потолков, м
0 19834 1.2
1 19863 2.0
2 14432 2.1
3 3326 2.3
4 7687 2.3

Top 10 самых высоких потолков:
index Высота потолков, м
0 7969 320.0
1 18076 320.0
2 9 318.0
3 15285 285.0
4 8624 265.0
5 3159 264.0
6 11323 264.0
7 2427 260.0
8 8176 260.0
9 18525 260.0
```

Мы нашли масштабы выбросов. Так что в дальнейшем нам нужно будет удалить из значений **1 самый низкий** потолок и **8 самых высоких**

Смотрим нечисловые данные

```
In [10]: formatted_r1st.describe(include=object)

Out[10]:
```

	Количество комнат	Тип	Метро	Адрес	Площадь, м2	Дом	Парковка	Цена	Телеф
count	19202	19737	19391	19737	19737	19737	8563	19737	-
unique	24	1	4871	10452	9585	2415	5	2219	-
top	1	Квартира	м. Водный станция пешком	Москва, Челявский бульвар, 3	40.0/20.0/10.0	3/5, Кирпичный	наземная	4000.0 руб./3а месца, Залог - 40000 руб., Ко...	-
freq	6646	19737	40	63	214	265	5283	1112	-

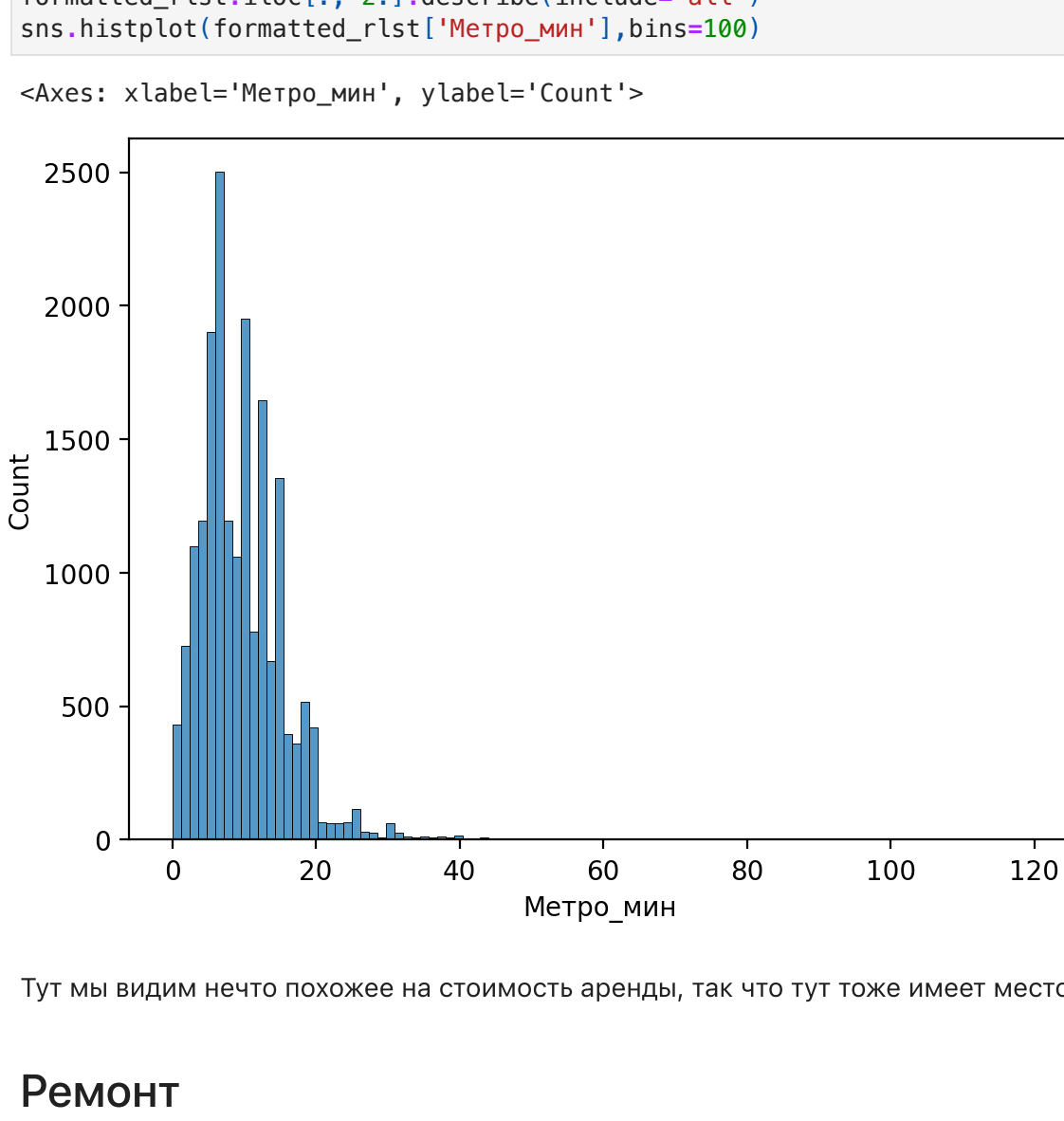
Что здесь?

- Видно, что у признака **Тип** всего одно уникальное значение, так что этим признаком мы смело пренебрегаем
- у **парковки, Ремонта, Окна, Можно с детьми/животными** достаточно небольшое число уникальных значений, так что эти данные можно назвать категориальными, и при дальнейшей обработке привести их к **типу категория или к числу**
- У **пола Мусоропровода** всего 2 уникальных значения, и одно из них "Да". Логично предположить, что другое значение в таком случае — это "Нет". В таком случае этот столбец можно привести к логическому типу данных
- Хотя **Ссылка на объявление** уникальна и так как данные не числовые, ничего полезного мы извлечь отсюда не сможем

Далее рассмотрим отдельные параметры

Парковка

```
In [11]: parking = pd.DataFrame(formatted_r1st['Парковка'].value_counts())
fig=plt.subplots(1,1,figsize=(10,5))
ax=sns.barplot(data = parking, x = parking.index, y = 'count')
ax.set_title('Распределение парковок')
ax.bar_label(ax.containers[0])
pass
```

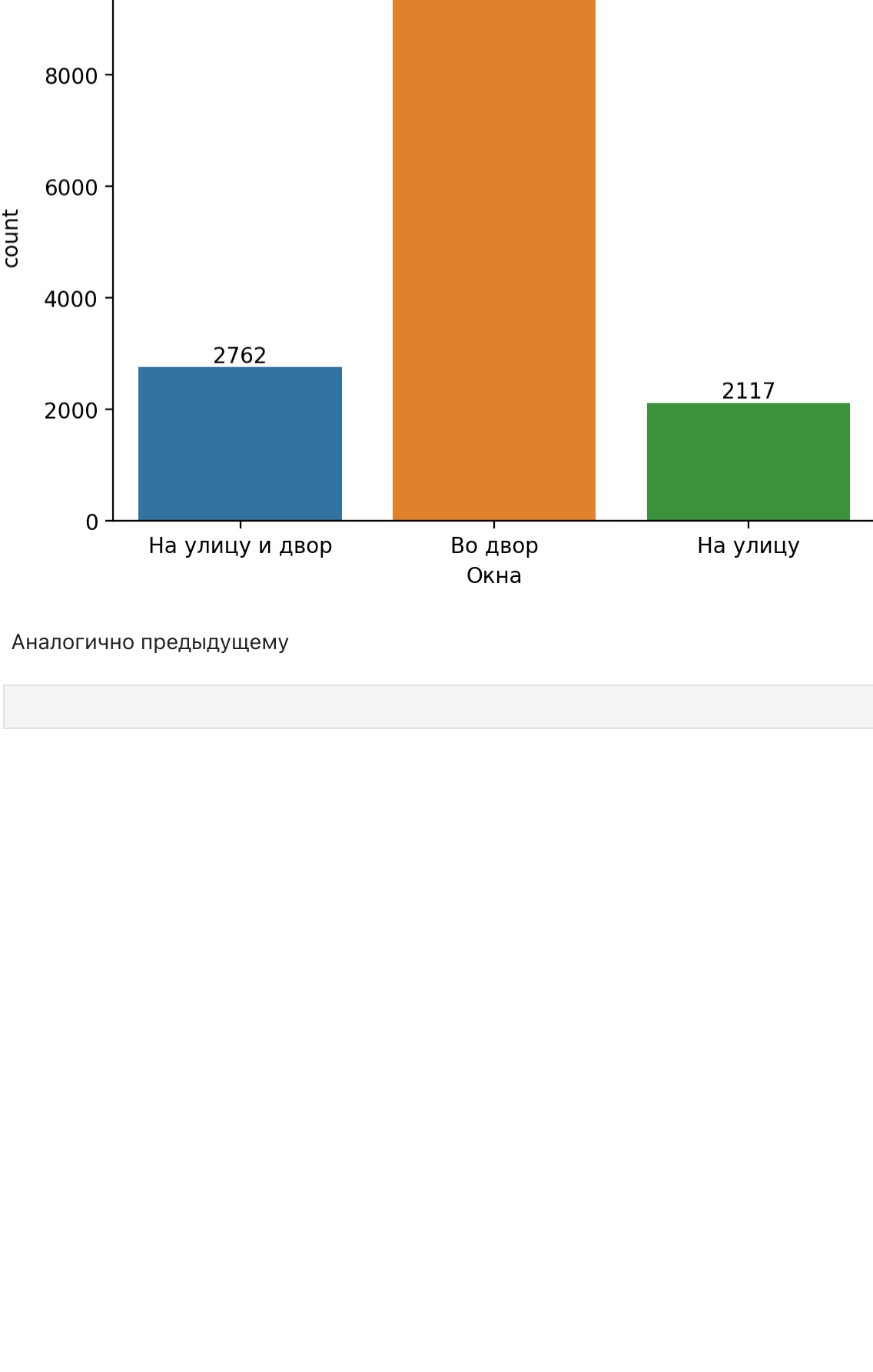


Здесь мы видим, что парковка на крыше всего одна, так что никаких зависимостей мы тут проследить не сможем, поэтому можем смело удалить это значение

Цена

```
In [12]: formatted_r1st['Цена аренды'] = formatted_r1st['Цена'].apply(lambda x: float(x.split('.')[0]))
plt.figure(figsize=(10,5))
sns.histplot(x='Цена аренды', data=formatted_r1st, bins=100, color='green')
# new_ticks = [40,000, 60,000, 80,000, 100,000, 120,000, 140,000, 160,000, 180,000, 200,000, 240,000, 280,000, 320,000, 360,000, 400,000, 440,000, 480,000, 520,000, 560,000, 600,000, 640,000, 680,000, 720,000, 760,000, 800,000, 840,000, 880,000, 920,000, 960,000, 1,000,000, 1,040,000, 1,080,000, 1,120,000, 1,160,000, 1,200,000, 1,240,000, 1,280,000, 1,320,000, 1,360,000, 1,400,000, 1,440,000, 1,480,000, 1,520,000, 1,560,000, 1,600,000, 1,640,000, 1,680,000, 1,720,000, 1,760,000, 1,800,000, 1,840,000, 1,880,000, 1,920,000, 1,960,000, 2,000,000, 2,040,000, 2,080,000, 2,120,000, 2,160,000, 2,200,000, 2,240,000, 2,280,000, 2,320,000, 2,360,000, 2,400,000, 2,440,000, 2,480,000, 2,520,000, 2,560,000, 2,600,000, 2,640,000, 2,680,000, 2,720,000, 2,760,000, 2,800,000, 2,840,000, 2,880,000, 2,920,000, 2,960,000, 3,000,000, 3,040,000, 3,080,000, 3,120,000, 3,160,000, 3,200,000, 3,240,000, 3,280,000, 3,320,000, 3,360,000, 3,400,000, 3,440,000, 3,480,000, 3,520,000, 3,560,000, 3,600,000, 3,640,000, 3,680,000, 3,720,000, 3,760,000, 3,800,000, 3,840,000, 3,880,000, 3,920,000, 3,960,000, 4,000,000, 4,040,000, 4,080,000, 4,120,000, 4,160,000, 4,200,000, 4,240,000, 4,280,000, 4,320,000, 4,360,000, 4,400,000, 4,440,000, 4,480,000, 4,520,000, 4,560,000, 4,600,000, 4,640,000, 4,680,000, 4,720,000, 4,760,000, 4,800,000, 4,840,000, 4,880,000, 4,920,000, 4,960,000, 5,000,000, 5,040,000, 5,080,000, 5,120,000, 5,160,000, 5,200,000, 5,240,000, 5,280,000, 5,320,000, 5,360,000, 5,400,000, 5,440,000, 5,480,000, 5,520,000, 5,560,000, 5,600,000, 5,640,000, 5,680,000, 5,720,000, 5,760,000, 5,800,000, 5,840,000, 5,880,000, 5,920,000, 5,960,000, 6,000,000, 6,040,000, 6,080,000, 6,120,000, 6,160,000, 6,200,000, 6,240,000, 6,280,000, 6,320,000, 6,360,000, 6,400,000, 6,440,000, 6,480,000, 6,520,000, 6,560,000, 6,600,000, 6,640,000, 6,680,000, 6,720,000, 6,760,000, 6,800,000, 6,840,000, 6,880,000, 6,920,000, 6,960,000, 7,000,000, 7,040,000, 7,080,000, 7,120,000, 7,160,000, 7,200,000, 7,240,000, 7,280,000, 7,320,000, 7,360,000, 7,400,000, 7,440,000, 7,480,000, 7,520,000, 7,560,000, 7,600,000, 7,640,000, 7,680,000, 7,720,000, 7,760,000, 7,800,000, 7,840,000, 7,880,000, 7,920,000, 7,960,000, 8,000,000, 8,040,000, 8,080,000, 8,120,000, 8,160,000, 8,200,000, 8,240,000, 8,280,000, 8,320,000, 8,360,000, 8,400,000, 8,440,000, 8,480,000, 8,520,000, 8,560,000, 8,600,000, 8,640,000, 8,680,000, 8,720,000, 8,760,000, 8,800,000, 8,840,000, 8,880,000, 8,920,000, 8,960,000, 9,000,000, 9,040,000, 9,080,000, 9,120,000, 9,160,000, 9,200,000, 9,240,000, 9,280,000, 9,320,000, 9,360,000, 9,400,000, 9,440,000, 9,480,000, 9,520,000, 9,560,000, 9,600,000, 9,640,000, 9,680,000, 9,720,000, 9,760,000, 9,800,000, 9,840,000, 9,880,000, 9,920,000, 9,960,000, 10,000,000, 10,040,000, 10,080,000, 10,120,000, 10,160,000, 10,200,000, 10,240,000, 10,280,000, 10,320,000, 10,360,000, 10,400,000, 10,440,000, 10,480,000, 10,520,000, 10,560,000, 10,600,000, 10,640,000, 10,680,000, 10,720,000, 10,760,000, 10,800,000, 10,840,000, 10,880,000, 10,920,000, 10,960,000, 11,000,000, 11,040,000, 11,080,000, 11,120,000, 11,160,000, 11,200,000, 11,240,000, 11,280,000, 11,320,000, 11,360,000, 11,400,000, 11,440,000, 11,480,000, 11,520,000, 11,560,000, 11,600,000, 11,640,000, 11,680,000, 11,720,000, 11,760,000, 11,800,000, 11,840,000, 11,880,000, 11,920,000, 11,960,000, 12,000,000, 12,040,000, 12,080,000, 12,120,000, 12,160,000, 12,200,000, 12,240,000, 12,280,000, 12,320,000, 12,360,000, 12,400,000, 12,440,000, 12,480,000, 12,520,000, 12,560,000, 12,600,000, 12,640,000, 12,680,000, 12,720,000, 12,760,000, 12,800,000, 12,840,000, 12,880,000, 12,920,000, 12,960,000, 13,000,000, 13,040,000, 13,080,000, 13,120,000, 13,160,000, 13,200,000, 13,240,000, 13,280,000, 13,320,000, 13,360,000, 13,400,000, 13,440,000, 13,480,000, 13,520,000, 13,560,000, 13,600,000, 13,640,000, 13,680,000, 13,720,000, 13,760,000, 13,800,000, 13,840,000, 13,880,000, 13,920,000, 13,960,000, 14,000,000, 14,040,000, 14,080,000, 14,120,000, 14,160,000, 14,200,000, 14,240,000, 14,280,000, 14,320,000, 14,360,000, 14,400,000, 14,440,000, 14,480,000, 14,520,000, 14,560,000, 14,600,000, 14,640,000, 14,680,000, 14,720,000, 14,760,000, 14,800,000, 14,840,000, 14,880,000, 14,920,000, 14,960,000, 15,000,000, 15,040,000, 15,080,000, 15,120,000, 15,160,000, 15,200,000, 15,240,000, 15,280,000, 15,320,000, 15,360,000, 15,400,000, 15,440,000, 15,480,000, 15,520,000, 15,560,000, 15,600,000, 15,640,000, 15,680,000, 15,720,000, 15,760,000, 15,800,000, 15,840,000, 15,880,000, 15,920,000, 15,960,000, 16,000,000, 16,040,000, 16,080,000, 16,120,000, 16,160,000, 16,200,000, 16,240,000, 16,280,000, 16,320,000, 16,360,000, 16,400,000, 16,440,000, 16,480,000, 16,520,000, 16,560,000, 16,600,000, 16,640,000, 16,680,000, 16,720,000, 16,760,000, 16,800,000, 16,840,000, 16,880,000, 16,920,000, 16,960,000, 17,000,000, 17,040,000, 17,080,000, 17,120,000, 17,160,000, 17,200,000, 17,240,000, 17,280,000, 17,320,000, 17,360,000, 17,400,000, 17,440,000, 17,480,000, 17,520,000, 17,560,000, 17,600,000, 17,640,000, 17,680,000, 17,720,000, 17,760,000, 17,800,000, 17,840,000, 17,880,000, 17,920,000, 17,960,000, 18,000,000, 18,040,000, 18,080,000, 18,120,000, 18,160,000, 18,200,000, 18,240,000, 18,280,000, 18,320,000, 18,360,000, 18,400,000, 18,440,000, 18,480,000, 18,520,000, 18,560,000, 18,600,000,
```


Out[19]: count 14587
unique 3
top 9708
freq 9708
Name: Окна, dtype: object



Аналогично предыдущему

In [1]: