

Подключаем необходимые для обработки данных библиотеки

```
In [619]: 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 [620]: path = 'data.csv'
pd.set_option('display.max_columns', 30)
#pd.set_option('max_info_columns', 30)
#pd.set_option('max_colwidth', 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. Смоленская (9 мин ...	Москва, улица Новый А...	200.0/20.0	5/16, Монолитный	подзем	
1	271634126	4	Квартира	М. Смоленская (8 мин ...	Москва, улица Новый А...	198.0/95.0/18.0	5/16, Монолитно-кирпичный	подзем	
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, перулеул Плот...	170.0/95.0/17.0	5/6	подзем	
4	273614615	2	Квартира	М. Арбатская (7 мин п...	Москва, улица Богданов...	58.0/38.0/5.0	12/26, Панельный		N
...	
23363	215565511	NaN	Квартира	М. Говерово (8 мин п...	Москва, Боровское шос...	35.0/16.4/8.0	10/14		N
23364	274654844	1	Квартира	М. Солнцево (7 мин п...	Москва, Производствен...	38.7/16.5/11.0	5/18, Монолитный		N
23365	268679909	2, 0ба варианта	Квартира	М. Солнцево (6 мин п...	Москва, Боровский про...	43.1	5/5, Кирпичный		N
23366	274807525	2	Квартира	М. Солнцево (11 мин п...	Москва, улица Богданов...	52.5/10.0	8/23, назем		
23367	274672243	2, 0ба варианта	Квартира	М. Говерово (10 мин п...	Москва, улица Богданов...	90.0/48.2/15.5	6/23, Панельный		N

23368 rows x 24 columns

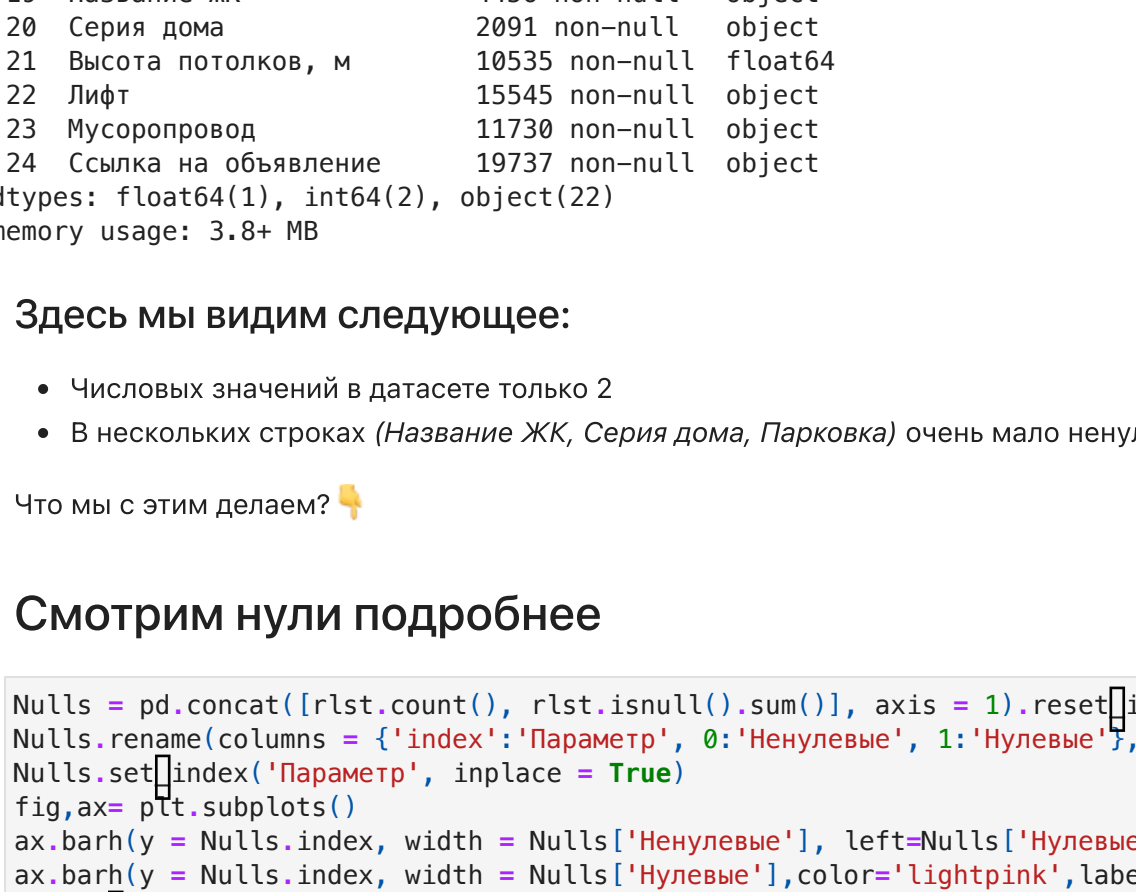
В датареиме мы видим **24 столбца** с данными. Можно сразу обратить внимание на следующее:

- 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 пилотного региона выбрана **Москва**, поэтому посмотрим, есть ли у нас квартиры **не в Москве**

```
In [621]: config InlineBackend.figure_format='retina'
bars=sns.countplot(x=r1st['Адрес'], format='retina')
bars.set(title='Города')
for i in bars.containers:
    bars.bar_label(i)
```



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

Удалим всё кроме Москвы

```
In [622]: r1st.drop(r1st.loc(r1st['Адрес']!=r1st['Адрес'].apply(lambda x: x.split(sep=',')[0])!='Москва').index, inplace=True)
r1st.reset_index(inplace=True)
print(f'Размер датареимы только с Москвой стан {r1st.shape}')
Размер датареимы 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 [623]: 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   ID                   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              13107 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  18535 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: 3.8+ MB
```

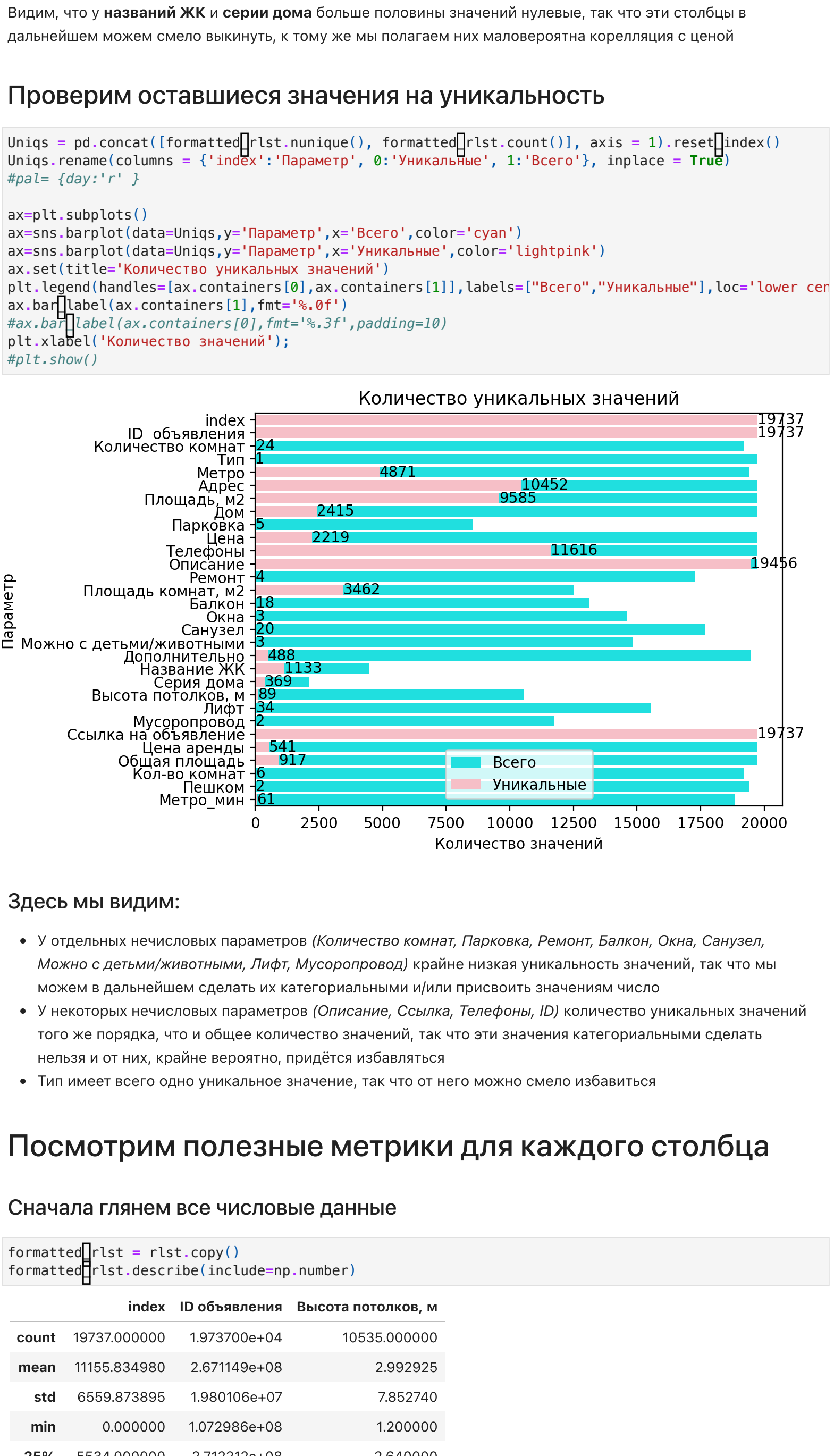
Здесь мы видим следующее:

- Числовые значения 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 [624]: Nulls = pd.concat([r1st.count(), r1st.isnull().sum()], axis=1).reset_index()
Nulls.rename(columns = {'index': 'Параметр', 0: 'Ненулевые', 1: 'Нулевые', 1: 'Нулевые', 1: 'Нулевые'})
Nulls.set_index('Параметр', inplace = True)
fig,ax=plt.subplots()
ax.bar(y = Nulls.index, width = Nulls['Ненулевые'], left=Nulls['Ненулевые'], color='cyan', label='Ненулевые')
ax.bar(y = Nulls.index, width = Nulls['Нулевые'], left=Nulls['Ненулевые'], color='lightpink', label='Нулевые')
ax.bar_label(ax.containers[1],fmt='%0.1f')
ax.set_title('Соотношение нулевых и ненулевых значений',xlabel='Кол-во значений')
plt.legend()
```

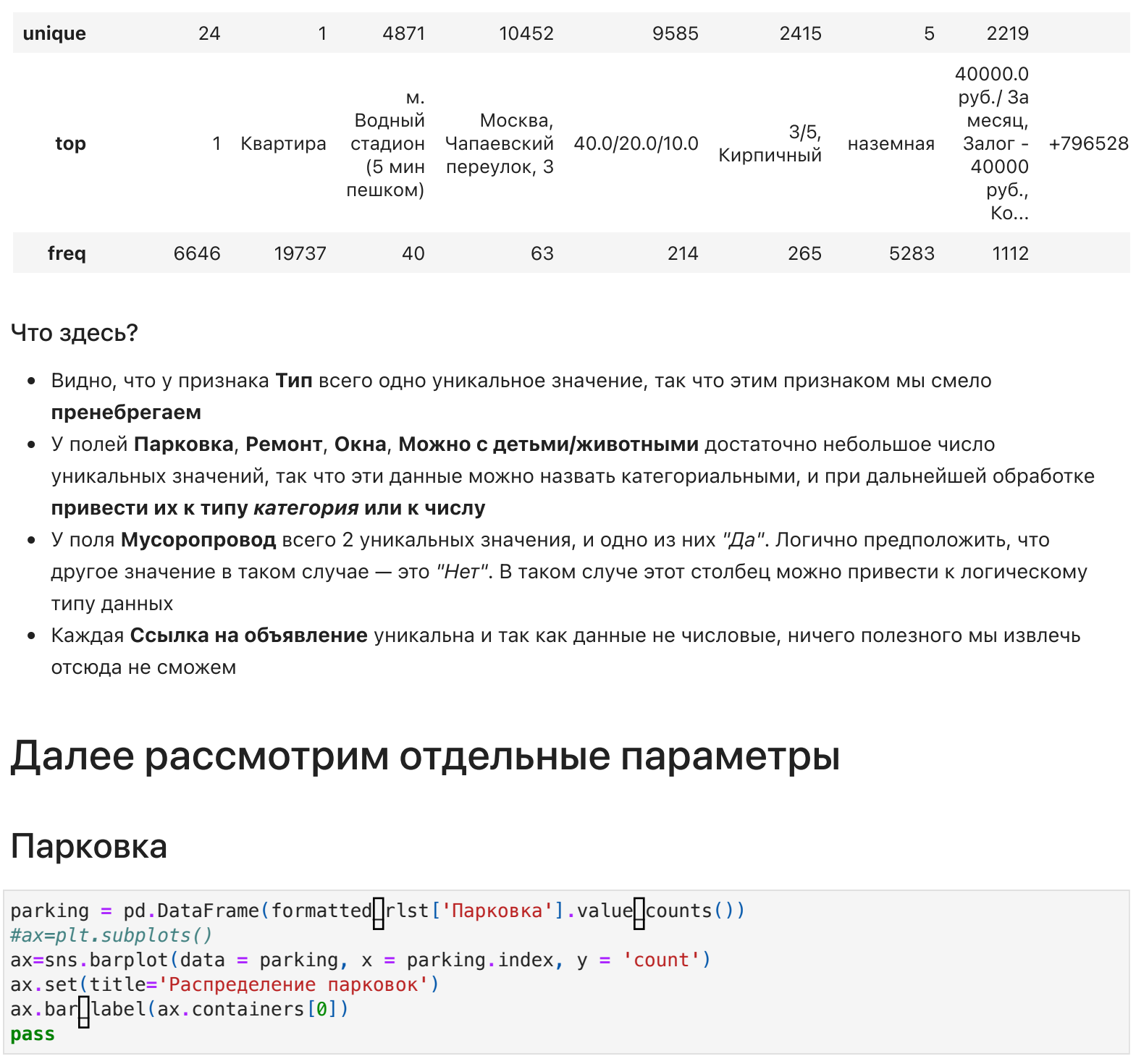


Что мы видим?

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

Проверим оставшиеся значения на уникальность

```
In [625]: Uniqs = pd.concat([formatted(r1st.nunique(), formatted(r1st.count(), axis=1).reset_index()
Uniqs.rename(columns = {'index': 'Параметр', 0: 'Уникальные', 1: 'Всего', 1: 'Всего', 1: 'Всего'})
#pal = day:'')
ax=plt.subplots()
ax=sns.barplot(data=Uniqs,y='Параметр',x='Всего',color='cyan')
ax=sns.barplot(data=Uniqs,y='Параметр',x='Уникальные',color='lightpink')
ax.set_title('Количество уникальных значений')
plt.legend(handles=[ax.containers[0],ax.containers[1]],labels=['Всего','Уникальные'],loc='lower center')
ax.bar_label(ax.containers[1],fmt='%0.1f')
ax.bar_label(ax.containers[0],fmt='%0.1f',padding=10)
ax.set_title('Количество значений')
plt.show()
```



Здесь мы видим:

- У отдельных нечисловых параметров (**Количество комнат, Парковка, Ремонт, Балкон, Окна, Санузел**) много значений, что объясняется тем, что эти значения категориальными, и при дальнейшей обработке можно с легкостью сделать их категориальными или присвоить значениям число
- У некоторых нечисловых параметров (**Описание, Ссылка, Телефоны, ID**) количество уникальных значений такое же, как и количество значений, что объясняется тем, что эти значения категориальными, и при дальнейшей обработке можно с легкостью сделать их категориальными или присвоить значениям число
- Тип имеет всего одно уникальное значение, так что от него можно смело избавиться

Посмотрим полезные метрики для каждого столбца

Сначала глянем все числовые данные

```
In [626]: formatted(r1st = r1st.copy()
formatted(r1st.describe(include=np.number))
Out[626]:
```

	index	ID объявления	Высота потолков, м
count	19737.000000	1973700e+04	10535.000000
mean	11155.834890	2.671149e+08	2.992925
std	6559.873895	1.980106e+07	7.852740
min	0.000000	1.02986e+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 не отличаются, поэтому глянем, насколько много таких выбросов

Потолок

```
In [627]: print(f'Top 5 самых низких потолков:\n{formatted(r1st[\"Высота потолков, м\"].sort(values().head(5)).reset_index()
print(f'Top 10 самых высоких потолков:\n{formatted(r1st[\"Высота потолков, м\"].sort(values().tail(10)).reset_index()
Top 5 самых низких потолков:
index Высота потолков, м
0 19034 1.2
1 9963 2.0
2 14432 2.1
3 3526 2.3
4 7687 2.3
Top 10 самых высоких потолков:
index Высота потолков, м
0 7969 320.0
1 10876 320.0
2 9 318.0
3 15285 280.0
4 8624 265.0
5 3159 264.0
6 11323 264.0
7 247 260.0
8 8176 26.0
9 18525 28.0
```

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

Смотрим нечисловые данные

```
In [628]: formatted(r1st.describe(include=object))
Out[628]:
```

	Количество комнат	Тип	Метро	Адрес	Площадь, м2	Дом	Парковка	Цена	Теле
count	19202	19737	19391	19737	19737	19737	8563	19737	
unique	24	1	40	10452	9585	2415	5	2219	
top	1	Квартира	М. Водный стан перулеул 3	Москва, Чкаловский перулеул 3	40.0/20.0/10.0	3/5, Кирпичный	наземная	40000.0 руб./3а месца, Залог - 40000 руб., Ко...	+795528
freq	66646	19737	40	63	214	265	5283	1112	

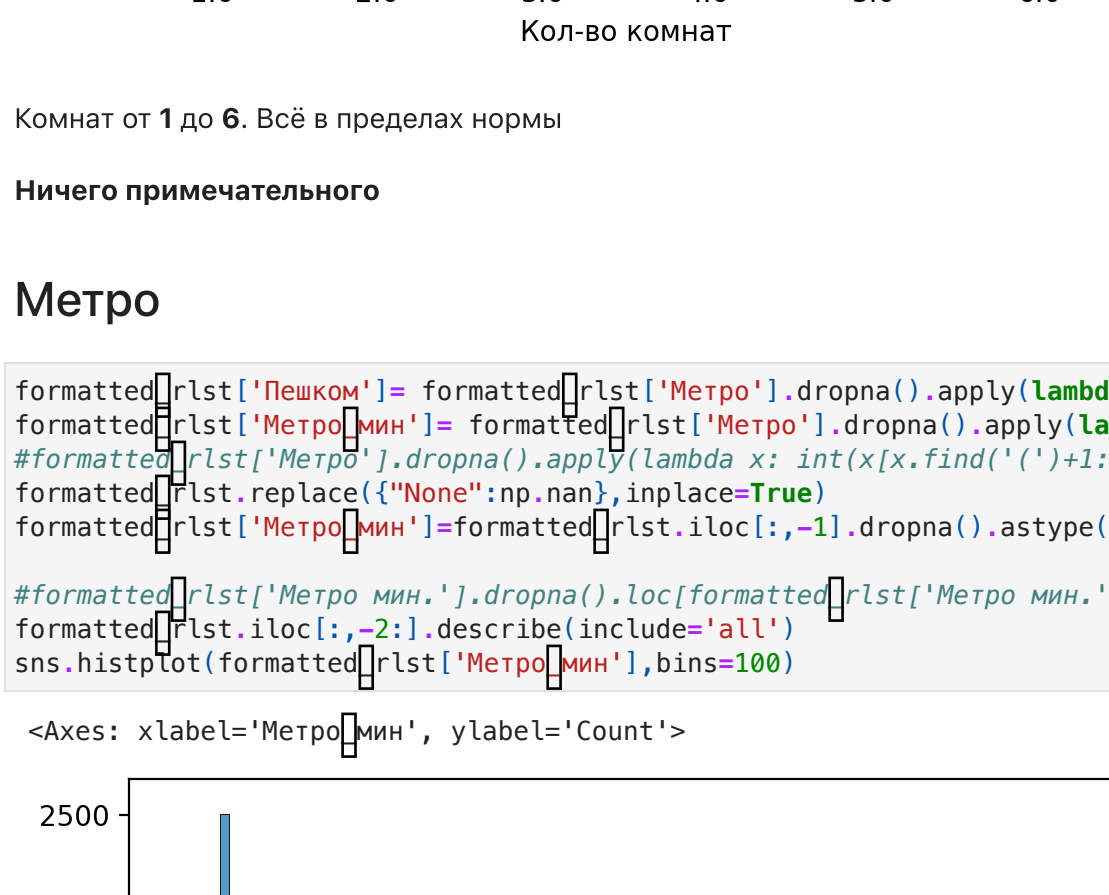
Что здесь?

- 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 **поля Мусоропровод** всего 2 уникальных значения, и одно из них "Да". Можно попробовать, что другое значение в таком случае — это "Нет". В таком случае этот столбец можно привести к логическому типу данных
- Каждая **Ссылка на объявление** уникальна и так как данные не числовые, ничего полезного мы извлечь отсюда не сможем

Далее рассмотрим отдельные параметры

Парковка

```
In [629]: parking = pd.DataFrame(formatted(r1st['Парковка'].value_counts()))
fig,ax=plt.subplots()
ax=sns.barplot(data=parking, x = parking.index, y = 'count')
ax.set_title('Распределение парковок')
ax.bar_label(ax.containers[0])
pass
```

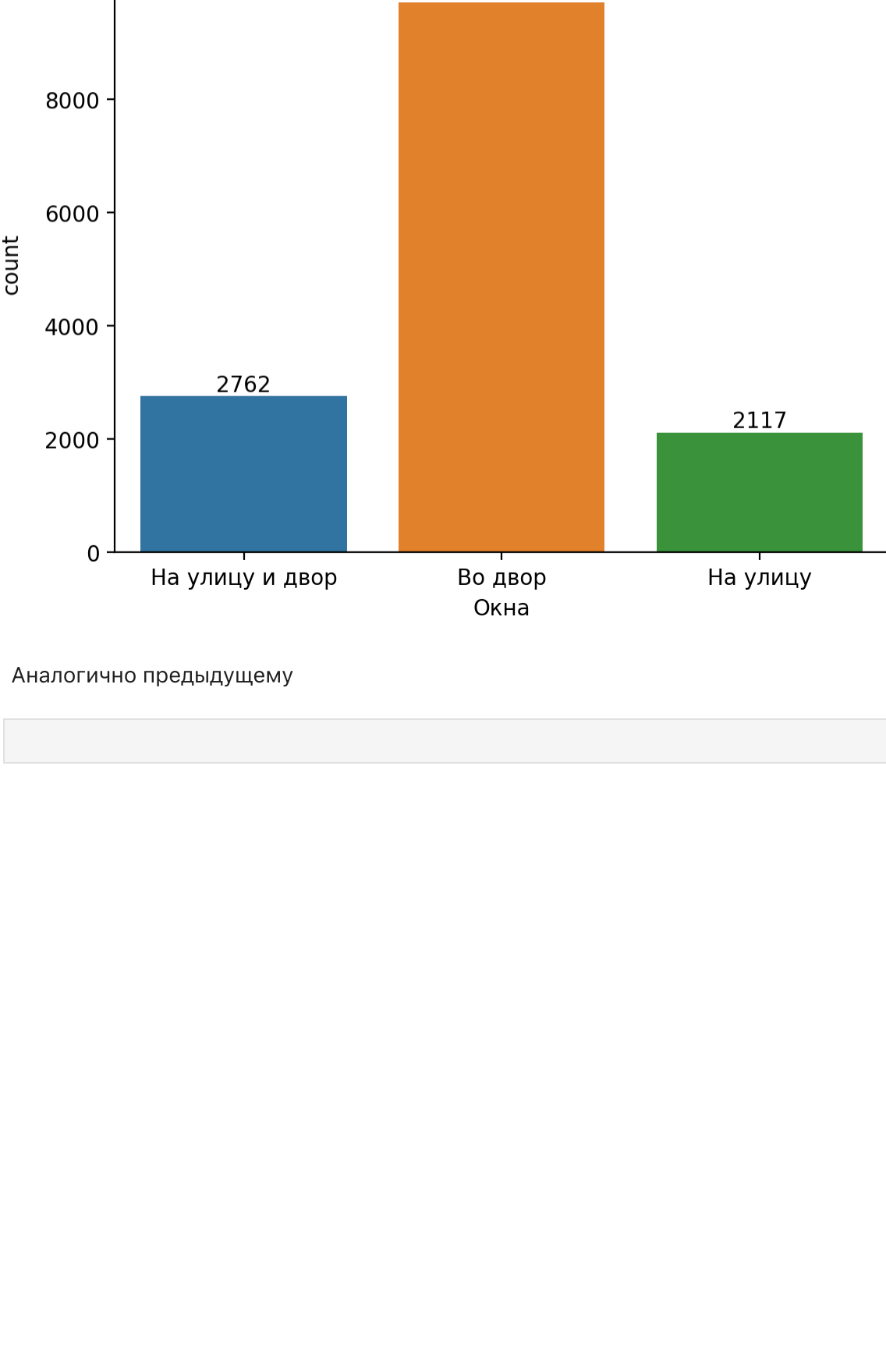


Здесь мы видим, что парковка на крыше всего одна, так что никаких зависимостей мы тут проследить не сможем, поэтому можем смело удалить это значение

Цена

```
In [630]: r1st[r1st['Цена аренды']!=formatted(r1st['Цена'].apply(lambda x: float(x.split(',')[0])))]
fig,figsize=(10,5)
sns.histplot(x='Цена аренды', data=formatted(r1st, bins=100, color='green'))
# new_ticks = [1000, 6000, 8000, 12000, 14000, 16000, 18000, 20000, 22000, 24000, 26000, 28000, 30000, 32000, 34000, 36000, 38000, 40000, 42000, 44000, 46000, 48000, 50000, 52000, 54000, 56000, 58000, 60000, 62000, 64000, 66000, 68000, 70000, 72000, 74000, 76000, 78000, 80000, 82000, 84000, 86000, 88000, 90000, 92000, 94000, 96000, 98000, 100000, 102000, 104000, 106000, 108000, 110000, 112000, 114000, 116000, 118000, 120000, 122000, 124000, 126000, 128000, 130000, 132000, 134000, 136000, 138000, 140000, 142000, 144000, 146000, 148000, 150000, 152000, 154000, 156000, 158000, 160000, 162000, 164000, 166000, 168000, 170000, 172000, 174000, 176000, 178000, 180000, 182000, 184000, 186000, 188000, 190000, 192000, 194000, 196000, 198000, 200000, 202000, 204000, 206000, 208000, 210000, 212000, 214000, 216000, 218000, 220000, 222000, 224000, 226000, 228000, 230000, 232000, 234000, 236000, 238000, 240000, 242000, 244000, 246000, 248000, 250000, 252000, 254000, 256000, 258000, 260000, 262000, 264000, 266000, 268000, 270000, 272000, 274000, 276000, 278000, 280000, 282000, 284000, 286000, 288000, 290000, 292000, 294000, 296000, 298000, 300000, 302000, 304000, 306000, 308000, 310000, 312000, 314000, 316000, 318000, 320000, 322000, 324000, 326000, 328000, 330000, 332000, 334000, 336000, 338000, 340000, 342000, 344000, 346000, 348000, 350000, 352000, 354000, 356000, 358000, 360000, 362000, 364000, 366000, 368000, 370000, 372000, 374000, 376000, 378000, 380000, 382000, 384000, 386000, 388000, 390000, 392000, 394000, 396000, 398000, 400000, 402000, 404000, 406000, 408000, 410000, 412000, 414000, 416000, 418000, 420000, 422000, 424000, 426000, 428000, 430000, 432000, 434000, 436000, 438000, 440000, 442000, 444000, 446000, 448000, 450000, 452000, 454000, 456000, 458000, 460000, 462000, 464000, 466000, 468000, 470000, 472000, 474000, 476000, 478000, 480000, 482000, 484000, 486000, 488000, 490000, 492000, 494000, 496000, 498000, 500000, 502000, 504000, 506000, 508000, 510000, 512000, 514000, 516000, 518000, 520000, 522000, 524000, 526000, 528000, 530000, 532000, 534000, 536000, 538000, 540000, 542000, 544000, 546000, 548000, 550000, 552000, 554000, 556000, 558000, 560000, 562000, 564000, 566000, 568000, 570000, 572000, 574000, 576000, 578000, 580000, 582000, 584000, 586000, 588000, 590000, 592000, 594000, 596000, 598000, 600000, 602000, 604000, 606000, 608000, 610000, 612000, 614000, 616000, 618000, 620000, 622000, 624000, 626000, 628000, 630000, 632000, 634000, 636000, 638000, 640000, 642000, 644000, 646000, 648000, 650000, 652000, 654000, 656000, 658000, 660000, 662000, 664000, 666000, 668000, 670000, 672000, 674000, 676000, 678000, 680000, 682000, 684000, 686000, 688000, 690000, 692000, 694000, 696000, 698000, 700000, 702000, 704000, 706000, 708000, 710000, 712000, 714000, 716000, 718000, 720000, 722000, 724000, 726000, 728000, 730000, 732000, 734000, 736000, 738000, 740000, 742000, 744000, 746000, 748000, 750000, 752000, 754000, 756000, 758000, 760000, 762000, 764000, 766000, 768000, 770000, 772000, 774000, 776000, 778000, 780000, 782000, 784000, 786000, 788000, 790000, 792000, 794000, 796000, 798000, 800000, 802000, 804000, 806000, 808000, 810000, 812000, 814000, 816000, 818000, 820000, 822000, 824000, 826000, 828000, 830000, 832000, 834000, 836000, 838000, 840000, 842000, 844000, 846000, 848000, 850000, 852000, 854000, 856000, 858000, 860000, 862000, 864000, 866000, 868000, 870000, 872000, 874000, 876000, 878000, 880000, 882000, 884000, 886000, 888000, 890000, 892000, 894000, 896000, 898000, 900000, 902000, 904000, 906000, 908000, 910000, 912000, 914000, 916000, 918000, 920000, 922000, 924000, 926000, 928000, 930000, 932000, 934000, 936000, 938000, 940000, 942000, 944000, 946000, 948000, 950000, 952000, 954000, 956000, 958000, 960000, 962000, 964000, 966000, 968000, 970000, 972000, 974000, 976000, 978000, 980000, 982000, 984000, 986000, 988000, 990000, 992000, 994000, 996000, 998000, 1000000, 1002000, 1004000, 1006000, 1008000, 1010000, 1012000, 1014000, 1016000, 1018000, 1020000, 1022000, 1024000, 1026000, 1028000, 1030000, 1032000, 1034000, 1036000, 1038000, 1040000, 1042000, 1044000, 1046000, 1048000, 1050000, 1052000, 1054000, 1056000, 1058000, 1060000, 1062000, 1064000, 1066000, 1068000, 1070000, 1072000, 1074000, 1076000, 1078000, 1080000, 1082000, 1084000, 1086000, 1088000, 1090000, 1092000, 1094000, 1096000, 1098000, 1100000, 1102000, 1104000, 1106000, 1108000, 1110000, 1112000, 1114000, 1116000, 1118000, 1120000, 1122000, 1124000, 1126000, 1128000, 1130000, 1132000, 1134000, 1136000, 1138000, 1140000, 1142000, 1144000, 1146000, 1148000, 1150000, 1152000, 1154000, 1156000, 1158000, 1160000, 1162000, 1164000, 1166000, 1168000, 1170000, 1172000, 1174000, 1176000, 1178000, 1180000,
```


Out[637]: count 14587
unique 3
top Во двор 9788
freq
Name: Окна, dtype: object



Аналогично предыдущему

In [1]: