

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('display.max_info_cols', 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 п...	Москва, улица Новый А...	200.0/20.0	5/16, Монолитный	подзем	
1	271634126	4	Квартира	м. Смоленский (8 мин п...	Москва, улица Новый А...	198.0/95.0/18.0	5/16, Монолитно-кирпичн...	подзем	
2	271173086	4, 0ба варианта	Квартира	м. Смоленская (7 мин п...	Москва, улица Новый А...	200.0/116.0/4.0	5/16	подзем	
3	272197456	4, 0ба варианта	Квартира	м. Смоленский (3 мин п...	Москва, перулеул Плот...	170.0/95.0/17.0	5/6	подзем	
4	273614415	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. Говорово (8 мин п...	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. Говорово (10 мин п...	Москва, улица Богданов...	90.0/48.2/15.5	6/23, Панельный		N

23368 rows x 24 columns

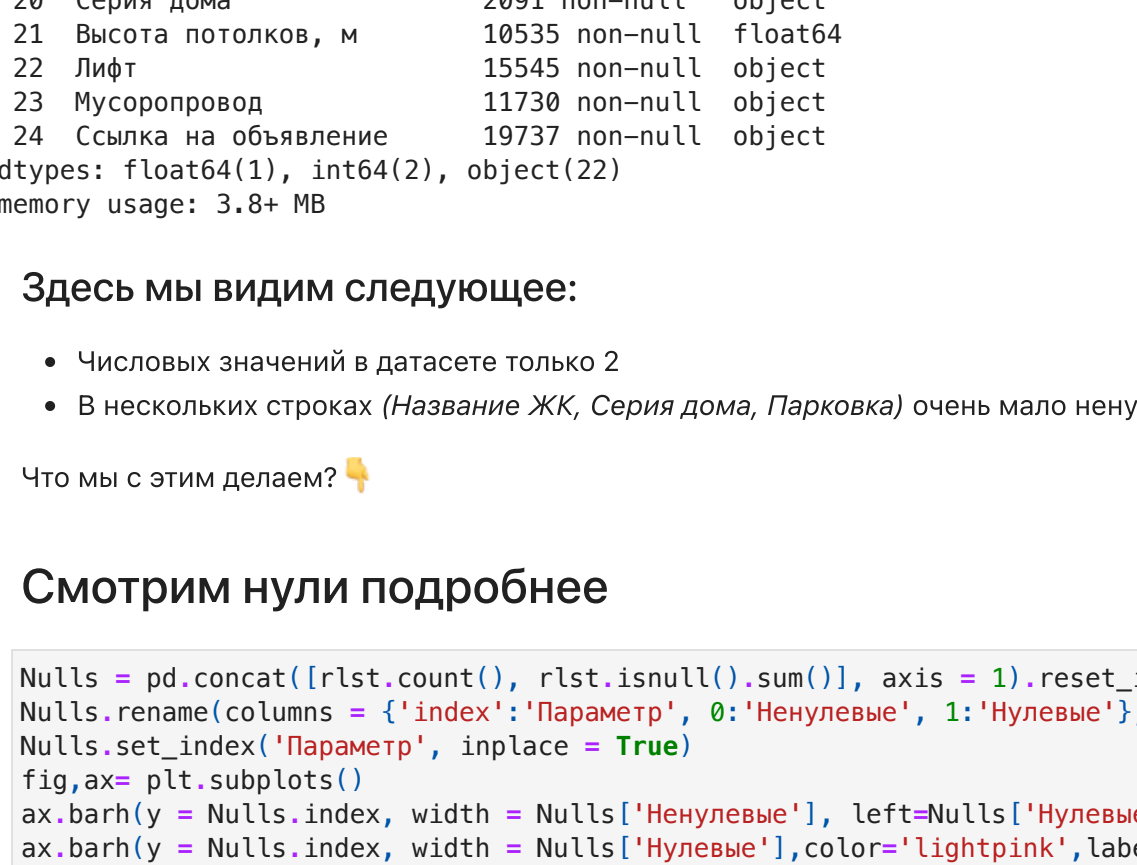
В датареиме мы видим **24 столбца** с данными. Можно сразу обратить внимание на следующее:

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

```
In [621]: conf = fig.BackendContext('figure', format='retina')
sns.set(rc={'figure.figsize': (10, 5)})
r1st.reset_index(inplace=True)
for i in r1st.columns:
    bars, bar_label = plt.subplots(1, 1, sharex=True)
    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'].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   Количество комнат    19737 non-null  int64
2   Тип                  19737 non-null  object
3   Метро                19391 non-null  object
4   Адрес                19737 non-null  object
5   Площадь, м2          19737 non-null  object
6   Дом                  19737 non-null  object
7   Парковка             8563 non-null   object
8   Цена                 19737 non-null  object
9   Телефоны             19737 non-null  object
10  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   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: 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 информацию подробнее

```
In [624]: Nulls = pd.concat([r1st.count(), r1st.isnull().sum()], axis=1).reset_index()
Nulls.rename(columns={'index': 'Параметр', 0: 'Ненулевые', 1: 'Нулевые', inplace=True)
Nulls.set_index('Параметр', inplace=True)

fig, ax = plt.subplots(1)
ax.barh(y = Nulls.index, width = Nulls['Ненулевые'], left=Nulls['Нулевые'], color='cyan', label='Ненулевые')
ax.barh(y = Nulls.index, width = Nulls['Нулевые'], color='lightpink', label='Нулевые')
ax.bar_label(ax.containers[1], fmt='%d')
ax.set_title('Соотношение нулевых и ненулевых значений', xlabel='Кол-во значений')
plt.legend()
```



Что мы видим?

Видно, что у **названия ЖК** и **серии дома** больше половины значений нулевые, так что эти столбцы в дальнейшем можем смело **закинуть**, 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', inplace=True)
#pal = [day: '']

ax=plt.subplots(1)
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='%d')
ax.set_title('Количество уникальных значений', xlabel='Кол-во значений')
plt.xlabel('Количество значений')
plt.show()
```



Здесь мы видим:

- У отдельных нечисловых параметров (*Количество комнат*, *Парковка*, *Ремонт*, *Балкон*, *Окна*, *Санузел*, *Можно с детьми/животными*, *Лифт*, *Мусоропровод*) крайне низкая уникальность значений, так что мы можем в дальнейшем сделать их категориальными и/или присвоить значениям число
- У некоторых нечисловых параметров (*Описание*, *Ссылка*, *Телефоны*, *ID*) количество уникальных значений почти равно количеству значений, так что эти значения категориальными сделать нельзя и от них крайне вероятно, придется избавиться
- 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)
```

	index	ID объявления	Количество комнат	Тип	Метро	Адрес	Площадь, м2	Дом	Парковка	Цена	Теле
count	19737.000000	19737.000e+04	10535.000000								
mean	11155.834980	2.671149e+08	2.9922925								
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(ascending=False).head(10).reset_index()})

Top 5 самых низких потолков:
index  Высота потолков, м
0      19834          2.0
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          328.0
1      10676          328.0
2      10676          328.0
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 выбросов. 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)
```

	Количество комнат	Тип	Метро	Адрес	Площадь, м2	Дом	Парковка	Цена	Теле
count	19202	19737	19391	19737	19737	19737	8563	19737	
unique	24	1	4871	10452	9585	2415	5	2219	
top	1	Квартира	м. Водный	Москва, станица Чкаловский	40.0/20.0/10.0	3/5, Кирпичный	наземная	40000.0 руб/3а месяц, Залог – 40000 руб., Ко...	+795528
freq	66646	19737	40	63	214	265	5283	1112	

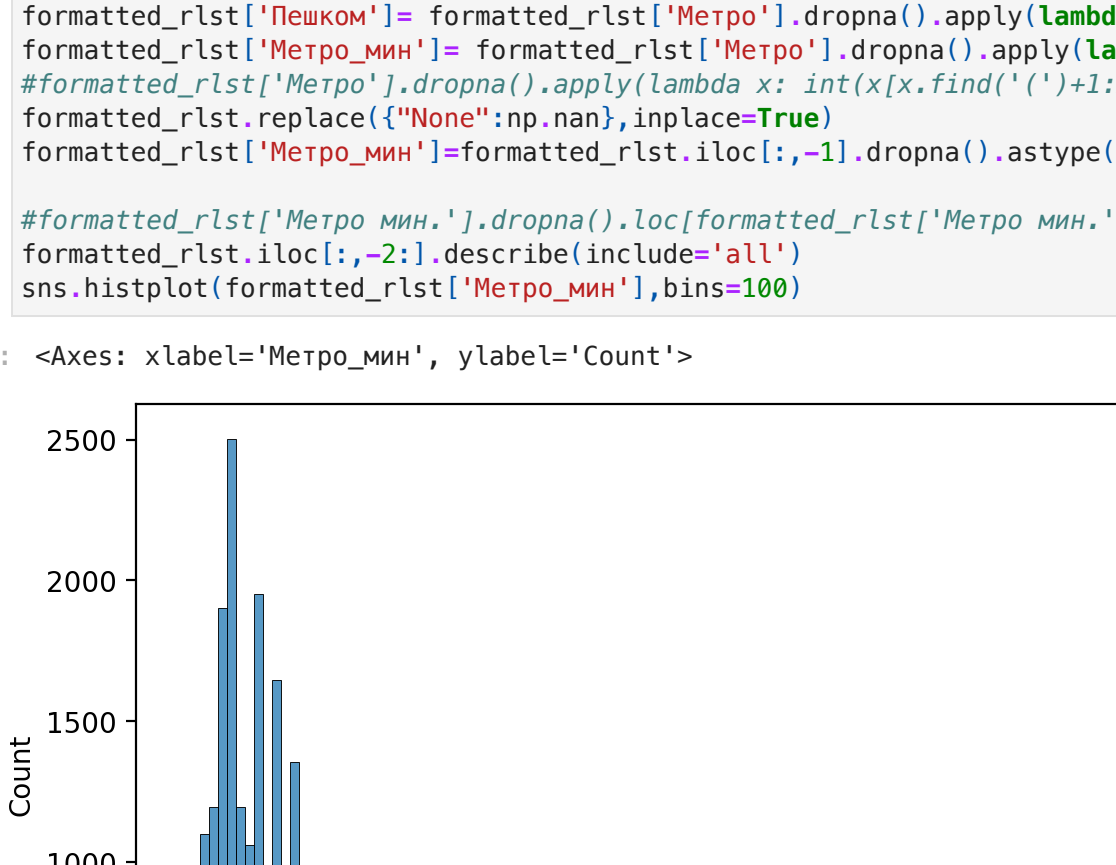
Что здесь?

- Видно, что у признака **Тип** всего одно уникальное значение, так что этим признаком мы смело **пренебрегаем**
- У **парковки**, **Ремонт**, **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 дополнительно сделать 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())
parking.reset_index(inplace=True)
ax=sns.barplot(x=parking.x, y=parking.y, color='cyan')
ax.set_title('Распределение парковок')
ax.bar_label(ax.containers[0])
pass
```



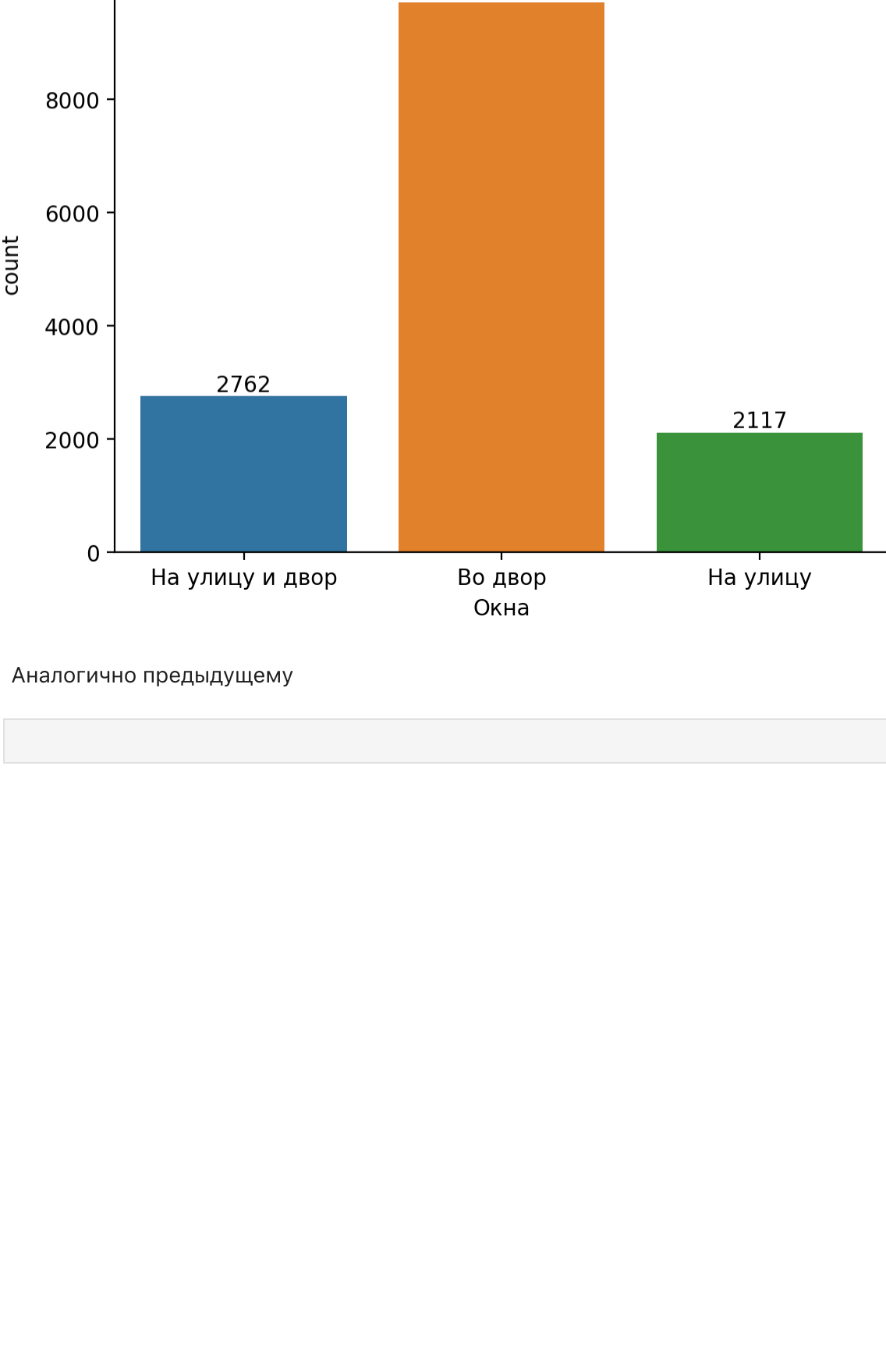
Здесь мы видим, что парковка на крыше всего одна, так что никаких зависимостей мы тут проследить не сможем, поэтому можем смело удалить это значение

Цена

```
In [630]: formatted_r1st[\"Цена аренды\"] = formatted_r1st[\"Цена\"].apply(lambda x: float(x.split('.')[0]))

plt.figure(figsize=(10,5))
sns.histplot(x=formatted_r1st[\"Цена аренды\"], data=formatted_r1st, bins=100, color='green')
# new_ticks = [100_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, 100_000, 104_000, 108_000, 112_000, 116_000, 120_000, 124_000, 128_000, 132_000, 136_000, 140_000, 144_000, 148_000, 152_000, 156_000, 160_000, 164_000, 168_000, 172_000, 176_000, 180_000, 184_000, 188_000, 192_000, 196_000, 200_000, 204_000, 208_000, 212_000, 216_000, 220_000, 224_000, 228_000, 232_000, 236_000, 240_000, 244_000, 248_000, 252_000, 256_000, 260_000, 264_000, 268_000, 272_000, 276_000, 280_000, 284_000, 288_000, 292_000, 296_000, 300_000, 304_000, 308_000, 312_000, 316_000, 320_000, 324_000, 328_000, 332_000, 336_000, 340_000, 344_000, 348_000, 352_000, 356_000, 360_000, 364_000, 368_000, 372_000, 376_000, 380_000, 384_000, 388_000, 392_000, 396_000, 400_000, 404_000, 408_000, 412_000, 416_000, 420_000, 424_000, 428_000, 432_000, 436_000, 440_000, 444_000, 448_000, 452_000, 456_000, 460_000, 464_000, 468_000, 472_000, 476_000, 480_000, 484_000, 488_000, 492_000, 496_000, 500_000, 504_000, 508_000, 512_000, 516_000, 520_000, 524_000, 528_000, 532_000, 536_000, 540_000, 544_000, 548_000, 552_000, 556_000, 560_000, 564_000, 568_000, 572_000, 576_000, 580_000, 584_000, 588_000, 592_000, 596_000, 600_000, 604_000, 608_000, 612_000, 616_000, 620_000, 624_000, 628_000, 632_000, 636_000, 640_000, 644_000, 648_000, 652_000, 656_000, 660_000, 664_000, 668_000, 672_000, 676_000, 680_000, 684_000, 688_000, 692_000, 696_000, 700_000, 704_000, 708_000, 712_000, 716_000, 720_000, 724_000, 728_000, 732_000, 736_000, 740_000, 744_000, 748_000, 752_000, 756_000, 760_000, 764_000, 768_000, 772_000, 776_000, 780_000, 784_000, 788_000, 792_000, 796_000, 800_000, 804_000, 808_000, 812_000, 816_000, 820_000, 824_000, 828_000, 832_000, 836_000, 840_000, 844_000, 848_000, 852_000, 856_000, 860_000, 864_000, 868_000, 872_000, 876_000, 880_000, 884_000, 888_000, 892_000, 896_000, 900_000, 904_000, 908_000, 912_000, 916_000, 920_000, 924_000, 928_000, 932_000, 936_000, 940_000, 944_000, 948_000, 952_000, 956_000, 960_000, 964_000, 968_000, 972_000, 976_000, 980_000, 984_000, 988_000, 992_000, 996_000, 1000_000, 1004_000, 1008_000, 1012_000, 1016_000, 1020_000, 1024_000, 1028_000, 1032_000, 1036_000, 1040_000, 1044_000, 1048_000, 1052_000, 1056_000, 1060_000, 1064_000, 1068_000, 1072_000, 1076_000, 1080_000, 1084_000, 1088_000, 1092_000, 1096_000, 1100_000, 1104_000, 1108_000, 1112_000, 1116_000, 1120_000, 1124_000, 1128_000, 1132_000, 1136_000, 1140_000, 1144_000, 1148_000, 1152_000, 1156_000, 1160_000, 1164_000, 1168_000, 1172_000, 1176_000, 1180_000, 1184_000, 1188_000, 1192_000, 1196_000, 1200_000, 1204_000, 1208_000, 1212_000, 1216_000, 1220_000, 1224_000, 1228_000, 1232_000, 1236_000, 1240_000, 1244_000, 1248_000, 1252_000, 1256_000, 1260_000, 1264_000, 1268_000, 1272_000, 1276_000, 1280_000, 1284_000, 1288_000, 1292_000, 1296_000, 1300_000, 1304_000, 1308_000, 1312_000, 1316_000, 1320_000, 1324_000, 1328_000, 1332_000, 1336_000, 1340_000, 1344_000, 1348_000, 1352_000, 1356_000, 1360_000, 1364_000, 1368_000, 1372_000, 1376_000, 1380_000, 1384_000, 1388_000, 1392_000, 1396_000, 1400_000, 1404_000, 1408_000, 1412_000, 1416_000, 1420_000, 1424_000, 1428_000, 1432_000, 1436_000, 1440_000, 1444_000, 1448_000, 1452_000, 1456_000, 1460_000, 1464_000, 1468_000, 1472_000, 1476_000, 1480_000, 1484_000, 1488_000, 1492_000, 1496_000, 1500_000, 1504_000, 1508_000, 1512_000, 1516_000, 1520_000, 1524_000, 1528_000, 1532_000, 1536_000, 1540_000, 1544_000, 1548_000, 1552_000, 1556_000, 1560_000, 1564_000, 1568_000, 1572_000, 1576_000, 1580_000, 1584_000, 1588_000, 1592_000, 1596_000, 1600_000, 1604_000, 1608_000, 1612_000, 1616_000, 1620_000, 1624_000, 1628_000, 1632_000, 1636_000, 1640_000, 1644_000, 1648_000, 1652_
```


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



Аналогично предыдущему

In [1]: