Государственное бюджетное профессиональное образовательное учреждение Московской области «Физико-технический колледж»

Отчёт по кейсу «Самолёт»:

Работу выполнил: Студент группы № ИСП-21 Бухаров Егор

Долгопрудный, 2024

Введение

В данном отчёте рассматриваются выводы, полученные с аналитической работы над данными в области «Квартиры в Московской области, Новой Москве и Москве».

Цель

Собрать данные и произвести аналитическую работу над ними для будущих работ, например, создание модели на основе выводов.

Задачи

- Используя открытые источники собрать список данных.
- На основе полученной информации произвести удаление ненужных данных, дополнение необходимых, выявление аномалий и их блокировка.
- Визуализация данных при помощи, как минимум, двух инструментов для подобных задач. Нахождение взаимосвязей между данными или их полное отсутствие, усреднённых показателей для уверенного отчёта.

Основная часть

Для выполнения основной задачи, существует небольшой выбор источников, откуда собирать данные, мною был выбрать интернет-ресурс «Циан». При помощи скриптов, написанных на языке Python и библиотеке CianParser было получено свыше десяти тысяч объявлений в нужных регионах.

К следующей задаче подходит такое начало, как соединение собранных данных в одну таблицу при помощи написанной функции с библиотекой Pandas.

После сбора всей информации воедино и уборки дубликатов, нужно узнать, какого типа наши данные (рис.1), так как отталкиваясь от типа данных, мы будем применять разные методы к их сортировке.

Теперь мы убираем -1 как базовое не собранное значение и смотрим, какие данные у нас смогли собраться(рис.2) при помощи библиотеки missingno. Как мы можем увидеть, object_type и heating_type практически нигде не указываются, значит мы вынуждены их удалить, так как строить анализ будет невозможно.

Далее отсеиваем ненужные данные и форматируем некоторые столбцы, чтобы их было легче анализировать с помощью сторонних инструментов(рис.3). После выполнения объёмной чистки данных, нужно проверить их состояние — смотрим внутрь файла и бегло проверяем на аномалии, в случае их отсутствия приступаем к кодовой проверке данных(рис.4-7).

После полной очистки данных вручную и программно можем приступать к сбору графиков/аналитической работе при помощи библиотеки matplotlib для вывода графических изображений. Например будет 4 графика:

- 1. Цена за м^2 по городам.
- 2. Цена за м^2 в зависимости от материала, используемого при постройке злания.
- 3. Количество объявлений по городам.
- 4. Количество объявлений по годам постройки здания.

Для первых и последних двух графиков будем использовать один метод подсчёта данных.

Первый метод — отбор цены за квадратный метр по трафаретному коду и запись в csv-файл(рис.8).

Второй метод – использование встроенной функции .value_counts() и запись в csv-файл для дальнейшей обработки(рис.9).

В итоге получаем графики(рис. 10), на основе которых уже можно проводить анализ, но мы перейдём к составлению графиков на Power BI.

Power BI – максимально удобный инструмент для составления графиков и аналитики данных.

Для первого графика мы выбираем данные price, меняем сумму на «среднее» и включаем их в график, на другую ось ставим «year_of_construction». Выбираем тип графика и получается примерный график со средней ценой квартиры, в зависимости от года его постройки(рис.11).

Второй график будет содержать в себе среднюю цену квартиры, в зависимости от материала здания(рис.12), просто вместо года постройки ставим тип материала. Добавлю к этим данным среднюю цену по виду отделки(рис.13). UPD: все графики были заменены на один dashboard(рис.14).

Аналитика данных

Благодаря выведенным графикам, можно сделать выводы, что цена в основном зависит от типа отделки, материала дома, города. От года постройки зданий зависит лишь их количество на рынке и количество комнат во время СССР, а на цену никак не влияет.

Заключение

В результате аналитической работы были собраны, отсортированы, почищены данные, простроены удобные для анализа графики, благодаря которым получилось выявить не маловажные критерии в оценивании стоимости недвижимости в Московской Области, Москве и Новой Москве. Основными факторами, оказывающими влияние на стоимость, выявились тип отделки, материала здания и расположение. Полученные данные могут быть использованы для дальнейшей разработки прогностических моделей.

```
Открываем наши драгоценные данные
   df = pd.read_csv('11K.csv')
Смотрим форму и колонки нашых данных
   print(df.dtypes)
   df.shape
 author
                        object
 author_type
                       obiect
 url
                       object
 location
                       object
 deal_type
                       object
                       object
 accommodation_type
                       float64
 floor
                       float64
 floors count
 rooms_count
                       float64
 total meters
                       obiect
 price
                       float64
year of construction
                      obiect
 object_type
                       float64
 house material type
                      obiect
 heating_type
                       float64
                       object
 finish_type
                       object
 living_meters
 kitchen meters
                       object
                       float64
 phone
 district
                       object
                       object
 street
                       object
 house number
 underground
                       object
 residential complex
                      object
 dtype: object
 (10582, 24)
                                         (рис.1)
```

```
Чистим от отрицательных значений
     df = df.replace("-1",np.nan) # Убираем все виды плохих или не собранных данных
df = df.replace(-1.0,np.nan) # Убираем все виды плохих или не собранных данных
      df = df.replace("���",0) # Убираем все виды плохих или не собранных данных df.to_csv("half_11K.csv", index=False) # Тут запись в отдельный файл, дабы поэтапно отслеживать, на каком моменте появляются аномальные действии # Ещё мы заходим в сохранённый файл и меняем ��� на 0, дабы избежать ошибок (сделал отдельно)
Чекаем корректность данных
             1.0
                                                                                                                 10582
     0.8
                                                                                                                 8465
                                                                                                                 6349
     0.6
                                                                                                                 4232
     0.4
                                                                                                                 2116
     0.0
                                 Toons counteer
                                     year of consti
                                              Police Mare
```

(рис.2)

```
df = pd.read_csv('half_11K.csv')
list = ["Напишите автору","Залоговая недвижимость","Аукцион","Позвоните автору","Подписаться на дом"]
for obj in list: df = df.replace(obj, np.nan) # Стираем
list1 = ["total_meters","living_meters","kitchen_meters"] # Данные, которые мы очистим от метров в ква
del df['author'] # Не надо
del df['author_type'] # Не надо
del df['accommodation_type'] # Это тип хаты
del df['deal_type'] # Тип сделки, у нас только продажа
del df['residential_complex'] # Много пропущено и не особо надо
del df['heating type'] # Подогрев, нигде не указан
del df['object_type'] # Чет ваще не нужная штука
del df['phone'] # звонить для аналитики не нужно
list = ["price","year_of_construction","floor","floors_count","rooms_count"] # Создаём лист с данными
   df[obj] = df[obj].astype(int) # Оформляем INT для чистоты разума
for obj in list1: df[obj] = pd.to_numeric(
   df[obj].str.replace(',', '.').apply(lambda x: x[:-3] if pd.notna(x) else np.nan),
    ).astype('float64')
<mark>df.to_csv("tret'_11K.csv", index=False)</mark> # Различия между cleaned_11К и tret'_11К - это ручной отброс
```

(рис.3)

```
# Связи числовых данных
 plt.figure(figsize=(8,4))
 plt.xticks(rotation=30, ha="right")
 plt.show()
                                                                                                   1.00
                floor -
                                0.71
                                         0.12
                                                  0.21
                                                          0.14
                                                                            0.15
                                                                                     0.2
                                                                                                   0.75
        floors_count -
                        0.71
                                         0.073
                                                 0.16
                                                          0.11
                                                                   0.14
                                                                           0.087
                                                                                    0.24
                                                                                                  - 0.50
                        0.12
        rooms_count -
                                0.073
                                                  0.76
                                                          0.41
                                                                            0.77
                                                                                    0.43
                                                                                                   0.25
        total_meters -
                        0.21
                                0.16
                                         0.76
                                                          0.72
                                                                  -0.0052
                                                                            0.93
                                                                                    0.65
                                                                                                   0.00
               price -
                        0.14
                                0.11
                                         0.41
                                                  0.72
                                                                  -0.004
                                                                            0.59
                                                                                    0.52
                                                                                                   -0.25
year_of_construction -
                       0.095
                                0.14
                                        -0.031
                                                -0.0052 -0.004
                                                                     1
                                                                                    0.041
                                                                                                   -0.50
       living_meters -
                        0.15
                                         0.77
                                                  0.93
                                                          0.59
                                                                                    0.51
                                                                                                   -0.75
                        0.2
                                0.24
                                                  0.65
     kitchen_meters -
                                         0.43
                                                          0.52
                                                                  0.041
                                                                            0.51
                                                  year of construction
                                                                                                   -1.00
                                                                        kitchen meters
                       Roors_count
                               rooms_count
                                       total_meters
                                                                 living meters
                     ROOT
```

(рис.4)

```
# Проверка на отсутствующие данные в процентах
   for col in df.columns:
       pct_missing = np.mean(df[col].isnull())
       print('{} - {}%'.format(col, round(pct_missing*100))
url - 0%
location - 4%
floor - 0%
floors_count - 0%
rooms count - 0%
total meters - 0%
price - 0%
year_of_construction - 0%
house material type - 88%
finish type - 84%
living meters - 23%
kitchen meters - 15%
district - 43%
street - 14%
house number - 11%
underground - 37%
                                                             (рис.5)
```

```
# Проверка на отсутствующие данные в цифрах
 msn.bar(df, figsize=(6,3), fontsize=10, color=(1, 0.75, 0.8))
 plt.show()
           ૢૺઌ૽ૺૢઌ૱ૢ૽ઌ૱ૢ૽ઌ૱ૢ૽ઌ૱ૢ૽ઌ૱ૢ૽ઌ૱ૢ૽ઌ૱ૢ૽ઌ૱ૢ૽ઌ૱ૢઌ૱ઌઌૢૹઌઌૹૹઌૢૹઌઌૹઌૣૡ૽૽૾ૢૹઌૣ૽ૹ
1.0
                                                                             10581
0.8
                                                                            8464
0.6
                                                                            6348
0.4
                                                                            4232
0.2
                                                                            2116
                                    wind neter neters
           of count count the set of
0.0
                  Pear of Constitution
                                                        residential complet
                      A house material rich
                                                   house humber
        ROOFS COUNT
                                                 district
                                 and type
```

(рис.6)

print(df.dtypes) df.shape url object location object floor int32 floors_count int32 rooms count int32 total_meters float64 int32 price year of construction int32 house_material_type object finish type object living_meters float64 kitchen meters float64 district object street object house_number object underground object dtype: object (10582, 16)

(рис.7)

```
df = pd.read_csv('cleaned_11K.csv')
list_of_cities = df['location'].unique()

def price_for_meter(location):
    city = df[df['location']==location]
    city_price = city['price'].sum();    cleaned_data = city['total_meters'].sum()
    return round(city_price/cleaned_data,2)

with open("dash_info_fifth.csv", 'w', newline='', encoding='UTF-8') as csvfile: # Создаем таблицу, чтобы y
    fieldnames = ['city', 'price_for_meter'];    writer = csv.DictWriter(csvfile, fieldnames=fieldnames)
    writer.writeheader()
    for city in list_of_cities: writer.writerow({'city': city, 'price_for_meter': price_for_meter(city)})
# Потом вручную удаляем строчку "nan,nan"

df = pd.read_csv('cleaned_11K.csv')

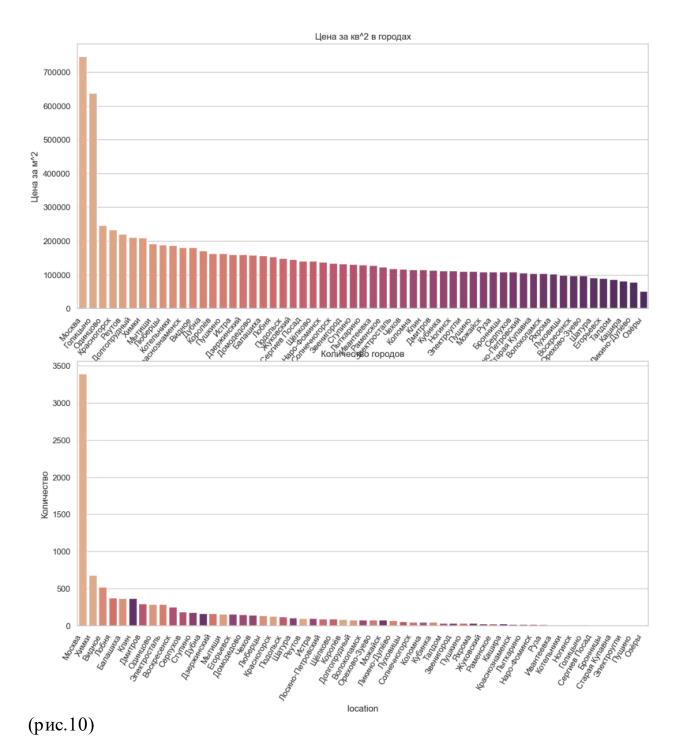
list_of_cities = df['location'].unique()

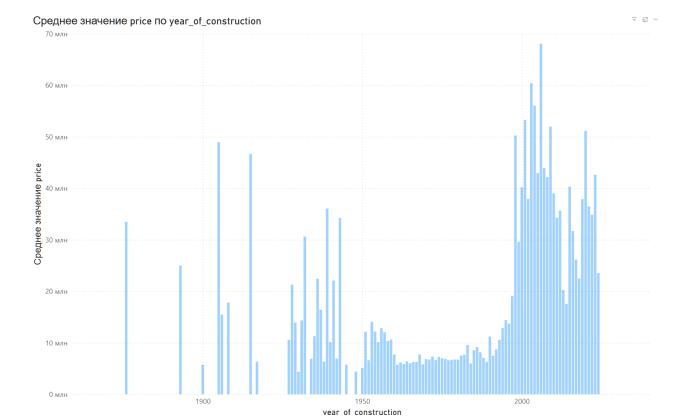
def price_for_meter(location):
    city = df[df['location']==location]
    city_price = city['price'].sum();    cleaned_data = city['total_meters'].sum()
    return round(city_price/cleaned_data,2)

with open("dash_info_fifth.csv", 'w', newline='', encoding='UTF-8') as csvfile: # Создаем таблицу, чтобы удобнее было вынимать данные
    fieldnames = ['city', 'price_for_meter'];    writer = csv.DictWriter(csvfile, fieldnames=fieldnames)
    writer.writeheader()
    for city in list_of_cities: writer.writerow({'city': city, 'price_for_meter': price_for_meter(city)})
# Потом вручную удаляем строчку "nan,nan"
```

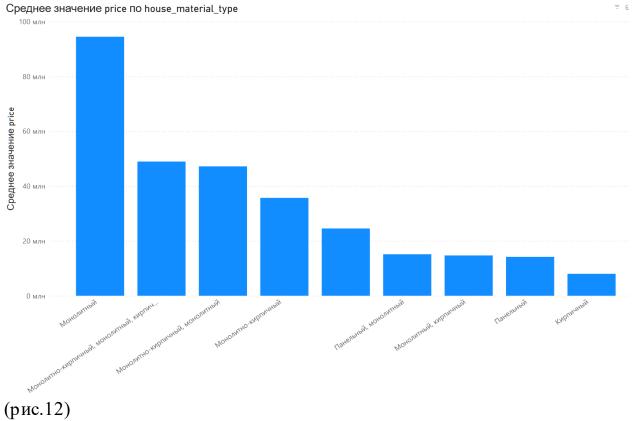
(рис.8)

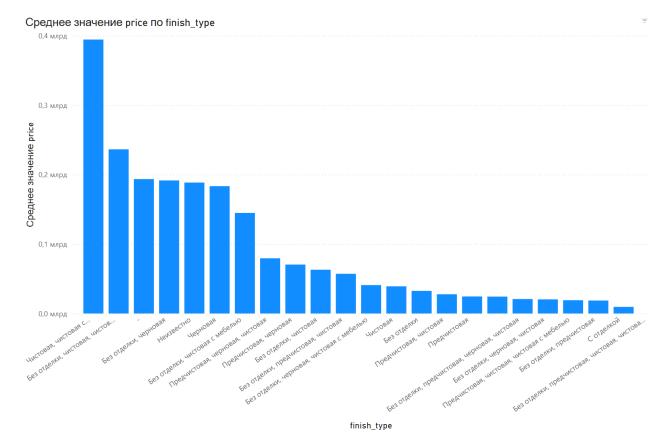
```
df['year_of_construction'] = df['year_of_construction'].fillna(-1) =
df['year_of_construction'] = df['year_of_construction'].astype(int)
df['year_of_construction'].value_counts().to_csv('years_count.csv')
(puc.9)
```



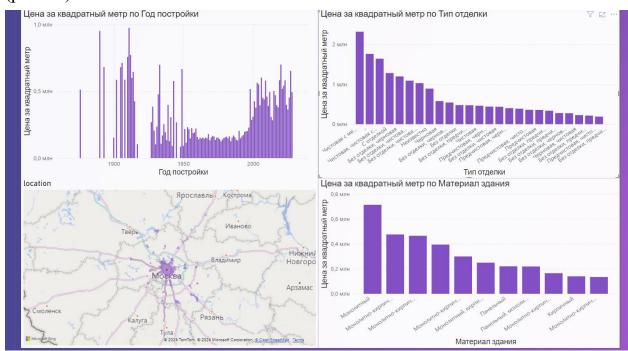












(рис.14)