

# Министерство образования и науки РФ

Федеральное государственное автономное образовательное учреждение высшего образования «Казанский (Приволжский) федеральный университет»

Институт математики и механики им. Н.И.Лобачевского Кафедра  
математического анализа

Направление: 02.03.01 – Математика и компьютерные науки

## ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА «СБОР ДАННЫХ И АНАЛИЗ ДАННЫХ С HTML-СТРАНИЦ»

### Работа завершена:

Студент гр. 05-503

"\_\_" \_\_\_\_\_ 2019 г. \_\_\_\_\_ Березин Н.С.

### Работа допущена к защите:

Научный руководитель

Кандидат физико-математических наук,

ассистент кафедры математического анализа

"\_\_" \_\_\_\_\_ 2019 г. \_\_\_\_\_ Новиков А.А.

Заведующий кафедрой

математического анализа

доктор физико-математических наук,

профессор

"\_\_" \_\_\_\_\_ 2019 г. \_\_\_\_\_ Насыров С.Р.

# Оглавление.

1. Введение. ....	3
2. Предварительные сведения.....	5
2.1. Язык HTML.....	5
2.2. Структура HTML-документа. ....	5
2.3. Система управления базами данных MySQL.....	6
2.4. Язык PHP.....	7
2.5. Язык Python.....	7
3. Материалы и методы.....	8
4. Основная часть. ....	9
4.1. Сбор данных с HTML-страницы.....	9
4.2. Построение регрессионной модели. ....	22
5. Заключение и выводы. ....	26
6. Список литературы. ....	27
7. Приложение. ....	28
7.1. Приложение 1. Код алгоритма на языке PHP. ....	28
7.2. Приложение 2. Часть кода базы данных на языке SQL.....	30
7.3. Приложение 3. Код алгоритма на языке Python.....	33

# 1. Введение.

Данная работа посвящена сбору данных с html страниц и последующим их анализом. Целью было собрать данные с объявлений на сайте недвижимости и построить регрессионную модель для их анализа. Регрессионная модель должна определить зависимость цены от других параметров указанных в объявлении, такие как район где расположена квартира, ближайшая станция метро, площадь квартиры. Также задачей было обучить данную модель с целью дальнейшего предсказания стоимости квартиры на основе собранных признаков.

Эта работа является актуальной, так как в настоящее время рынок недвижимости один из крупнейших рынков в мире. Только в России по данным Росреестра за первое полугодие 2018 года было продано более 4,2 миллионов объектов недвижимости. Из них более 3,9 миллионов было продано на вторичном рынке. [1] Исходя из этих данных можно сделать вывод о том, что система для рекомендации стоимости недвижимости была бы очень востребована на этом рынке.

В связи с этим был реализован алгоритм по сбору данных с объявлений на сайте по продаже недвижимости на языке PHP с последующей записью этих данных в базу данных с помощью системы управления базами данных MySQL и дальнейшим анализом этих данных на языке Python.

В пункте 2 приведены предварительные сведения в которых рассказывается, что такое язык HTML, как устроен HTML-документ, что такое язык PHP, что такое система управления базами данных MySQL и что такое язык Python.

В пункте 3 рассказано с помощью чего была сделана данная работа.

В пункте 4 расписано как был реализован алгоритм на языке PHP по сбору данных и записи их в базу данных используя СУБД MySQL, а также алгоритм анализа этих данных на языке Python.

В пункте 5 приведены итоги по проделанной работе.

В пункте 7 представлены код алгоритма на языке PHP, код базы данных на языке SQL и код алгоритма на языке Python.

## 2. Предварительные сведения.

### 2.1. Язык HTML.

HTML — язык разметки веб-документов. Позволяет создавать веб-страницы, которые потом отображаются в окне браузера компьютера или мобильного устройства. HTML 5 – последняя версия языка. В этой версии он имеет большое количество возможностей, таких как создание таблиц, отображение картинок, музыки или видео. [2]

### 2.2. Структура HTML-документа.

HTML-документ состоит из открывающихся и закрывающихся тегов. Каждый HTML-документ должен начинаться с тега DOCTYPE который определяет версию используемого HTML и также определяет соответствующий DTD-файл в интернет. DTD-файл – это XML-документ, который определяет какие теги, атрибуты и их значения действительны для конкретной версии HTML. После тега DOCTYPE, идет тег <html>. Элементы, которые находятся между тегами <html> и </html> образуют дерево элементов, где элемент <html> является корневым.

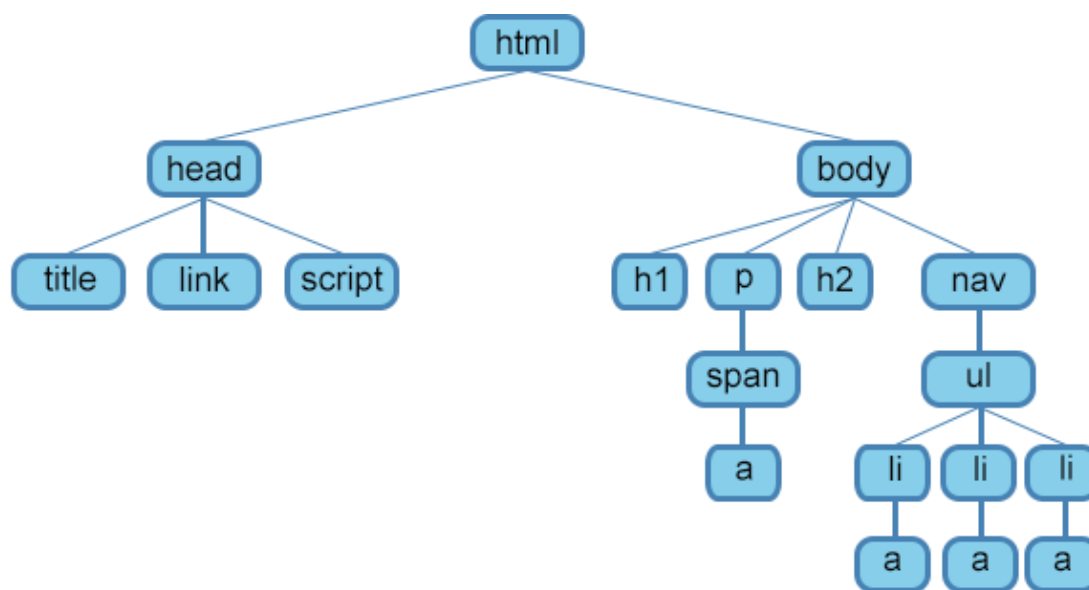


Рисунок 1. Простейшая структура HTML-документа.

Внутри тега <html> находятся два элемента <head> и <body>.

Внутри тега <head> содержится техническая информация, такая как кодировка, заголовки, описание, ключевые слова для поисковых систем, а также в этом теге происходит подключение файла отвечающего за оформление страницы CSS, и подключение файла, отвечающего за какую-либо динамику на веб-странице JavaScript. Вся информация, находящаяся внутри тега <head> не отображается в браузере, но содержит указания, как правильно браузеру отображать эту страницу.

Тег <body> содержит в себе данные которые должны отобразиться в браузере.

У каждого элемента могут быть различные атрибуты, например, id= " " или class= " ". Эти атрибуты помогают обратиться к этим элементам для того чтобы задать им стили с помощью CSS, либо придать им какую-то динамику с помощью JavaScript.

## 2.3. Система управления базами данных MySQL.

База данных — совокупность данных, организованных в соответствии с концептуальной структурой, описывающей характеристики этих данных и взаимоотношения между ними, которая поддерживает одну или более областей применения. [3]

MySQL — свободная реляционная система управления базами данных. Разработку и поддержку MySQL осуществляет корпорация Oracle. MySQL является решением для малых и средних приложений. Обычно MySQL используется в качестве сервера, к которому обращаются локальные или удалённые клиенты, однако в дистрибутив входит библиотека внутреннего сервера, позволяющая включать MySQL в автономные программы. Гибкость СУБД MySQL обеспечивается поддержкой большого количества типов таблиц: пользователи могут выбрать как таблицы типа MyISAM, поддерживающие

полнотекстовый поиск, так и таблицы InnoDB, поддерживающие транзакции на уровне отдельных записей. Более того, СУБД MySQL поставляется со специальным типом таблиц EXAMPLE, демонстрирующим принципы создания новых типов таблиц. Благодаря открытой архитектуре и GPL-лицензированию, в СУБД MySQL постоянно появляются новые типы таблиц. [4]

## 2.4. Язык PHP.

PHP — скриптовый язык, часто применяемый для разработки веб-приложений. Является самым популярным языком для создания динамических веб-сайтов. Чаще всего выполняет роль сервера, то есть обрабатывает команды, приходящие от html-страниц, и отправляет обработанную информацию. Также выводит и записывает данные из базы данных. Имеет большую популярность из-за простоты и удобства. Включает в себя много различных функций для простого и удобного создания веб-приложений. [5]

## 2.5. Язык Python.

Python — высокоуровневый язык программирования общего назначения, ориентированный на повышение производительности разработчика и читаемости кода. Синтаксис ядра Python минималистичен. В то же время стандартная библиотека включает большой объем полезных функций.

Python поддерживает структурное, объектно-ориентированное, функциональное, императивное и аспектно-ориентированное программирование. Основные архитектурные черты — динамическая типизация, автоматическое управление памятью, полная интроспекция, механизм обработки исключений, поддержка многопоточных вычислений, высокоуровневые структуры данных. Поддерживается разбиение программ на модули, которые, в свою очередь, могут объединяться в пакеты. [6]

### 3. Материалы и методы.

Алгоритм по сбору данных был написан на языке PHP версии 7.2. Выбор пал на язык PHP, потому что он наиболее удобен в работе с веб-страницами. Для написания программы была использована среда разработки PhpStorm версии 2018.1.2 от компании JetBrains. Эта среда разработки была выбрана в связи с ее удобством, так как она имеет приятный интерфейс и делает подсказки в написании кода. Так же эта среда разработки позволяет просматривать базы данных различных СУБД. Полученные в ходе работы алгоритма данные записывались в базу данных с помощью СУБД MySQL версии 8.0. Использование именно этой СУБД обусловлено тем, что она является одной из самых быстрых, в моей работе это важно, так как собираются и записываются в базу большие объемы данных. Для запуска алгоритма была использована платформа Open Server версии 5.2.2. Эта платформа представляет из себя локальный веб-сервер для операционной системы Windows. Open Server включает в себя большой набор серверного программного обеспечения. Для своей работы из платформы Open Server я использовал веб-сервер Apache версии 2.2 и веб-интерфейс для администрирования СУБД MySQL – phpMyAdmin версии 4.8. Алгоритм анализа данных был написан на языке Python версии 3.7. Python был выбран как наиболее удобный язык для математических вычислений. Для написания алгоритма на Python была использована среда разработки PyCharm версии 2018.2 от компании JetBrains. Была выбрана именно эта среда разработки так как она очень удобна в использовании, имеет приятный интерфейс и дает подсказки в написании кода.



## 4. Основная часть.

### 4.1. Сбор данных с HTML-страницы.

Сбор данных осуществлялся с сайта по продаже недвижимости с помощью написанного алгоритма на языке PHP. Помимо стандартного набора функций языка PHP для написания алгоритма также использовалась библиотека phpQuery.

Предварительно скачаем и установим локальный веб-сервер. Загрузим с официального сайта платформу OpenServer [7], включающую в себя сервер Apache и интерфейс для администрирования СУБД MySQL phpMyAdmin. Далее произведем установку этой платформы. После установки, в папке domain, создаем папку, ей можно дать любое название, назовём ее diplom. Потом внутри этой папки создадим php файл с названием index, в котором будет весь код алгоритма для сбора данных. Теперь запустим сервер и откроем phpMyAdmin. Далее нужно создать базу данных. Для этого нужно нажать на кнопку «Создать БД» (Рисунок 2).

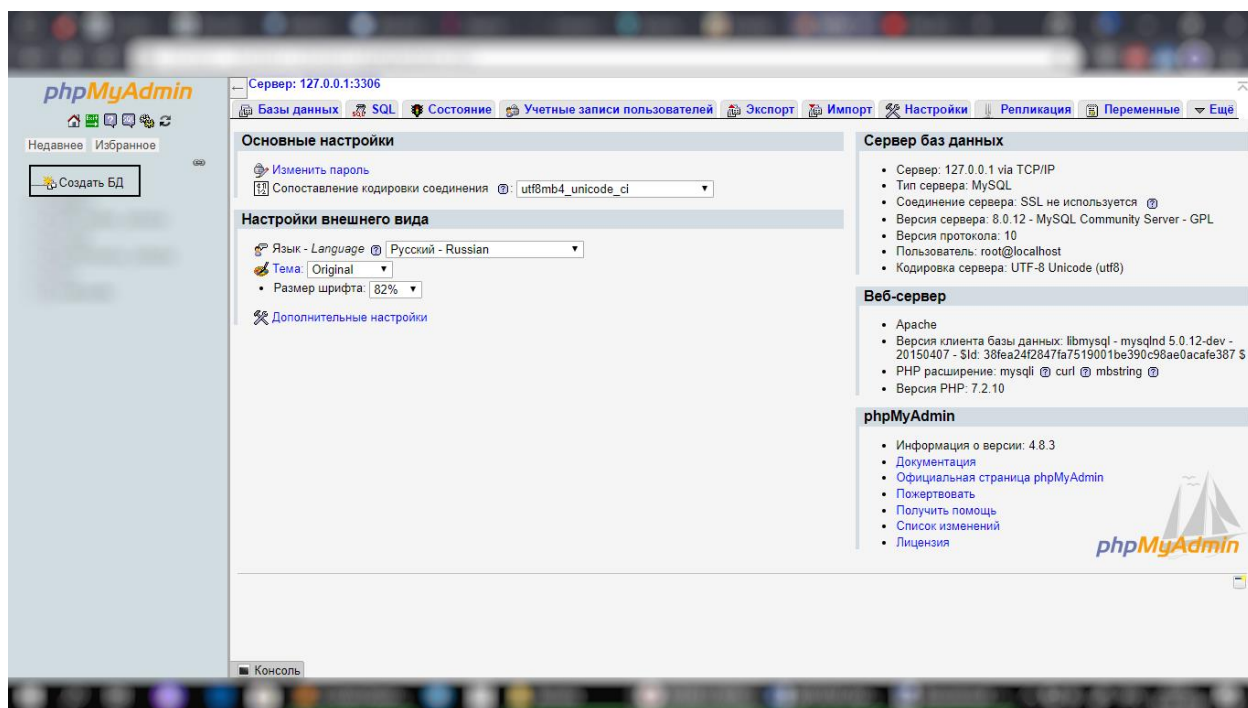


Рисунок 2. Создание БД в phpMyAdmin.

После этого откроется окно, нужно написать имя базы данных, назовем ее `diplom`, и нажмем на кнопку создать (Рисунок 3).

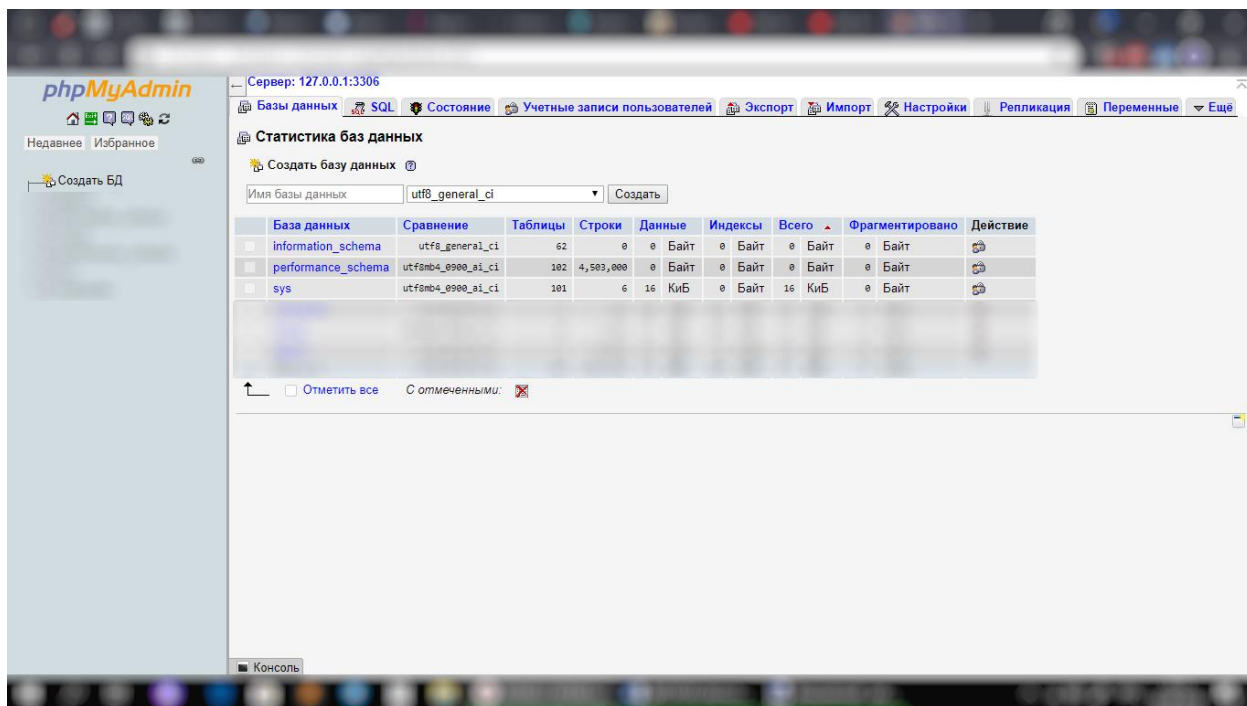


Рисунок 3. Создание БД в phpMyAdmin.

Далее создаем таблицу, у которой будет 11 столбцов и название `apartaments` (Рисунок 4).

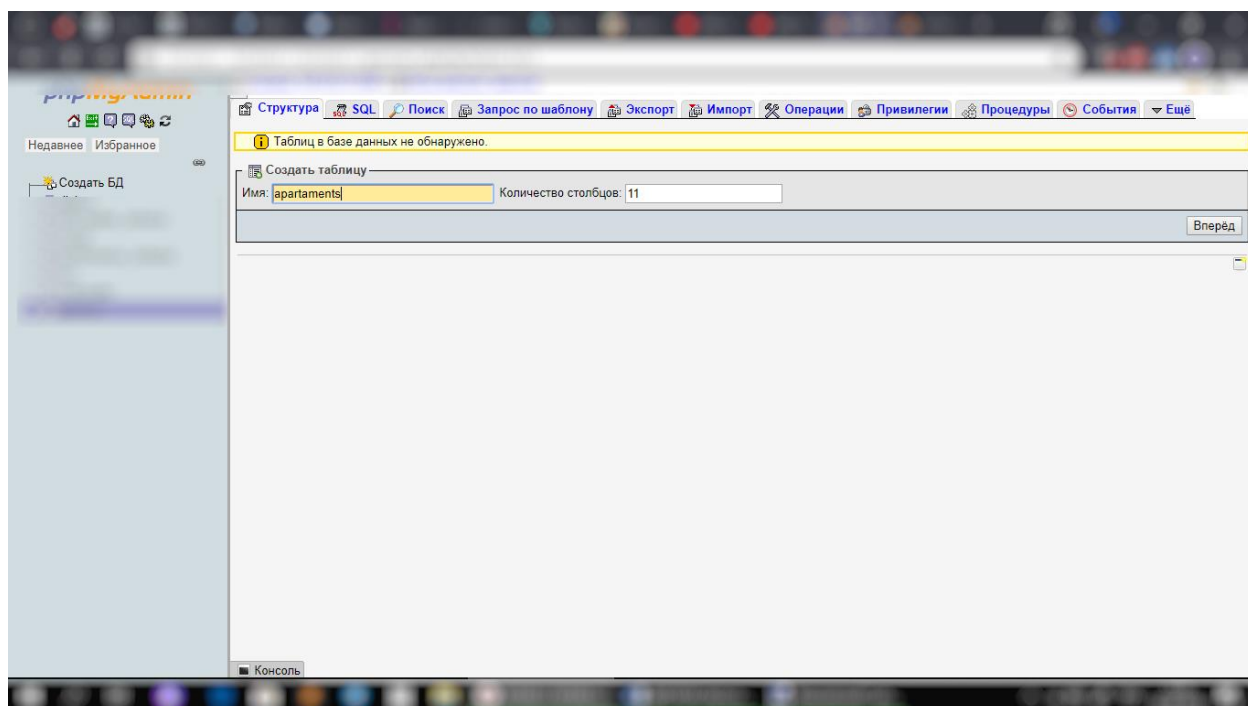


Рисунок 4. Создание таблицы.

Таблица `apartaments` содержит 11 столбцов (Рисунок 5):

1. Столбец `id` имеет тип `int` – в этот столбец записывается номер строки с данными. Пример записи: 1.
2. Столбец `quantity_of_rooms` имеет тип `int` – в этот столбец записывается количество комнат в квартире. Пример записи: 1.
3. Столбец `area_of_apartaments` имеет тип `int` – в этот столбец записывается площадь квартиры. Пример записи: 23.
4. Столбец `floor` имеет тип `varchar` – в этот столбец записывается на каком из скольких этажей находится квартира. Пример записи: 17/23 этаж.
5. Столбец `metro` имеет тип `varchar` – в этот столбец записывается название ближайшей к квартире станции метро. Пример записи: Охотный ряд.
6. Столбец `time_to_metro` имеет тип `int` – в этот столбец записывается время ходьбы или проезда до ближайшей станции метро. Пример записи: 5.
7. Столбец `name_of_district` имеет тип `varchar` – в этот столбец записывается название района, в котором находится квартира. Пример записи: Дмитровский.
8. Столбец `price` имеет тип `int` – в этот столбец записывается цена квартиры. Пример записи 8564987.
9. Столбец `price_per_square_metr` имеет тип `int` – в этот столбец записывается цена квартиры за квадратный метр. Пример записи: 372390.
10. Столбец `date` имеет тип `date` – в этот столбец записывается дата сбора данных. Пример записи: 2019-05-18.
11. Столбец `collection number` имеет тип `int` – в этот столбец записывается номер сбора данных. Пример записи: 1.

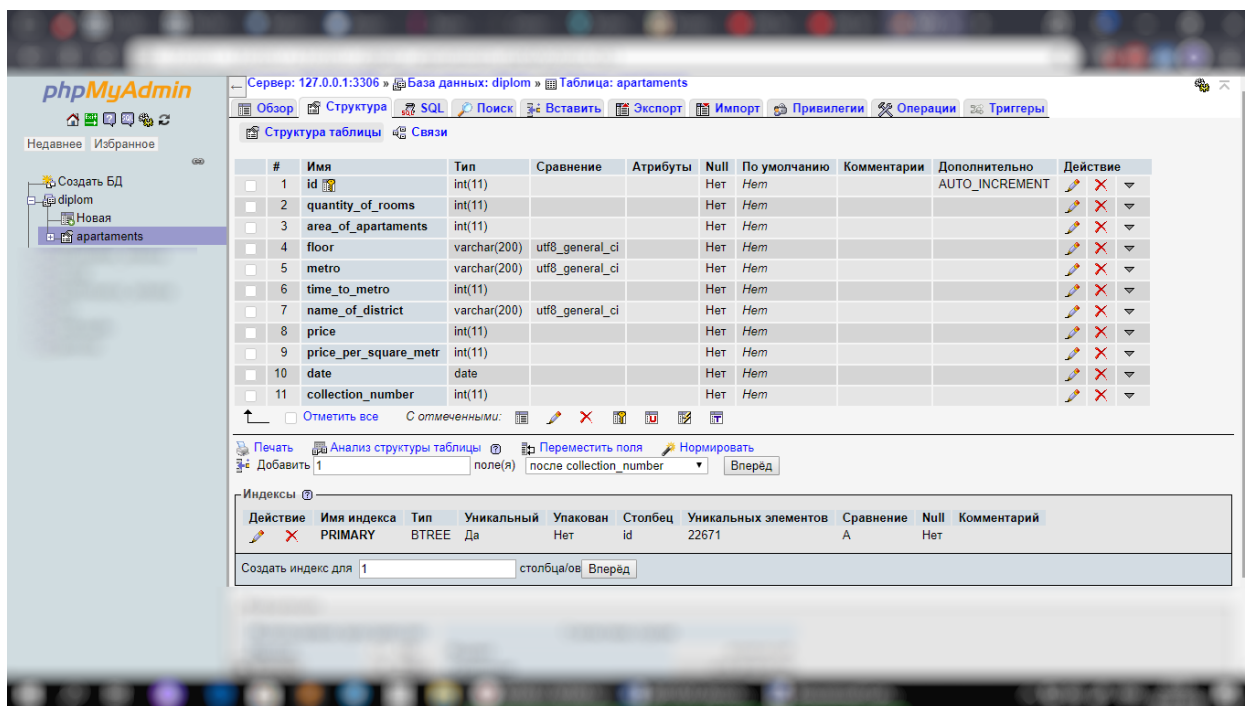


Рисунок 5. Структура таблицы apartaments.

Платформа OpenServer установлена, база данных создана. Далее подробно разберем алгоритм по сбору данных.

Для начала подключим библиотеку phpQuery, для этого нужно добавить файл «phpQuery-onfile.php» в папку с проектом, далее с помощью функции include() или с помощью функции require() подключаем библиотеку к проекту. Отличие этих функций заключается в том, что при неудачном подключении файла функция include() выдает предупреждение, но продолжает работу, а функция require() выдает ошибку и полностью останавливает работу программы. Так как, библиотека phpQuery играет важную роль в работе алгоритма, подключим ее с помощью функции require(), то есть вот так require("phpQuery-onfile.php").

После подключение библиотеки перейдем к написанию алгоритма. Первым делом получим весь HTML-код со страницы с помощью функции \$file = file\_get\_contents(\$url), где \$url это ссылка на нужный нам сайт, в данном случае это сайт по продаже недвижимости. Далее используя библиотеку

phpQuery, мы создаем документ с помощью \$doc = phpQuery::newDocument(\$file), где \$file это HTML-код, который мы получили с помощью функции file\_get\_contents(\$url). В результате мы получим объект, к которому можно будет применять определенные методы. Далее на нужной нам странице нужно выделить нужную информацию, для начала отделим каждое объявление. На одной странице располагается 25 объявлений. Объявления выглядят как блоки располагающиеся друг за другом. В каждом объявлении располагается основная информация о квартире. Такая как, фотографии дома и квартиры изнутри, количество комнат в квартире, площадь квартиры в квадратных метрах, на каком этаже располагается квартира и сколько всего этажей в доме, ближайшее метро, расстояние до метро в минутах, стоимость квартиры в рублях, стоимость за квадратный метр в рублях, адрес дома, небольшое описание и контакты для связи с владельцем квартиры (Рисунок 6).

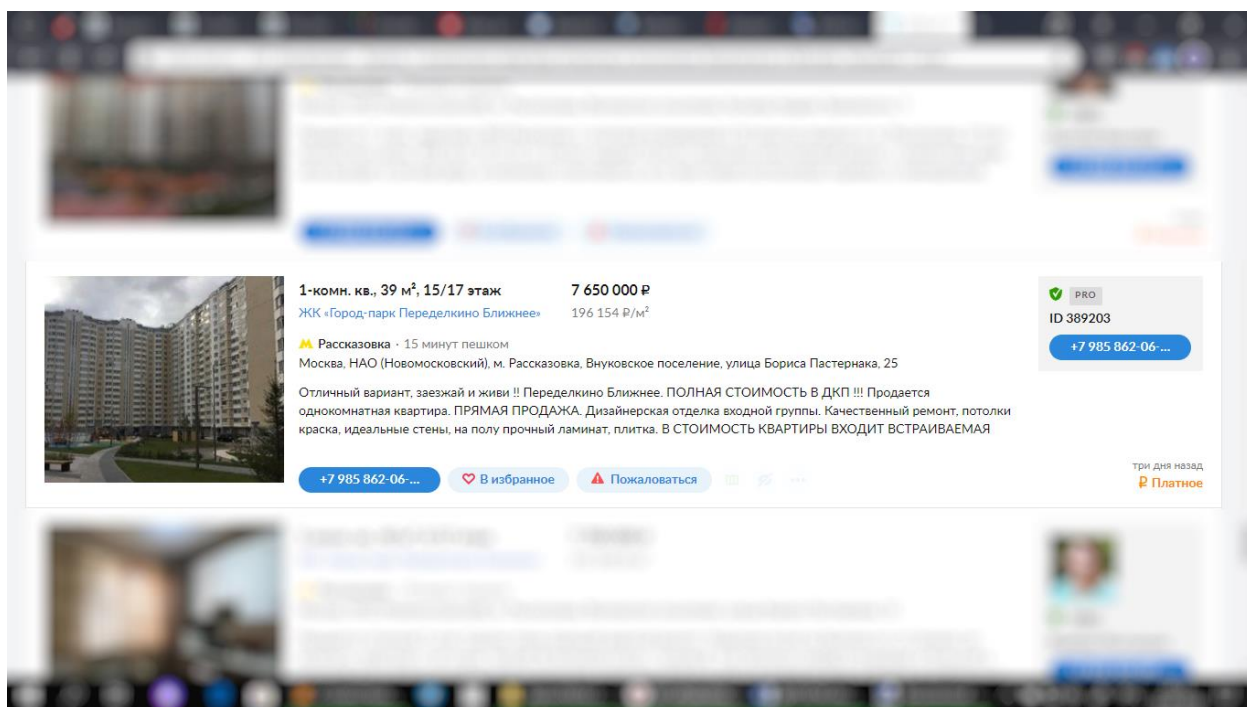


Рисунок 6. Вид объявления на сайте по продаже недвижимости.

Каждый блок объявления на сайте имеет одинаковый атрибут class. Здесь

каждое объявление имеет значение атрибута `class = "_93444fe79c--card--_yguQ"`. В данном случае объявления имеют одинаковый атрибут для того чтобы можно было привести все объявления к одному стилю и задать одинаковое оформление.

Чтобы узнать значение атрибута у требуемого нам блока, необходимо зайти на выбранный нами сайт, нажать на нужный нам блок правой кнопкой мыши и выбрать пункт исследовать элемент (Рисунок 7).

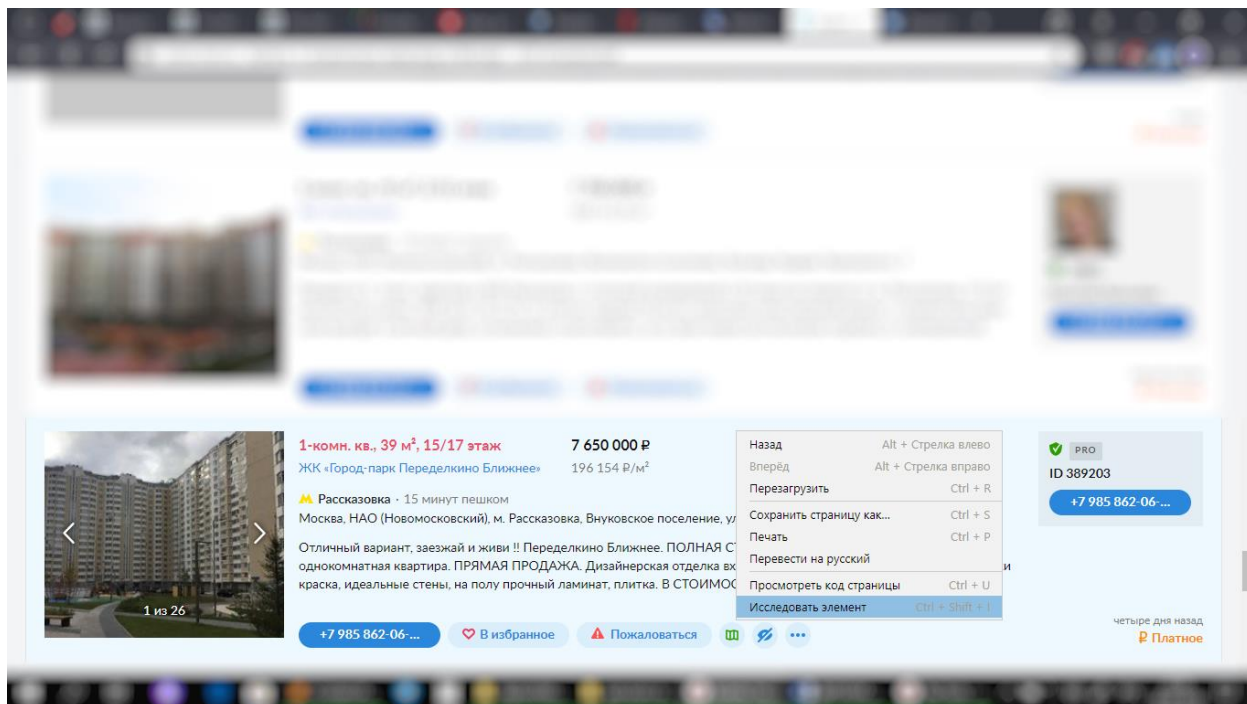


Рисунок 7. Исследование элемента.

После нажатие на пункт «Исследовать элемент» откроется окно DevTools (Инструменты разработчика), в котором во вкладке Elements можно посмотреть HTML-код всей страницы. Во вкладке Console показываются ошибки JavaScript кода, а также там можно писать JavaScript код, например, чтобы узнать значение переменной. Во вкладке Sources можно посмотреть все подключенные файлы к этой странице, а также открыв JavaScript файл можно провести его пошаговое выполнение и сделать отладку. У нас по умолчанию откроется вкладка Elements и будет выделен нужный нам блок. Как мы видим (Рисунок 8), нужный нам блок имеет тег `<div>` и атрибут `class = "_93444fe79c--card--_yguQ"`. Также (Рисунок 8,



рисунок 9) видно, что другие объявления имеют такой же тег `<div>` с таким же значением атрибута `class = "_93444fe79c--card--_yguQ"`.

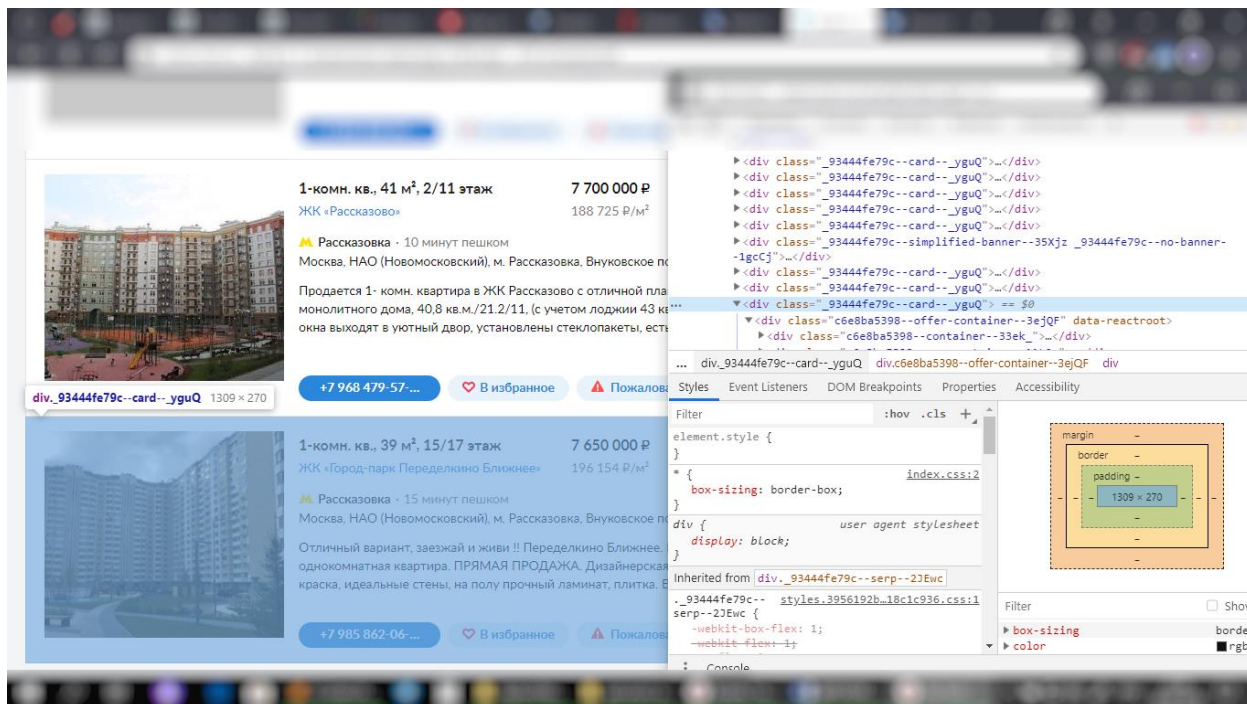


Рисунок 8. Исследование элемента.

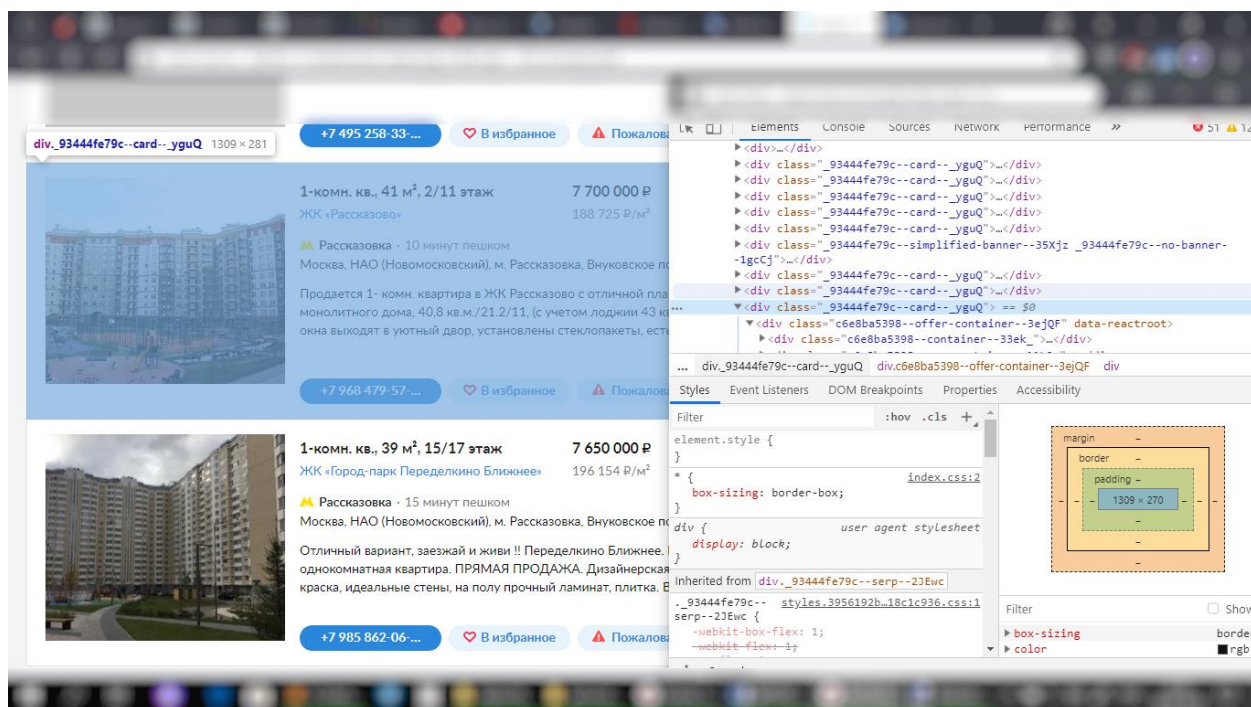


Рисунок 9. Исследование элемента.

Итак, так как мы узнали значение атрибута, у блока, который содержит информацию об объявлении, теперь его можно найти в нашем документе. Для

поиска по объекту, который вернулся с помощью `phpQuery::NewDocument($file)`, мы используем метод `$doc→find()`, где `$doc` это сам документ. Чтобы найти нужное нам значение атрибута в документе, в скобках метода `find()`, мы пишем значение атрибута в кавычках и ставим точку в начале, чтобы показать, что это атрибут `class`, то есть так `$dv1 = $doc→find("._93444fe79c--card--_yguQ")`. В переменную `$dv1` записываются только данные с объявлениями. С этого момента поиск нужных нам данных будет производиться по переменной `$dv1`. С помощью алгоритма собираются следующие данные, количество комнат в квартире, площадь квартиры в квадратных метрах, на каком этаже располагается квартира и сколько всего этажей в доме, ближайшее метро, расстояние до метро в минутах, стоимость квартиры в рублях, стоимость за квадратный метр в рублях (Рисунок 10).

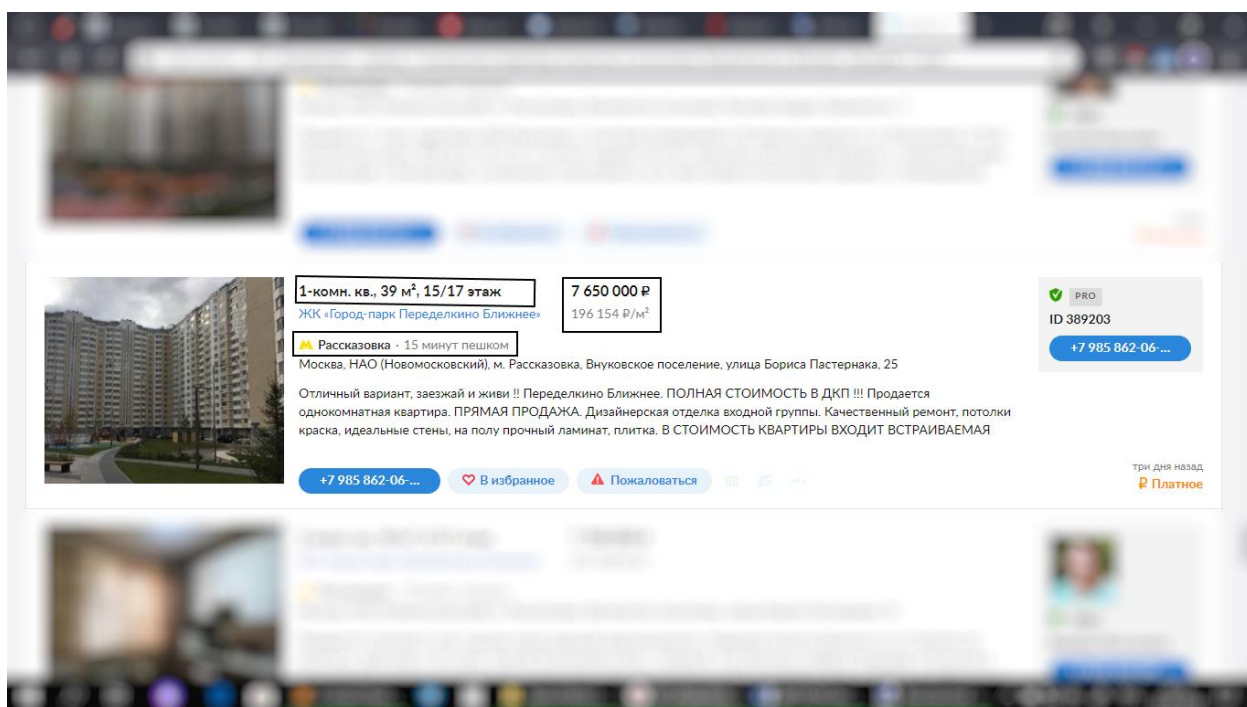


Рисунок 10. Собираемые данные (выделены).

Сперва с помощью окна DevTools узнаем значение атрибута у блока, который содержит в себе информацию о количестве комнат, площади квартиры и этажность. Нужный нам блок также имеет тег `<div>` и обладает значением атрибута `class="сбе8ba5398--title--2CW78"` (Рисунок 11). Блок с этой



информацией содержится в каждом объявлении на странице, и имеет одинаковое значение атрибута class.

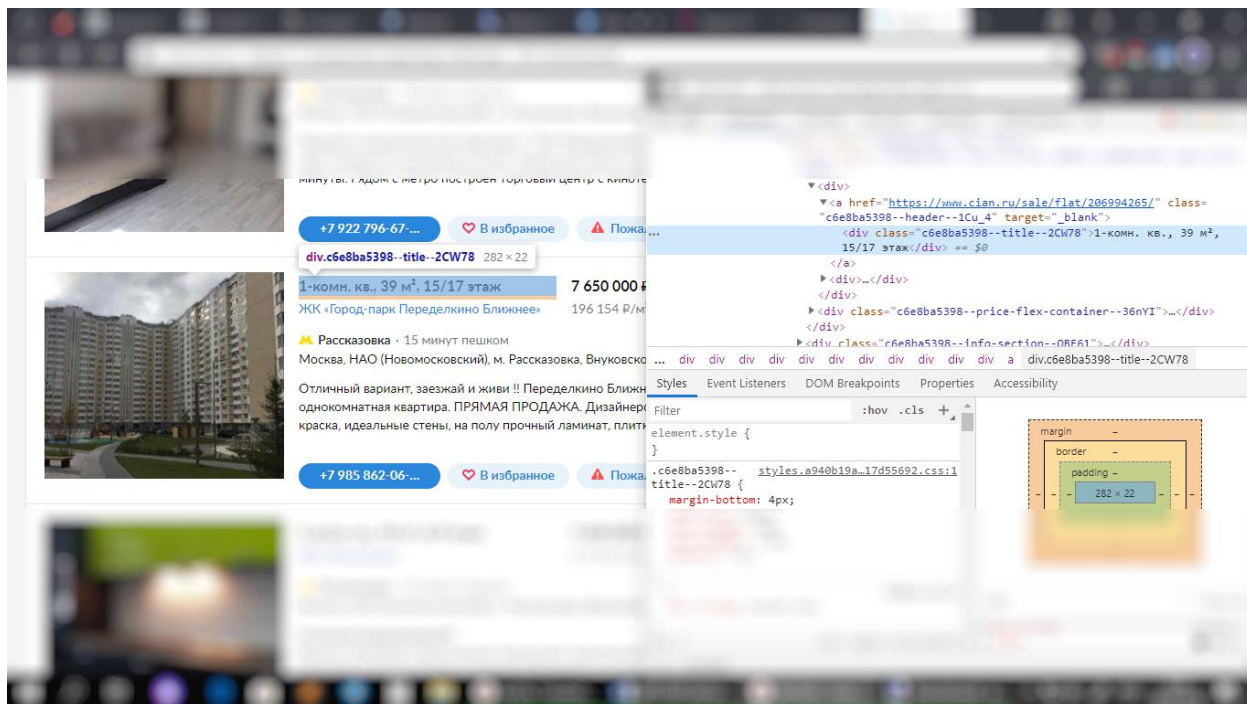


Рисунок 11. Исследование элемента с основными данными о квартире.

Теперь, зная значение атрибута, используем метод `find()`, чтобы найти его в переменной `$dv1`. Запишем найденные данные в переменную `$dv2`. `$dv2 = $dv1 → find('.c6e8ba5398--title--2CW78')`. Если вывести на экран переменную `$dv2` внутри тега `<html>`, то видно что данные там записаны вместе с тегами и атрибутами (Рисунок 12).

```
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 18 м², 1/2 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 20 м², 1/3 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 22 м², 1/11 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. апарт., 40 м², 8/19 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. апарт., 37 м², 12/13 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. апарт., 42 м², 11/19 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. апарт., 41 м², 8/13 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. апарт., 35 м², 2/13 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. апарт., 42 м², 7/13 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 28 м², 6/23 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 27 м², 9/23 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 30 м², 1/17 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 34 м², 1/17 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 26 м², 3/17 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 38 м², 2/3 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. апарт., 39 м², 12/17 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 35 м², 2/17 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 34 м², 2/3 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 30 м², 14/17 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 38 м², 13/23 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 39 м², 1/17 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 35 м², 13/23 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 35 м², 9/23 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 30 м², 14/17 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 35 м², 14/14 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 39 м², 2/12 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. кв., 41 м², 1/17 этаж</div>
<div class="c6e8ba5398--title--2CW78">1-комн. апарт., 45 м², 3/8 этаж</div>
```

Рисунок 12. Значение переменной `$dv2`.

Для удобного анализа данных, необходимо убрать лишнюю информацию. Поэтому с помощью функции `str_replace()` уберем всё, кроме

нужной для записи в базу данных информации. После применения этой функции, тип переменной `$dv2` поменяется с объекта на строку. Эта функция позволяет заменить определенную информацию на нужную. В данном случае заменим сначала открывающийся тег `<div>`. `$dv2 = str_replace ('<div class="сбе8ba5398--title--2CW78">', '', $dv2)`. Эта запись означает, что все что в переменной `$dv2` равно `<div class="сбе8ba5398--title--2CW78">` заменяется на пустоту, то есть по сути удаляется. Далее заменим закрывающий тег `</div>`. `$dv2 = str_replace('</div>', '~', $dv2)`. Здесь закрывающий тег мы заменяем на тильду, чтобы в дальнейшем разделить строку по этому символу и создать массив для последующей записи в базу данных. Функция `str_replace` помогла избавиться от лишних данных. И теперь в переменной `$dv2` содержится только нужная информация (Рисунок 13).



```
1-комн. кв., 18 м², 1/2 этаж~ 1-комн. кв., 22 м², 1/11 этаж~ 1-комн. апарт., 40 м², 8/19 этаж~ 1-комн. апарт., 37 м², 12/13 этаж~ 1-комн. апарт., 42 м², 11/19 этаж~ 1-комн. апарт., 41 м², 9/13 этаж~ 1-комн. апарт., 35 м², 2/13 этаж~ 1-комн. апарт., 42 м², 7/13 этаж~ 1-комн. кв., 28 м², 6/23 этаж~ 1-комн. кв., 27 м², 9/23 этаж~ 1-комн. кв., 30 м², 1/17 этаж~ 1-комн. кв., 34 м², 1/17 этаж~ 1-комн. кв., 26 м², 3/17 этаж~ 1-комн. кв., 38 м², 2/3 этаж~ 1-комн. апарт., 39 м², 12/17 этаж~ 1-комн. кв., 35 м², 2/17 этаж~ 1-комн. кв., 34 м², 2/3 этаж~ 1-комн. кв., 30 м², 14/17 этаж~ 1-комн. кв., 38 м², 13/23 этаж~ 1-комн. кв., 39 м², 1/17 этаж~ 1-комн. кв., 35 м², 13/23 этаж~ 1-комн. кв., 35 м², 9/23 этаж~ 1-комн. кв., 30 м², 14/17 этаж~ 1-комн. кв., 35 м², 14/14 этаж~ 1-комн. кв., 39 м², 2/12 этаж~ 1-комн. кв., 41 м², 1/17 этаж~ 1-комн. апарт., 45 м², 3/8 этаж~ 1-комн. кв., 40 м², 2/12 этаж~
```



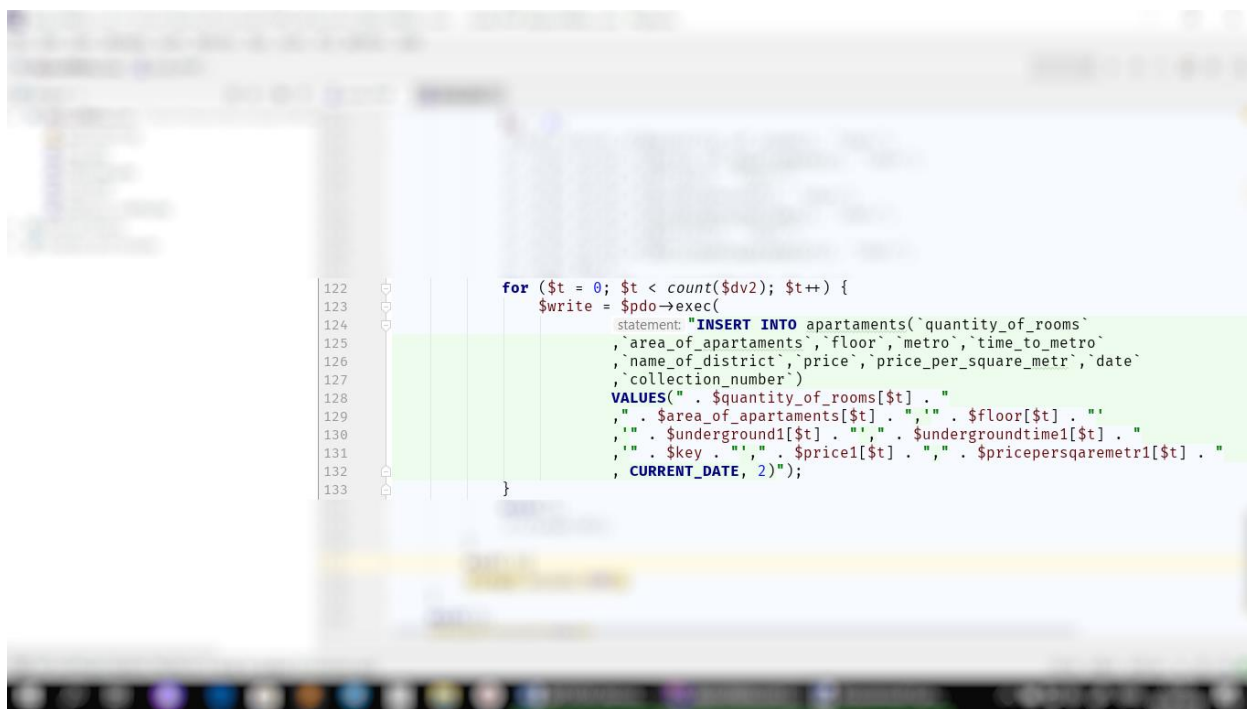
Рисунок 13. Значение переменной `$dv2`.

Теперь применим функцию `preg_split()`. `$dv2 = preg_split('~', $dv2)`. С помощью неё мы поделили строку `$dv2` по символу `~`. Теперь `$dv2` имеет тип переменной массив, каждый элемент которого имеет такой вид: 1-комн. кв., 18 м², 1/2 этаж. Также применим функцию `array_pop($dv2)`. Эта функция позволяет удалить последний элемент массива. В данном случае нужно применить эту

функцию, так как последний элемент массива \$dv2 пустой, и для последующего корректного использования массив не должен содержать пустых элементов.

Все остальные данные собираются по точно такому же алгоритму. То есть сначала с помощью окна DevTools нужно узнать значение атрибута class, далее применяем нужные методы и функции, в итоге получаем массив с данными.

Теперь, все данные нужно записать в базу данных. С помощью строчки `$pdo = new PDO('mysql:host=localhost;dbname=diplom', 'root', ' ')` выполним подключение к базе данных, в скобках указывается, вид СУБД, в этом проекте это MySQL, адрес сервера, здесь это локальный сервер, то есть localhost, а также название базы данных, логин и пароль от БД. Подключение выполнено, запись полученных данных будет происходить в цикле, который будет проходить по всем массивам и выполнять SQL запрос (Рисунок 14).



```
122     for ($t = 0; $t < count($dv2); $t++) {
123         $write = $pdo->exec(
124             statement: "INSERT INTO apartaments(`quantity_of_rooms`
125             ,`area_of_apartaments`, `floor`, `metro`, `time_to_metro`
126             ,`name_of_district`, `price`, `price_per_square_metr`, `date`
127             ,`collection_number`)
128             VALUES(" . $quantity_of_rooms[$t] . "
129             , " . $area_of_apartaments[$t] . " , " . $floor[$t] . "
130             , " . $underground1[$t] . " , " . $undergroundtime1[$t] . "
131             , " . $key . " , " . $price1[$t] . " , " . $pricepersquaremetr1[$t] . "
132             , CURRENT_DATE, 2)");
133     }
```

Рисунок 14. Цикл исполнения запроса.

Здесь (Рисунок 14) `count($dv2)` это количество объявлений на странице. в столбцы базы данных записываются значения массивов, в котором хранятся все данные по объявлению. Для перехода к следующей странице на сайте, для

последующего дальнейшего сбора данных находим с помощью окна DevTools значение атрибута class у страницы, которая открыта сейчас (Рисунок 15).

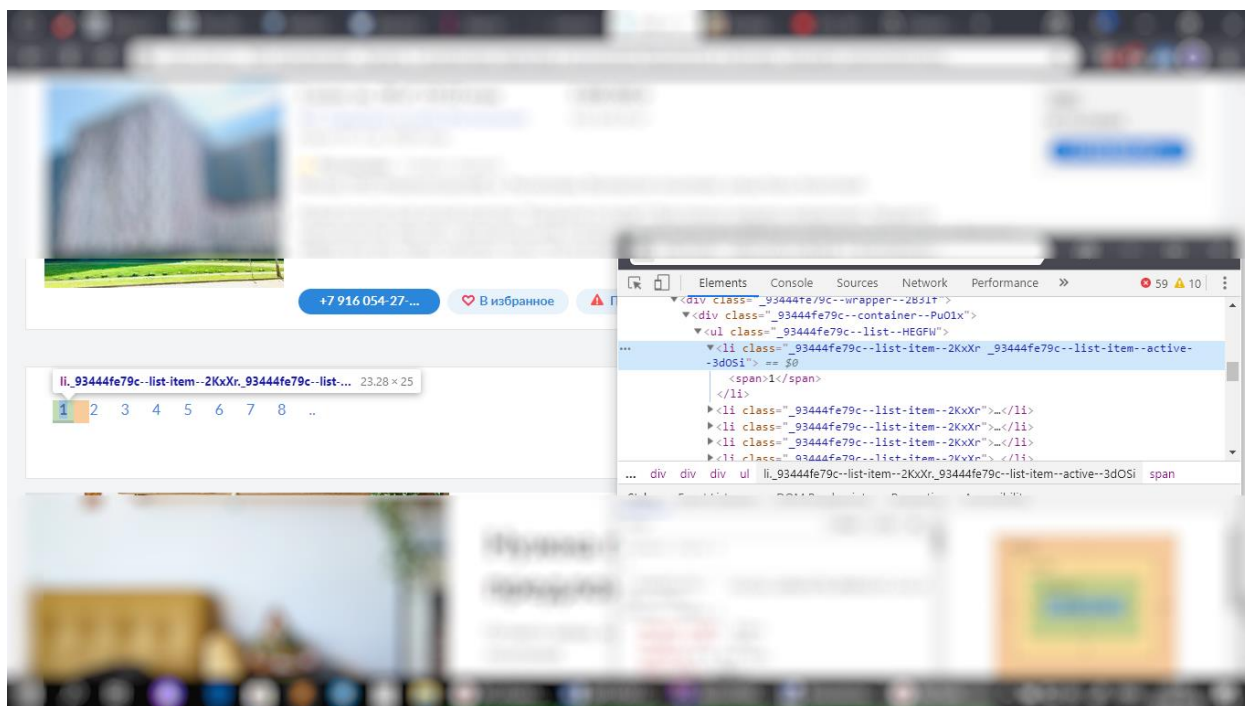


Рисунок 15. Исследование элемента.

У активной страницы значение атрибута равно `class="_93444fe79c--list-item--2KxXr_93444fe79c--list-item--active--3dOSi"`. Далее с помощью методов `jQuery` найдем следующую за этой страницу. Это делается так: `$next = $doc->find('._93444fe79c--list-item--active--3dOSi')->next()->children()->attr('href')`.

Теперь в переменной `$next` хранится ссылка, на следующую страницу. Далее мы записываем значение переменной `$next` в переменную `$url`. Теперь мы оборачиваем весь алгоритм в цикл `while`. В итоге алгоритм будет производить сбор данных со страниц, пока они не закончатся.

Чтобы запустить алгоритм, включаем локальный веб-сервис, заходим в браузер и вводим в адресную строку название папки, которая была создана в папке `domains`, здесь она названа `diplom`. Далее нажимаем `enter`, теперь алгоритм заработал и собирает данные. Все данные будут записываться в таблицу в базу данных (Рисунок 16).

Сервер: 127.0.0.1:3306 » База данных: diplom » Таблица: apartments

Обзор Структура SQL Поиск Вставить Экспорт Импорт Привилегии Операции Триггеры

+ Параметры

	id	quantity_of_rooms	area_of_apartaments	floor	metro	time_to_metro	name_of_district	price	price_per_square_metr
<input type="checkbox"/>	1	1	26	3/3 этаж	Рассказовка	20	Внуковское	2700000	152273
<input type="checkbox"/>	2	1	20	1/3 этаж	Новопеределкино	29	Внуковское	2700000	92868
<input type="checkbox"/>	3	1	23	8/15 этаж	Рассказовка	15	Внуковское	3200000	106400
<input type="checkbox"/>	4	1	22	1/9 этаж	Рассказовка	18	Внуковское	3350000	126000
<input type="checkbox"/>	5	1	40	8/19 этаж	Рассказовка	10	Внуковское	3750000	96145
<input type="checkbox"/>	6	1	37	12/13 этаж	Рассказовка	7	Внуковское	3926160	98900
<input type="checkbox"/>	7	1	32	5/23 этаж	Рассказовка	12	Внуковское	3969000	96500
<input type="checkbox"/>	8	1	42	11/19 этаж	Рассказовка	20	Внуковское	3990000	117700
<input type="checkbox"/>	9	1	41	8/13 этаж	Рассказовка	7	Внуковское	4005450	90508
<input type="checkbox"/>	10	1	42	2/13 этаж	Рассказовка	7	Внуковское	4014400	153846
<input type="checkbox"/>	11	1	35	2/13 этаж	Рассказовка	7	Внуковское	4095960	153571
<input type="checkbox"/>	12	1	45	6/8 этаж	Рассказовка	15	Внуковское	4100000	144328
<input type="checkbox"/>	13	1	27	8/23 этаж	Рассказовка	1	Внуковское	4200000	123500
<input type="checkbox"/>	14	1	28	13/23 этаж	Рассказовка	1	Внуковское	4300000	102326
<input type="checkbox"/>	15	1	30	1/17 этаж	Рассказовка	5	Внуковское	4373142	131255
<input type="checkbox"/>	16	1	36	3/23 этаж	Рассказовка	12	Внуковское	4396600	170233

Консоль

Рисунок 16. Данные в БД.

## 4.2. Построение регрессионной модели.

Для построения регрессионной модели решим задачу минимизации.

$$\sum_{i=1}^n (y_i - (\vec{a}, \vec{x}_i) - b)^2 \xrightarrow{a,b} \min$$

Поскольку функция является выпуклой всюду в  $\mathbb{R}^{n+1}$ , то она обладает единственным экстремумом. Для его нахождения возьмем производные по каждой из переменных  $a_1, a_2, \dots, a_n$  и  $b$  и приравняем к нулю. Тогда получим систему уравнений:

$$\begin{cases} \sum_{i=1}^n 2x_i(y_i - (a, x_i) - b) = 0 \\ \sum_{i=1}^n 2x_i^2(y_i - (a, x_i) - b) = 0 \\ \dots \\ \sum_{i=1}^n 2x_i^n(y_i - (a, x_i) - b) = 0 \\ \sum_{i=1}^n 2(y_i - (a, x_i) - b) = 0 \end{cases}$$

Эта система сводится к матричному уравнению:

$$\begin{pmatrix} Cov(x_1, x_1) & \dots & Cov(x_1, x_n) \\ \vdots & \ddots & \vdots \\ Cov(x_n, x_1) & \dots & Cov(x_n, x_n) \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} Cov(x_1, y) \\ \vdots \\ Cov(x_n, y) \end{pmatrix}$$

с дополнительным условием

$$b = \bar{y} - (\vec{a}, \vec{x}) = \bar{y} - a_1\bar{x}_1 - a_2\bar{x}_2 - \dots - a_n\bar{x}_n$$

Находя  $a_1, a_2, \dots, a_n$  получаем линейную регрессионную модель.

$$y \sim a_1x_1 + a_2x_2 + \dots + a_nx_n + b$$

Для определения значимости параметров сделаем центрирование и нормировку, то есть преобразование

$$x \rightarrow x' = \frac{x - \bar{x}}{sd(x)}$$

Где

$$sd(x) = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n - 1}}$$

Делаем нормировку данных и строим регрессионную модель по  $x'_i$  и находим  $a'$ . Пропорция между компонентами  $a'_i$  позволяет определить значимость параметров для модели. В качестве измерения качества модели возьмем индекс детерминированности  $R^2$ , который равен:

$$1 - \frac{var(y(x))}{var(x)}$$

Далее разберем алгоритм построения регрессионной модели на языке Python. Изначально считаем все собранные записи с базы данных, для этого нужно подключиться к БД. Далее после запроса получим все эти данные. Для большей видимости результатов, возьмем однокомнатные квартиры, стоимостью от 1000000Р до 50000000Р. Для анализа и визуализации данных воспользуемся библиотекой matplotlib. Так как в БД большинство полей – строки, нужно произвести классификацию данных. Для начала определим корреляцию признаков от цены, то есть узнаем насколько сильно взаимосвязаны признаки с ценой. Сделаем это с помощью функции corrcoef из библиотеки numpy (Рисунок 17, Рисунок 18).

```
72 print("Зависимость цены от площади")
73 print(np.corrcoef(prices,area)[0][1])
74 print("Зависимость цены от времени до метро")
75 print(np.corrcoef(prices,time)[0][1])
76 print("Зависимость цены от района")
77 print(np.corrcoef(prices,classification)[0][1])
78 print("Зависимость цены от этажа квартиры")
79 print(np.corrcoef(prices,floor0)[0][1])
80 print("Зависимость цены от этажности в доме")
81 print(np.corrcoef(prices,floor1)[0][1])
82 print("Зависимость цены от цены за квадратный метр")
83 print(np.corrcoef(prices,pricePerMeter)[0][1])
```

Рисунок 17. Вычисление коэффициента корреляции.



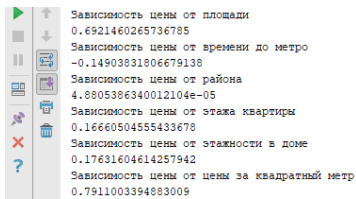


Рисунок 18. Значение коэффициента корреляции.

Видно, что цена наиболее взаимосвязана с такими признаками как площадь и цена за квадратный метр. В этих двух случаях коэффициент корреляции положительный и наиболее близок к единице. Это означает, что чем больше площадь квартиры или цена за квадратный метр, тем больше цена за квартиру. У признака время до метро отрицательный коэффициент корреляции, это означает, что чем ближе метро, тем выше цена. Построим графики зависимости цены от этих трех признаков (Рисунок 19).

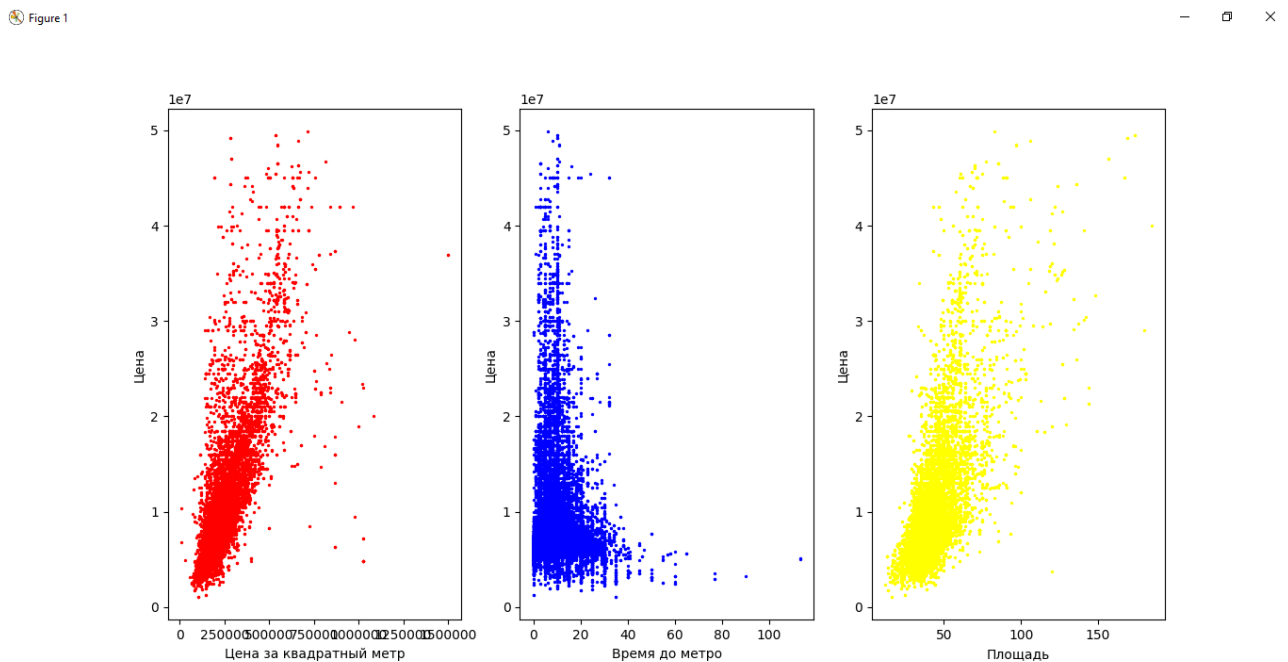


Рисунок 19. График зависимости цены от признаков.

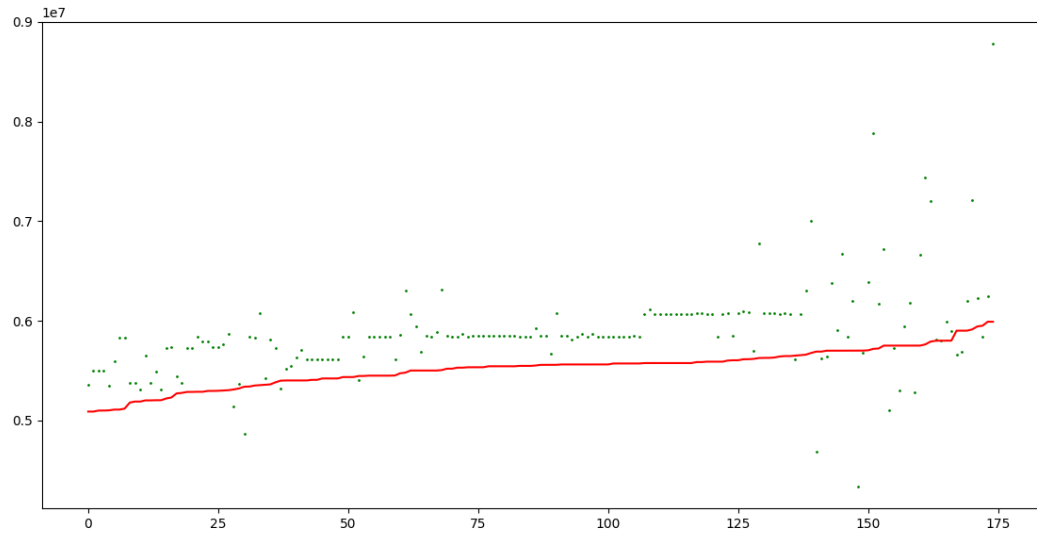
На графиках также видно эти зависимости. Далее используя эти три признака обучим модель. В качестве модели для обучения была выбрана модель `LinearRegression()`. Обучаем используя функцию `fit()`. После обучения модели, сделаем предсказание по тем же трем признакам. Для этого используем функцию `predict()`. В итоге получаем точность предсказания равную 0,82 (Рисунок 20).



Рисунок 20. Точность предсказания.

Также на графике сопоставим значения предсказанной цены и реальной, зеленые точки — это предсказанная цена, а красные это реальная (Рисунок 21).

Figure 1



x=116.522 y=8.92396e+06

Рисунок 21. Сопоставление предсказанной цены и реальной.

## 5. Заключение и выводы.

В данной работе, успешно был осуществлен сбор данных с html-страниц с помощью языка PHP и библиотеки jQuery. Все собранные данные были упорядочены и успешно записаны в базу данных. В дальнейшем производился анализ этих данных. В ходе анализа были выведены зависимости признаков от результативного признака. На основе этих признаков была обучена линейная регрессионная модель. В итоге после обучения на основе трех признаков, а именно площадь квартиры, цена за квадратный метр и расстояние до метро, получилось сделать предсказание цены основываясь на эти же признаки с точность 0,82.

## 6. Список литературы.

- 1) Ситуация на рынке жилья по регионам к середине 2018 года. [Электронный ресурс] URL: <https://kazan.cian.ru/stati-situatsija-na-rynke-zhilja-po-regionam-k-seredine-2018-goda-283917/>
- 2) Стивен Шафер. HTML, XHTML и CSS. Библия пользователя, 5-е издание = HTML, XHTML, and CSS Bible, 5th Edition. — М.: «Диалектика», 2010. — 656 с. — ISBN 978-5-8459-1676-1.
- 3) Кузнецов С. Д. Основы баз данных. — 2-е изд. — М.: Интернет-университет информационных технологий; БИНОМ. Лаборатория знаний, 2007. — 484 с. — [ISBN 978-5-94774-736-2](#).
- 4) В. Васвани. MySQL: использование и администрирование = MySQL Database Usage & Administration. — М.: «Питер», 2011. — 368 с. — [ISBN 978-5-459-00264-5](#).
- 5) Кузнецов Максим, Симдянов Игорь. PHP на примерах. — 2-е изд. перераб. и доп. — СПб.: «БХВ-Петербург», 2011. — С. 400. — [ISBN 978-5-9775-0445-4](#).
- 6) Фёдоров Д. Ю. [Основы программирования на примере языка Python](#) / Учебное пособие. — СПб.: Юрайт, 2018. — 167 с. — [ISBN 978-5-534-04479-9](#).

# 7. Приложение.

## 7.1. Приложение 1. Код алгоритма на языке PHP.

Файл index.php.

```
1. <?php
2. ini_set('max_execution_time', 30000);
3. header('Content-type: text/html; charset=utf-8');
4. include("array.php");
5. require 'phpQuery-onefile.php';
6. $pdo = new PDO('mysql:host=localhost;dbname=diplom', 'root', '');
7. for ($i = 0; $i < count($districts); $i++) {
8.     for ($j = 0; $j < count($districts[$i]); $j++) {
9.         $flag = true;
10.        $url = array_values($districts[$i])[$j];
11.        $key = array_search(array_values($districts[$i])[$j], $districts[$i]);
12.        while ($flag) {
13.            $file = file_get_contents($url);
14.            $doc = phpQuery::newDocument($file);
15.            $dv1 = $doc->find('._93444fe79c--card--yguQ');
16.            $dv2 = $dv1->find('.c6e8ba5398--title--2CW78');
17.            $dv2 = str_replace('<div class="c6e8ba5398--title--2CW78">', '', $dv2);
18.            $dv2 = str_replace('</div>', '~', $dv2);
19.            $dv2 = preg_split('#~#', $dv2);
20.            $elem = array_pop($dv2);
21.            $underground = $dv1->find('.c6e8ba5398--underground-name--3YjAi');
22.            $underground = str_replace('<div class="c6e8ba5398--underground-name--3YjAi">', '', $underground);
23.            $underground = str_replace('</div>', '~', $underground);
24.            $underground = preg_split('#~#', $underground);
25.            while (count($underground) > count($dv2)) {
26.                $elem = array_pop($underground);
27.            }
28.            $undergroundtime = $dv1->find('.c6e8ba5398--remoteness--1MdE6');
29.            $undergroundtime = str_replace('<div class="c6e8ba5398--remoteness--1MdE6">', '', $undergroundtime);
30.            $undergroundtime = str_replace('</div>', '~', $undergroundtime);
31.            $undergroundtime = preg_split('#~#', $undergroundtime);
32.            while (count($underground) > count($dv2)) {
33.                $elem = array_pop($undergroundtime);
34.            }
35.            $price = $dv1->find('.c6e8ba5398--header--1df-X');
36.            $price = str_replace('<div class="c6e8ba5398--header--1df-X">', '', $price);
37.            $price = str_replace('</div>', '~', $price);
38.            $price = preg_split('#~#', $price);
39.            $elem = array_pop($price);
40.            $pricepersquaremetr = $dv1->find('.c6e8ba5398--term--3kvtJ');
41.            $pricepersquaremetr = str_replace('<div class="c6e8ba5398--term--3kvtJ">', '', $pricepersquaremetr);
42.            $pricepersquaremetr = str_replace('</div>', '~', $pricepersquaremetr);
43.            $pricepersquaremetr = preg_split('#~#', $pricepersquaremetr);
44.            for ($z = 0; $z <= 2; $z++) {
45.                $elem = array_shift($pricepersquaremetr);
46.                $elem = array_pop($pricepersquaremetr);
47.            }
48.            $elem = array_pop($pricepersquaremetr);
49.            $next = $doc->find('._93444fe79c--list-item--active--3d0Si')->next()->children()->attr('href');
50.            $next = str_replace('https://www.!!!.ru', '', $next);
51.            if (!empty($next)) {
```

```

52.         $url = 'https://www.!!!.ru'.$next;
53.     } else {
54.         $flag = false;
55.     }
56.     $g = 0;
57.     $quantity_of_rooms = array();
58.     $area_of_apartaments = array();
59.     $floor = array();
60.     $underground1 = array();
61.     $undergroundtime1 = array();
62.     $price1 = array();
63.     $pricepersquaremetr1 = array();
64.     foreach ($dv2 as $value) {
65.         $value = preg_split('#,#', $value);
66.         foreach ($value as $value2) {
67.             if ($g == 0) {
68.                 array_push($quantity_of_rooms, intval($value2));
69.                 $g = 1;
70.             } else {
71.                 if ($g == 1) {
72.                     array_push($area_of_apartaments, intval($value2));
73.                     $g = 2;
74.                 } else {
75.                     array_push($floor, $value2);
76.                     $g = 0;
77.                 }
78.             }
79.         }
80.     }
81.     $g = 0;
82.     foreach ($underground as $value) {
83.         array_push($underground1, $value);
84.     }
85.     $g = 0;
86.     foreach ($undergroundtime as $value) {
87.         array_push($undergroundtime1, intval($value));
88.     }
89.     $g = 0;
90.     foreach ($price as $value) {
91.         $value = str_replace(" ", "", $value);
92.         array_push($price1, intval($value));
93.     }
94.     $g = 0;
95.     foreach ($pricepersquaremetr as $value) {
96.         $value = str_replace(" ", "", $value);
97.         array_push($pricepersquaremetr1, intval($value));
98.     }
99.     $g = 0;
100.     for ($t = 0; $t < count($dv2); $t++) {
101.         $write = $pdo-
>exec("INSERT INTO apartaments(`quantity_of_rooms`,`area_of_apartaments`,`floor`,`metro`,`
time_to_metro`,`name_of_district`,`price`,`price_per_square_metr`,`date`,`collection_numbe
r`) VALUES(" . $quantity_of_rooms[$t] . "," . $area_of_apartaments[$t] . "," . $floor[$t]
. "," . $underground1[$t] . "," . $undergroundtime1[$t] . "," . $key . "," . $price1
[$t] . "," . $pricepersquaremetr1[$t] . ", CURRENT_DATE, 2)");
102.     }
103.     }
104.     }
105.     }

```

## 7.2. Приложение 2. Часть кода базы данных на языке SQL.

Файл apartments.sql.

```
1. -- phpMyAdmin SQL Dump
2. -- version 4.8.3
3. -- https://www.phpmyadmin.net/
4. --
5. -- Хост: 127.0.0.1:3306
6. -- Время создания: Май 22 2019 г., 16:13
7. -- Версия сервера: 8.0.12
8. -- Версия PHP: 7.2.10
9.
10. SET SQL_MODE = "NO_AUTO_VALUE_ON_ZERO";
11. SET AUTOCOMMIT = 0;
12. START TRANSACTION;
13. SET time_zone = "+00:00";
14.
15.
16. /*!40101 SET @OLD_CHARACTER_SET_CLIENT=@@CHARACTER_SET_CLIENT */;
17. /*!40101 SET @OLD_CHARACTER_SET_RESULTS=@@CHARACTER_SET_RESULTS */;
18. /*!40101 SET @OLD_COLLATION_CONNECTION=@@COLLATION_CONNECTION */;
19. /*!40101 SET NAMES utf8mb4 */;
20.
21. --
22. -- База данных: `diplom`
23. --
24.
25. -- -----
26.
27. --
28. -- Структура таблицы `apartaments`
29. --
30.
31. CREATE TABLE `apartaments` (
32.   `id` int(11) NOT NULL,
33.   `quantity_of_rooms` int(11) NOT NULL,
34.   `area_of_apartaments` int(11) NOT NULL,
35.   `floor` varchar(200) NOT NULL,
36.   `metro` varchar(200) NOT NULL,
37.   `time_to_metro` int(11) NOT NULL,
38.   `name_of_district` varchar(200) NOT NULL,
39.   `price` int(11) NOT NULL,
40.   `price_per_square_metr` int(11) NOT NULL,
41.   `date` date NOT NULL,
42.   `collection_number` int(11) NOT NULL
43. ) ENGINE=InnoDB DEFAULT CHARSET=utf8;
44.
45. --
46. -- Дамп данных таблицы `apartaments`
47. --
48.
49. INSERT INTO `apartaments` (`id`, `quantity_of_rooms`, `area_of_apartaments`, `floor`, `metro`, `time_to_metro`, `name_of_district`, `price`, `price_per_square_metr`, `date`, `collection_number`) VALUES
50. (1, 1, 26, ' 3/3 этаж', '\nРассказовка', 20, 'Внуковское', 2700000, 152273, '2019-05-18', 1),
51. (2, 1, 20, ' 1/3 этаж', '\nНовопеределкино', 29, 'Внуковское', 2700000, 92868, '2019-05-18', 1),
52. (3, 1, 23, ' 8/15 этаж', '\nРассказовка', 15, 'Внуковское', 3200000, 106400, '2019-05-18', 1),
```

53. (4, 1, 22, ' 1/9 этаж', '\nРассказовка', 18, 'Внуковское', 3350000, 126000, '2019-05-18', 1),
54. (5, 1, 40, ' 8/19 этаж', '\nРассказовка', 10, 'Внуковское', 3750000, 96145, '2019-05-18', 1),
55. (6, 1, 37, ' 12/13 этаж', '\nРассказовка', 7, 'Внуковское', 3926160, 98900, '2019-05-18', 1),
56. (7, 1, 32, ' 5/23 этаж', '\nРассказовка', 12, 'Внуковское', 3969000, 96500, '2019-05-18', 1),
57. (8, 1, 42, ' 11/19 этаж', '\nРассказовка', 20, 'Внуковское', 3990000, 117700, '2019-05-18', 1),
58. (9, 1, 41, ' 8/13 этаж', '\nРассказовка', 7, 'Внуковское', 4005450, 90508, '2019-05-18', 1),
59. (10, 1, 42, ' 2/13 этаж', '\nРассказовка', 7, 'Внуковское', 4014400, 153846, '2019-05-18', 1),
60. (11, 1, 35, ' 2/13 этаж', '\nРассказовка', 7, 'Внуковское', 4095960, 153571, '2019-05-18', 1),
61. (12, 1, 45, ' 6/8 этаж', '\nРассказовка', 15, 'Внуковское', 4100000, 144328, '2019-05-18', 1),
62. (13, 1, 27, ' 8/23 этаж', '\nРассказовка', 1, 'Внуковское', 4200000, 123500, '2019-05-18', 1),
63. (14, 1, 28, ' 13/23 этаж', '\nРассказовка', 1, 'Внуковское', 4300000, 102326, '2019-05-18', 1),
64. (15, 1, 30, ' 1/17 этаж', '\nРассказовка', 5, 'Внуковское', 4373142, 131255, '2019-05-18', 1),
65. (16, 1, 36, ' 3/23 этаж', '\nРассказовка', 12, 'Внуковское', 4396600, 170233, '2019-05-18', 1),
66. (17, 1, 43, ' 19/19 этаж', '\nРассказовка', 3, 'Внуковское', 4400000, 118421, '2019-05-18', 1),
67. (18, 1, 34, ' 1/17 этаж', '\nРассказовка', 5, 'Внуковское', 4423291, 135294, '2019-05-18', 1),
68. (19, 1, 26, ' 3/17 этаж', '\nРассказовка', 5, 'Внуковское', 4426056, 117207, '2019-05-18', 1),
69. (20, 1, 38, ' 2/3 этаж', '\nНовоперedelкино', 21, 'Внуковское', 4500000, 138701, '2019-05-18', 1),
70. (21, 1, 34, ' 2/3 этаж', '\nРассказовка', 15, 'Внуковское', 4600000, 129500, '2019-05-18', 1),
71. (22, 1, 40, ' 12/16 этаж', '\nРассказовка', 23, 'Внуковское', 4700000, 132713, '2019-05-18', 1),
72. (23, 1, 30, ' 14/17 этаж', '\nРассказовка', 5, 'Внуковское', 5100000, 157576, '2019-05-18', 1),
73. (24, 1, 35, ' 13/23 этаж', '\nРассказовка', 2, 'Внуковское', 5100000, 134127, '2019-05-18', 1),
74. (25, 1, 35, ' 9/23 этаж', '\nРассказовка', 10, 'Внуковское', 5195000, 138588, '2019-05-18', 1),
75. (26, 1, 33, ' 13/17 этаж', '\nРассказовка', 5, 'Внуковское', 5200000, 129051, '2019-05-18', 1),
76. (27, 1, 39, ' 2/12 этаж', '\nРассказовка', 5, 'Внуковское', 5257792, 132500, '2019-05-18', 1),
77. (28, 1, 38, ' 4/23 этаж', '\nРассказовка', 5, 'Внуковское', 5266344, 134127, '2019-05-18', 1),
78. (29, 1, 41, ' 1/17 этаж', '\nРассказовка', 26, 'Внуковское', 5278202, 147383, '2019-05-18', 1),
79. (30, 1, 40, ' 19/19 этаж', '\nРассказовка', 5, 'Внуковское', 5300000, 127381, '2019-05-18', 1),
80. (31, 1, 40, ' 2/12 этаж', '\nРассказовка', 5, 'Внуковское', 5324855, 150000, '2019-05-18', 1),
81. (32, 1, 36, ' 18/23 этаж', '\nРассказовка', 17, 'Внуковское', 5350000, 134127, '2019-05-18', 1),
82. (33, 1, 42, ' 14/15 этаж', '\nРассказовка', 5, 'Внуковское', 5350000, 126293, '2019-05-18', 1),
83. (34, 1, 36, ' 22/24 этаж', '\nРассказовка', 5, 'Внуковское', 5400000, 146793, '2019-05-18', 1),
84. (35, 1, 41, ' 2/12 этаж', '\nРассказовка', 5, 'Внуковское', 5432157, 147521, '2019-05-18', 1),
85. (36, 1, 43, ' 1/12 этаж', '\nРассказовка', 5, 'Внуковское', 5455860, 165165, '2019-05-18', 1),

```

86. (37, 1, 37, ' 6/17 этаж', '\nРассказовка', 7, 'Внуковское', 5460700, 141388, '2019-05-18', 1),
87. (38, 1, 37, ' 8/17 этаж', '\nРассказовка', 20, 'Внуковское', 5487787, 130952, '2019-05-18', 1),
88. (39, 1, 33, ' 14/14 этаж', '\nРассказовка', 17, 'Внуковское', 5500000, 141026, '2019-05-18', 1),
89. (40, 1, 39, ' 17/19 этаж', '\nРассказовка', 5, 'Внуковское', 5500000, 137744, '2019-05-18', 1),
90. (41, 1, 42, ' 10/16 этаж', '\nРассказовка', 5, 'Внуковское', 5500000, 124570, '2019-05-18', 1),
91. (42, 1, 39, ' 17/19 этаж', '\nРассказовка', 5, 'Внуковское', 5500000, 144492, '2019-05-18', 1),
92. (43, 1, 40, ' 2/17 этаж', '\nРассказовка', 5, 'Внуковское', 5551099, 156239, '2019-05-18', 1),
93. (44, 1, 45, ' 1/17 этаж', '\nРассказовка', 5, 'Внуковское', 5568299, 144172, '2019-05-18', 1),
94. (45, 1, 39, ' 11/17 этаж', '\nРассказовка', 5, 'Внуковское', 5684268, 136081, '2019-05-18', 1),
95. (46, 1, 49, ' 12/23 этаж', '\nРассказовка', 5, 'Внуковское', 5695000, 143510, '2019-05-18', 1),
96. (47, 1, 42, ' 21/25 этаж', '\nРассказовка', 24, 'Внуковское', 5700000, 148718, '2019-05-18', 1),
97. (48, 1, 42, ' 3/17 этаж', '\nРассказовка', 5, 'Внуковское', 5715404, 152632, '2019-05-18', 1),
98. (49, 1, 40, ' 1/17 этаж', '\nРассказовка', 5, 'Внуковское', 5726053, 153034, '2019-05-18', 1),
99. (50, 1, 39, ' 6/17 этаж', '\nРассказовка', 27, 'Внуковское', 5800000, 161111, '2019-05-18', 1),
100. (51, 1, 38, ' 2/17 этаж', '\nРассказовка', 27, 'Внуковское', 5800000, 133028, '2019-05-18', 1),
101. (52, 1, 38, ' 9/9 этаж', '\nРассказовка', 24, 'Внуковское', 5800000, 159341, '2019-05-18', 1),
102. (53, 1, 36, ' 3/20 этаж', '\nРассказовка', 8, 'Внуковское', 5800000, 136081, '2019-05-18', 1),
103. (54, 1, 44, ' 7/24 этаж', '\nРассказовка', 28, 'Внуковское', 5800000, 155263, '2019-05-18', 1),
104. --
105. -- Индексы сохранённых таблиц
106. --
107.
108. --
109. -- Индексы таблицы `apartaments`
110. --
111. ALTER TABLE `apartaments`
112.     ADD PRIMARY KEY (`id`);
113.
114. --
115. -- AUTO_INCREMENT для сохранённых таблиц
116. --
117.
118. --
119. -- AUTO_INCREMENT для таблицы `apartaments`
120. --
121. ALTER TABLE `apartaments`
122.     MODIFY `id` int(11) NOT NULL AUTO_INCREMENT, AUTO_INCREMENT=23522;
123. COMMIT;
124.
125. /*!40101 SET CHARACTER_SET_CLIENT=@OLD_CHARACTER_SET_CLIENT */;
126. /*!40101 SET CHARACTER_SET_RESULTS=@OLD_CHARACTER_SET_RESULTS */;
127. /*!40101 SET COLLATION_CONNECTION=@OLD_COLLATION_CONNECTION */;

```



## 7.3. Приложение 3. Код алгоритма на языке Python.

```
1. import numpy as np
2. from sklearn.linear_model import LinearRegression
3. import pymysql
4. connection = pymysql.connect(
5.     host='localhost',
6.     user='root',
7.     password='',
8.     db='diplom',
9.     charset='utf8mb4',
10.    cursorclass=pymysql.cursors.DictCursor
11.)
12. class Apartament():
13.     def __init__(self,classification = None,pricePerMeter = None,time = None,floor0 = None,
14.                 floor1 = None,price = None,area = None):
15.         self.price = price
16.         self.floor0 = floor0
17.         self.floor1 = floor1
18.         self.classification = classification
19.         self.time = time
20.         self.pricePerMeter = pricePerMeter
21.         self.area = area
22.
23. apartaments = []
24. query = connection.cursor()
25. query.execute('SELECT * FROM apartaments WHERE price < 50000000 AND price > 1000000 AND id
26. < 23000')
27. result = query.fetchall()
28. x = []
29. y = []
30. uniqueNames = {}
31. for res in result:
32.     if res['name_of_district'] not in uniqueNames.keys():
33.         uniqueNames[res['name_of_district']] = len(uniqueNames)
34.     def distToClass(dist):
35.         return -1 if dist not in uniqueNames.keys() else uniqueNames[dist]
36. # print(uniqueNames)
37. for res in result[:200]:
38.     floor = res['floor'].replace(' этаж', '').split('/')
39.     appendToX = [res['price_per_square_metr'],
40.                 res['time_to_metro'],
41.                 res['area_of_apartaments']
42.                 ]
43.     appendToY = res['price']
44.     x.append(appendToX)
45.     y.append(appendToY)
46.     apartaments.append(Apartament(classification=distToClass(res['name_of_district']),
47.                                   pricePerMeter=res['price_per_square_metr'],
48.                                   time=res['time_to_metro'],area=res['area_of_apartaments'],
49.                                   floor0=int(floor[0]),floor1=(int)(floor[1]),price=res['price']))
50.
51. reg = LinearRegression()
52. print(len(y))
53. print(len(x))
54.
55. reg.fit(x, y)
56.
57.
58. test = result[-175:]
59. print(len(test))
60. print(reg.score(x,y))
```

```

61. import matplotlib.pyplot as plt
62. apartments.sort(key= lambda apartment: apartment.price)
63. prices = [apartment.price for apartment in apartments]
64. classification = [apartment.classification for apartment in apartments]
65. pricePerMeter = [apartment.pricePerMeter for apartment in apartments]
66. time = [apartment.time for apartment in apartments]
67. area = [apartment.area for apartment in apartments]
68. floor0 = [apartment.floor0 for apartment in apartments]
69. floor1 = [apartment.floor1 for apartment in apartments]
70. print("Зависимость цены от площади")
71. print(np.corrcoef(prices,area)[0][1])
72. print("Зависимость цены от времени до метро")
73. print(np.corrcoef(prices,time)[0][1])
74. print("Зависимость цены от района")
75. print(np.corrcoef(prices,classification)[0][1])
76. print("Зависимость цены от этажа квартиры")
77. print(np.corrcoef(prices,floor0)[0][1])
78. print("Зависимость цены от этажности в доме")
79. print(np.corrcoef(prices,floor1)[0][1])
80. print("Зависимость цены от цены за квадратный метр")
81. print(np.corrcoef(prices,pricePerMeter)[0][1])
82. fig, f = plt.subplots(nrows=1,ncols=3,figsize=(40,15))
83. f[0].scatter(pricePerMeter,prices,color='red',s=2)
84. f[0].set_xlabel('Цена за квадратный метр')
85. f[0].set_ylabel('Цена')
86. f[1].scatter(time,prices,color='blue',s=2)
87. f[1].set_xlabel('Время до метро')
88. f[1].set_ylabel('Цена')
89. f[2].scatter(area,prices,color='yellow',s=2)
90. f[2].set_xlabel('Площадь')
91. f[2].set_ylabel('Цена')
92. plt.show()
93.
94. testList = [Apartament(classification=distToClass(res['name_of_district']),
95.                        pricePerMeter=res['price_per_square_metr'],
96.                        time=res['time_to_metro'],area=res['area_of_apartaments'],
97.                        floor0=int(floor[0]),floor1=(int)(floor[1]),price=res['price']) for res in
n test]
98. testList.sort(key = lambda apartment: apartment.price)
99. x_test = []
100.     y_test = []
101.     y1 = []
102.     x1 = []
103.     for apartment in testList[:]:
104.         x_test = [np.array([
105.                     apartment.pricePerMeter,
106.                     apartment.time,
107.                     apartment.area
108.                 ])]
109.         y_test = reg.predict(x_test)
110.         if y_test > 0:
111.             y1.append(y_test)
112.             x1.append(apartment.price)
113.             print(y_test,apartment.price)
114.         print(np.nanmean(prices))
115.         plt.scatter([i for i in range(len(y1))],x1, color = "red", s = 1)
116.         plt.scatter([i for i in range(len(y1))],y1, color = "green", s = 1)
117.         plt.show()
118.         allTrue = len(x1)
119.         x1 = np.array(x1)
120.         y1 = np.array(y1)
121.         trueCount = 0
122.         for i in range(len(y1)):
123.             if abs(y1[i]-x1[i]) < 500000:
124.                 trueCount += 1
125.         print('Accuracy:', np.round(trueCount/allTrue,2))

```