### 1. 项目介绍

# 项目背景与目的

随着互联网技术的发展和数据的普及,数据驱动决策在各行各业中变得越来越重要。对于房地产市场而言,获取并分析房源信息对买房者、卖房者以及房地产投资者都有重要的参考价值。本项目旨在通过技术手段实现对厦门二手房信息的自动化采集、分析与展示,帮助用户更好地了解市场动态,做出更加明智的决策。

#### 项目概述

本项目设计并实现了一个综合系统,涵盖了数据采集、数据处理与分析、以及 Web 应用展示三个主要部分。项目使用 Scrapy 爬取厦门二手房网站

(https://xm.lianjia.com/ershoufang/)的房源信息,并将数据存储到 SQLite 数据库中。随后,使用 Pandas 进行数据清洗和整理,利用 Matplotlib 生成各种数据可视化图表。最终,通过 Flask 框架开发一个 Web 应用,展示采集到的房源信息和市场分析结果,用户可以通过该应用进行数据搜索、过滤与查看分析结果。

### 项目目标

- 1. **实现数据采集自动化**:使用 Scrapy 爬取厦门二手房网站的房源数据,包括房源名称、地址、价格、面积、房型、发布日期等信息,并存储到 SQLite数据库中。
- 2. **数据处理与分析**:使用 Pandas 进行数据清洗和整理,对房源数据进行统计分析,如房价分布、面积分布、各区域房源数量等,生成相应的可视化图表。
- 3. **Web 应用展示**:使用 Flask 开发一个用户友好的 Web 应用,提供房源信息展示和市场分析结果查看功能,用户可以通过搜索和过滤功能查找特定房源信息,并查看数据分析结果的可视化图表。
- 4. **提升用户体验**:通过美观的界面设计和便捷的功能操作,提升用户体验,帮助用户更加直观和高效地获取所需信息。

## 2. 技术方案

项目架构概述 项目主要分为三个模块:数据采集模块、数据处理与分析模块、Web 应用展示模块。这三个模块相互独立又相互配合,形成一个完整的数据采集、分析与展示系统。

#### 使用的技术栈

- Scrapy
  - 用于爬取厦门二手房(https://xm.lianjia.com/ershoufang/)网站的 房源数据。

- Scrapy 是一个用于提取网站数据的强大且高效的 Python 框架,支持 多种数据解析方式和数据存储方式。

#### Pandas

- 用于数据的加载、清洗和处理。
- Pandas 是一个数据处理和分析的 Python 库,提供了高效的数据结构和数据分析工具。

### Matplotlib

- 用于数据可视化,生成图表如房价分布图、面积分布图和各区域房源数量图等。
- Matplotlib 是一个 Python 2D 绘图库,能够生成高质量的图表,支持 多种输出格式。

#### Flask

- 用于开发 Web 应用,展示爬取到的房源信息和数据分析结果。
- Flask 是一个轻量级的 Python Web 框架,具有灵活性高、易于扩展的特点,非常适合快速开发 Web 应用。

#### SQLite

- 用于存储爬取到的房源数据。
- SQLite 是一个轻量级的嵌入式关系型数据库,具有自包含、无服务器、零配置等特点,适合中小规模的数据存储。

### 系统模块划分

#### 1. 数据采集模块

- **目标网站分析**:分析厦门二手房网站的网页结构,确定需要爬取的信息,包括房源名称、地址、价格、面积、房型、发布日期等。
- 爬虫设计与实现:使用 Scrapy 框架编写爬虫,配置爬虫的各项参数,如爬取的页面范围、爬取频率、数据解析方式等。
- **数据存储**: 将爬取到的数据存储到 SQLite 数据库中,设计适合的数据表结构。

# 2. 数据处理与分析模块

- **数据加载与清洗**: 使用 Pandas 加载 SQLite 数据库中的数据,对数据进行清洗和转换,包括数据类型转换、缺失值处理、冗余数据清理等。
- **数据分析:** 对房源数据进行统计分析,如房价分布、面积分布、各区域房源数量等,计算相关统计指标。

- **数据可视化**:使用 Matplotlib 生成数据可视化图表,包括但不限于房价分布图、面积分布图、各区域房源数量图等,为数据展示提供直观的图形支持。

#### 3. **Web** 应用展示模块

- **Web 应用架构设计**:基于 Flask 框架设计 Web 应用的架构,规划项目结构,配置必要的依赖和参数。
- **路由设计与实现**: 定义 Web 应用的各个路由,编写对应的视图函数,处理房源信息和数据分析结果的展示和搜索功能。
- **前端页面设计**:使用 HTML、CSS 和 Bootstrap 设计前端页面布局, 创建用户友好的界面,提供数据展示和交互功能。
- **数据展示与交互**:在前端页面中嵌入数据可视化图表,提供房源信息的搜索和过滤功能,让用户能够方便地查看和分析房源数据。

### 3. 实现步骤

# 3.1 数据采集

# 3.1.1 目标网站选择及分析

选择厦门链家二手房(https://xm.lianjia.com/ershoufang/)作为数据源,分析网页结构,确定需要爬取的信息字段,包括:

- 房源名称(标题)
- 地址
- 价格
- 面积
- 房型
- 发布日期

# 3.1.2 Scrapy 爬虫的设计与实现

使用 Scrapy 框架编写爬虫,爬取目标网站的房源信息。

# 爬虫配置

- 1. 创建 Scrapy 项目
- 2. 在 settings.py 中配置 User-Agent 等必要设置
- 3. 设置下载延迟和并发请求数,避免被封禁

### 数据解析与提取

- 1. 在 spiders 目录下创建爬虫文件,例如 xm lianjia spider.py
- 2. 编写解析函数,提取房源信息字段

```
import scrapy
from real_estate.items import RealEstateItem
class XiamenLianjiaSpider(scrapy.Spider):
    name = 'xm lianjia'
    allowed domains = ['xm.lianjia.com']
    start_urls = ['https://xm.lianjia.com/ershoufang/']
   def parse(self, response):
        for listing in response.css('div.info'):
            item = RealEstateItem()
            item['name'] = listing.css('div.title a::text').get()
            item['address'] = listing.css('div.address a::text').get()
            item['price'] = listing.css('div.totalPrice span::text').ge
t()
            item['area'] = listing.css('div.houseInfo::text').re first
(r'\d+')
            item['house type'] = listing.css('div.houseInfo::text').re
first(r'\d 室\d 厅')
            item['publish_date'] = listing.css('div.followInfo::text').
re first(r'\d+-\d+-\d+')
           yield item
        next_page = response.css('div.house-lst-page-box a:last-child::
attr(href)').get()
        if next_page is not None:
           yield response.follow(next page, self.parse)
数据存储
      配置 Item Pipeline,将数据存储到 SQLite 数据库中
  1.
      在 pipelines.py 中实现数据存储逻辑
import sqlite3
class RealEstatePipeline:
    def open_spider(self, spider):
        self.conn = sqlite3.connect('real estate.db')
        self.cursor = self.conn.cursor()
        self.cursor.execute('''
            CREATE TABLE IF NOT EXISTS listings (
                id INTEGER PRIMARY KEY AUTOINCREMENT,
               name TEXT,
               address TEXT,
                price REAL,
                area REAL,
               house type TEXT,
               publish_date TEXT
        ''')
```

```
def close_spider(self, spider):
       self.conn.commit()
       self.conn.close()
   def process_item(self, item, spider):
       self.cursor.execute('
           INSERT INTO listings (name, address, price, area, house_typ
e, publish_date) VALUES (?, ?, ?, ?, ?)
       ''', (item['name'], item['address'], item['price'], item['area
'], item['house_type'], item['publish_date']))
       return item
3.2 数据处理与分析
3.2.1 数据加载与清洗
使用 Pandas 加载 SOLite 数据库中的数据,进行必要的清洗和转换。
import sqlite3
import pandas as pd
DATABASE_PATH = 'real_estate.db'
def load_data():
   conn = sqlite3.connect(DATABASE PATH)
   df = pd.read_sql_query("SELECT * FROM listings", conn)
   conn.close()
   return df
def clean data(df):
   df['price'] = df['price'].astype(float)
   df['area'] = df['area'].astype(float)
   df['publish_date'] = pd.to_datetime(df['publish_date'])
   df.dropna(inplace=True)
   return df
df = load data()
df = clean_data(df)
3.2.2 数据分析
对房源数据进行统计分析,例如房价分布、面积分布、各区域房源数量等。
def analyze_data(df):
   price_stats = df['price'].describe()
   area_stats = df['area'].describe()
   region counts = df['address'].value counts()
   return price_stats, area_stats, region_counts
```

```
price_stats, area_stats, region_counts = analyze_data(df)
print("Price Statistics:\n", price_stats)
print("Area Statistics:\n", area_stats)
print("Region Counts:\n", region_counts)
```

#### 3.2.3 数据可视化

使用 Matplotlib 生成数据可视化图表,例如房价分布图、面积分布图、各区域房源数量图等。

```
import matplotlib.pyplot as plt

def plot_histogram(data, title, xlabel, ylabel, output_path):
    plt.figure(figsize=(10, 6))
    plt.hist(data, bins=30, color='blue', alpha=0.7)
    plt.title(title)
    plt.xlabel(xlabel)
    plt.ylabel(ylabel)
    plt.savefig(output_path)
    plt.close()
```

```
def plot_bar_chart(data, title, xlabel, ylabel, output_path):
    plt.figure(figsize=(10, 6))
    data.plot(kind='bar', color='purple', alpha=0.7)
    plt.title(title)
    plt.xlabel(xlabel)
    plt.ylabel(ylabel)
    plt.savefig(output_path)
    plt.close()
```

```
plot_histogram(df['price'], 'Price Distribution', 'Price (in ten thousa
nds)', 'Frequency', 'price_distribution.png')
plot_histogram(df['area'], 'Area Distribution', 'Area (in square meter
s)', 'Frequency', 'area_distribution.png')
plot_bar_chart(region_counts, 'Number of Listings per Region', 'Region
', 'Number of Listings', 'region_counts.png')
```

#### 3.3 Web 应用开发

#### 3.3.1 Flask 应用的设计与实现

项目结构 创建 Flask 项目目录结构,包括模板和静态文件目录。

```
data_process_analysis/
    ├─ data_process_analysis.py
   real estate/
      - spiders/
      - items.py
      middlewares.py
      pipelines.py
      - settings.py
   run.py

    requirements.txt

   scrapy.cfg
路由设计 定义 Flask 应用的路由,处理房源信息和市场分析结果的展示和搜索功
能。
from flask import Flask, render_template, request
import sqlite3
import pandas as pd
from data process analysis import data process analysis
app = Flask(__name__)
@app.route('/')
def index():
   filters = {
        'name': request.args.get('name'),
        'address': request.args.get('address'),
        'min_price': request.args.get('min_price'),
        'max_price': request.args.get('max_price'),
        'min area': request.args.get('min area'),
        'max_area': request.args.get('max_area'),
        'house_type': request.args.get('house_type')
    }
    filters = {k: v for k, v in filters.items() if v}
    listings = get listings(filters)
   house_types = get_unique_house_types()
    return render_template('index.html', listings=listings, filters=fil
ters, house_types=house_types)
def get_listings(filters):
    conn = sqlite3.connect('real_estate.db')
    query = "SELECT * FROM listings"
    params = []
    conditions = []
   if 'name' in filters:
        conditions.append("name LIKE ?")
        params.append(f"%{filters['name']}%")
```

```
if 'address' in filters:
        conditions.append("address LIKE ?")
        params.append(f"%{filters['address']}%")
    if 'min_price' in filters:
        conditions.append("price >= ?")
       params.append(filters['min_price'])
    if 'max_price' in filters:
        conditions.append("price <= ?")</pre>
        params.append(filters['max_price'])
    if 'min_area' in filters:
        conditions.append("area >= ?")
        params.append(filters['min_area'])
    if 'max_area' in filters:
        conditions.append("area <= ?")</pre>
       params.append(filters['max_area'])
    if 'house_type' in filters:
        conditions.append("house type = ?")
        params.append(filters['house type'])
    if conditions:
        query += " WHERE " + " AND ".join(conditions)
    df = pd.read_sql_query(query, conn, params=params)
    conn.close()
    return df.to_dict(orient='records')
def get_unique_house_types():
    conn = sqlite3.connect('real estate.db')
    query = "SELECT DISTINCT house_type FROM listings"
    df = pd.read_sql_query(query, conn)
    conn.close()
    return df['house_type'].tolist()
@app.route('/analysis')
def analysis():
    data process analysis.generate plots()
    return render template('analysis.html')
if __name__ == '__main__':
    app.run(debug=True)
前端页面设计设计 HTML 模板,使用 Bootstrap 进行页面布局,包含搜索表单、
数据展示表格和分析图表展示。
<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <title>Real Estate Listings</title>
```

```
<link rel="stylesheet" href="https://maxcdn.bootstrapcdn.com/bootst</pre>
rap/4.5.2/css/bootstrap.min.css">
</head>
<body>
    <div class="container">
        <h1>Real Estate Listings</h1>
        <form method="get" action="/" class="mb-4">
            <div class="row">
                <div class="col-md-3 mb-2">
                    <input type="text" class="form-control" name="name"</pre>
 placeholder="Name" value="{{ filters.get('name', '') }}">
                </div>
                <div class="col-md-3 mb-2">
                     <input type="text" class="form-control" name="addre</pre>
ss" placeholder="Address" value="{{ filters.get('address', '') }}">
                </div>
                <div class="col-md-3 mb-2">
                    <input type="number" class="form-control" name="min</pre>
_price" placeholder="Min Price" value="{{ filters.get('min_price', '')
}}">
                </div>
                <div class="col-md-3 mb-2">
                     <input type="number" class="form-control" name="max</pre>
_price" placeholder="Max Price" value="{{ filters.get('max_price', '')
}}">
                </div>
            </div>
            <div class="row">
                <div class="col-md-3 mb-2">
                     <input type="number" class="form-control" name="min</pre>
_area" placeholder="Min Area" value="{{ filters.get('min_area', '') }}
                </div>
                <div class="col-md-3 mb-2">
                     <input type="number" class="form-control" name="max</pre>
_area" placeholder="Max Area" value="{{ filters.get('max_area', '') }}
                </div>
                <div class="col-md-3 mb-2">
                    <select class="form-control" name="house_type">
                         <option value="">All House Types</option>
                         {% for house_type in house_types %}
                         <option value="{{ house_type }}" {% if filters.</pre>
get('house_type') == house_type %} selected {% endif %}>{{ house_type}
 }}</option>
                         {% endfor %}
                     </select>
                </div>
                <div class="col-md-3 mb-2">
                     <button type="submit" class="btn btn-primary">Searc
```

```
h</button>
            </div>
            <div class="col-md-3 mb-2">
               <a href="/" class="btn btn-secondary">Reset</a>
            </div>
         </div>
      </form>
      <a href="/analysis" class="btn btn-info mb-4">View Analysis</a>
      <thead>
            Name
               Address
               Price
               Area
               House Type
               Publish Date
            </thead>
         {% for listing in listings %}
            <a href="/listing/{{ listing.id }}">{{ listing.
name }}</a>
               {{ listing.address }}
               {{ listing.price }}
               {{ listing.area }}
               {{ listing.house_type }}
               {{ listing.publish_date }}
            {% endfor %}
         </div>
</body>
</html>
数据可视化嵌入将生成的图表嵌入到Web页面中,提供数据分析结果的展示。
<!-- templates/analysis.html -->
<!DOCTYPE html>
<html lang="en">
 <head>
   <meta charset="UTF-8" />
   <title>Market Analysis</title>
   k
    rel="stylesheet"
    href="https://maxcdn.bootstrapcdn.com/bootstrap/4.5.2/css/bootstr
ap.min.css"
   />
```

```
</head>
  <body>
    <div class="container">
      <h1>Market Analysis</h1>
        src="/static/analysis_output/price_distribution.png"
        alt="Price Distribution"
      />
      <img
        src="/static/analysis_output/area_distribution.png"
        alt="Area Distribution"
      />
      <img
        src="/static/analysis_output/region_counts.png"
        alt="Region Counts"
      />
    </div>
  </body>
</html>
```

这份详细的实现步骤可以帮助你完整地实现项目的各个部分,包括数据采集、数据处理与分析以及 Web 应用的开发。

### 4. 使用说明

#### 环境配置

# Python 及依赖包安装

- 1. 安装 Python:
  - 请确保安装 Python 3.6 或以上版本。可以从 Python 官方网站下载并 安装。
- 2. 创建并激活虚拟环境:

```
python -m venv venv source venv/bin/activate # 在Windows 上使用 `venv\Scripts\activate`
```

- 安装项目依赖包:
  - 在项目根目录下创建一个名为 requirements.txt 的文件,并添加以下内容:

```
Flask==2.1.1
pandas==1.4.2
matplotlib==3.5.1
Scrapy==2.6.1
sqlite3
```

- 运行以下命令安装依赖包:

pip install -r requirements.txt

### 数据库配置

- 1. 确保 SQLite 已安装并配置好。
- 2. 确保项目根目录下有一个名为 real\_estate.db 的 SQLite 数据库文件。如果 没有,请按照下面的数据采集部分进行数据爬取并生成数据库文件。

#### 项目运行步骤

# 数据采集

- 1. 确保已在项目目录中创建并配置好 Scrapy 爬虫。
- 2. 运行 Scrapy 爬虫:

scrapy crawl your\_spider\_name # 将 your\_spider\_name 替换为实际爬虫 名称

3. 爬虫运行完成后,确保爬取的数据已存储到 real\_estate.db 数据库中。

# 数据处理与分析

1. 在项目目录中运行数据处理脚本:

python data\_process\_analysis.py

- 该脚本会加载数据库中的数据,进行清洗和分析,并生成数据可视化图表,保存到 app/static/analysis\_output 目录下。

#### Web 应用启动

- 1. 确保已完成 Flask 应用的配置和实现。
- 2. 在项目根目录下运行 Flask 应用:

export FLASK\_APP=run.py # 在Windows 上使用 `set FLASK\_APP=run.py` flask run

3. 打开浏览器,访问 http://127.0.0.1:5000,查看 Web 应用。

#### 功能使用指南

# 搜索与过滤功能

1. 在首页搜索栏中输入房源名称、地址、价格范围、面积范围或选择房型,然 后点击"Search"按钮进行搜索。

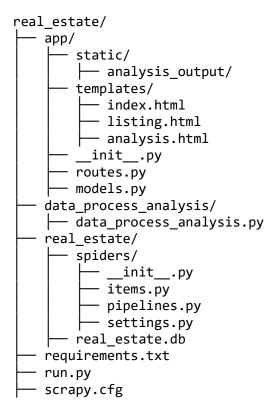
- 2. 页面将根据输入的条件显示符合条件的房源列表。
- 3. 使用表头的排序链接,可以按名称、地址、价格、面积、房型进行升序或降序排序。

# 数据分析结果查看

- 1. 在首页点击"View Analysis"按钮,进入数据分析结果页面。
- 2. 页面将展示房价分布图、面积分布图和各区域房源数量图。
- 3. 点击各图表,可以查看详细的市场分析结果。

# 项目目录结构

确保项目目录结构如下:



此使用说明提供了如何配置环境、运行项目和使用功能的具体步骤,帮助用户顺利完成项目的搭建和使用。