# python 中的数据爬取

杨晨

学号 2021212171

北京邮电大学计算机学院

日期：2024 年 2 月 26 日

## 1 概述

### 1.1 实验内容

1. 爬取学堂在线的计算机类课程页面内容
   https://www.xuetangx.com/search?query=&org=&classify=1&type=&status=&page=1
   要求将课程名称、老师、所属学校和选课人数信息，保存到一个 csv 文件中。
2. 爬取链家官网二手房的数据
   https://bj.lianjia.com/ershoufang/
   要求爬取北京市东城、西城、海淀和朝阳四个城区的数据（每个区爬取 5 页），将楼盘名称、总价、平米数、单价保存到 json 文件中。

### 1.2 开发环境

- Windows10
- PyCharm 2023.2.4 (Professional Edition)

## 2 实验过程

### 2.1 学堂在线计算机类课程内容爬取

#### 2.1.1 介绍

首先，按照 PPT 上的方法，用静态页面的方式进行爬取，但是发现爬虫关闭后，也没有爬取到任何数据
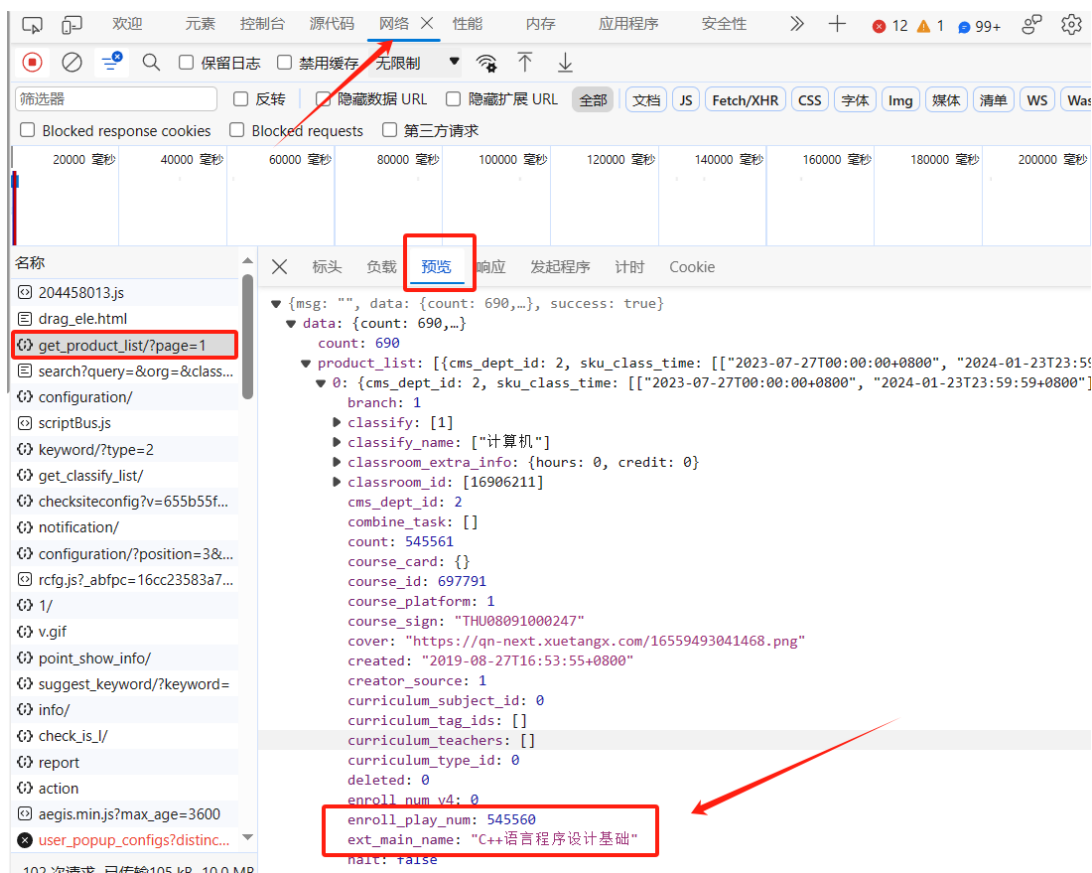
考虑网页是动态加载的，进入网页审查，发现以下结果

现在明确确实是动态页面。

然后进一步发现是 Ajax 类型，以下为截图
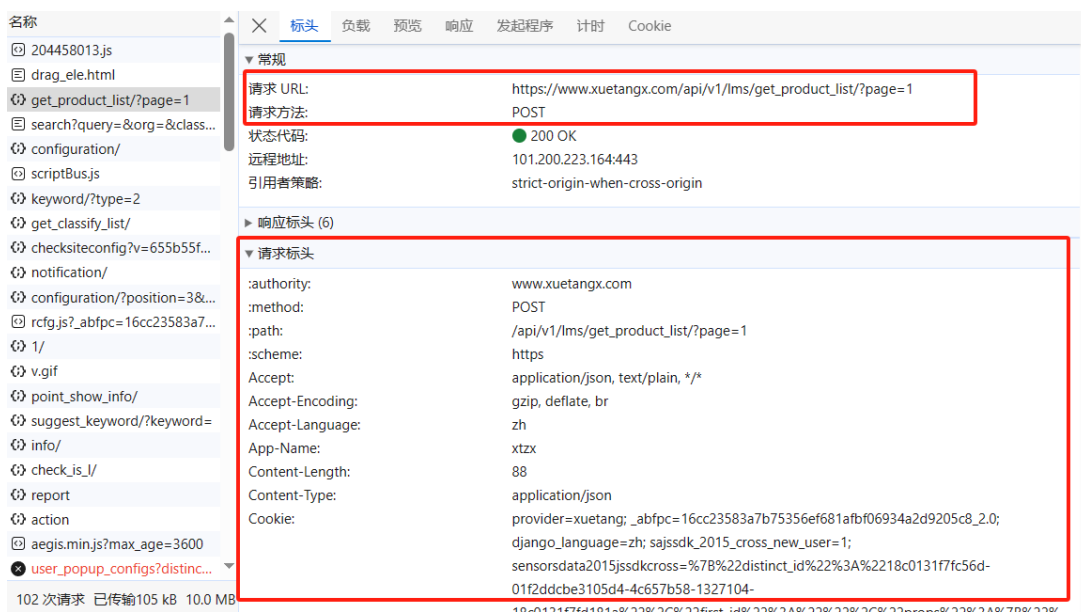


　　点击打开 page 里的 product_lists 可以看见课程的相关信息，正是我们需要爬取的内容，如下所示

基于以上的发现明确进行动态页面的爬取

### 2.1.2 爬取方法

1. 找到爬取的链接以及请求的方式



2. 构造 header，找到 Headers 选项，查看请求头
3. 将此代码赋值下来，填充 headers 的构造，代码如下所示

```
# 请求头
headers = {
    "authority": "www.xuetangx.com",
    "method": "POST",
    "scheme": "https",
    "accept": "application/json, text/plain, */*",
    "accept-encoding": "gzip, deflate, br",
    "accept-language": "zh",
    "content-type": "application/json",
    "cookie": "provider=xuetang; django_language=zh",
    "django-language": "zh",
    "origin": "https://www.xuetangx.com",
    "Referer": "https://www.xuetangx.com/search?query=&org=&classify=1&type=&
        status=&page=1",
    "sec-fetch-dest": "empty",
    "sec-fetch-mode": "cors",
    "sec-fetch-site": "same-origin",
    "user-agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit
        /537.36 "
    "(KHTML, like Gecko) Chrome/119.0.0.0 Safari/537.36",
    "X-Client": "web",
    "Xtbz": "xt",
}
```

4. 除了请求头此外，还需要 post 的提交表格数据



代码如下

```
data = {
    "query": "",
    "chief_org": [],
    "classify": ["1"],
    "selling_type": [],
    "status": [],
    "appid": 10000,
}
```

5. 在前面第一步已经分析，采取的是 post 的页面的请求方法，那么借助FormRequest函数即可实现 post 请求。重写tart_request(self)函数，然后经过这个函数不断的循环发送请

求，该函数代码实现如下

```python
def start_requests(self):
    for page_num in range(1, 70):
        yield scrapy.FormRequest(
            url="https://www.xuetangx.com/api/v1/lms/get_product_list/?page={}"
                ".format(
                page_num
            ),
            headers=self.headers,
            method="POST",
            body=json.dumps(self.data),
            callback=self.parse,
        )
```

6. 接着就是解析内容的阶段根据我们的需求，需要爬取课程名称，老师，所属学校以及课程人数，那么将 items.py 文件如下实现：

```python
class XuetangxItem(scrapy.Item):
    # 课程名称、教师、学校、学生人数
    course_name = scrapy.Field()
    teacher = scrapy.Field()
    school = scrapy.Field()
    student_num = scrapy.Field()
```

7. 那么对于 response 返回的内容，使用 json.loads 处理获取到的 json 文件，json.loads() 函数是将 json 格式数据转换为字典。网页 product_list 部分显示如下：

```
▶ classify_name: ["计算机"]
▶ classroom_extra_info: {hours: 0, credit: 0}
▶ classroom_id: [16906211]
  cms_dept_id: 2
  combine_task: []
  count: 545561
  course_card: {}
  course_id: 697791
  course_platform: 1
  course_sign: "THU08091000247"
  cover: "https://qn-next.xuetangx.com/16559493041468.png"
  created: "2019-08-27T16:53:55+0800"
  creator_source: 1
  curriculum_subject_id: 0
  curriculum_tag_ids: []
  curriculum_teachers: []
  curriculum_type_id: 0
  deleted: 0
  enroll_num_v4: 0
  enroll_play_num: 545560
  ext_main_name: "C++语言程序设计基础"
  halt: false
  id: 4333
  is_train_valid: 0
  lecture: []
  live_status: {}
  mobile_cover: ""
  modified: "2023-11-24T02:02:26+0800"
  name: "C++语言程序设计基础"
▼ org: {big_logo_name: "https://qn-next.xuetangx.com/15676613232012.png",…}
    big_logo: "https://qn-next.xuetangx.com/15675957781755.png"
```

那么对进行 product_list 内容如下方式进行提取

```python
def parse(self, response):
    for product in json.loads(response.body)["data"]["product_list"]:
        item = XuetangxItem()
        item["course_name"] = "课程名称: " + product["name"]
        item["teacher"] = "课程讲师: "
        for teacher in product["teacher"]:  # 一个课程可能有多个讲师
            item["teacher"] += teacher["name"] + " "
        item["school"] = "开课学校: " + product["org"]["name"]
        item["student_num"] = "学生人数: " + str(product["count"])
```

8. 然后 pipelines 进行将数据写入 csv 文件，那么 pipelines.py 文件如下所示

```python
import csv

class XuetangxPipeline:
    def __init__(self):
        self.file = open('XuetangxData.csv', 'w+', encoding='utf-8', newline='
            ')
        self.writer = csv.writer(self.file)

    def process_item(self, item, spider):
        self.writer.writerow([item['course_name'], item['teacher'], item['
            school'], item['student_num']])
        return item

    def close_spider(self, spider):
        self.file.close()
```

9. 最后当然需要将管道打开，那么 settings.py 文件如下所示

```python
ROBOTSTXT_OBEY = False
# ITEM_PIPELINES = {'myproject.pipelines.MyPipeline': 300}


ITEM_PIPELINES = {'myproject.pipelines.XuetangxPipeline': 300}
```

### 2.1.3　爬取结果

**图 1:** 学堂在线爬取结果

## 2.2 链家官网北京二手房数据爬取

### 2.2.1 介绍

链家的官网是静态页面，相比之下要容易处理

审查页面，找到想提取的楼盘名称、总价、平米数、单价的 xpath，进行提取即可

此外，对于第 1 页，第 2 页，网页的 url 具有相似特征/pg1，/pg2；这使得爬取较为方便

### 2.2.2 爬取方法

1. 构造 spider

```python
class LianjiaSpider(scrapy.spiders.Spider):
    name = "lianjia"
    allowed_domains = ["lianjia.com"]
    start_urls = [
        "https://bj.lianjia.com/ershoufang/dongcheng/pg1/",
        "https://bj.lianjia.com/ershoufang/xicheng/pg1/",
        "https://bj.lianjia.com/ershoufang/chaoyang/pg1/",
        "https://bj.lianjia.com/ershoufang/haidian/pg1/",
    ]
```

2. 需要爬取的信息有：楼盘名称、总价、平米数、单价。那么 items.py 中的代码如下

```python
class LianjiaItem(scrapy.Item):
    # 楼盘名称、总价、平米数、单价
    name = scrapy.Field()
    price = scrapy.Field()
    area = scrapy.Field()
```

```
        unit_price = scrapy.Field()
```

3. 审查页面，筛选想提取的信息，如下



在每个 sellListContent 类下面，每个 li 代表一个房子的信息。进行提取的代码如下

```python
def parse(self, response):
    item = LianjiaItem()
    distinct = response.url.split("/")[4]
    page = response.url.split("/")[5]
    for each in response.xpath('//ul[@class="sellListContent"]/li'):
        item["name"] = "楼盘名称: " + each.xpath("div/div/a/text()").get()
        price_value = each.xpath(
            "div/div[@class='priceInfo']/div[@class='totalPrice totalPrice2']/
                span/text()"
        ).get()
        price_unit = each.xpath(
            "div/div[@class='priceInfo']/div[@class='totalPrice totalPrice2']/
                i[last()]/text()"
        ).get()
        item["price"] = "总价: " + f"{price_value}{price_unit}"
        area_text = each.xpath(
            ".//div[@class='address']/div[@class='houseInfo']/text()"
        ).get()
        match = re.search(r"(\d+(\.\d+)?)平米", area_text)
```

```
        if match:
            item["area"] = "平米数: " + match.group(1) + "平米"
        else:
            item["area"] = "平米数: " + "unknown"
        item["unit_price"] = "单价: " + each.xpath(
            "div/div[@class='priceInfo']/div[@class='unitPrice']/span/text()"
        ).get()
        if item["name"] and item["price"] and item["area"] and item[
            "unit_price"]:
            yield item
```

为了便于提取平米数，使用了正则表达式 re 库

4. 因为需要提取前 5 页的内容，所以需要构造下一页的 URL，并发送一个新的请求，回调函数为自身的parse方法，实现翻页功能

```
if page != "pg5":
    next_page = int(page[2]) + 1
    next_url = "https://bj.lianjia.com/ershoufang/{}/pg{}/".format(distinct,
        next_page)
    yield scrapy.Request(next_url, callback=self.parse)
```

5. 然后 pipelines 进行将数据写入 json 文件，那么 pipelines.py 文件如下所示

```
class LianjiaPipeline:
    def __init__(self):
        self.file = open("LianjiaData.json", "w+", encoding="utf-8")
        self.writer = csv.writer(self.file)

    def process_item(self, item, spider):
        dict_item = dict(item)
        json_str = json.dumps(dict_item, ensure_ascii=False) + "\n"
        self.file.write(json_str)
        return item

    def close_spider(self, spider):
        self.file.close()
```

6. 最后需要将管道打开，那么 settings.py 文件如下所示

```
ROBOTSTXT_OBEY = False
# ITEM_PIPELINES = {'myproject.pipelines.MyPipeline': 300}


# ITEM_PIPELINES = {'myproject.pipelines.XuetangxPipeline': 300}


ITEM_PIPELINES = {'myproject.pipelines.LianjiaPipeline': 300}
```

### 2.2.3 爬取结果



**图 2:** 链家爬取结果

# 3 附录：spider.py 完整代码

## 3.1 爬取学堂在线

```python
import scrapy
from ..items import XuetangxItem
import json


class XuetangxSpider(scrapy.spiders.Spider):
    name = "xuetangx"
    allowed_domains = ["xuetangx.com"]
    # 请求头
    headers = {
        "authority": "www.xuetangx.com",
        "method": "POST",
        "scheme": "https",
        "accept": "application/json, text/plain, */*",
        "accept-encoding": "gzip, deflate, br",
        "accept-language": "zh",
        "content-type": "application/json",
        "cookie": "provider=xuetang; django_language=zh",
        "django-language": "zh",
        "origin": "https://www.xuetangx.com",
```

```python
        "Referer": "https://www.xuetangx.com/search?query=&org=&classify=1&type=&
            status=&page=1",
        "sec-fetch-dest": "empty",
        "sec-fetch-mode": "cors",
        "sec-fetch-site": "same-origin",
        "user-agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36
            "
        "(KHTML, like Gecko) Chrome/119.0.0.0 Safari/537.36",
        "X-Client": "web",
        "Xtbz": "xt",
    }
    data = {
        "query": "",
        "chief_org": [],
        "classify": ["1"],
        "selling_type": [],
        "status": [],
        "appid": 10000,
    }
    download_delay = 1

    def start_requests(self):
        for page_num in range(1, 70):
            yield scrapy.FormRequest(
                url="https://www.xuetangx.com/api/v1/lms/get_product_list/?page={}"
                    .format(
                    page_num
                ),
                headers=self.headers,
                method="POST",
                body=json.dumps(self.data),
                callback=self.parse,
            )

    def parse(self, response):
        for product in json.loads(response.body)["data"]["product_list"]:
            item = XuetangxItem()
            item["course_name"] = "课程名称: " + product["name"]
            item["teacher"] = "课程讲师: "
            for teacher in product["teacher"]:  # 一个课程可能有多个讲师
                item["teacher"] += teacher["name"] + " "
            item["school"] = "开课学校: " + product["org"]["name"]
            item["student_num"] = "学生人数: " + str(product["count"])

            if (
                item["course_name"]
                and item["teacher"]
```

```
            and item["school"]
            and item["student_num"]
        ):
            yield item
```

## 3.2 爬取链家

```python
import scrapy
from ..items import LianjiaItem
import re


class LianjiaSpider(scrapy.spiders.Spider):
    name = "lianjia"
    allowed_domains = ["lianjia.com"]
    start_urls = [
        "https://bj.lianjia.com/ershoufang/dongcheng/pg1/",
        "https://bj.lianjia.com/ershoufang/xicheng/pg1/",
        "https://bj.lianjia.com/ershoufang/chaoyang/pg1/",
        "https://bj.lianjia.com/ershoufang/haidian/pg1/",
    ]

    def parse(self, response):
        item = LianjiaItem()
        distinct = response.url.split("/")[4]
        page = response.url.split("/")[5]
        for each in response.xpath('//ul[@class="sellListContent"]/li'):
            item["name"] = "楼盘名称: " + each.xpath("div/div/a/text()").get()
            price_value = each.xpath(
                "div/div[@class='priceInfo']/div[@class='totalPrice totalPrice2']/
                    span/text()"
            ).get()
            price_unit = each.xpath(
                "div/div[@class='priceInfo']/div[@class='totalPrice totalPrice2']/i
                    [last()]/text()"
            ).get()
            item["price"] = "总价: " + f"{price_value}{price_unit}"
            area_text = each.xpath(
                ".//div[@class='address']/div[@class='houseInfo']/text()"
            ).get()
            match = re.search(r"(\d+(\.\d+)?)平米", area_text)
            if match:
                item["area"] = "平米数: " + match.group(1) + "平米"
            else:
                item["area"] = "平米数: " + "unknown"
            item["unit_price"] = "单价: " + each.xpath(
```

```python
                "div/div[@class='priceInfo']/div[@class='unitPrice']/span/text()"
            ).get()
            if item["name"] and item["price"] and item["area"] and item["unit_price
                "]:
                yield item

        if page != "pg5":
            next_page = int(page[2]) + 1
            next_url = "https://bj.lianjia.com/ershoufang/{}/pg{}/".format(distinct
                , next_page)
            yield scrapy.Request(next_url, callback=self.parse)
```