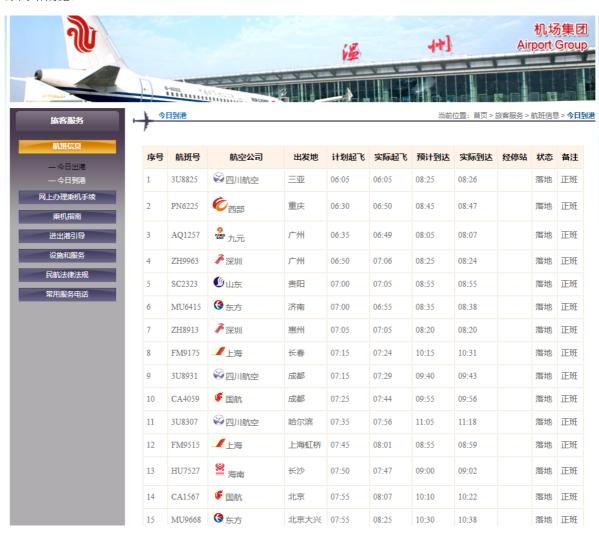
# 计网作业: 爬虫

#### 温州机场数据地址

#### 原网站概览:



# 使用说明

# 环境配置方法

- Scrapy
- 1 pip install Scrapy
- 2 scrapy startproject airport
- 3 cd airport # 进入项目airport目录
- 4 scrapy genspider wenzhou http://www.wzair.cn/lkfw/hbxx/jrdg/index.html? v=1606439135598&rxLoad=1&\_rand=1606586459057
- selenium

见教程: 爬取动态渲染网站

### 代码运行说明

在外层airport文件夹下,输入以下命令,即可运行(每24h爬取一次数据)

```
1 | python main.py # 使用python 3版本
```

# 实验过程

一开始,我使用scrapy框架进行爬取,发现自xpath中自//table后,便无法获得向下节点的对应信息。于是,我考虑了用以下两种方式来爬:

- 获取API
- selenium模拟网页行为爬取

### 通过API

```
import requests
    import time
3
   import re
4
    import datetime
    import json
 6
    import os
7
8
    se = requests.session()
9
    post_url = "http://121.40.33.186:12800/services/wzairAuto/arrival?
10
    callback=jsonp1606647111005&_=1606647111010"
11
    while True:
12
       ti = datetime.datetime.now().strftime('%Y-%m-%d')
13
        filename = '[温州机场]'+ti+'.json'
14
        data = se.get(post_url).text.replace("'", '"').replace('/ ', '/')
        with open(filename, "w", encoding="utf-8") as f:
15
16
            f.write(data)
        # os.system("scrapy crawl wenzhou")
17
18
        print("load data success! ",ti)
19
        time.sleep(3600*24) # 每24小时爬一次
```

### 通过selenium

在 airport/app.py 中,代码如下:

```
from selenium import webdriver
 2
    from selenium.webdriver.chrome.options import Options
 3
    import time
    chrome_options = Options()
6
    chrome_options.add_argument("--headless")
7
    driver = webdriver.Chrome(
8
        executable_path='chromedriver',
9
        options=chrome_options)
10
    driver.get('http://www.wzair.cn/lkfw/hbxx/jrdg/index.html?v=1606439135598/')
    time.sleep(5)
11
12
    driver.find_elements_by_xpath("//td[@id='hangbanInfoBox']/table/tbody/tr")
    for i in li:
13
```

## 数据相关

### API法

数据截图

```
{} [温州机场]2020-11-29.json
       jsonp1606647111005([
               "air attr": "落地",
               "air_company": "四川航空",
               "air etime1": "08:25",
               "air etime2": "08:26",
               "air etime3": "08:25",
               "air from": "三亚",
               "air id": 1734569,
               "air memo": "正班",
               "air number": "3U8825",
               "air_stay": "",
  12
               "air_stime1": "06:05",
               "air stime2": "06:05",
               "air time": 1606610700000,
               "air to": "温州",
               "air type": 0,
               "via3": "",
               "via4": ""
               "viac": ""
               "air_attr": "落地",
               "air_company": "西部",
               "air etime1": "08:45"
               "air etime2": "08:47",
               "air etime3": "08:45",
               "air from": "重庆",
               "air id": 1734404,
               "air memo": "正班",
               "air number": "PN6225",
               "air stav": ""
```

#### 字段说明

以爬得的第一条数据为例:

```
1
        {
 2
           "air_attr": "落地", // 状态
 3
           "air_company": "四川航空", // 航空公司
           "air_etime1": "08:25", // 预计到达
 4
 5
           "air_etime2": "08:26", // 实际到达
           "air_etime3": "08:25", // 多余字段
 6
 7
           "air_from": "三亚", // 出发地
 8
           "air_id": 1734569, // 航班编号
9
           "air_memo": "正班", // 备注
10
           "air_number": "3U8825", // 航班号
           "air_stay": "", // 经停站
11
           "air_stime1": "06:05", // 计划起飞
12
13
           "air_stime2": "06:05", // 实际起飞
14
           "air_time": 1606610700000, // 时间
15
           "air_to": "温州", // 到达地
           "air_type": 0, // 类型
16
           "via3": "", // 冗余字段
17
           "via4": "", // 冗余字段
18
           "viac": "" // 冗余字段
19
20
        }
```

### selenium法

#### 数据截图

```
序号 航班号 航空公司 出发地 计划起飞 实际起飞 预计到达 实际到达 经停站 状态 备注
1 3U8825 四川航空 三亚 06:05 06:07 08:30 08:29 落地 正班
2 PN6225 西部 重庆 06:30 06:37 08:30 08:36 落地 正班
3 AQ1257 九元 广州 06:35 06:50 08:10 08:05 落地 正班
4 ZH9963 深圳 广州 06:50 07:01 08:20 08:17 落地 正班
5 MU6415 东方 济南 07:00 06:55 08:35 08:43 落地 正班
6 SC2323 山东 贵阳 07:00 07:21 09:10 09:08 落地 正班
7 ZH8913 深圳 惠州 07:05 07:09 08:25 08:32 落地 正班
8 3U8931 四川航空 成都 07:15 07:31 09:40 09:45 落地 正班
9 FM9175 上海 长春 07:15 07:22 10:20 起飞 正班
10 CA4059 国航 成都 07:25 07:46 10:00 09:55 落地 正班
11 8L9517 祥鹏航空 郑州 07:30 07:27 09:15 09:13 落地 正班
12 308307 四川航空 哈尔滨 07:35 07:48 11:05 起飞 正班
13 FM9515 上海 上海虹桥 07:55 08:01 08:55 08:53 落地 正班
14 CZ3811 南方 广州 07:55 08:16 09:35 09:30 落地 正班
15 CA1567 国航 北京 07:55 07:56 10:10 起飞 正班
16 MU9668 东方 北京大兴 08:00 08:10 10:20 起飞 正班
17 CZ3833 南方 贵阳 08:15 取消 正班
18 MU5861 东方 泸州 08:15 08:09 10:25 起飞 正班
19 KY8223 昆航 昆明 08:15 08:34 10:55 起飞 正班
20 MU2655 东方 成都 08:20 08:49 11:05 起飞 正班
21 CZ8295 南方 郑州 08:25 08:21 10:05 10:05 落地 正班
22 QW6090 青岛航空 长沙 08:40 08:54 10:05 起飞 正班
23 ZH9925 深圳 深圳 08:55 09:22 10:45 起飞 正班
24 HU7481 海南 重庆 08:55 09:53 11:45 起飞 正班
25 9H8327 长安航空 西安 09:00 09:24 10:50 宜昌 起飞 正班
```

### 字段说明

与网页中的表格项一一对应。分别是:

序号 航班号 航空公司 出发地 计划起飞 实际起飞 预计到达 实际到达 经停站 状态 备注