## 爬取知乎关注者信息

---

### 爬虫的一般思路

抓取网页

解析网页

储存数据

### 实战环节

分析知乎

动手写爬虫

# 爬虫的一般思路

**1** 抓取网页、分析请求。

**2** 解析网页、寻找数据。
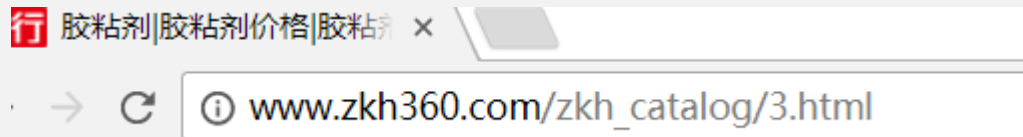
**3** 储存数据、多页处理。

# 分析具体网页

http://www.zkh360.com/zkh_catalog/3.html

胶粘剂|胶粘剂价格|胶粘 ×

→ C ⓘ www.zkh360.com/zkh_catalog/3.html

翻页后URL不变
如何寻找请求？

哪些网站也是这样？

# 使用Chrome



翻页后出现的新变化

# 如何寻找真实的请求

分析

真实的
请求

测试

重复

尝试寻找
淘宝的请求

# 实战环节

https://www.zhihu.com/people/excited-vczh/following



寻找请求

# 动手写爬虫

知乎请求头信息

常见的请求头

▼ Request Headers    view source
accept: application/json, text,
Accept-Encoding: gzip, deflate,
Accept-Language: zh-CN,zh;q=0.
authorization: Bearer Mi4xUUhNc(
7|b43e672333e0fe298111f2b18ed!
Connection: keep-alive
Cookie: d_c0="AJAClbp8MQyPTvIc
7962.1502441667; _xsrf=cf9ef5
0; aliyungf_tc=AQAAAPRDWmS6EA(
5ea2e42c1a1d5ff217b1c0fa09c9b
d2d2db42baf3c904a9"; __utma=5
=51854390; __utmz=51854390.15(
e=20170809=1; z_c0=Mi4xUUhNc0[
|b43e672333e0fe298111f2b18ed9(
Host: www.zhihu.com
Referer: https://www.zhihu.com/
User-Agent: Mozilla/5.0 (Windo
X-UDID: AJAClbp8MQyPTvIcjXY7p-

找不同

Accept:
Accept-Encoding:
Accept-Language:
Connection:
Cookie:
Host:
Referer:
User-Agent:

# 保存单网页信息

```
# -*- coding:utf-8 -*-

import requests
import pandas as pd

headers = {
    'authorization':'Bearer Mi4xUUhNc0FnQUFBQUFBa0FLVnVud3hEQmNBQUFFCaEFsVk5oeW5qV1FBcmVhY2
    F0VUhWenJTc1hVcUlycW1tRzAtMXpB|1505467527|b43e672333e0fe298111f2b18ed9e52c9f8f23d7',
    'User-Agent':'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like
        Gecko) Chrome/62.0.3192.0 Safari/537.36',
}

url = 'https://www.zhihu.com/api/v4/members/excited-vczh/followees?offset=0&limit=20'
r = requests.get(url,headers=headers).json()

df = pd.DataFrame.from_dict(r['data'])
df.to_csv('users.csv')
```

设置请求头

获得数据

保存数据

DC学院 | 造数
class.pkbigdata.com

# 保存多页

```python
# -*- coding:utf-8 -*-

import requests
import pandas as pd
import time

headers = {
    'authorization':'Bearer Mi4xUUhNc0FnQUFBQUFBa0FFVnVud3hEQmNBQUFFaEFFsVk5oeW5qV1FBcmVhY2F0F0VUhWenJTc1hVcU1ycW1tRzAtMXpB|'\
    '1505467527|b43e672333e0fe298111f2b18ed9e52c9f8f23d7',
    'User-Agent':'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/62.0.3192.0 Safari/537.36',
}

user_data = []

def get_user_data(page):
    for i in range(page):
        url = 'https://www.zhihu.com/api/v4/members/excited-vczh/followees?offset={}&limit=20'.format(i*20)
        response = requests.get(url, headers=headers).json()
        user_data.extend(response['data'])
        print("成功爬取第%s页" % str(i+1))
        time.sleep(1)

if __name__ == '__main__':
    get_user_data(8)
    df = pd.DataFrame.from_dict(user_data)
    df.to_csv('user.csv')
```

翻页处理

暂停爬虫

调用函数

# 简单数据可视化



vczh关注男女比

# 完成作业

- **查看淘宝的翻页**

- **将vczh改为其他人**

- **尝试爬取知乎粉丝信息**