# 爬虫案例3(模拟登陆,验证码登陆,cookie登陆)

http://example.webscraping.com/places/default/user/login

模拟登陆的网站
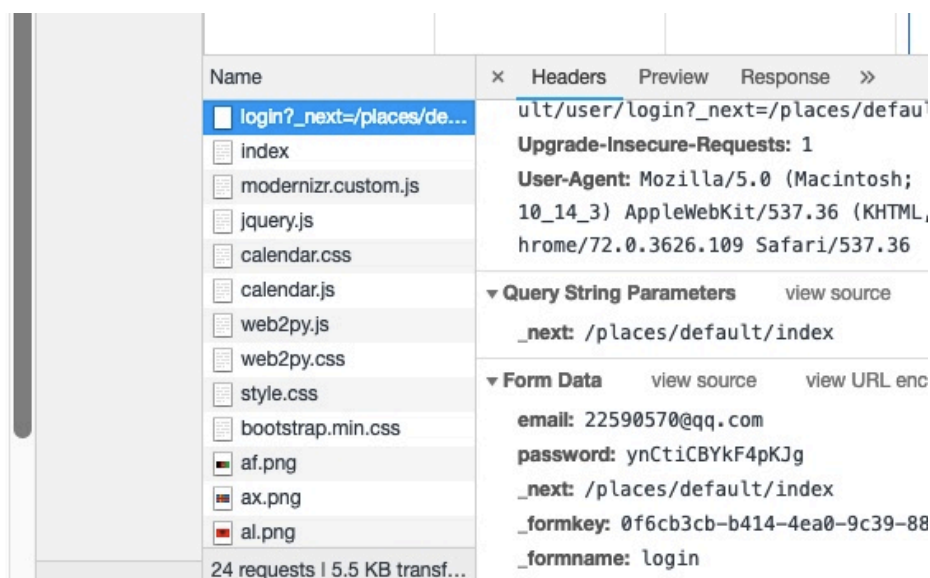
form表单说明下:

method 决定了http请求方法

action 决定了请求的url

enctype 决定了表单数据编码类型

input 决定了表单数据的内容



以上是请求信息

—————接下来分析响应信息

▼ **Response Headers**     view source

Connection: keep-alive

Content-Type: text/html; charset=UTF-8

Date: Tue, 19 Feb 2019 01:36:25 GMT

Location: /places/default/index

Server: nginx

Set-Cookie: session_id_places=True; httponly; Path=/

Set-Cookie: session_data_places="2d908112833a47176ff21f1cbd95e5a9:g7Z
1G3xOgOJuA1r-HzjGZXD4zRcvjN8-ek5LVge28ZEpglqfGogBEhzSpvMkjpYhTDEV2a
6gBAZ2i__WlstAyPwvN8d_YeTfxq8xMSJ2gLKJHI9QkX4xG0S-9HVIStNJnFIyD_SvT
aEdxfG_rgNKTFSstfZYJPVyk8zvap2RSF0tLTu7Ud195Kh6ozDjkBYNj6tNmLWZh1lW
MUQsFmXIrS4DG2J040x7Z4MQBYxA9BTyGqz2s3Q2F4YbEnr7ApVdPT8UeXPS162kFDS
dpYSOYvuzpyXIy9tF1bpkyaorBRbvCDr4CIxlgkSv638Cgp6u1wLaPhJloyR3bDzMWY
Rk1FRFEDDkM1wA3Zg4hU-0R-u7W74rL_8VkUnaVkmFPLhdiUPHoZ6vhuvNdImo5oBsh

——控制台操作

```
1  scrapy shell http://example.webscraping.com/places/default/user/login?
   _next=/places/default/index
```

接下来模拟

```
In [1]: from scrapy.http import FormRequest

In [2]: fd ={'email':'22590570@qq.com','password':'ynCtiCBYkF4pKJg'}

In [3]: fd
Out[3]: {'email': '22590570@qq.com', 'password': 'ynCtiCBYkF4pKJg'}

In [4]: request = FormRequest.from_response(response,formdata=fd)
```

```
In [8]: 'Welcome 伟杰' in response.text
Out[8]: True
```

```
1  ----代码如下
```

```
1  # -*- coding: utf-8 -*-
2  import scrapy
3  from ..items import LogindemoItem
4  from scrapy.http import  Request,FormRequest
5  class LoginSpider(scrapy.Spider):
6      name = 'login'
7      allowed_domains = ['example.webscraping.com']
8      #这个是登陆之后请求的详情页面
9      #目的就是主要是为了能够访问登陆之后能够看到的页面
10     start_urls = ['http://example.webscraping.com/places/default/user/profile']
11
12     def parse(self, response):
13         value = response.css('#auth_user_email__row > td.w2p_fw::text').extract_first()
14         loginMM = LogindemoItem()
15         loginMM['email_value'] = value
16         yield loginMM
17
18     #设置登陆页面的url
```

```
19        login_url = 'http://example.webscraping.com/places/default/user/login'
20
21    #这个方法是爬虫进入后第一次执行的方法
22    def start_requests(self):
23        yield Request(self.login_url,callback=self.login)
24    #调用login请求登陆
25    def login(self,response):
26        fd = {'email':'22590570@qq.com','password':'ynCtiCBYkF4pKJg'}
27        yield  FormRequest.from_response(response,formdata=fd,callback=self.parse_login)
28
29    #检查是否成功通过判断response里面的内容,调用父类的start_requests方法
30    #父类的方法请求的是start_urls的地址
31    def parse_login(self,response):
32        if 'Welcome 伟杰' in response.text:
33            yield  from  super().start_requests()
34
35
36
37
```

———识别验证码-------

需要用到的东西

OCR 验证码识别库

pytesseract 第三方库

```
brew  tesseract
```

```
1 sudo pip3 install pillow
```

```
1 sudo pip3 install pytesseract
```

等等，我找个网站

https://my.vultr.com/

以这个网站为基础

[https://my.vultr.com/billing/](https://my.vultr.com/billing/) 这个是查看个人余额的页面

来走起，新建一个项目 核心代码如下

```python
# -*- coding: utf-8 -*-
import scrapy
import requests
from scrapy.http import Request,FormRequest
from PIL import Image
from io import BytesIO
import  pytesseract

class VultrSpider(scrapy.Spider):
    name = 'vultr'
    allowed_domains = ['my.vultr.com']
    #解析余额的页面
    start_urls = ['https://my.vultr.com/billing/']

    def parse(self, response):
        print('突突突突突突')
        print(response.css('#header2_right > div:nth-child(2) > a').extract_first())


    #设置登录的网址
    login_url = 'https://my.vultr.com/'

    #请求登录网站调用login方法
    def start_requests(self):
        yield  Request(self.login_url,callback=self.login,dont_filter=True)

```

```python
27      #login方法解析验证码和构建表单
28      def login(self, response):
29          #获取验证码图片url
30          captchaUrl = response.css('#loggedout_module > form > div.captcha_container >
    img::attr(src)').extract_first()
31          captchaUrl = response.urljoin(captchaUrl)
32          #识别一下咯
33          imgresponse = requests.get(captchaUrl)
34          BytesIOObj = BytesIO(imgresponse.content)
35          #BytesIOObj.write(imgresponse.content)
36          img = Image.open(BytesIOObj)
37          img = img.convert('L')
38          img.show()
39          #captcha = pytesseract.image_to_string(img)
40          #改为手动输入
41          captcha = input('输入图片中的验证码')
42          print('验证码'+captcha)
43          #构建表单数据
44          fd = {'username': '22590570@qq.com', 'password': 'Wwjxly0215', 'captcha': captcha}
45          print(fd)
46          #因为vultr登录和访问进去后都是用一个网址
47          #scrapy默认会过滤重复网页，发起Request添加dont_filter=True，则可以重复请求
48          yield FormRequest.from_response(response, formdata=fd,
    callback=self.parse_login,dont_filter=True)
49      #检查是否成功
50      def parse_login(self,response):
51          if 'Welcome to Vultr.com!' in response.text:
52              yield from super().start_requests()
53
54
55
56
57
58
59
60
61
```



```
                                        'scrapy.spidermiddlewares.offsite.OffsiteMiddleware',
                              .Spi    'scrapy.spidermiddlewares.referer.RefererMiddleware',
                                        'scrapy.spidermiddlewares.urllength.UrlLengthMiddleware',
                                        'scrapy.spidermiddlewares.depth.DepthMiddleware']
                              my.v    2019-02-19 17:01:53 [scrapy.middleware] INFO: Enabled item pipelines:
                                        []
                              ://m    2019-02-19 17:01:53 [scrapy.core.engine] INFO: Spider opened
                                        2019-02-19 17:01:53 [scrapy.extensions.logstats] INFO: Crawled 0 pages (at 0 pa
                                        es/min), scraped 0 items (at 0 items/min)
                              onse    2019-02-19 17:01:53 [scrapy.extensions.telnet] DEBUG: Telnet console listening
                              ')      n 127.0.0.1:6023
                                        2019-02-19 17:01:54 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://my.v
   elines.py        (response.css('  ltr.com/robots.txt> (referer: None)
   tings.py      17                   2019-02-19 17:01:54 [scrapy.core.engine] DEBUG: Crawled (200) <GET https://my.v
   y.cfg         18                   ltr.com/> (referer: None)
   sv            19     网址          2019-02-19 17:01:54 [urllib3.connectionpool] DEBUG: Starting new HTTPS connect
   braries       20          = 'https://my.  n (1): my.vultr.com:443
   and Consoles  21                   2019-02-19 17:01:56 [urllib3.connectionpool] DEBUG: https://my.vultr.com:443 "
                 22                   T /_images/captcha.php?s=5c6bc60296f17 HTTP/1.1" 200 None
                                        2019-02-19 17:01:56 [PIL.PngImagePlugin] DEBUG: STREAM b'IHDR' 16 13
                 VultrSpider › parse()  2019-02-19 17:01:56 [PIL.PngImagePlugin] DEBUG: STREAM b'IDAT' 41 2724
                                        输入图片中的验证码[]
```

```
ltr.com/billing/> (referer: https://my.vultr.com/)
突突突突突突
<a href="/billing/#billinghistory" style="color:#7cb342;" onclick="changeTabSub
enu('billinghistory'); hidePostMessages();">

7.40
                                        </a>
```

———cookie登录----

有些验证码特别难读，另外一种思路就是：

　　登录网站后，用户信息的cookie会被保存在本地，利用爬虫直接使用cookie发起http请求，就可以绕过表单登录

1.先安装第三方库，获取谷歌和火狐的cookie

pip3 install browsercookie

```
>>> import browsercookie
>>> chrome_cookiejar = browsercookie.chrome()
>>> for cookie in chrome_cookiejar:
...     print(cookie)
...
```

打是打出来了很多cookie信息

（这个思路也可以参考）

https://blog.csdn.net/fox64194167/article/details/81055327

———实现一个自定义的使用cookie中间件

在middlewares中

加入

```
 1  import browsercookie
 2  from scrapy.downloadermiddlewares.cookies import CookiesMiddleware
 3  #实现cookie的中间件
 4  class BrowserCookiesMiddleware(CookiesMiddleware):
 5      def __init__(self,debug=False):
 6          super().__init__(debug)
 7          self.load_browser_cookies()
 8
 9      #自定义加载浏览器的cookie的方法
10      def load_browser_cookies(self):
11          jar = self.jars['chrome']
12          chrome_cookiejar = browsercookie.chrome()
13          for cookie in chrome_cookiejar:
14              jar.set_cookie(cookie)
15          jar = self.jars['firefox']
16          firefox_cookiejar = browsercookie.firefox()
17          for cookie in firefox_cookiejar:
18              jar.set_cookie(cookie)
```

----接下来试一下模拟知乎中的cookie访问

1.创建项目：

scrapy startproject browser_cookie

```python
1  # -*- coding: utf-8 -*-
2
3  # Define here the models for your spider middleware
4  #
5  # See documentation in:
6  # https://doc.scrapy.org/en/latest/topics/spider-middleware.html
7
8  from scrapy import signals
9  import browsercookie
10 from scrapy.downloadermiddlewares.cookies import CookiesMiddleware
11
12
13 class BrowserCookieSpiderMiddleware(object):
14     # Not all methods need to be defined. If a method is not defined,
15     # scrapy acts as if the spider middleware does not modify the
16     # passed objects.
17
18     @classmethod
19     def from_crawler(cls, crawler):
20         # This method is used by Scrapy to create your spiders.
21         s = cls()
22         crawler.signals.connect(s.spider_opened, signal=signals.spider_opened)
23         return s
24
25     def process_spider_input(self, response, spider):
26         # Called for each response that goes through the spider
27         # middleware and into the spider.
28
29         # Should return None or raise an exception.
30         return None
31
32     def process_spider_output(self, response, result, spider):
33         # Called with the results returned from the Spider, after
34         # it has processed the response.
35
36         # Must return an iterable of Request, dict or Item objects.
37         for i in result:
38             yield i
39
40     def process_spider_exception(self, response, exception, spider):
41         # Called when a spider or process_spider_input() method
42         # (from other spider middleware) raises an exception.
43
44         # Should return either None or an iterable of Response, dict
45         # or Item objects.
46         pass
47
48     def process_start_requests(self, start_requests, spider):
49         # Called with the start requests of the spider, and works
50         # similarly to the process_spider_output() method, except
51         # that it doesn't have a response associated.
52
```

```python
        # Must return only requests (not items).
        for r in start_requests:
            yield r

    def spider_opened(self, spider):
        spider.logger.info('Spider opened: %s' % spider.name)


class BrowserCookieDownloaderMiddleware(object):
    # Not all methods need to be defined. If a method is not defined,
    # scrapy acts as if the downloader middleware does not modify the
    # passed objects.

    @classmethod
    def from_crawler(cls, crawler):
        # This method is used by Scrapy to create your spiders.
        s = cls()
        crawler.signals.connect(s.spider_opened, signal=signals.spider_opened)
        return s

    def process_request(self, request, spider):
        # Called for each request that goes through the downloader
        # middleware.

        # Must either:
        # - return None: continue processing this request
        # - or return a Response object
        # - or return a Request object
        # - or raise IgnoreRequest: process_exception() methods of
        #   installed downloader middleware will be called
        return None

    def process_response(self, request, response, spider):
        # Called with the response returned from the downloader.

        # Must either;
        # - return a Response object
        # - return a Request object
        # - or raise IgnoreRequest
        return response

    def process_exception(self, request, exception, spider):
        # Called when a download handler or a process_request()
        # (from other downloader middleware) raises an exception.

        # Must either:
        # - return None: continue processing this exception
        # - return a Response object: stops process_exception() chain
        # - return a Request object: stops process_exception() chain
        pass

    def spider_opened(self, spider):
        spider.logger.info('Spider opened: %s' % spider.name)


#实现cookie的中间件
```

```
109  class BrowserCookiesMiddleware(CookiesMiddleware):
110      def __init__(self,debug=False):
111          super().__init__(debug)
112          self.load_browser_cookies()
113
114      #自定义加载浏览器的cookie的方法
115      def load_browser_cookies(self):
116          jar = self.jars['chrome']
117          chrome_cookiejar = browsercookie.chrome()
118          for cookie in chrome_cookiejar:
119              jar.set_cookie(cookie)
120          jar = self.jars['firefox']
121          firefox_cookiejar = browsercookie.firefox()
122          for cookie in firefox_cookiejar:
123              jar.set_cookie(cookie)
```

2.在setting.py中添加一下配置

伪装浏览器

http://blog.51cto.com/laoyinga/2046970 随机

```
1  ROBOTSTXT_OBEY = False
```

https://blog.csdn.net/you_are_my_dream/article/details/60479699 关于上面设置false的说明

```
1  USER_AGENT = 'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_14_3) AppleWebKit/537.36 (KHTML,
   like Gecko) Chrome/72.0.3626.109 Safari/537.36'
2  DOWNLOADER_MIDDLEWARES = {
3      'scrapy.downloadermiddlewares.cookies.CookiesMiddleware':None,
4      'browser_cookie.middlewares.BrowserCookiesMiddleware': 5,
5  }
```

3.构建代码

```
1  # -*- coding: utf-8 -*-
2  import scrapy
3
4
5  class ZhihuSpider(scrapy.Spider):
6      name = 'zhihu'
7      allowed_domains = ['www.zhihu.com']
8      start_urls = ['https://www.zhihu.com/settings/account']
9
10     def parse(self, response):
11         print(response.css('#SettingsMain > div > div:nth-child(1) > div > div > div > div >
   h2'))
12
13     #重写请求
14     def start_requests(self):
15         url = 'https://www.zhihu.com/settings/account'
16         yield scrapy.Request(url,meta={'cookiejar':'chrome'},callback=self.parse)
17
```