

requests模块

Requests: 让 HTTP 服务人类 虽然Python的标准库中 urllib 模块已经包含了平常我们使用的大多数功能，但是它的 API 使用起来让人感觉不太好，而 Requests 自称 "HTTP for Humans"，说明使用更简洁方便。

Requests 唯一的一个非转基因的 Python HTTP 库，人类可以安全享用：)

Requests 继承了urllib的所有特性。Requests支持HTTP连接保持和连接池，支持使用cookie保持会话，支持文件上传，支持自动确定响应内容的编码，支持国际化的 URL 和 POST 数据自动编码。

Requests的文档非常完备，中文文档也相当不错。Requests能完全满足当前网络的需求，支持Python 2.6--3.5，而且能在PyPy下完美运行。

开源地址：<https://github.com/kennethreitz/requests>

中文文档 API：http://docs.pythonrequests.org/zh_CN/latest/index.html

安装方式

利用pip安装

```
1 $ pip install requests -i https://pypi.douban.com/simple
```

基本GET请求（headers参数 和 parmas参数）

1. 最基本的GET请求可以直接用get方法

```
1 response = requests.get("http://www.baidu.com/")
2
3 # 也可以这么写
4 # response = requests.request("get", "http://www.baidu.com/")
```

2. 添加 headers 和 查询参数

如果想添加 headers，可以传入 `headers` 参数来增加请求头中的headers信息。如果要将参数放在url中传递，可以利用 `params` 参数。

```
1 import requests
2
3 kw = {'wd': '长城'}
4
5 headers = {
6     "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64)
7     AppleWebKit/537.36 (KHTML, like Gecko) Chrome/54.0.2840.99
8     Safari/537.36"
9 }
10
11 # params 接收一个字典或者字符串的查询参数，字典类型自动转换为url
12 # 编码，不需要urlencode()
13 response = requests.get("http://www.baidu.com/s?", params =
14 kw, headers = headers)
15
16 # 查看响应内容，response.text 返回的是str格式的数据
17 print(response.text)
18
19 # 查看响应内容，response.content返回的字节流数据
20 print(response.content)
21
22 # 查看完整url地址
23 print(response.url)
24
25 # 查看响应头部字符编码
26 print(response.encoding)
```

```
23
24 # 查看响应码
25 print(response.status_code)
```

运行结果

```
1 .....
2
3 .....
4
5 'http://www.baidu.com/s?wd=%E9%95%BF%E5%9F%8E '
6
7 'utf-8'
8
9 200
```

- 使用response.text 时，Requests 会基于 HTTP 响应的文本编码自动解码响应内容，大多数 Unicode 字符集都能被无缝地解码。
- 使用response.content 时，返回的是服务器响应数据的原始二进制字节流，可以用来保存图片等二进制文件。

为什么要使用User-Agent?

```
1 In [12]: resp1 = requests.get("https://www.baidu.com")
2
3 In [13]: resp1.status_code
4 Out[13]: 200
5
6 In [14]: len(resp1.text)
7 Out[14]: 2443
8
9 In [15]: headers
10 Out[15]: {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/54.0.2840.99 Safari/537.36'}
11
```

```
12 In [16]: resp2 = requests.get("https://www.baidu.com",
13     headers=headers)
14
14 In [17]: resp2.status_code
15 Out[17]: 200
16
17 In [18]: len(resp2.text)
18 Out[18]: 156612
```

3. POST请求

github登录

```
1 import requests
2
3 login_url = "https://github.com/session"
4
5 headers = {
6     "User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X
7     10_13_4) AppleWebKit/537.36 (KHTML, like Gecko)
8     Chrome/66.0.3359.170 Safari/537.36",
9     "Cookie": "_ga=GA1.2.1855430798.1461857641;
10     _octo=GH1.1.783519559.1525492869; tz=Asia%2FShanghai;
11     logged_in=no; _gat=1;
12     _gh_sess=TTZZNnVPOTRwTlI2YW1UOWFqTUZwa1JHVENFOVpjYWxoekJLbS9
13     GTGFQqkR40DZVSTIrV1NTWUFQS0dETUhURmg4YVR3bGxjV1hJT2dMZ2NIOEp
14     tZlREUUXrOXV0eG1EcG5kUi9adUZQamFpWmxmMHVhY04vckJVWkxXbkNVa3Z
15     HNCtjekF2WEVnMGVzaElqMnBpVUIvVGVSZjmdjNQMFxYwPONE1HMks2eDh
16     pVzZ0Wk9ZQUZLMVJ0OTR5YXJWUjV6VUNmRFhyaH1YczUzdUozSWR5M210akh
17     3dkcvaXRhY2ZmanRNRC9IbElUMm5OSmkzVDhtbEwvSEdGWFmvd0xySWIxcFR
18     rbGZ0RDQwQit5eGUvaGZkdWp5U1dYSnZ0VzRRb1FqUThZZXlreTBNU1RicUh
19     heGJGQjFRcDlnN1N2b2RXRXRsT21lRFB5Q3RFSHQ3V1FpR05QSlo2TVBic2o
20     3R0hDaUZOSmhST0l2ZXdabHVzVEdBcGdMRWpiS2lXME5jOWV2WTFJckFGUXI
21     2WHpjjeTZ5ZXPXMUJvSXlpdMzZlZjZlZjZlZjZlZjZlZjZlZjZlZjZlZjZl
22     NU0QrWDFQ1E9PQ%3D%3D--
23     ab7943f35d872df42bb86b8086b5cec50d3ef0ce"
24 }
```

```
9
10 post_data = {
11     "commit": "Sign in",
12     "utf8": "✓",
13     "authenticity_token":
14     "UFKuL5RE8DTUXAc0cddcawtX3gWADuVQInNPqIacBfESUsUZwZ8jNQ24sYp
15     QVHS6vqFX1ci9FqeTV9aZ+wqa+Q==",
16     "login": "czwspider",
17     "password": "qwer1234"
18 }
19
20 response = requests.post(login_url, post_data,
21 headers=headers)
22
23 print(response.status_code)
24
25 with open("github_login.html", "wb") as f:
26     f.write(response.content)
```

小栗子

通过requests获取新浪首页

```
1 import requests
2 response = requests.get("http://www.sina.com")
3 print(response.request.headers)
4 print(response.content.decode())
```

结果

```
1 {'User-Agent': 'python-requests/2.12.4', 'Accept-Encoding':  
  'gzip, deflate', 'Accept': '/*/*', 'Connection': 'keep-  
  alive'}  
2 <!DOCTYPE html>  
3 <!-- [ published at 2017-06-09 15:15:23 ] -->  
4 <html>  
5 <head>  
6     <meta http-equiv="Content-type" content="text/html;  
  charset=utf-8" />  
7     <meta http-equiv="X-UA-Compatible" content="IE=edge" />  
8     <title>新浪首页</title>  
9     <meta name="keywords" content="新浪,新浪  
  网,SINA,sina,sina.com.cn,新浪首页,门户,资讯" />  
10     ...
```

```
1 import requests  
2 response = requests.get("http://www.sina.com")  
3 print(response.request.headers)  
4 print(response.text)
```

结果

```

1 {'User-Agent': 'python-requests/2.12.4', 'Accept-Encoding':
  'gzip, deflate', 'Accept': '/*/*', 'Connection': 'keep-
  alive'}
2 <!DOCTYPE html>
3 <!-- [ published at 2017-06-09 15:18:10 ] -->
4 <html>
5 <head>
6     <meta http-equiv="Content-type" content="text/html;
  charset=utf-8" />
7     <meta http-equiv="X-UA-Compatible" content="IE=edge" />
8     <title>æ-°æµ³é|-éµ</title>
9     <meta name="keywords" content="æ-°æµ³,æ-
  °æµ³ç½‘,SINA,sina,sina.com.cn,æ-°æµ³é|-éµ,é-“æ^·,èµ,,è®-” />
10    <meta name="description" content="æ-°æµ³ç½‘ä,°å...
  “ç·fç”“æ^·24å°·æ-¶æ··ä¾>å...“é·çå·šæ-¶çš,,ä,·æ-¶èµ,,è®-“i¼åå†...
  å®¹è|†ç>-å>½å†...åæ-ç··å·‘æ-°é
  -»ä°<ä»¶ä€·ä½““å·>èµ>ä°<ä€·å·±ä¹·æ-
  ¶å°šä€·ä°šä,šèµ,,è®-ä€·å®žç”“ä¿µæ·ç·%i¼åè¾æææ-°é
  -»ä€·ä½“è,²ä€·å·±ä¹·ä€·è
  ´çç»·ä€·çš‘æšæä€·æ^¿ä°šä€·æ±½è½|ç·%30åæšä,åå†...
  å®¹éç‘é·“i¼åå·ææ-
  ¶å¼åè®¼å·šå®çä€·èš†éç‘ä€·è®å·>ç·%è†æç”±ä°’åš“ä°ææµ·ç®é-´ä€
  ,” />
11    <link rel="mask-icon" sizes="any"
  href="//www.sina.com.cn/favicon.svg" color="red">

```

产生问题的原因分析

1. requests默认自带的Accept-Encoding导致或者新浪默认发送的就是压缩之后的网页
2. 但是为什么content.decode()没有问题，因为requests，自带解压压缩网页的功能
3. 当收到一个响应时，Requests 会猜测响应的编码方式，由于在你调用 response.text 方法时对响应进行解码。Requests 首先在 HTTP 头部检测是否存在指定的编码方式，如果不存在，则会使用 chardet.detect 来尝试猜测编码方式（存在误差）

4. 更推荐使用response.content.deocde()

Cookies 和 Session

Cookies

如果一个响应中包含了cookie，那么我们可以利用 cookies参数拿到：

```
1
2 import requests
3
4 response = requests.get("http://www.baidu.com/")
5
6 # 返回CookieJar对象:
7 cookiejar = response.cookies
8
9 # 将CookieJar转为字典:
10 cookiedict = requests.utils.dict_from_cookiejar(cookiejar)
11
12 print (cookiejar)
13
14 print (cookiedict)
```

运行结果：

```
1 <RequestsCookieJar[<Cookie BDORZ=27315 for .baidu.com/>]>
2
3 {'BDORZ': '27315'}
```

实现github登录并跳转

```
1
2 #!/usr/bin/python
3 # -*- coding: utf-8 -*-
4
```



```
5 # 模拟 github 登录
6
7 # 0. 分析爬虫
8 '''
9 登录地址
10 确定请求地址 https://github.com/session
11 确定请求方式 POST
12 确定请求的内容
13 commit: Sign in
14 utf8: ✓
15 authenticity_token:
    FKPt8/jlSD6VqqKbJqQUylCZaArCLMEhyIYWtA12LSzK47nyaPOs8IoIZ04o
    5AGJiQIc04jX9b0lsWuETzc8+g==
16 login: czwspider
17 password: qwer1234
18 当请求成功 获取 页面的 Cookie 值并且，保存下来以后发送的请求都携
    带这个Cookie值
19 '''
20
21 import requests
22 import requests.utils
23
24 # 1. 定义请求参数
25 login_url = "https://github.com/session"
26
27 headers = {
28     "User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X
    10_13_4) AppleWebKit/537.36 (KHTML, like Gecko)
    Chrome/66.0.3359.170 Safari/537.36",
```

```

29     "Cookie": "_ga=GA1.2.1855430798.1461857641;
    _octo=GH1.1.783519559.1525492869; tz=Asia%2FShanghai;
    logged_in=no; _gat=1;
    _gh_sess=TTZZNnVPOTRwTlI2YW1UOWFqTUZwa1JHVENFOVpjYWxoekJLbS9
    GTGFQqkR40DZVSTIrV1NTWUFQS0dETUhURmg4YVR3bGxjV1hJT2dMZ2NIOEp
    tZlREUUXrOXV0eG1EcG5kUi9adUZQamFpWmxmMHVhY04vckJVWkxXbkNVa3Z
    HNCtjekF2WEVnMGVzaElqMnBpVUIvVGVSZJmdjNQMFxYWpONE1HMks2eDh
    pVzZ0Wk9ZQUZLMVJ0OTRsYXJWUjV6VUNmRFhyaH1YczUzdUozSWR5M210akh
    3dkcvaXRhY2ZmanRNRC9IbElUMm50SmkzVDhtbEwwSEdGWFMvd0xySWIxcFR
    rbGZ0RDQwQit5eGUvaGZKdWp5U1dYSnZ0VzRRb1FqUTHZZXlreTBNU1RicUh
    heGJGQjFRcDlnN1N2b2RXRXRsT21lRFB5Q3RF5HQ3V1FpR05QSlo2TVBic2o
    3R0hDaUZOSmhST0l2ZXdaHVzVEdBcGdMRWpiS2lXME5jOWV2WTFJckFGUXI
    2WHpjeTZ5ZXpXMUJvSXlpdMRpZWNONUhFejFUMD0tLVhyaTZJdlZBMudNwnB
    NU0QrWDFQ1E9PQ%3D%3D--
    ab7943f35d872df42bb86b8086b5cec50d3ef0ce"
30 }
31 post_data = {
32     "commit": "Sign in",
33     "utf8": "✓",
34     "authenticity_token":
    "UFKuL5RE8DTUXAc0cddcawtX3gWADuVQInNPqIacBfESUsUZwZ8jNQ24sYp
    QVHS6vqFXlci9FqeTV9aZ+wqa+Q==",
35     "login": "czwspider",
36     "password": "qwer1234"
37 }
38
39
40 # 获取session对象通过session对象进行请求，session对象的作用就是
    自动记录Cookies值，代码中不需要关心
41 session = requests.session()
42 # 2. 发送请求获取响应
43 response = session.post(login_url, post_data,
    headers=headers)
44
45
46 # 3. 从响应中获取Cookies值并且在以后的通讯中都使用这个Cookies值

```

```
47 # requests.utils.dict_from_cookiejar 用于把Cookie对象转换成词
   # 典对象
48 #
   print(requests.utils.dict_from_cookiejar(response.cookies))
49
50 # 4. 爬取
51 setting_headers = {
52     "User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X
   10_13_4) AppleWebKit/537.36 (KHTML, like Gecko)
   Chrome/66.0.3359.170 Safari/537.36",
53 }
54 setting_url = "https://github.com/settings/profile"
55 settting_resp = session.get(setting_url,
   headers=setting_headers)
56
57 # print(response.content.decode('utf-8'))
58 print(settting_resp.status_code)
59
60 with open('github_login_03.html', 'wb') as f:
61     f.write(settting_resp.content)
62
63
```

练一练

- 使用requests爬取任意输入的百度贴吧的网页，并保存到本地