

# 第 129 天：爬取微信公众号文章内容

原创 某某白米饭 Python技术 2月18日

有时候我们遇到一个好的公众号，里面的每篇都是值得反复阅读的，这时就可以使用公众号爬虫将内容抓取保存下来慢慢赏析。

## 安装 Fiddler

Fiddler 的下载地址为：<https://www.telerik.com/download/fiddler>，安装好之后，确保手机和电脑的网络为同一个局域网。

## Finddler 的配置

点击 Tools >> Options >> Connections 面板，参考下图配置，Fiddler 的默认端口为 8888，如果 8888 端口被占用了，可修改为其他端口。

点击 Tools >> Options >> HTTPS 面板，参考下图配置

## Android 手机配置

进入 WLAN 设置，选择当前所在局域网的 WIFI 设置，代理设置为手动，代理服务器主机名为 Finddler 中 右上角 Online 点击显示，端口号为 8888。

在手机浏览器中访问配置的地址：<http://ip:8888>，当显示 Fiddler Echo Service，则配置手机成功。

Finddler 为了拦截 HTTPS 请求，手机中必须安装 CA 证书，在 <http://ip:8888> 也中点击 FiddlerRoot certificate，下载并安装证书。此时配置工作全部完成。

## 微信历史页面

以【腾旭大申网】为例，点击【上海新闻】菜单的二级菜单【历史消息】。

观察 Fiddler 的变化，此时在左侧窗口中会陆续出现多个 URL 连接地址，这个就是 Fiddler 拦截的 Android 请求。

1. Result: 服务器的响应结果
2. Protocol: 请求协议, 微信协议都是 HTTPS 所以需要在手机端和PC端安装证书
3. HOST: 主机名
4. URL: URL 地址

其中有一条以 `https://mp.weixin.qq.com/mp/profile_ext?action=home...` 开头的URL就是我们需要的。点击 右侧 Inspectors 面板, 再点击下面的 Headers 和 WebView 面板, 会出现如下图样

#### Headers 面板

1. Request Headers: 请求行, 里面有请求方式、请求地址、请求协议等待
2. Client、Cookies: 请求头

#### WebView 面板

WebView 面板显示的是服务器返回的 HTML 代码渲染后的结果, Textview 面板则显示的为服务器返回的 HTML 源代码

### 抓取历史页面

在上一节中公众号消息历史页面已经可以显示在 Fiddler 的 WebView 面板了, 这一节则使用 Python 抓取历史页面。创建一个名为 `wxcrawler.py` 的脚本, 抓取页面我们需要 URL 地址和 HEADER 请求头, 直接从 Fiddler 中拷贝

#### 把 header 转换为 Json

```
1 # coding:utf-8
2 import requests
3
4 class WxCrawler(object):
5
6     # 复制出来的 Headers, 注意这个 x-wechat-key, 有时间限制, 会过期。当返回的内容出现 验证 的情况, 就需要换 x-
7     headers = """Connection: keep-alive
8         x-wechat-uin: MTY4MTI3NDIxNg%3D%3D
9         x-wechat-key: 5ab2dd82e79bc5343ac5fb7fd20d72509db0ee1772b1043c894b24d441af288ae942feb4c4fb4d2
10        Upgrade-Insecure-Requests: 1
11        User-Agent: Mozilla/5.0 (Linux; Android 10; GM1900 Build/QKQ1.190716.003; wv) AppleWebKit/53
12        Accept: text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,image/wxpica,image/a
13        Accept-Encoding: gzip, deflate
14        Accept-Language: zh-CN,en-US;q=0.9
15        Cookie: wxuin=1681274216; devicetype=android-29; version=27000933; lang=zh_CN; pass_ticket=J
```

```

15         X-Requested-With: com.tencent.mm""
16
17     url = "https://mp.weixin.qq.com/mp/profile_ext?action=home&__biz=MjEwNDI4NTA2MQ==&scene=123&dev
18
19
20     # 将 Headers 转换为 字典
21     def header_to_dict(self):
22         headers = self.headers.split("\n")
23         headers_dict = dict()
24         for h in headers:
25             k,v = h.split(":")
26             headers_dict[k.strip()] = v.strip()
27         return headers_dict;
28
29
30     def run(self):
31         headers = self.header_to_dict()
32         response = requests.get(self.url, headers=headers, verify=False)
33
34         print(response.text)
35
36
37 if __name__ == "__main__":
38
39     wx = WxCrawler()
40     wx.run()
41

```

下图就是打印在控制台的内容，其中在 JavaScript 中 变量 msgList 的值就是需要的内容

接下来就是提取 msgList 内容，使用正则表达式提取内容，返回一个文章列表

```

1 import re
2 import html
3 import json
4
5 def article_list(self, context):
6     rex = "msgList = '({.*?})'"
7     pattern = re.compile(pattern=rex, flags=re.S)
8     match = pattern.search(context)
9     if match:
10         data = match.group(1)
11         data = html.unescape(data)
12         data = json.loads(data)
13         articles = data.get("list")
14         return articles

```

下面就是解析 msgList 的结果

1. title: 文章标题
2. content\_url: 文章链接
3. source\_url: 原文链接, 有可能为空
4. digest: 摘要
5. cover: 封面图
6. datetime: 推送时间

其他的内容保存在 multi\_app\_msg\_item\_list 中

```
1  {'comm_msg_info':
2    {
3      'id': 1000033457,
4      'type': 49,
5      'datetime': 1575101627,
6      'fakeid': '2104285061',
7      'status': 2,
8      'content': ''
9    },
10   'app_msg_ext_info':
11     {
12       'title': '快查手机! 5000多张人脸照正被贱卖, 数据曝光令人触目惊心!',
13       'digest': '谁有权收集人脸信息?',
14       'content': '',
15       'fileid': 0,
16       'content_url': 'http:\\\\mp.weixin.qq.com\\s?__biz=MjEwNDI4NTA2MQ==&mid=2651824634&',
17       'source_url': '',
18       'cover': 'http:\\\\mmbiz.qpic.cn\\mmbiz_jpg\\G8vkERUJibkstwkIvXB960sM0yQdYF2x2qibTxAIq2e',
19       'subtype': 9,
20       'is_multi': 1,
21       'multi_app_msg_item_list':
22         [{
23           'title': '先有鸡还是先有蛋? 6.1亿年前的胚胎化石揭晓了',
24           'digest': '解决了困扰大申君20多年的问题',
25           'content': '',
26           'fileid': 0,
27           'content_url': 'http:\\\\mp.weixin.qq.com\\s?__biz=MjEwNDI4NTA2MQ==&mid=26518',
28           'source_url': '',
29           'cover': 'http:\\\\mmbiz.qpic.cn\\mmbiz_jpg\\yl6JkZAE3S92BESibpZgTPE1BcBhSLiaG0g',
30           'author': '',
31           'copyright_stat': 100,
32           'del_flag': 1,
33           'item_show_type': 0,
34           'audio_fileid': 0,
35           'duration': 0,
36           'play_url': '',
37           'malicious_title_reason_id': 0
```

```

37         'malicious_title_reason_id': 0,
38         'malicious_content_type': 0
39     },
40     {
41         'title': '外交部惊现“李佳琦”！网友直呼：“OMG被种草了！”',
42         'digest': '种草了！',
43         'content': '', ...}
44     ...]

```

## 抓取单个页面

在上节中我们可以得到 `app_msg_ext_info` 中的 `content_url` 地址了，这是需要从 `comm_msg_info` 这个不规则的 Json 中取出。这是需要使用 `demjson` 模块补全不规则的 `comm_msg_info`。

安装 `demjson` 模块

```
1 pip3 install demjson
```

```

1 import demjson
2
3 # 获取单个文章的URL
4 content_url_array = []
5
6 def content_url(self, articles):
7     content_url = []
8     for a in articles:
9         a = str(a).replace("\\/", "/")
10        a = demjson.decode(a)
11        content_url_array.append(a['app_msg_ext_info']['content_url'])
12        # 取更多的
13        for multi in a['app_msg_ext_info']['multi_app_msg_item_list']:
14            self.content_url_array.append(multi['content_url'])
15    return content_url

```

获取到单个文章的地址之后，使用 `requests.get()` 函数取得 HTML 页面并解析

```

1
2 # 解析单个文章
3 def parse_article(self, headers, content_url):
4     for i in content_url:
5         content_response = requests.get(i, headers=headers, verify=False)
6         with open("wx.html", "wb") as f:
7             f.write(content_response.content)
8         html = open("wx.html", encoding="utf-8").read()
9         soup_body = BeautifulSoup(html, "html.parser")

```

```
10     context = soup_body.find('div', id = 'js_content').text.strip()
11     print(context)
```

## 所有历史文章

把历史消息往下滑动时出现了正在加载中..., 这是公众号的历史消息正在翻页, 在 Fiddler 中查看得知, 公众号请求的地址为 [https://mp.weixin.qq.com/mp/profile\\_ext?action=getmsg&\\_biz...](https://mp.weixin.qq.com/mp/profile_ext?action=getmsg&_biz...)

翻页请求地址返回结果, 一般可以分析出

1. ret: 是否成功, 0为成功
2. msg\_count: 每页的条数
3. can\_msg\_continue: 是否继续翻页, 1为继续翻页
4. general\_msg\_list: 数据, 包含了标题、文章地址等信息

```
1 def page(self, headers):
2     response = requests.get(self.page_url, headers=headers, verify=False)
3     result = response.json()
4     if result.get("ret") == 0:
5         msg_list = result.get("general_msg_list")
6         msg_list = demjson.decode(msg_list)
7         self.content_url(msg_list["list"])
8         #递归
9         self.page(headers)
10    else:
11        print("无法获取内容")
```

## 总结

到这里已经爬取到了公众号的内容, 但是单个文章的阅读数和在看数还未爬取。思考一下, 这些内容该如何爬取?

示例代码: Python-100-days

## 系列文章

第128天: Seaborn-可视化数据集的分布

第127天: Seaborn-可视化分类数据

第126天: Seaborn-可视化统计关系

第125天: Flask 项目结构

第124天: Web 开发 Django 模板

第123天: Web 开发 Django 管理工具

第122天: Flask 单元测试

第121天: 机器学习之决策树

从 0 学习 Python 0 - 120 大合集总结

**PS:** 公号内回复: Python, 即可进入Python 新手学习交流群, 一起**100天计划!**

-END-

**Python 技术**

关于 Python 都在这里

---