

面向知乎 QAWEB 的网络爬虫设计与实现

游浩然
黎张帆
郭金城

u201515429 u201514574 u201511174 @hust.edu.cn

January 2, 2018

Table of Contents

I. Introduction

II. Method

III. Experiment

项目描述

- 爬虫控制端：启动爬虫，停止爬虫，监视爬虫的运行情况
- 爬虫运行模块：包含三个小模块，URL 管理器、网页下载器、网页解析器
 - URL 管理器：对需要爬取的 URL 和已经爬取过的 URL 进行管理，可以从其中取出待爬取的 URL 传递给网页下载器。
 - 网页下载器：网页下载器将 URL 指定的网页下载下来，存储成一个字符串，传递给网页解析器。
 - 网页解析器：网页解析器解析传递的字符串，解析器不仅可以解析出需要爬取的数据，而且还可以解析出每一个网页指向其他网页的 URL，这些 URL 被解析出来会补充进 URL 管理器。
- 数据输出模块：存储爬取的数据

开发平台

- 硬件平台：个人 PC，可连接 WEB 网络
- 软件环境：
 - 操作系统：WINDOWS 10 & Ubuntu
 - IDE：PyCharm Community Edition 2017.2.3
 - Python 版本：Python 3.6.3 |Anaconda, Inc.| (default, Oct 15 2017, 03:27:45)
- 项目所需库
 - Wxpython：Python GUI 库
 - Selenium：Python 模拟浏览器操作库
 - Phantom：基于 selenium 的 Python 工业化模拟浏览器操作库，需设置运行路径或将 bin 文件夹放入系统路径中。

任务分工

姓名	负责模块	主要任务
黎张帆	爬虫控制器	启动、停止、监视爬虫的运行情况；实现多线程爬取。GUI 设计
游浩然	爬虫运行模块	知乎登陆；URL 管理器、网页下载器、网页解析器；动态加载
郭金城	数据输出模块	通过网页标签抓取问题 URL 下的相关信息；动态加载

Table of Contents

I. Introduction

II. Method

- 整体框架设计
- 爬虫控制模块
- 爬虫运行模块
- 数据输出模块

III. Experiment

框架

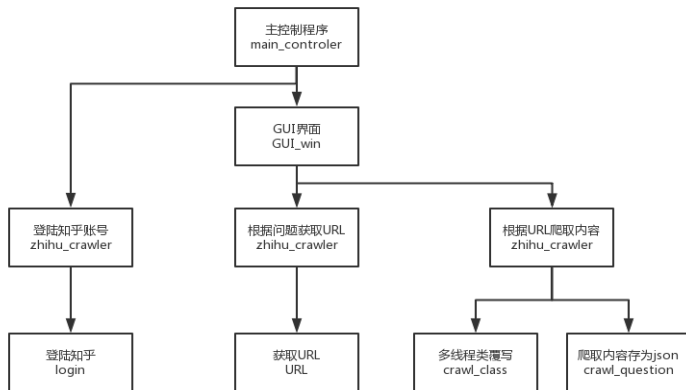


Figure 1: 整体框架设计

GUI 使用说明

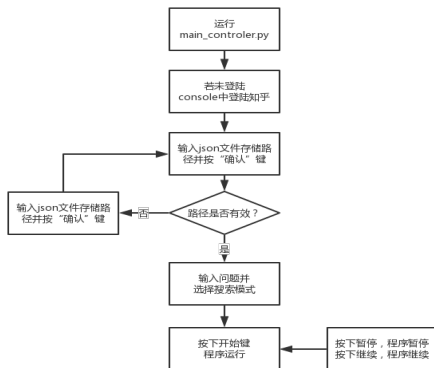


Figure 2: GUI 使用说明



Figure 3: GUI 界面

GUI 设计

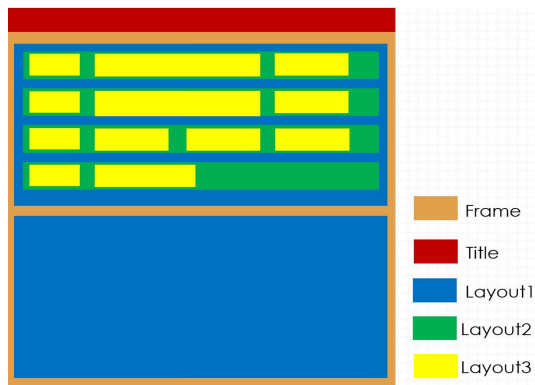


Figure 4: GUI 设计框图



Figure 5: GUI 界面

主程序流程

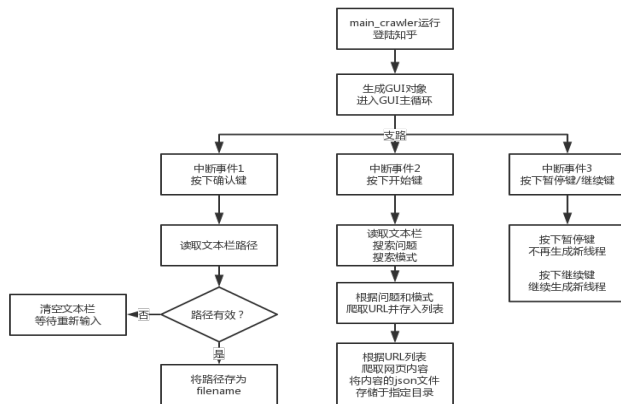


Figure 6: 主程序流程图

多线程控制

Thread 子类 Crawler

Class Crawler

Parent: threading.Thread

Method:

__init__(self, funcName, args=)#覆写父类方法

run(self)#覆写父类方法

_run(self)

Figure 7: Crawl

多线程控制

```
def run(self):
    try:
        self._run()
    except Exception as e:
        self.exitcode = 1
        self.exception = e
        self.exc_traceback = ''.join(traceback.format_exception(*sys.exc_info()))
        print("URL Crawler Exception: ", self.args[0])
        print(" self.exc_traceback")

def _run(self):
    try:
        self.funcName(self.args[0],self.args[1])
    except Exception as e:
        raise e
```

Figure 8: 覆写 run 方法

爬虫运行模块

本模块主要实现了爬虫的基本运行，分为两个主要部分：

- 知乎登陆
- 建立问题 URL 库
 - 针对具体问题 (Question) 搜索知乎的 “综合” 模块
 - 针对具体话题 (Topic) 搜索知乎的 “话题” 模块

其中对问题 (Question) 的爬取采用 PhantomJS 库来实现动态加载，对话题 (Topic) 的爬取采用模拟翻页来实现动态加载。

爬虫运行框图

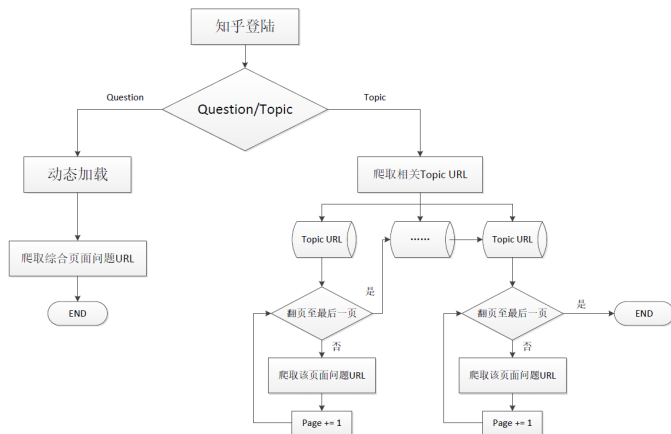


Figure 9: 爬虫运行框图

知乎登陆

使用手机的代理 (User-Agent) 来模拟登陆

```
1 headers = {  
2     "Host": "www.zhihu.com",  
3     "Referer" : url,  
4     'X-Requested-With' : 'XMLHttpRequest',  
5     'User-Agent' : 'Mozilla/5.0 (Linux; Android 6.0; Nexus 5 Build/MRA58N)  
6                   AppleWebKit/537.36 (KHTML, like Gecko) Chrome/56.0.2924.87 Mobile Safari/537.36'  
7 }
```

Figure 10: headers 头文件信息

知乎登陆

采用验证码方式登陆，测试结果如下：

```
1 C:\Users\dell-pc\Anaconda3\python.exe C:/Users/dell-pc/Desktop/zhihu_crawler_py3/login.py
2 Cookie cannot be load!
3 your name: None
4 please input your account:*****
5 please input your secret:*****
6 please input the captcha:f8dk
7 the status code returned by server: 200
8 {'r': 0, 'msg': '登录成功'}
9 your name: <span class="name">游浩然</span>
```

Figure 11: login

网页下载和解析

网页下载函数即为下代码所示，使用 requests 库来抓取指定 URL 的内容：

```
1 def getHTMLText(url):  
2     try:  
3         r = requests.get(url, headers=headers)  
4         # r.status_code 200  
5         # r.encoding 'utf-8'  
6         return r.content  
7     except:  
8         return "Error"
```

Figure 12: 网页下载

网页下载和解析

网页解析函数即为下代码所示，使用 BeautifulSoup 库来抓取解析网页：

```
1 demo = demo_temp.pop()
2 soup = BeautifulSoup(demo, 'html.parser')
3 for meta in soup.find_all('meta', attrs={'itemprop': 'url'}):
4     url_temp = meta.get('content')
5     url_list.append(url_temp)
```

Figure 13: 网页解析

URL 爬取和管理

- 对知乎“综合”模块进行爬取
- 对知乎“话题”模块进行爬取

存储为 `urllist[]` 列表格式供数据爬取模块调用

Question 页面动态加载

- 问题的完整描述需要执行点击操作，使用 selenium 模拟
- 通过鼠标的下滑不断地动态加载，使用 selenium 模拟

Question 页面动态加载

```
1 # 点击操作
2 try:
3     click_btn = driver.find_element_by_xpath('//button[@class="Button QuestionRichText-more Button--plain"]')
4     ActionChains(driver).click(click_btn).perform()
5     time.sleep(0.5)
6     html_text = driver.page_source
7 except:
8     html_text = driver.page_source
9 # 下滑操作
10 while(True):
11     driver.execute_script('window.scrollTo(0,document.body.scrollHeight)')
12     time.sleep(0.8)
13     if(html_text == driver.page_source):
14         break
15     html_text = driver.page_source
```

Figure 14: selenium 模拟点击与下滑

Question 数据存储

效果展示

```

{"question_text": "", "question_comment": "19 条评论", "question_id": "62968639", "question_followers": "427", "question_title": "文字是否是抵制 AI 的最后一个阵地?", "answer_number": "95 个回答", "b",
{"answer_id": "252689132", "answer_comment": "7 条评论", "answer_author": "知乎用户", "answer_votes": "27", "answer_text": "文字不是, 但文学批评 (literacy criticism) 估计是 AI 第一个取代的领域"}
{"answer_id": "205125138", "answer_comment": "3 条评论", "answer_author": "吕学平", "answer_votes": "6", "answer_text": "文字的意义在于“记录”和“表达”AI 的记录功能不用怀疑, 强悍的信息检索与总结能力是远远在人"}
{"answer_id": "252710799", "answer_comment": "9 条评论", "answer_author": "Karl Ernst", "answer_votes": "27", "answer_text": "我知道 AI 搞文学创作肯定是有可能的。但是, 我想说的不是这个问题, 也不是这个“抵制"}
{"answer_id": "252231747", "answer_comment": "66 条评论", "answer_author": "墨略集团", "answer_votes": "89", "answer_text": "<img src='\"https://pic2.zhimg.com/50/v2-ca629e3baa845f880c24f7a6539e76"}
{"answer_id": "204802622", "answer_comment": "4 条评论", "answer_author": "Farhi AI", "answer_votes": "28", "answer_text": "AI 的最后一个障碍只会是我编译"}
{"answer_id": "204566871", "answer_comment": "83 条评论", "answer_author": "廖颖", "answer_votes": "314", "answer_text": "谢邀。比起下围棋, 作者的生活变量太多, 无法预料。曹雪芹早上多吃个馒头, 今晚就就"}
{"answer_id": "252739979", "answer_comment": "李怡", "answer_votes": "4", "answer_text": "我觉得单纯从文学水平上去聊 AI 文学毫无意义, 就算 AI 将来能一字不差写出一个能得诺贝尔的文字"}
{"answer_id": "204789891", "answer_comment": "8 条评论", "answer_author": "Sirius", "answer_votes": "17", "answer_text": "结论: 几乎不能创作纯粹的 AI 文学, 不是机器的原因, 是人的原因。先看文学是怎么来的。"}
{"answer_id": "204659547", "answer_comment": "36 条评论", "answer_author": "五伯击", "answer_votes": "101", "answer_text": "AI 早就开始文学创作, 而且还差点荣获了 2016 年“星新一文学奖”(以日本著名科幻小说家"}
{"answer_id": "252725645", "answer_comment": "12 条评论", "answer_author": "理想的蘑菇", "answer_votes": "45", "answer_text": "别的不说, 现代诗肯定挡不住而且不是在来就聊四在已经被干掉了甚至都没有使用 ai"}
{"answer_id": "204796511", "answer_comment": "9 条评论", "answer_author": "冬黑", "answer_votes": "25", "answer_text": "在 2014 年翻一下知乎里围棋相关的回答, 基本上没有认为 AI 可以达到职业顶尖水平的。而现在"}
{"answer_id": "252748396", "answer_comment": "添加评论", "answer_author": "王汪汪", "answer_votes": "0", "answer_text": "文学也严重文学和通俗文学, 那些流布网上的网文甚至是畅销一时的通俗小说, 有朝一日 ai"}
{"answer_id": "204678305", "answer_comment": "23 条评论", "answer_author": "三当家的", "answer_votes": "64", "answer_text": "阿尔法。先生, 我在这, 我是不是暴死了? 我很遗憾, 是的, 先生。我死之前不可不给"}
{"answer_id": "253198697", "answer_comment": "添加评论", "answer_author": "文艺范托样作死君", "answer_votes": "0", "answer_text": "ai 可能永远搞不懂文学, 但当 ai 成为主流以后, 文学可能就没有存在的价值了"}
{"answer_id": "204934130", "answer_comment": "5 条评论", "answer_author": "谢丹", "answer_votes": "4", "answer_text": "不能。因为图灵判断, 一个 AI 写出的文学, 不需要超越人类才能存在, 只要大家读起来觉得不"}
{"answer_id": "204046842", "answer_comment": "12 条评论", "answer_author": "W3 小蘑菇", "answer_votes": "9", "answer_text": "<img src='\"https://pic4.zhimg.com/50/v2-cf6d9054c234c428c75d39bb788cd0e2e"}
{"answer_id": "204895251", "answer_comment": "1 条评论", "answer_author": "王玄豪", "answer_votes": "31", "answer_text": "<img data-rwheight='\"222\", src='\"https://pic3.zhimg.com/50/v2-348b256c"}
{"answer_id": "205080032", "answer_comment": "添加评论", "answer_author": "稻黍人", "answer_votes": "1", "answer_text": "回答这个问题前, 必须先回答另一个问题: 意识到底是什么? 这里的意识指的是物质的某种"}
{"answer_id": "204645686", "answer_comment": "1 条评论", "answer_author": "华说旧雨", "answer_votes": "7", "answer_text": "文学不只会由语言组成。写作的愉悦也是文学的一部分。"}
{"answer_id": "204989494", "answer_comment": "3 条评论", "answer_author": "咕咕咕", "answer_votes": "4", "answer_text": "以现代人的阅读习惯, 说不定这一块会比你想象的更快呢。这里有一位著名作家回答了, 我知"}
{"answer_id": "204259204", "answer_comment": "2 条评论", "answer_author": "Grotruprida", "answer_votes": "4", "answer_text": "微软小冰已经开始写诗了"}
{"answer_id": "205176195", "answer_comment": "添加评论", "answer_author": "Jawesome", "answer_votes": "3", "answer_text": "https://arxiv.org/pdf/1705.08807.pdf 这篇文里作者调查了 AI 研究者对未来自 AI 取代"}
{"answer_id": "204673094", "answer_comment": "1 条评论", "answer_author": "AABb", "answer_votes": "3", "answer_text": "还有打从东 50 一晚的房子 <img src='\"https://pic1.zhimg.com/50/v2-fc96b4a113da"}
{"answer_id": "204573387", "answer_comment": "添加评论", "answer_author": "黑果大树对", "answer_votes": "3", "answer_text": "不是还有音乐、绘画、雕塑等艺术另外, 我见过有人把艺术要的是人创造出来的东西, 而"}
{"answer_id": "204600313", "answer_comment": "添加评论", "answer_author": "匿名用户", "answer_votes": "4", "answer_text": "基础学科研究才是"}
{"answer_id": "205445454", "answer_comment": "1 条评论", "answer_author": "blank-shiki", "answer_votes": "2", "answer_text": "不好说啊。比如: 脑洞一下。根据人类的思维与行为也是按照有逻辑性的算法形成, 那 AI"}
{"answer_id": "204701484", "answer_comment": "4 条评论", "answer_author": "匿名用户", "answer_votes": "2", "answer_text": "还记得《锦绣未央》(原名《庶女有毒》) 是如何成书的吗? 而那不是 AI, 不过是区区书作"}
{"answer_id": "205012767", "answer_comment": "添加评论", "answer_author": "劳哈若", "answer_votes": "2", "answer_text": "以艺术文学家的视角, 怕不是早就被 AI 爆了十条街。至于应声文学家(哲学家, 科学家), 早就已"}
{"answer_id": "253874917", "answer_comment": "添加评论", "answer_author": "周别", "answer_votes": "1", "answer_text": "你们平时看到的部分新闻稿已经是 AI 自动撰写成自动变来的了。平时有感觉到吗? 国外好像多媒体都"}
{"answer_id": "252628580", "answer_comment": "添加评论", "answer_author": "匿名用户", "answer_votes": "1", "answer_text": "最后一个阵地。。。。把基础科学研究搞?"}
{"answer_id": "205591314", "answer_comment": "添加评论", "answer_author": "匿名用户", "answer_votes": "1", "answer_text": "人类写作的意义应该不只在吸引别人来看, 还可以自我愉悦, AI 应该不能时常靠自己从前"}
{"answer_id": "205168293", "answer_comment": "添加评论", "answer_author": "如是", "answer_votes": "1", "answer_text": "01. 这问题把文学和人工智能联合在一起了, 文学是经久不衰的, 人工智能又是时代的宠儿, 挺有前"}
{"answer_id": "204848234", "answer_comment": "添加评论", "answer_author": "泽鸿翔", "answer_votes": "1", "answer_text": "微软小冰不是一部诗吗? <img src='\"https://pic2.zhimg.com/50/v2-3232de7828"}
{"answer_id": "204648795", "answer_comment": "5 条评论", "answer_author": "Evangeliion", "answer_votes": "1", "answer_text": "文学算什么来抵制..."}

```

Figure 15: json 效果展示

Question 数据存储

```
{
  "question_followers": "15,436",
  "question_title": "如何在四小时内学会用 Ai 做 UI?",
  "answer_number": "86 个回答",
  "question_text": "整个问题的初衷在于看到了@黎敏 的回答, 产品经理新人是否有必要学习Photoshop? 我这个问题主要在于希望一些新人能速度的进入AI操作, 而不必花时间去找教程, 需要这个问题回答的人不少。再次谢谢 @黎敏。",
  "question_commet": "15 条评论",
  "brower_number": "989,872",
  "question_id": "21378838"
}
```

Figure 16: json 效果展示

```
{
  "answer_votes": "0",
  "answer_author": "mac mico",
  "answer_text": "工欲善其事必先利其器, 效率一定是首要的。UI设计之效率为王 - 设计与开发之效率 - 知乎专栏",
  "answer_id": "61575458",
  "answer_comment": "添加评论"
}
```

Figure 17: json 效果展示

Table of Contents

I. Introduction

II. Method

III. Experiment

- Thank You.