

# 一、系统要求

python3.4以上版本, 不支持python2.x

# 二、准备工具

requests、beautifulsoup 是爬虫两大神器，requests 用于网络请求，beautifulsoup 用于操作 html 数据。有了这两把梭子，干起活来利索。scrapy 这样的爬虫框架我们就不用了，这样的小程序派上它有点杀鸡用牛刀的意思。此外，既然是把 html 文件转为 pdf，那么也要有相应的库支持，wkhtmltopdf 就是一个非常的工具，它可以用适用于多平台的 html 到 pdf 的转换，pdfkit 是 wkhtmltopdf 的Python封装包。首先安装好下面的依赖包

```
1 pip install requests
2 pip install beautifulsoup4
3 pip install pdfkit
```

# 三、安装 wkhtmltopdf

Windows平台直接在 <http://wkhtmltopdf.org/downloads.html> 下载稳定版的 wkhtmltopdf 进行安装，安装完成之后把该程序的执行路径加入到系统环境 \$PATH 变量中，否则 pdfkit 找不到 wkhtmltopdf 就出现错误 "No wkhtmltopdf executable found"。Ubuntu 和 CentOS 可以直接用命令行进行安装

```
1 $ sudo apt-get install wkhtmltopdf # ubuntu
2 $ sudo yum install wkhtmltopdf    # centos
```

# 四、运行

```
1 python crawler.py
```

# 五、常见问题

1. SyntaxError: Missing parentheses in call to 'print'

beautifulsoup3不支持python2,所以下载beautifulsoup是要指定 beautifulsoup4

2. 如果是使用PyCharm开发, 那么运行的时候要在shell/cmd 窗口执行脚本, 直接在Pycharm中运行会找不

到 wkhtmltopdf命令

## 六、完整代码如下

```
1  # coding=utf-8
2  from __future__ import unicode_literals
3
4  import logging
5  import os
6  import re
7  import time
8
9  try:
10     from urllib.parse import urlparse # py3
11 except:
12     from urlparse import urlparse # py2
13
14 import pdfkit
15 import requests
16 from bs4 import BeautifulSoup
17
18 html_template = """
19 <!DOCTYPE html>
20 <html lang="en">
21 <head>
22     <meta charset="UTF-8">
23 </head>
24 <body>
25 {content}
26 </body>
27 </html>
28
29 """
30
31
32 class Crawler(object):
33     """
34     爬虫基类，所有爬虫都应该继承此类
35     """
36     name = None
37
```

```
38     def __init__(self, name, start_url):
39         """
40         初始化
41         :param name: 将要被保存为PDF的文件名称
42         :param start_url: 爬虫入口URL
43         """
44         self.name = name
45         self.start_url = start_url
46         self.domain = '{uri.scheme}://{uri.netloc}'.format(uri=urlparse(self.start_url))
47
48     @staticmethod
49     def request(url, **kwargs):
50         """
51         网络请求,返回response对象
52         :return:
53         """
54         response = requests.get(url, **kwargs)
55         return response
56
57     def parse_menu(self, response):
58         """
59         从response中解析出所有目录的URL链接
60         """
61         raise NotImplementedError
62
63     def parse_body(self, response):
64         """
65         解析正文,由子类实现
66         :param response: 爬虫返回的response对象
67         :return: 返回经过处理的html正文文本
68         """
69         raise NotImplementedError
70
71     def run(self):
72         start = time.time()
73         options = {
74             'page-size': 'Letter',
75             'margin-top': '0.75in',
76             'margin-right': '0.75in',
77             'margin-bottom': '0.75in',
78             'margin-left': '0.75in',
79             'encoding': "UTF-8",
```

```

80         'custom-header': [
81             ('Accept-Encoding', 'gzip')
82         ],
83         'cookie': [
84             ('cookie-name1', 'cookie-value1'),
85             ('cookie-name2', 'cookie-value2'),
86         ],
87         'outline-depth': 10,
88     }
89     htmls = []
90     for index, url in enumerate(self.parse_menu(self.request(self.start_url))):
91         html = self.parse_body(self.request(url))
92         f_name = ".".join([str(index), "html"])
93         with open(f_name, 'wb') as f:
94             f.write(html)
95         htmls.append(f_name)
96
97     pdfkit.from_file(htmls, self.name + ".pdf", options=options)
98     for html in htmls:
99         os.remove(html)
100     total_time = time.time() - start
101     print(u"总共耗时: %f 秒" % total_time)
102
103
104 class LiaoxuefengPythonCrawler(Crawler):
105     """
106     廖雪峰Python3教程
107     """
108
109     def parse_menu(self, response):
110         """
111         解析目录结构, 获取所有URL目录列表
112         :param response 爬虫返回的response对象
113         :return: url生成器
114         """
115         soup = BeautifulSoup(response.content, "html.parser")
116         menu_tag = soup.find_all(class_="uk-nav uk-nav-side")[1]
117         for li in menu_tag.find_all("li"):
118             url = li.a.get("href")
119             if not url.startswith("http"):
120                 url = "".join([self.domain, url]) # 补全为全路径
121             yield url

```

```
122
123     def parse_body(self, response):
124         """
125         解析正文
126         :param response: 爬虫返回的response对象
127         :return: 返回处理后的html文本
128         """
129         try:
130             soup = BeautifulSoup(response.content, 'html.parser')
131             body = soup.find_all(class_='x-wiki-content')[0]
132
133             # 加入标题，居中显示
134             title = soup.find('h4').get_text()
135             center_tag = soup.new_tag("center")
136             title_tag = soup.new_tag('h1')
137             title_tag.string = title
138             center_tag.insert(1, title_tag)
139             body.insert(1, center_tag)
140
141             html = str(body)
142             # body中的img标签的src相对路径的改成绝对路径
143             pattern = "<img .*?src=\"(.*)\">"
144
145             def func(m):
146                 if not m.group(2).startswith("http"):
147                     rtn = "".join([m.group(1), self.domain, m.group(2), m.group(3)])
148                     return rtn
149                 else:
150                     return "".join([m.group(1), m.group(2), m.group(3)])
151
152             html = re.compile(pattern).sub(func, html)
153             html = html_template.format(content=html)
154             html = html.encode("utf-8")
155             return html
156         except Exception as e:
157             logging.error("解析错误", exc_info=True)
158
159
160 if __name__ == '__main__':
161     start_url = "http://www.liaoxuefeng.com/wiki/0013739516305929606ddd18361248578c67b8067c8c017b000"
162
163     crawler = LiaoxuefengPythonCrawler("廖雪峰Git", start_url)
```

## 七、结果显示

