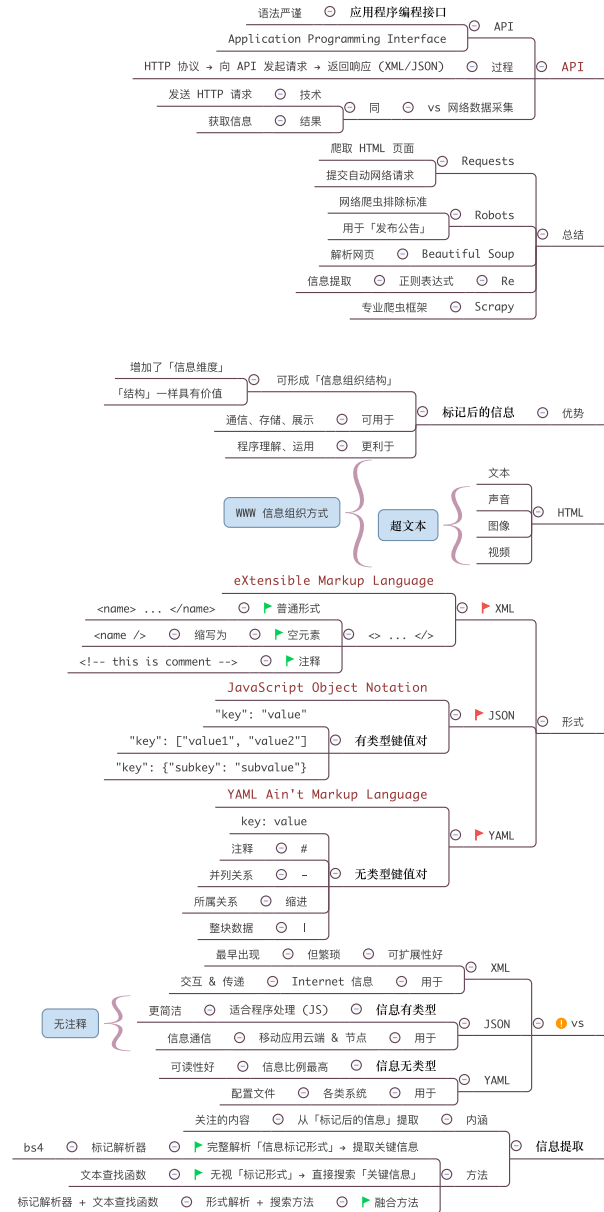
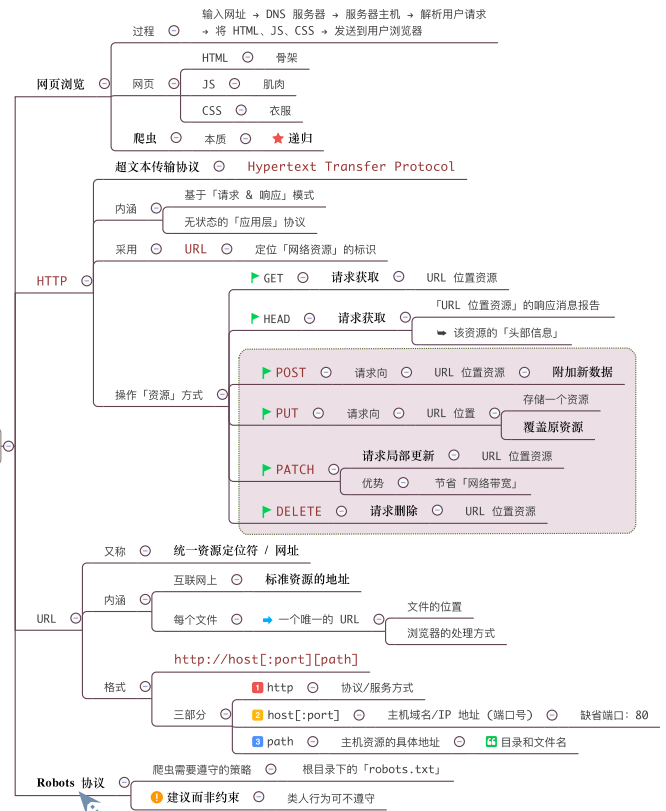


ZOE's MindMap
1-4 | v2.0 | 20180426

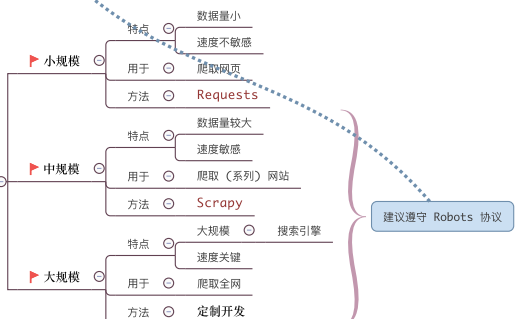


延伸知识

背景知识



网络爬虫尺寸



公众号: 数林觅风

ZOE's MindMap

2-4 | v2.0 | 20180426

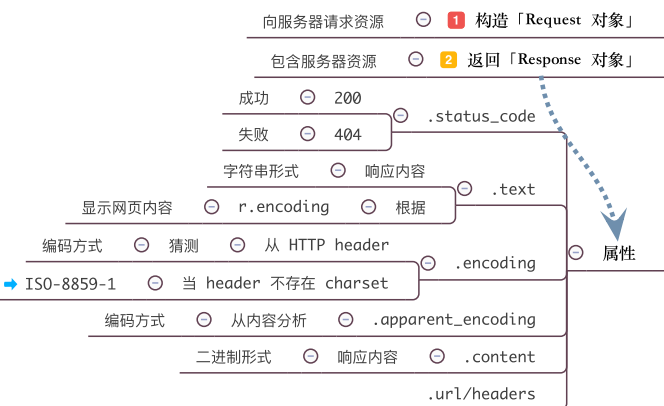
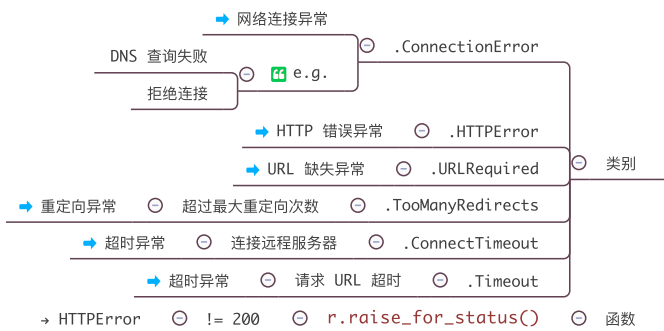
公众号: 数林觅风
<https://woaielf.github.io/>
2 Requests 库

异常

步骤

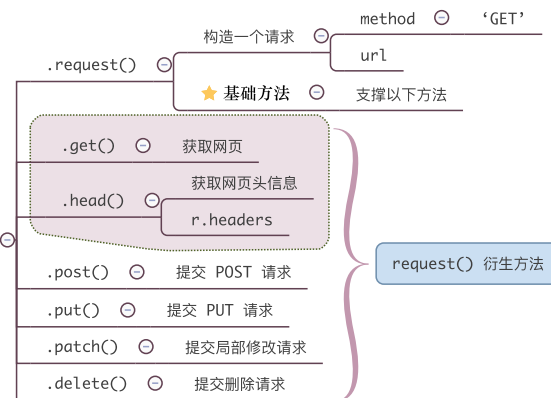
主要方法

控制访问参数

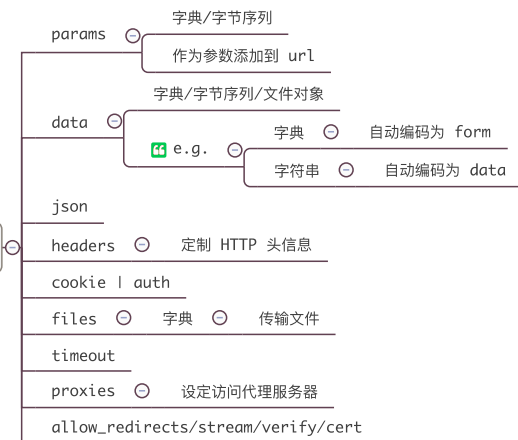


★ 通用代码框架

```
def getHTMLText(url):
    try:
        r = requests.get(url, timeout=30)
        r.raise_for_status()
        r.encoding = r.apparent_encoding
        return r.text
    except:
        return '产生异常'
```



request() 衍生方法



作为 Request 的内容



公众号: 数林觅风

ZOE's MindMap

3-4 | v2.0 | 20180426

公众号: 数林觅风

<https://woaielf.github.io/>

3 BeautifulSoup & urllib

查找方式

名称 & 属性 ○ find & find_all

基于「文档位置」

导航树

内容查找

基本方法

BeautifulSoup

HTML 解析库

- 安装 ○ pip install beautifulsoup4
- 引入 ○ from bs4 import BeautifulSoup
- 功能 ○ 解析、遍历、维护「标签树」
- 编码 ○ 任何 HTML 输入 + utf-8 编码

对象

- BeautifulSoup 对象
 - find & find_all
 - 调用「子标签」
- Tag 标签对象 ○ 获取
 - 标签里的文字
- NavigableString 对象 ○ 注释标签
- Comment 对象 ○ 注释标签

★ 优化显示

- 内涵 ○ 为 HTML 文本 <> ○ 增加 '\n'
- 方法
 - soup.prettify()
 - <tag>.prettify()
- print() 输出

★ 中文对齐

- 描述 ○ 中文字符宽度不够 ○ 用「西文字符」填充
- 解决 ○ 中文字符的空格 ○ chr(12288)
- e.g. ○ string.format(a, b, c, chr(12288))

```
soup = BeautifulSoup(demo, 'html.parser')
```

理解

- 三者等价 ○ ★ HTML 文档 ⇌ 标签树 ⇌ BeautifulSoup 类
- BeautifulSoup 类型 ○ 标签树的根节点

1 Python 自带解析器

- 'html.parser' ○ HTML
- 'lxml' ○ lxml 的 HTML
- 'xml' ○ lxml 的 XML

2 lxml

3 html5lib

- 'html5lib' ○ html5lib

Tag ○ soup.<tag> ○ 标签 ○ 基本信息组织单元

Name ○ <tag>.name ○ 名字 ○ 字符串

Attributes ○ <tag>.attrs ○ 属性 ○ 字典

NavigableString ○ <tag>.string ○ 非属性字符串

Comment ○ 清除所有标签

获取字符串 ○ .get_text() ○ 仅返回字符串

最后才用

soup.a ○ 只返回第一个

e.g. ○ soup.a ○ .name/.attrs/.string

soup.a.parent.name



公众号: 数林觅风

urllib2 ○ Python 2.x ○ 原为

标准库

urllib ○ 背景知识

.request

.parse

.error

from urllib.request import urlopen

from bs4 import BeautifulSoup

html = urlopen('...')

解析 ○ bsObj = BeautifulSoup(html.read(), 'html.parser')

links = bsObj.find_all('a')

for link in links:

if 'href' in link.attrs:

print('link.attrs[href]')

soup.body.h1 ○ 标题

名称 & 属性 ○ find & find_all

基于「文档位置」

导航树

内容查找

子节点的「列表」类型 ○ .contents

子节点的「迭代」类型 ○ .children

子孙节点的「迭代」类型 ○ .descendants

节点的父亲标签 ○ .parent

节点的先辈标签的「迭代」类型 ○ .parents

下一个平行节点 ○ .next_sibling

上一个平行节点 ○ .previous_sibling

后续所有平行节点 ○ .next_siblings

前续所有平行节点 ○ .previous_siblings

只能用于「循环遍历」

返回所有标签 ○ True

e.g. ○ 标签名称 ○ name=

re.compile('b')

{'class': 'green'} ○ 标签属性值 ○ attrs={}

< ... </> 中字符串 ○ 标签文本内容 ○ text=None

是否对子孙全部检索 ○ 默认 True ○ recursive=True

limit=None

keyword

<tag>(...) | soup(...) ○ ⇔

返回一个结果 ○ <>.find()

先辈 ○ <>.find_parents()/parent()

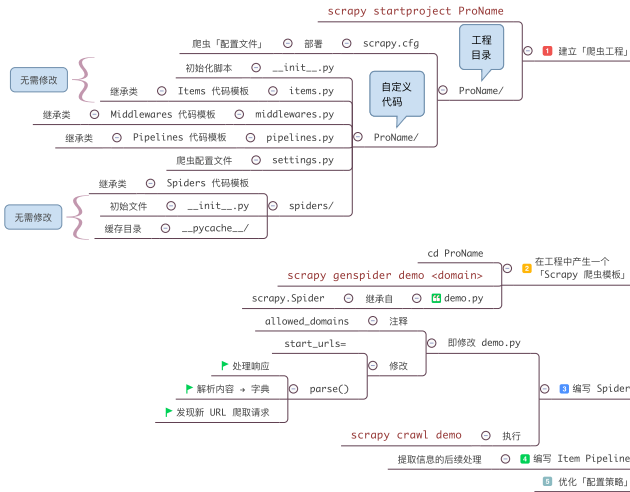
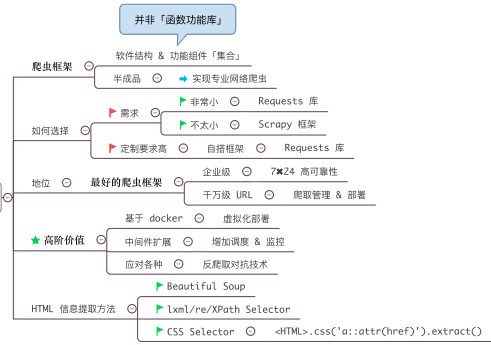
后续 ○ <>.find_next_siblings()/sibling()

前续 ○ <>.find_previous_siblings()/sibling()

Requests	Scrapy	同! !
页面级 功能库 并发性不足 性能差 重点: 页面下载 定制灵活	网站级 框架 并发性好 性能高 重点: 爬虫结构 一般定制 OK, 深度定制困难	页面请求 & 爬取 同为两条重要技术路线 requests-bs4-re scrapy 处理 js 提交表单 应对验证码 PhantomJS

Requests vs Scrapy

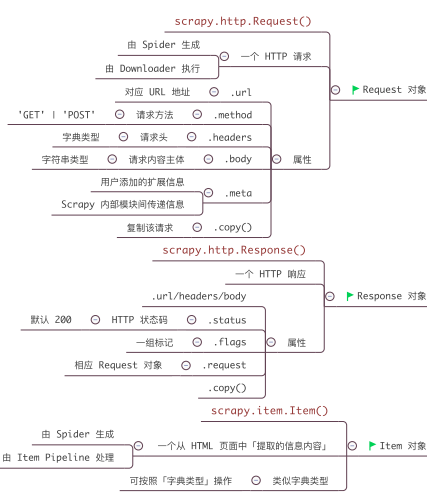
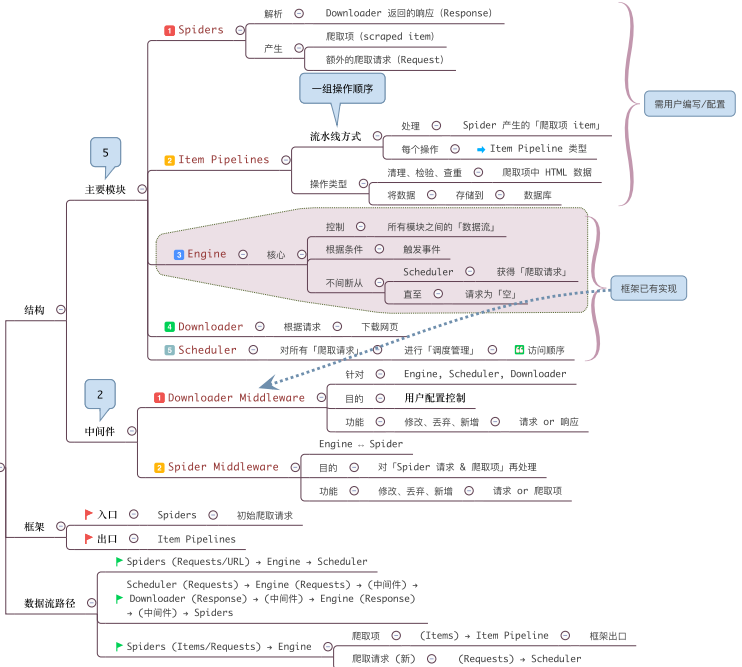
简介



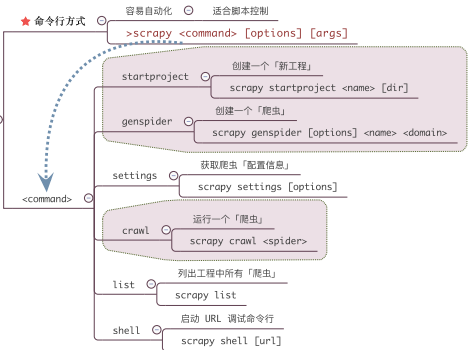
使用步骤

公众号: 数林梵风
<https://woaief.github.io/>
4 Scrapy 爬虫框架

“5+2” 结构



常用命令



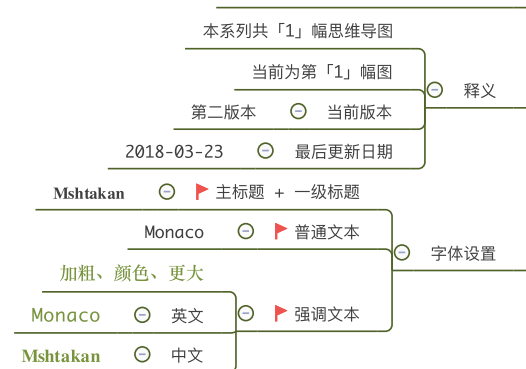
ZOE's MindMap

1-1 | v2.1 | 20180402



编号-系列总数 | 版本号 | 更新日期

1-1 | v2.0 | 20180323

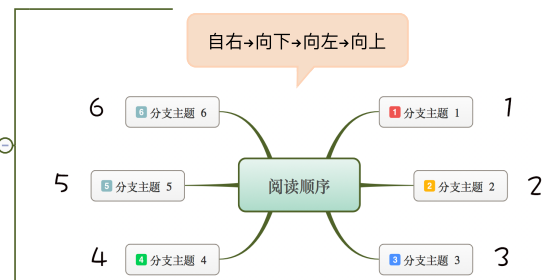


4 关于作者

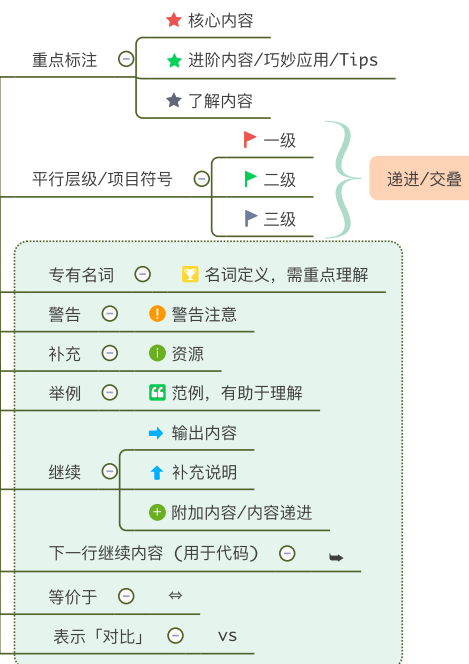
公众号：数林觅风
<https://woaielf.github.io/>
❤️ ZOE 思维导图规则

1 阅读顺序

顺时针方向！



2 ★ 图标含义



图例

- 1 编号1
- 2 编号2
- 3 编号3
- 4 编号4
- 一级
- 二级
- 三级
- 核心理解
- 进阶内容
- 简单了解
- 输出内容
- 补充前面
- 注意
- 附加/递进
- 资源
- 举例说明
- 爱心
- 名词定义



公众号：数林觅风



获取更多 ZOE 思维导图

❤️ <https://woaielf.github.io/>

👤 数林觅风 | zoemindmap

📺 @数林觅风

📖 ZOE 酱



思维导图 | 数据科学 | 编程语言 | 读书笔记 | 精进思考

