

目录

前言	1.1
PyQuery概述	1.2
PyQuery常见操作	1.3
find() vs items()	1.3.1
PyQuery语法	1.4
PyQuery实例	1.5
汽车之家品牌车型车系	1.5.1
附录	1.6
参考资料	1.6.1

HTML解析库Python版jQuery: PyQuery

- 最新版本: `v1.0`
- 更新时间: `20210414`

简介

介绍HTML解析领域的Python版的jQuery: PyQuery。先对PyQuery概述, 再介绍常见的基本操作, 包括如何获取元素的文本值, 如何获取带子元素的元素的文本值; 获取元素的属性, 如何用CSS选择器去选择定位查找元素; 对比了find和items的返回结果; 整理了CSS选择器的基本语法; 给出了具体的实际例子, 比如用PySpider抓取汽车之家之家的品牌和车型车系数据中如何使用PyQuery。

源码+浏览+下载

本书的各种源码、在线浏览地址、多种格式文件下载如下:

Gitbook源码

- [crifan/python_html_parse_pyquery: HTML解析库Python版jQuery: PyQuery](#)

如何使用此Gitbook源码去生成发布为电子书

详见: [crifan/gitbook_template: demo how to use crifan gitbook template and demo](#)

在线浏览

- [HTML解析库Python版jQuery: PyQuery book.crifan.com](#)
- [HTML解析库Python版jQuery: PyQuery crifan.github.io](#)

离线下载阅读

- [HTML解析库Python版jQuery: PyQuery PDF](#)
- [HTML解析库Python版jQuery: PyQuery ePub](#)
- [HTML解析库Python版jQuery: PyQuery Mobi](#)

版权说明

此电子书教程的全部内容，如无特别说明，均为本人原创和整理。其中部分内容参考自网络，均已备注了出处。如有发现侵犯您的版权，请通过邮箱联系我 [admin 艾特 crifan.com](mailto:admin@crifan.com)，我会尽快删除。谢谢合作。

鸣谢

感谢我的老婆陈雪的包容理解和悉心照料，才使得我 [crifan](#) 有更多精力去专注技术专研和整理归纳出这些电子书和技术教程，特此鸣谢。

更多其他电子书

本人 [crifan](#) 还写了其他 100+ 本电子书教程，感兴趣可移步至：

[crifan/crifan_ebook_readme](#): Crifan的电子书的使用说明

[crifan.com](#)，使用[署名4.0国际\(CC BY 4.0\)协议](#)发布 all right reserved,
powered by Gitbook最后更新：2021-04-14 19:50:55

PyQuery概述

- 概述
 - 作用
 - Python 版的 jQuery
 - 注: jQuery是 JavaScript 中的库, 可以 (很方便的) 操作 HTML (中的元素)
 - 使用
 - PySpider内部也用到了PyQuery
- 官网文档
 - 入口
 - <https://pythonhosted.org/pyquery/index.html>
 - Quickstart = 快速开始
 - <https://pythonhosted.org/pyquery/index.html#quickstart>
 - Attributes = 元素属性
 - <https://pythonhosted.org/pyquery/attributes.html>
 - CSS = 用CSS选择器选择元素
 - <https://pythonhosted.org/pyquery/css.html>
 - Manipulating = 操作元素
 - <https://pythonhosted.org/pyquery/manipulating.html>
 - Traversing = 元素遍历
 - <https://pythonhosted.org/pyquery/traversing.html>
 - API
 - <https://pythonhosted.org/pyquery/api.html>
 - Scraping = 抓取 (页面)
 - <https://pythonhosted.org/pyquery/scrap.html>
 - pyquery.ajax – PyQuery AJAX extension = AJAX扩展
 - <https://pythonhosted.org/pyquery/ajax.html>
 - Tips = 提示
 - <https://pythonhosted.org/pyquery/tips.html>
 - Testing = 测试
 - <https://pythonhosted.org/pyquery/testing.html>

crifan.com, 使用署名4.0国际(CC BY 4.0)协议发布 all right reserved,
powered by Gitbook最后更新: 2021-04-14 19:36:05

PyQuery常见操作

获取元素的文本值

举例：html

```
<p class="contributors">
  By
  <a href="/teachers/authors/judy-love.html" target="_self"
  ,
  <a href="/teachers/authors/julie-danneberg.html" target="
</p>
```

写法：

```
curHtmlElement.text()
```

举例：

```
authors = []
for eachAuthor in response.doc('p[class="contributors"] a[]
    print("eachAuthor=%s", eachAuthor)
    authorText = eachAuthor.text()
```

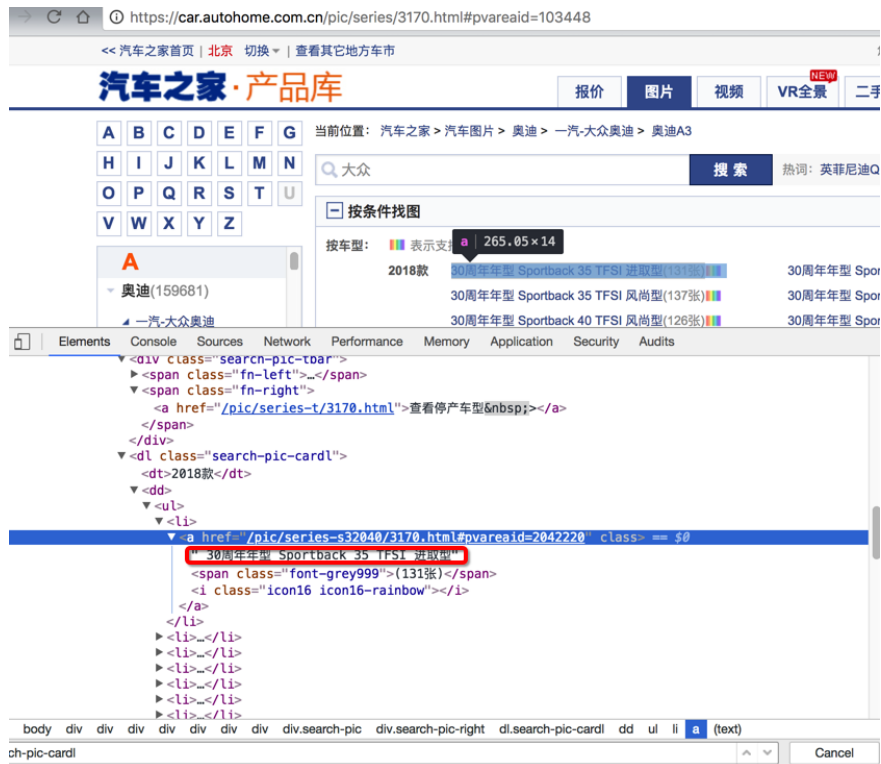
详见官网文档：

<https://pythonhosted.org/pyquery/api.html#pyquery.pyquery.PyQuery.text>

获取含有子元素的 a 的本身的 text

页面：

find() vs items()



html

```
<a href="/pic/series-s32040/3170.html#pvareaid=2042220" class="font-grey999">30周年年型 Sportback 35 TFSI 进取型</a>
  <span class="font-grey999">(131张)</span>
  <i class="icon16 icon16-rainbow"></i>
</a>
```

想要：获取a的本身的text

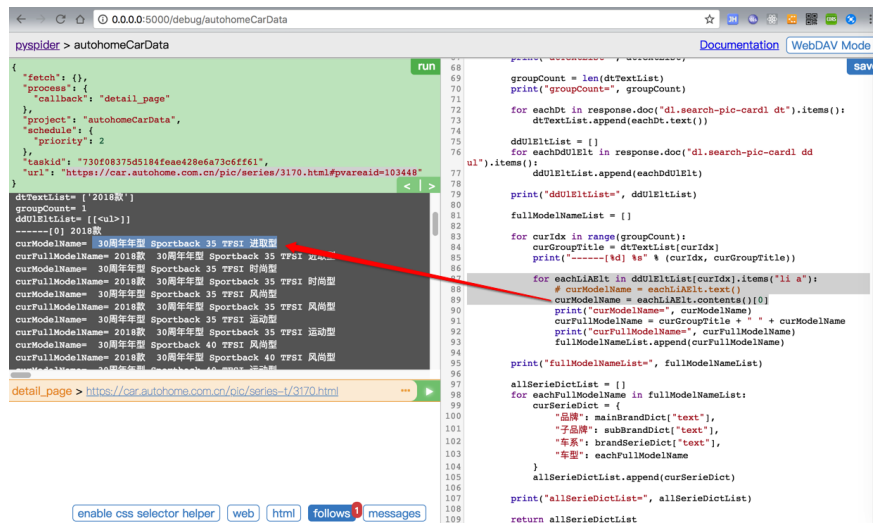
代码：

```
for eachLiAEl in ddUleltList[curIdx].items("li a"):
    curModelName = eachLiAEl.contents()[0]
```

输出值： 30周年年型 Sportback 35 TFSI 进取型

效果：

find() vs items()



- 官网文档
 - `pyquery.pyquery.PyQuery.contents`
 - Return contents (with text nodes)
 - 和之前用过的 `BeautifulSoup` 的 `contents` 是类似的
 - `BeautifulSoup.contents`
 - tag的 `.contents` 属性可以将tag的子节点以列表的方式输出

直接用text()会包含子元素的text值

如果直接用 `text()` :

```
for eachLiAElt in ddUlEltList[curIdx].items("li a"):
    curModelName = eachLiAElt.text()
```

是错误的写法。

因为会把 span中的 text，此处的（131张）也包括在内

获取元素的属性的值

语法：

- `attr`

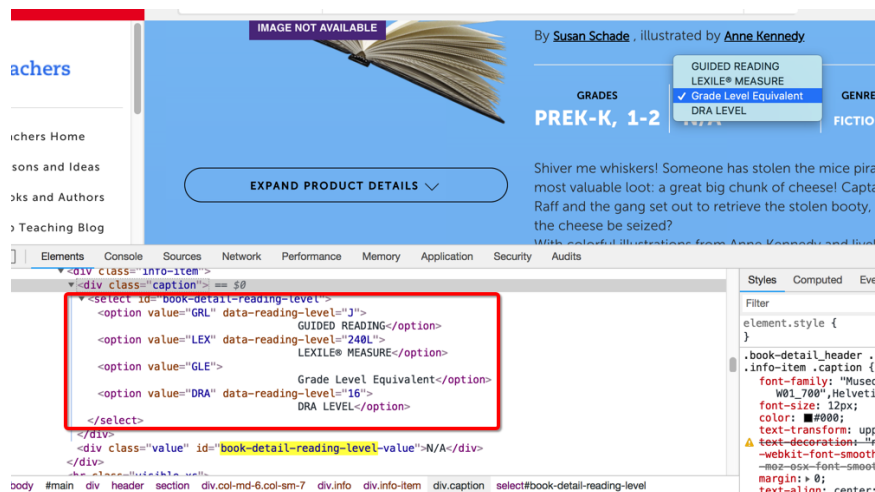
```
htmlElement.attr("property name")
```
- 如果 `property name` 是连接在一起的单词这种，比如 `xxx`，`xxx_yyy`，那么可以写成：

htmlElement.attr.property_name

find() vs items()

举例1

```
<div class="caption">
  <select id="book-detail-reading-level">
    <option value="GRL" data-reading-level="J">
      GUIDED READING</option>
    <option value="LEX" data-reading-level="240L">
      LEXILE® MEASURE</option>
    <option value="GLE">
      Grade Level Equivalent</option>
    <option value="DRA" data-reading-level="16">
      DRA LEVEL</option>
  </select>
</div>
```



对应代码:

```
readingLevelSelectElement = readingLevelElement.find('select')
print("readingLevelSelectElement=%s" % readingLevelSelectElement)
guidedReading = readingLevelSelectElement.find('option[value="GRL"]')
lexileMeasure = readingLevelSelectElement.find('option[value="LEX"]')
gradeLevelEquivalent = readingLevelSelectElement.find('option[value="GLE"]')
draLevel = readingLevelSelectElement.find('option[value="DRA"]')
print("guidedReading=%s" % guidedReading)
print("lexileMeasure=%s" % lexileMeasure)
print("gradeLevelEquivalent=%s" % gradeLevelEquivalent)
print("draLevel=%s" % draLevel)
```

即可提取到内容:

find() vs items()

```
readingLevelSelectElement=<select id="book-detail-reading-level"
    <option value="GRL" data-read="GUIDED READING"></option>
    <option value="LEX" data-read="LEXILE® MEASURE"></option>
    <option value="GLE">Grade Level Equivalent</option>
    <option value="DRA" data-read="DRA LEVEL"></option>
</select>

guidedReading=J
lexileMeasure=240L
gradeLevelEquivalent=None
draLevel=16
```

举例2

页面：

The screenshot shows a web page for a car listing on Autohome. The car is a white sedan. The page includes a car image, specifications, and a user review. The browser's developer tools are open, showing the HTML structure. A specific span element with data-price='14.59万元' and class='fn-left' is highlighted.

html:

```
<span data-price="14.59万元" class="fn-left">
    厂商指导价：14.59万元</span>
```

想要：获取 span 的 data-price 的属性的值

- 正确写法

find() vs items()

```
msrpPriceStr = msrpPriceElt.attr("data-price")
```

- 错误写法:

```
msrpPriceStr = msrpPriceElt.attr.data-price
```

不是用PyQuery.val去获取元素属性值

注意此处不是

[pyquery – PyQuery complete API — pyquery 1.2.4 documentation](#)

中的:

```
PyQuery.val(value=<NoDefault>)
```

去获取属性的值的

CSS 选择器的写法

attribute = value 类的元素定位

例子1

调试看到的 html代码是:

```
<ul class="list-user list-user-1" id="list-user-1">
  <li>
    <a href="/index.php?m=home&c=match_new&a=video&show_id=
```

PySpider 内部会自动加上当前host, 即http前缀的 http://xxx , 实际上是:

```
<ul class="list-user list-user-1" id="list-user-1">
  <li>
    <a href="http://xxx/index.php?m=home&c=match_new&a=vic
```

(PySpider 中的) PyQuery 的写法:

```
response.doc('ul[id^="list-user"] a[href^="http"]')
```

或:

```
response.doc('ul[id^="list-user"] li a[href^="http"]')
```

(PySpider中的) 效果：可以找到元素



```
<video controls="" class="video-box" poster="https://img.q
    <source src="https://cdn2.qupeiyin.cn/2017-12-15/id151
</video>
```

11

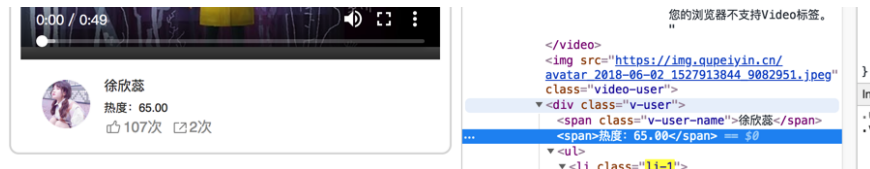
find() vs items()

代码:

```
# videoUrl = response.doc('video source[src$=".mp4"]')
videoUrl = response.doc('video source[src^="http"]').attr('src')
```

nth-child 类的元素定位

页面:



```
<div class="v-user">
  <span class="v-user-name">徐欣蕊</span>
  <span>热度: 65.00</span>
```

想要定位: `div` 下面第二个 `span`

css选择器的写法:

```
div[class="v-user"] span:nth-child(2)
```

crifan.com, 使用署名4.0国际(CC BY 4.0)协议发布 all right reserved,
powered by Gitbook最后更新: 2021-04-14 19:49:46

find() vs items()

概述

- `PyQuery.find()` : 返回的是 lxml 的 element 的 generator
- `PyQuery.items()` : 返回的是 PyQuery 的 generator

详解

- find

- 官网文档: [pyquery.pyquery.PyQuery.find](#)

`PyQuery.find(selector)`

Find elements using selector traversing down from self:

```
>> m = '<p><span><em>Whoah!</em></span></p><p>'
>> d = PyQuery(m)
>> d('p').find('em')
[<em>, <em>]
>> d('p').eq(1).find('em')
[<em>]
```

- -> `find()` 返回的是 element 元素
= `lxml.html.HtmlElement` 的 generator

- items

- 官网文档: [pyquery.pyquery.PyQuery.items](#)

`PyQuery.items(selector=None)`

Iter over elements. Return PyQuery objects:

```
>> d = PyQuery('<div><span>foo</span><span>bar</span>')
>> [i.text() for i in d.items('span')]
['foo', 'bar']
>> [i.text() for i in d('span').items()]
['foo', 'bar']
```

- -> `items()` 返回的
是 `PyQuery = pyquery.pyquery.PyQuery` 的 generator

举例

find() vs items()

```
<ul class="rank-list-ul" 0>
    <li id="s3170">
...
    </li>
    <li id="s692">
...
    </li>
```

find()

如果是 `find()` :

```
carSeriesDocGenerator = merchantRankDoc.find("li[id*='s']")
print("type(carSeriesDocGenerator)=%s" % type(carSeriesDocGenerator))
```

则返回的是 `PyQuery`

```
type(carSeriesDocGenerator)=<class 'pyquery.pyquery.PyQuery'>
```

然后 `generator` 转换成 `list` 后:

```
carSeriesDocList = list(carSeriesDocGenerator)
for curSeriesIdx, eachCarSeriesDoc in enumerate(carSeriesDocList):
    print("type(eachCarSeriesDoc)=%s" % type(eachCarSeriesDoc))
```

每个元素是: `lxml.html.HtmlElement`

```
type(eachCarSeriesDoc)=<class 'lxml.html.HtmlElement'>
```

items()

如果换成 `items()`

```
carSeriesDocGenerator = merchantRankDoc.items("li[id*='s']")
print("type(carSeriesDocGenerator)=%s" % type(carSeriesDocGenerator))
```

则返回的是 `generator` :

```
type(carSeriesDocGenerator)=<class 'generator'>
```

然后 `generator` 转换成 `list` 后:

find() vs items()

```
carSeriesDocList = list(carSeriesDocGenerator)
for curSeriesIdx, eachCarSeriesDoc in enumerate(carSeriesDocList):
    print("type(eachCarSeriesDoc)=%s" % type(eachCarSeriesDoc))
```

每个元素是: `pyquery.pyquery.PyQuery`

```
type(eachCarSeriesDoc)=<class 'pyquery.pyquery.PyQuery'>
```

crifan.com, 使用[署名4.0国际\(CC BY 4.0\)协议](#)发布 all right reserved,
powered by Gitbook最后更新: 2021-04-14 19:49:27

PyQuery语法

CSS Selector = CSS选择器

PyQuery中选择=定位=查找元素的话，主要是用 CSS选择器 。

此处贴出 CSS选择器 的常用的基本的语法：

选择器	举例	解释
<code>.class</code>	<code>.intro</code>	选择 <code>class="intro"</code> 所有元素
<code>#id</code>	<code>#firstname</code>	选择 <code>id="firstname"</code> 所有元素
元素相关		
<code>*</code>	<code>*</code>	选择所有元素
<code>element</code>	<code>p</code>	选择所有 <code><p></code> 元素
<code>element,element</code>	<code>div,p</code>	选择所有 <code><div></code> 元素和所有 <code><p></code> 元素
<code>element element</code>	<code>div p</code>	选择 <code><div></code> 元素的所有 <code><p></code> 元素
<code>element>element</code>	<code>div>p</code>	选择父元素为 <code><div></code> 元素的所有 <code><p></code> 元素
<code>element+element</code>	<code>div+p</code>	选择紧接在 <code><div></code> 元素之后的所有 <code><p></code> 元素
<code>element1~element2</code>	<code>p~ul</code>	选择前面有 <code><p></code> 元素的每个 <code></code> 元素
属性相关		
<code>[attribute]</code>	<code>[target]</code>	选择带有 <code>target</code> 属性所有元素
<code>[attribute=value]</code>	<code>[target=_blank]</code>	选择 <code>target="_blank"</code> 所有元素
<code>[attribute~=value]</code>	<code>[title~=flower]</code>	选择 <code>title</code> 属性单词 <code>"flower"</code> 有元素
<code>[attribute =value]</code>	<code>[lang=en]</code>	选择 <code>lang</code> 属性 <code>"en"</code> 开头的元素
<code>[attribute^=value]</code>	<code>a[src^="https"]</code>	选择其 <code>src</code> 属性 <code>"https"</code> 开头的 <code><a></code> 元素

选择器	举例	解释
<code>[attribute\$=value]</code>	<code>a[src\$=".pdf"]</code>	选择其 <code>src</code> 属性以 <code>".pdf"</code> 结尾的 <code><a></code> 元素
<code>[attribute*=value]</code>	<code>a[src*="abc"]</code>	选择其 <code>src</code> 属性包含 <code>"abc"</code> 子串的 <code><a></code> 元素
元素的修饰? 相关		
<code>:link</code>	<code>a:link</code>	选择所有未被访问的链接
<code>:visited</code>	<code>a:visited</code>	选择所有已被访问的链接
<code>:active</code>	<code>a:active</code>	选择活动链接
<code>:hover</code>	<code>a:hover</code>	选择鼠标指针位于其上的链接
<code>:focus</code>	<code>input:focus</code>	选择获得焦点的 <code>input</code> 元素
<code>:before</code>	<code>p:before</code>	在每个 <code><p></code> 元素内容之前插入内容
<code>:after</code>	<code>p:after</code>	在每个 <code><p></code> 元素内容之后插入内容
<code>:first-letter</code>	<code>p:first-letter</code>	选择每个 <code><p></code> 元素的首字母
<code>:first-line</code>	<code>p:first-line</code>	选择每个 <code><p></code> 元素的首行
<code>:first-child</code>	<code>p:first-child</code>	选择属于父元素的第一个子元素的每个元素
<code>:first-of-type</code>	<code>p:first-of-type</code>	选择属于其父元素的第一个 <code><p></code> 元素的每个 <code><p></code> 元素
<code>:last-of-type</code>	<code>p:last-of-type</code>	选择属于其父元素的最后一个 <code><p></code> 元素的每个 <code><p></code> 元素
<code>:only-of-type</code>	<code>p:only-of-type</code>	选择属于其父元素的唯一的 <code><p></code> 元素的每个 <code><p></code> 元素

选择器	举例	解释
<code>:only-child</code>	<code>p:only-child</code>	选择属于其父元素唯一子元素的每个元素
<code>:nth-child(n)</code>	<code>p:nth-child(2)</code>	选择属于其父元素第n个子元素的每个元素 <p> 元素
<code>:nth-last-child(n)</code>	<code>p:nth-last-child(2)</code>	同上，从最后一个元素开始计数
<code>:nth-of-type(n)</code>	<code>p:nth-of-type(2)</code>	选择属于其父元素第n个 <p> 元素的每个 <p> 元素
<code>:nth-last-of-type(n)</code>	<code>p:nth-last-of-type(2)</code>	同上，但是从最后一个子元素开始计数
<code>:last-child</code>	<code>p:last-child</code>	选择属于其父元素最后一个子元素的每个元素
<code>:lang(language)</code>	<code>p:lang(it)</code>	选择带有以 "it" 开头的 lang 属性的每个 <p> 元素
<code>:root</code>	<code>:root</code>	选择文档的根元素
<code>:empty</code>	<code>p:empty</code>	选择没有子元素的 <p> 元素（包括空节点）
<code>:target</code>	<code>#news:target</code>	选择当前活动的 #news 元素
<code>:enabled</code>	<code>input:enabled</code>	选择每个启用的 <input> 元素
<code>:disabled</code>	<code>input:disabled</code>	选择每个禁用的 <input> 元素
<code>:checked</code>	<code>input:checked</code>	选择每个被选中的 <input> 元素
<code>:not(selector)</code>	<code>:not(p)</code>	选择非 <p> 元素的所有元素
<code>::selection</code>	<code>::selection</code>	选择被用户选取的部分

find() vs items()

crifan.com, 使用署名4.0国际(CC BY 4.0)协议发布 all right reserved,
powered by Gitbook最后更新: 2021-04-14 19:42:30

PyQuery实例

下面记录一些实际案例，供参考。

crifan.com，使用署名4.0国际(CC BY 4.0)协议发布 all right reserved,
powered by Gitbook最后更新： 2021-04-14 19:50:20

汽车之家品牌车型车系

下面贴出，用PySpider爬取汽车之家车型的代码，其中包含 PyQuery 相关代码，供参考。

```
#!/usr/bin/env python
# -*- encoding: utf-8 -*-
# Created on 2018-04-27 21:53:02
# Project: autohomeBrandData

from pyspider.libs.base_handler import *
import string
import re

class Handler(BaseHandler):
    crawl_config = {
    }

    # @every(minutes=24 * 60)
    def on_start(self):
        for eachLetter in list(string.ascii_lowercase):
            self.crawl("https://www.autohome.com.cn/grade/("

    @catch_status_code_error
    def gradCarHtmlPage(self, response):
        print("gradCarHtmlPage: response=", response)
        picSeriesItemList = response.doc('.rank-list-ul li')
        print("picSeriesItemList=", picSeriesItemList)
        # print("len(picSeriesItemList)=%s"%(len(picSeriesItemList)))
        for each in picSeriesItemList:
            self.crawl(each.attr.href, callback=self.picSeriesPage)

    @config(priority=2)
    def picSeriesPage(self, response):
        # <a href="/pic/series-t/66.html">查看停产车型&
        # <a class="ckmore" href="/pic/series/588.html">
        # <span class="fn-right">&nbsp;&nbsp;&nbsp;</span>
        fnRightPicSeries = response.doc('.search-pic-tbar')
        print("fnRightPicSeries=", fnRightPicSeries)
        if fnRightPicSeries:
            # hrefValue = fnRightPicSeries.attr.href
            # print("hrefValue=", hrefValue)
            # fullPicSeriesUrl = "https://car.autohome.com.cn/pic/" + hrefValue
            fullPicSeriesUrl = fnRightPicSeries.attr.href
            print("fullPicSeriesUrl=", fullPicSeriesUrl)
            self.crawl(fullPicSeriesUrl, callback=self.picSeriesPage)

        # continue parse brand data
        aDictList = []
        # for eachA in response.doc('.breadnav a[href^="/"]'):
        for eachA in response.doc('.breadnav a[href*="/pic/"]'):
            eachADict = {
                "text": eachA.text(),
                "href": eachA.attr.href
            }
            aDictList.append(eachADict)
```

```

        print("eachADict=", eachADict)
        aDictList.append(eachADict)

    print("aDictList=", aDictList)

    mainBrandDict = aDictList[-3]
    subBrandDict = aDictList[-2]
    brandSerieDict = aDictList[-1]
    print("mainBrandDict=%s, subBrandDict=%s, brandSerieDict=%s" % (mainBrandDict, subBrandDict, brandSerieDict))

    dtTextList = []
    for eachDt in response.doc("dl.search-pic-cardl dt"):
        dtTextList.append(eachDt.text())

    print("dtTextList=", dtTextList)

    groupCount = len(dtTextList)
    print("groupCount=", groupCount)

    for eachDt in response.doc("dl.search-pic-cardl dt"):
        dtTextList.append(eachDt.text())

    ddUleltList = []
    for eachDdUlelt in response.doc("dl.search-pic-cardl ddUlelt"):
        ddUleltList.append(eachDdUlelt)

    print("ddUleltList=", ddUleltList)

    modelDetailDictList = []

    for curIdx in range(groupCount):
        curGroupTitle = dtTextList[curIdx]
        print("-----[%d] %s" % (curIdx, curGroupTitle))

        for eachLiAElt in ddUleltList[curIdx].items("li"):
            # 1. model name
            curModelName = eachLiAElt.text()
            curModelName = eachLiAElt.contents()[0]
            curModelName = curModelName.strip()
            print("curModelName=", curModelName)
            curFullModelName = curGroupTitle + " " + curModelName
            print("curFullModelName=", curFullModelName)

            # 2. model id + carSeriesId + spec url
            curModelId = ""
            curSeriesId = ""
            curModelSpecUrl = ""
            modelSpecUrlTemplate = "https://www.autohome.com.cn/pic/series-s3/"
            curModelPicUrl = eachLiAElt.attr.href
            print("curModelPicUrl=", curModelPicUrl)
            # https://car.autohome.com.cn/pic/series-s3/

```



```

        foundModelSeriesId = re.search("pic/series-
print("foundModelSeriesId=", foundModelSer:
if foundModelSeriesId:
    curModelId = foundModelSeriesId.group(
    curSeriesId = foundModelSeriesId.group(
    print("curModelId=%s, curSeriesId=%s",
    curModelSpecUrl = (modelSpecUrlTemplate
    print("curModelSpecUrl=", curModelSpec

# 3. model status
modelStatus = "在售"
foundStopSale = eachLiAElt.find('i[class*=
if foundStopSale:
    modelStatus = "停售"
else:
    foundWseason = eachLiAElt.find('i[class
    if foundWseason:
        modelStatus = "未上市"

modelDetailDictList.append({
    "url": curModelSpecUrl,
    "车系ID": curSeriesId,
    "车型ID": curModelId,
    "车型": curFullModelName,
    "状态": modelStatus
})
print("modelDetailDictList=", modelDetailDictList)

allSerieDictList = []
for curIdx, eachModelDetailDict in enumerate(model
    curSerieDict = {
        "品牌": mainBrandDict["text"],
        "子品牌": subBrandDict["text"],
        "车系": brandSerieDict["text"],
        "车系ID": eachModelDetailDict["车系ID"],
        "车型": eachModelDetailDict["车型"],
        "车型ID": eachModelDetailDict["车型ID"],
        "状态": eachModelDetailDict["状态"]
    }
    allSerieDictList.append(curSerieDict)
    # print("before send_message: [%d] curSerieDict
    # self.send_message(self.project_name, curSerie
    print("[%d] curSerieDict=%s" % (curIdx, curSer
    self.crawl(eachModelDetailDict["url"], callback

# print("allSerieDictList=", allSerieDictList)
# return allSerieDictList

#def on_message(self, project, msg):
#    print("on_message: msg=", msg)
#    return msg

```

find() vs items()

```
@catch_status_code_error
def carModelSpecPage(self, response):
    print("carModelSpecPage: response=", response)
    # https://www.autohome.com.cn/spec/32708/#pvareaid=
    curSerieDict = response.save
    print("curSerieDict", curSerieDict)

    # cityDealerPriceInt = 0
    # cityDealerPriceElt = response.doc('.cardetail-infor-price')
    # print("cityDealerPriceElt=%s" % cityDealerPriceElt)
    # if cityDealerPriceElt:
    #     cityDealerPriceFloatStr = cityDealerPriceElt.text
    #     print("cityDealerPriceFloatStr=", cityDealerPriceFloatStr)
    #     cityDealerPriceFloat = float(cityDealerPriceFloatStr)
    #     print("cityDealerPriceFloat=", cityDealerPriceFloat)
    #     cityDealerPriceInt = int(cityDealerPriceFloat * 10000)
    #     print("cityDealerPriceInt=", cityDealerPriceInt)

    msrpPriceInt = 0
    # body > div.content > div.row > div.col
    # 厂商指导价=厂商建议零售价格=MSRP=Manufacturer's suggested retail price
    msrpPriceElt = response.doc('.cardetail-infor-price')
    print("msrpPriceElt=", msrpPriceElt)
    if msrpPriceElt:
        msrpPriceStr = msrpPriceElt.attr("data-price")
        print("msrpPriceStr=", msrpPriceStr)
        foundMsrpPrice = re.search("(?P<msrpPrice>[0-9]+)")
        print("foundMsrpPrice=", foundMsrpPrice)
        if foundMsrpPrice:
            msrpPrice = foundMsrpPrice.group("msrpPrice")
            print("msrpPrice=", msrpPrice)
            msrpPriceFloat = float(msrpPrice)
            print("msrpPriceFloat=", msrpPriceFloat)
            msrpPriceInt = int(msrpPriceFloat * 10000)
            print("msrpPriceInt=", msrpPriceInt)

    # curSerieDict["经销商参考价"] = cityDealerPriceInt
    curSerieDict["厂商指导价"] = msrpPriceInt

    return curSerieDict
```

详见：

【已解决】写Python爬虫爬取汽车之家品牌车系车型数据

crifan.com, 使用署名4.0国际(CC BY 4.0)协议发布 all right reserved,
powered by Gitbook最后更新: 2021-04-14 19:50:09

附录

下面列出相关参考资料。

crifan.com, 使用署名4.0国际(CC BY 4.0)协议发布 all right reserved,
powered by Gitbook最后更新: 2021-04-14 19:42:48

参考资料

- [【已解决】PySpider中获取PyQuery获取到节点的子元素](#)
- [【已解决】写Python爬虫爬取汽车之家品牌车系车型数据](#)
- [【已解决】PySpider中PyQuery选择div后面的第二个span元素](#)
-
- [CSS Selectors Reference](#)
- [CSS selectors - CSS: Cascading Style Sheets | MDN](#)
- [CSS 选择器参考手册](#)
-

crifan.com, 使用署名4.0国际(CC BY 4.0)协议发布 all right reserved,
powered by Gitbook最后更新: 2021-04-14 19:42:55