



链滴

# python 爬虫 (SGMLParser)

作者: [cactus0509](#)

原文链接: <https://ld246.com/article/1470105137836>

来源网站: [链滴](#)

许可协议: [署名-相同方式共享 4.0 国际 \(CC BY-SA 4.0\)](#)

```

<p>`# -<em>- coding: utf-8 -</em>-<br>
import requests<br>
import urllib2<br>
import sys<br>
from sgmlib import SGMLParser<br>
from cgitb import text<br>
reload(sys)<br>
sys.setdefaultencoding('utf8')</p>
<p>class CLAS_EXPERT_LIST(SGMLParser):<br>
def <strong>init</strong>(self):<br>
reload(sys)<br>
SGMLParser.<strong>init</strong>(self)<br>
self.is_a = ""<br>
self.name = []<br>
self.urls = []</p>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">def start_a(self, attrs):
</span></span><span class="highlight-line"><span class="highlight-cl">    for k, v in attrs :
</span></span><span class="highlight-line"><span class="highlight-cl">        if k=='href'
</span></span><span class="highlight-line"><span class="highlight-cl">            nd v.count('detail') &gt; 0 :
</span></span><span class="highlight-line"><span class="highlight-cl">                self.is_a =
</span></span><span class="highlight-line"><span class="highlight-cl">                    self.urls.ap
end(v)
</span></span><span class="highlight-line"><span class="highlight-cl">def end_a(self):
</span></span><span class="highlight-line"><span class="highlight-cl">    self.is_a = 0
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">def handle_data(s
lf, text):
</span></span><span class="highlight-line"><span class="highlight-cl">    if self.is_a == 1:
</span></span><span class="highlight-line"><span class="highlight-cl">        self.name.ap
end(text)
</span></span></code></pre>
<p>class EXPERT(SGMLParser):</p>
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">def __init__(self):
</span></span><span class="highlight-line"><span class="highlight-cl">    SGMLParser.__in
t__(self)
</span></span><span class="highlight-line"><span class="highlight-cl">    self.is_div = 0
</span></span><span class="highlight-line"><span class="highlight-cl">    self.is_h3 = 0
</span></span><span class="highlight-line"><span class="highlight-cl">    self.is_div_p = 0
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">    self.p_cnt = 0
</span></span><span class="highlight-line"><span class="highlight-cl">    self.image = {}
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">    self.is_div_exper
= 0
</span></span><span class="highlight-line"><span class="highlight-cl">    self.is_div_exper
_p = 0
</span></span><span class="highlight-line"><span class="highlight-cl">    self.is_div_exper
_p_cnt = 0
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">def start_div(self,

```

```

ttrs):
</span></span><span class="highlight-line"><span class="highlight-cl"> for k, v in attrs :
</span></span><span class="highlight-line"><span class="highlight-cl">     if k=='class'
nd v.count('name') &gt; 0 :
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_div
1
</span></span><span class="highlight-line"><span class="highlight-cl">     if k=='class'
nd v.count('expert_content') &gt; 0 :
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_div_
xpert = 1
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> def end_div(self):
</span></span><span class="highlight-line"><span class="highlight-cl">     if self.is_div ==
:
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_div = 0
</span></span><span class="highlight-line"><span class="highlight-cl">     elif self.is_div_e
pert == 1 :
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_div_ex
ert = 0
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> def start_h3(self, a
trs):
</span></span><span class="highlight-line"><span class="highlight-cl">     if self.is_div :
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_h3 = 1
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> def end_h3(self):
</span></span><span class="highlight-line"><span class="highlight-cl">     self.is_h3 = 0
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> def start_p(self, att
s):
</span></span><span class="highlight-line"><span class="highlight-cl">     if self.is_div ==
:
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_div_p =
1
</span></span><span class="highlight-line"><span class="highlight-cl">     elif self.is_div_e
pert == 1 :
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_div_ex
ert_p = 1
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_div_ex
ert_p_cnt = self.is_div_expert_p_cnt + 1
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> def end_p(self):
</span></span><span class="highlight-line"><span class="highlight-cl">     if self.is_div :
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_div_p =
0
</span></span><span class="highlight-line"><span class="highlight-cl">         self.p_cnt = s
lf.p_cnt + 1
</span></span><span class="highlight-line"><span class="highlight-cl">     elif self.is_div_e
pert == 1 :
</span></span><span class="highlight-line"><span class="highlight-cl">         self.is_div_ex
ert_p = 0
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">

```

```

</span></span><span class="highlight-line"><span class="highlight-cl">def handle_data(s
lf, text):
</span></span><span class="highlight-line"><span class="highlight-cl">    try:
</span></span><span class="highlight-line"><span class="highlight-cl">        if self.is_div
= 1:
</span></span><span class="highlight-line"><span class="highlight-cl">            if self.is_h3
== 1 :
</span></span><span class="highlight-line"><span class="highlight-cl">                self.ima
e["name"] = text
</span></span><span class="highlight-line"><span class="highlight-cl">            if self.is_di
_p == 1 :
</span></span><span class="highlight-line"><span class="highlight-cl">                if self.p_
nt == 0 :
</span></span><span class="highlight-line"><span class="highlight-cl">                    self.i
age["job"] = text
</span></span><span class="highlight-line"><span class="highlight-cl">                else:
</span></span><span class="highlight-line"><span class="highlight-cl">                    self.i
age["title"] = text
</span></span><span class="highlight-line"><span class="highlight-cl">            if self.is_div_e
pert == 1:
</span></span><span class="highlight-line"><span class="highlight-cl">                if self.is_di
_expert_p == 1 :
</span></span><span class="highlight-line"><span class="highlight-cl">                    #print se
f.is_div_expert_p_cnt,text
</span></span><span class="highlight-line"><span class="highlight-cl">                if self.is
div_expert_p_cnt == 2:
</span></span><span class="highlight-line"><span class="highlight-cl">                    self.i
age["employer"] = text
</span></span><span class="highlight-line"><span class="highlight-cl">                elif self.i
_div_expert_p_cnt == 6:
</span></span><span class="highlight-line"><span class="highlight-cl">                    self.i
age["filed"] = text
</span></span><span class="highlight-line"><span class="highlight-cl">                elif self.i
_div_expert_p_cnt == 16:
</span></span><span class="highlight-line"><span class="highlight-cl">                    self.i
age["conn_info"] = text
</span></span><span class="highlight-line"><span class="highlight-cl">                #print text
</span></span><span class="highlight-line"><span class="highlight-cl">                #print self.
mage["name"] , self.image["title"] , self.image["job"] , self.image["employer"] , self.image["file
"] , self.image["conn_info"]
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">    except Exceptio
,e:
</span></span><span class="highlight-line"><span class="highlight-cl">        print e
</span></span></code></pre>

```

<p>def list\_expert():<br>

headers = {<br>

"Connection": "keep-alive",<br>

"Cookie": "Ecp\_IpLoginFail=160726111.205.187.18; kc\_cnki\_net\_uid=ff38e944-e46c-2d76-349c24a97e03ded8; ASP.NET\_SessionId=ysbae4exnu0vkugigsdnknps; AutoIpLogin=; LID=; SID=12103; CNZZDATA4922505=cnzz\_eid%3D1343153553-1469773415-%26ntime%3D1469782211 FileNameM=cnki%3A; c\_m\_LinID=LinID=WEEvREcwSIJHSldTTGJhYIRtMVNwOTZ6Q1UzaHd

```

OFN2RzR2MEEyUkJPWmE=$9A4hF_YAuvQ5obgVAqNKPCYcEjKensW4IQMowwHtwkF4VYPoH
KxJw!!&ot=07/29/2016 18:16:04", <br>
"Host": "elib.cnki.net", <br>
"Accept": "text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8", <br>
"Accept-Encoding": "gzip, deflate", <br>
"Accept-Language": "zh-CN,zh;q=0.8,en-US;q=0.5,en;q=0.3", <br>
"Referer": "<a href='\"https://ld246.com/forward?goto=http%3A%2F%2Fwww.example.com%2
%2522\" target='\"_blank\" rel='\"nofollow ugc\">http://www.example.com/\" </a>,<br>
"User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10.11; rv:47.0) Gecko/20100101 Firefox
47.0" <br>
}<br>
r = requests.get('<a href='\"https://ld246.com/forward?goto=http%3A%2F%2Fwww.chinathink
anks.org.cn%2Fcontent%2Fexpert\" target='\"_blank\" rel='\"nofollow ugc\">http://www.chinathin
tanks.org.cn/content/expert</a>') <br>
content=r.text <br>
listname = CLAS_EXPERT_LIST() <br>
listname.feed(content) <br>
rn = zip(listname.urls,listname.name) <br>
return rn </p>
<p>def get_expert(url,name): <br>
headers = { <br>
"Connection": "keep-alive", <br>
"Cookie": "Ecp_IpLoginFail=160726111.205.187.18; kc_cnki_net_uid=ff38e944-e46c-2d76-349c
24a97e03ded8; ASP.NET_SessionId=ysbae4exnu0vkugigsdnkps; AutoIpLogin=; LID=; SID=1
2103; CNZZDATA4922505=cnzz_eid%3D1343153553-1469773415-%26ntime%3D1469782211
FileNameM=cnki%3A; c_m_LinID=LinID=WEEvREcwSIJHSlDTTGJhYIRtMVNwOTZ6Q1UzaHd
OFN2RzR2MEEyUkJPWmE=$9A4hF_YAuvQ5obgVAqNKPCYcEjKensW4IQMowwHtwkF4VYPoH
KxJw!!&ot=07/29/2016 18:16:04", <br>
"Host": "elib.cnki.net", <br>
"Accept": "text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8", <br>
"Accept-Encoding": "gzip, deflate", <br>
"Accept-Language": "utf-8,zh;q=0.8,en-US;q=0.5,en;q=0.3", <br>
"Referer": "<a href='\"https://ld246.com/forward?goto=http%3A%2F%2Fwww.example.com%2
%2522\" target='\"_blank\" rel='\"nofollow ugc\">http://www.example.com/\" </a>,<br>
"User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10.11; rv:47.0) Gecko/20100101 Firefox
47.0" <br>
}<br>
r = requests.get(url) <br>
content=r.text <br>
expert = EXPERT() <br>
expert.feed(content) <br>
return expert.image </p>
<p>if <strong>name</strong> == "<strong>main</strong>": <br>
try: <br>
i = 0 <br>
import chardet <br>
with open("/tmp/expert.txt","w") as f: <br>
exports = list_expert() <br>
for k,v in exports: <br>
v = get_expert(k,v) <br>
name,job,title,employer,filed,conn_info=None,None,None,None,None,None <br>
if "name" in v: <br>
name= v["name"] <br>
if "job" in v: <br>

```

```

job= v["job"]<br>
if "title" in v:<br>
title= v["title"]<br>
if "employer" in v:<br>
employer= v["employer"]<br>
if "filed" in v:<br>
filed= v["filed"]<br>
if "conn_info" in v:<br>
conn_info= v["conn_info"]</p>
<pre> <code class="highlight-chroma"> <span class="highlight-line"> <span class="highlight
cl">          print "{0}#{1}#{2}#{3}#{4}#{5}".format(name,job,title,employer,filed,conn_info)
</span> </span> <span class="highlight-line"> <span class="highlight-cl">          f.write("{0}
{1}#{2}#{3}#{4}#{5}\n".format(name,job,title,employer,filed,conn_info))
</span> </span> <span class="highlight-line"> <span class="highlight-cl">          f.flush()
</span> </span> <span class="highlight-line"> <span class="highlight-cl">          i = i + 1
</span> </span> <span class="highlight-line"> <span class="highlight-cl">          left = div
od(i,50)
</span> </span> <span class="highlight-line"> <span class="highlight-cl">          if left[1]
= 0 :
</span> </span> <span class="highlight-line"> <span class="highlight-cl">          print i
</span> </span> <span class="highlight-line"> <span class="highlight-cl">
</span> </span> <span class="highlight-line"> <span class="highlight-cl"> >except Exception,e
:
</span> </span> <span class="highlight-line"> <span class="highlight-cl"> print e
</span> </span> <span class="highlight-line"> <span class="highlight-cl"> `
</span> </span> </code> </pre>

```