

xpath练习:抓取猪八戒数据

不墨迹, 直接怼.

干猪八戒.

```
1 import requests
2 from lxml import etree
3
4 url = "https://beijing.zbj.com/search/f/?
    type=new&kw=saas"
5
6 headers = {
7     "User-Agent": "Mozilla/5.0 (Macintosh; Intel
    Mac OS X 10_15_4) AppleWebKit/537.36 (KHTML, like
    Gecko) Chrome/87.0.4280.88 Safari/537.36"
8 }
9 resp = requests.get(url, headers=headers)
10
11 # 丢给etree, 生成Element对象
12 tree = etree.HTML(resp.text)
13 # 拿数据吧
```

```
14 els = tree.xpath("//div[@class='witkey-list-grid  
j-service-provider-wrap ']/*/div[@class='witkey-  
item grid-box']")  
15  
16 for div in els:  
17     name = div.xpath("./div/div[@class='grid-top-  
right']/section/h4/a/text()")[0]  
18     intro = div.xpath("./div/div[@class='grid-  
top-right']/section/h4/div/div[last()]/text()")  
[0]  
19     shopdesc =  
"".join(div.xpath("./div/div[@class='grid-top-  
right']/section/h4/div/div[@class='witkey-  
shopdesc']//text()")).replace("\n", "")  
20  
21     firstline =  
div.xpath("./div/div[@class='grid-top-  
right']/section/h4/div/div[@class='witkey-  
shopdesc-firstline']//text()")  
22     if firstline:  
23         shopdesc +=  
"".join(firstline).replace("\n", "")  
24     fav =  
"".join(div.xpath("./div/div[@class='grid-top-  
right']/section/h4/div/div[@class='expert-  
tree']//text()")).replace("\n", "")
```

25

26

27 # 后面你想干嘛就干嘛吧。 记录数据库,写入文件,扔给全文检索工
具 都可以

28