

20200821信息安全实训

笔记本: 我的第一个笔记本

创建时间: 2020/8/21 18:19

更新时间: 2020/8/21 21:05

作者: 820410740@qq.com

全站爬虫 scrapy 爬虫框架
QuotesSpider类

```
1 import scrapy
2
3 class QuotesSpider(scrapy.Spider):
4     name = "quote"
5
6     def start_requests(self):
7         urls = [
8             "http://quotes.toscrape.com/page/1/",
9             "http://quotes.toscrape.com/page/2/"
10        ]
11        for url in urls:
12            yield scrapy.Request(url=url, callback=self.parse)
13
14    def parse(self, response):
15        page = response.url.split("/")[-2]
16        filename = 'quotes-%s.html' % page
17        with open(filename, 'wb') as f:
18            f.write(response.body)
19        self.log('Saved file %s' % filename)
```

scrapy shell

css

XPath

```
1 import scrapy
2
3 class LinkSpider(scrapy.Spider):
4     name = "links"
5
6     start_urls = [
7         'https://www.sziangyou.com',
8     ]
9
10    def parse(self, response):
11        for a in response.css('a'):
12            print(a.css('::attr(href)').get())
```

```

1 import scrapy
2
3 class LinkSpider(scrapy.Spider):
4
5     name = "links"
6
7     start_urls = [
8         'https://www.szlangyou.com',
9     ]
10
11     def parse(self, response):
12         for href in response.css('a::attr(href)').getall():
13             print(href)

```

```

1 import scrapy
2
3 class LinkSpider(scrapy.Spider):
4
5     name = "links"
6
7     start_urls = [
8         'https://www.szlangyou.com',
9     ]
10
11     def parse(self, response):
12         for href in response.css('a'):
13             yield {'link': href.css('::attr(href)').get(), }

```

```

1 import scrapy
2 from urllib.parse import urljoin
3
4 class LinkSpider(scrapy.Spider):
5
6     name = "links"
7
8     start_urls = [
9         'https://www.szlangyou.com',
10     ]
11
12     def parse(self, response):
13         for href in response.css('a'):
14             link = href.css('::attr(href)').get()
15             if "https://www.szlangyou.com" not in link and link[0] == '/':
16                 link = urljoin("https://www.szlangyou.com", link)
17             yield {'link': link, }
18             yield response.follow(link, callback=self.parse)

```

指纹识别

消息摘要算法

MD5 SHA-1

```

1 import json
2 import hashlib
3 import scrapy
4 from scrapy.crawler import CrawlerProcess
5
6 with open("./links.json", "r") as f:
7     raw = f.read()
8
9 raw = json.loads(raw)
10
11 links = []
12 for i in raw:
13     links.append(i["link"])
14
15 fgprts = []
16 class MySpider(scrapy.Spider):
17
18     name = "fgprpt"
19
20     start_urls = links
21
22     def parse(self, response):
23         md5 = hashlib.md5()
24         md5.update(response.body)
25         finger_print = md5.hexdigest()
26         fgprts.append(
27             {"url": response.url, "finger_print": finger_print}
28         )
29
30 process = CrawlerProcess()
31 process.crawl(MySpider)
32 process.start()
33
34 process.start()
35
36
37 with open("./fgprts.json", "w") as f:
38     f.write(json.dumps(fgprts))

```

mechanicalsp@osptech@UspMol:~/tutorial\$

管理 控制 视图 热键 设备 帮助

```

1 import json
2 import hashlib
3 import scrapy
4 from scrapy.crawler import CrawlerProcess
5
6 with open("./fgprts.json") as f:
7     fgprts = json.loads(f.read())
8
9 urls = [{"url": i["url"]} for i in fgprts]
10
11 results = []
12 class MySpider(scrapy.Spider):
13
14     name = "check"
15
16     start_urls = urls
17
18     def parse(self, response):
19         for i in fgprts:
20             if i["url"] == response.url:
21                 finger_print = i["finger_print"]
22
23         -

```

```

4 from scrapy.crawler import CrawlerProcess
5
6 with open('./fgprts.json', 'r') as f:
7     fgprts = json.loads(f.read())
8
9 urls = [i['url'] for i in fgprts]
10
11 results = []
12 class MySpider(scrapy.Spider):
13
14     name = 'check'
15
16     start_urls = urls
17
18     def parse(self, response):
19         for i in fgprts:
20             if i['url'] == response.url:
21                 finger_print = i['finger_print']
22                 md5 = hashlib.md5()
23                 md5.update(response.body)
24                 finger_print_now = md5.hexdigest()
25                 results.append(
26                     (
27                         'url': response.url,
28                         'isDistorted': finger_print == finger_print_now
29                     )
30                 )
31
32 process = CrawlerProcess()
33 process.crawl(MySpider)
34 process.start()
35
36 with open('./results', 'w') as f:
37     f.write('%s %s\n' % ())

```

```

34 process.start()
35
36 with open('./results', 'w') as f:
37     for i in results:
38         f.write('%s %s\n' % (i['url'], str(i['isDistorted'])))

```

```

30         )
31
32 process = CrawlerProcess()
33 process.crawl(MySpider)
34 process.start()
35

```