

爬虫scrapy框架练习答案

题目要求

- 根据已知网址爬取数据,将本网址的第一页数据使用scrapy框架爬取到本地并写入excel表格中,具体爬取数据位置见图片,红色框部分是需要爬取的数据部分。

<http://www.duozhi.com/>



浙江、山东、辽宁公布开学时间

浙江省校外培训机构面向中小学生的线下培训,不得早于当地中小学全部开学后一周开展。

K12

2020/04/03 | by 余甜

温馨提示:

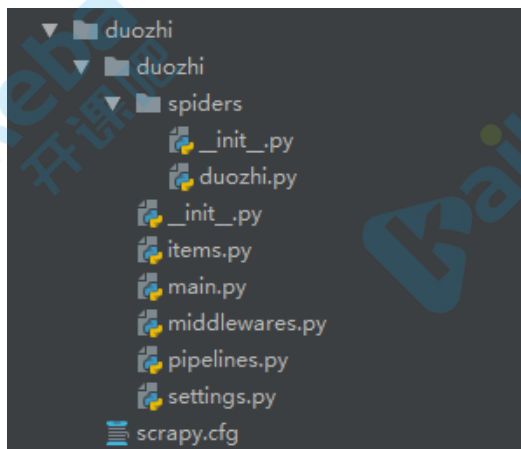
- 1.用命令来创建项目
- 2.settings里面要修改配置

解题步骤:

1、创建项目

使用命令创建即可

```
1 scrapy
```



settings.py文件

```
1 # -*- coding: utf-8 -*-
2 # Scrapy settings for duozhi project
```

```
4 #
5 # For simplicity, this file contains only settings considered important or
6 # commonly used. You can find more settings consulting the documentation:
7 #
8 #     https://docs.scrapy.org/en/latest/topics/settings.html
9 #     https://docs.scrapy.org/en/latest/topics/downloader-middleware.html
10 #     https://docs.scrapy.org/en/latest/topics/spider-middleware.html
11
12 BOT_NAME = 'duozhi'
13
14 SPIDER_MODULES = ['duozhi.spiders']
15 NEWSPIDER_MODULE = 'duozhi.spiders'
16
17 # Crawl responsibly by identifying yourself (and your website) on the
18 # user-agent
19 # USER_AGENT = 'duozhi (+http://www.yourdomain.com)'
20 # Obey robots.txt rules
21 ROBOTSTXT_OBEY = False
22
23 # Configure maximum concurrent requests performed by Scrapy (default: 16)
24 # CONCURRENT_REQUESTS = 32
25
26 # Configure a delay for requests for the same website (default: 0)
27 # See https://docs.scrapy.org/en/latest/topics/settings.html#download-
28 # delay
29 # See also autothrottle settings and docs
30 # DOWNLOAD_DELAY = 3
31 # The download delay setting will honor only one of:
32 # CONCURRENT_REQUESTS_PER_DOMAIN = 16
33 # CONCURRENT_REQUESTS_PER_IP = 16
34 # Disable cookies (enabled by default)
35 # COOKIES_ENABLED = False
36 # Disable Telnet Console (enabled by default)
37 # TELNETCONSOLE_ENABLED = False
38
39 # Override the default request headers:
40 # DEFAULT_REQUEST_HEADERS = {
41 #     'Accept':
42 #         'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
43 #     'Accept-Language': 'en',
44 # }
45
46 # Enable or disable spider middlewares
47 # See https://docs.scrapy.org/en/latest/topics/spider-middleware.html
48 # SPIDER_MIDDLEWARES = {
49 #     'duozhi.middlewares.DuozhiSpiderMiddleware': 543,
50 # }
51
52 # Enable or disable downloader middlewares
53 # See https://docs.scrapy.org/en/latest/topics/downloader-middleware.html
54 # DOWNLOADER_MIDDLEWARES = {
55 #     'duozhi.middlewares.DuozhiDownloaderMiddleware': 543,
56 # }
57
58 # Enable or disable extensions
59 # See https://docs.scrapy.org/en/latest/topics/extensions.html
60 # EXTENSIONS = {
61 #     'scrapy.extensions.telnet.TelnetConsole': None,
```

```

62 # }
63 # Configure item pipelines
64 # See https://docs.scrapy.org/en/latest/topics/item-pipeline.html
65 ITEM_PIPELINES = {
66     'duozhi.pipelines.DuozhiPipeline': 300,
67 }
68
69 # Enable and configure the AutoThrottle extension (disabled by default)
70 # See https://docs.scrapy.org/en/latest/topics/autothrottle.html
71 # AUTOTHROTTLE_ENABLED = True
72 # The initial download delay
73 # AUTOTHROTTLE_START_DELAY = 5
74 # The maximum download delay to be set in case of high latencies
75 # AUTOTHROTTLE_MAX_DELAY = 60
76 # The average number of requests Scrapy should be sending in parallel to
77 # each remote server
78 # AUTOTHROTTLE_TARGET_CONCURRENCY = 1.0
79 # Enable showing throttling stats for every response received:
80 # AUTOTHROTTLE_DEBUG = False
81 # Enable and configure HTTP caching (disabled by default)
82 # See https://docs.scrapy.org/en/latest/topics/downloader-middleware.html#httpcache-middleware-settings
83 # HTTPCACHE_ENABLED = True
84 # HTTPCACHE_EXPIRATION_SECS = 0
85 # HTTPCACHE_DIR = 'httpcache'
86 # HTTPCACHE_IGNORE_HTTP_CODES = []
87 # HTTPCACHE_STORAGE = 'scrapy.extensions.httpcache.FilesystemCacheStorage'

```

duozhi.py文件

```

1 # -*- coding: utf-8 -*-
2 import scrapy
3 import bs4
4 from ..items import DuozhiItem
5 class DangdangBookSpider(scrapy.Spider):
6     name = 'duozhi'
7     allowed_domains = ['duozhi.com']
8     start_urls = ['http://www.duozhi.com/']
9     def parse(self, response):
10         soup = bs4.BeautifulSoup(response.text, 'html.parser')
11         elements = soup.find_all('div', class_='post-item')
12         for element in elements:
13             title = element.find('a', class_='post-title').text
14             content = element.find('p', class_='post-desc').text
15             item = DuozhiItem()
16             item['title'] = title
17             item['content'] = content
18             yield item

```

items文件

```

1 # -*- coding: utf-8 -*-

```

```

3 # Define here the models for your scraped items
4 #
5 # See documentation in:
6 # https://docs.scrapy.org/en/latest/topics/items.html
7 import scrapy
8
9 class DuozhiItem(scrapy.Item):
10     # define the fields for your item here like:
11     # name = scrapy.Field()
12     title = scrapy.Field()
13     content = scrapy.Field()

```

pipelines.py文件

```

1 # -*- coding: utf-8 -*-
2 # Define your item pipelines here
3 #
4 # Don't forget to add your pipeline to the ITEM_PIPELINES setting
5 # See: https://docs.scrapy.org/en/latest/topics/item-pipeline.html
6 import openpyxl
7
8 class DuozhiPipeline(object):
9     def __init__(self):
10         self.wb = openpyxl.Workbook()
11         self.ws = self.wb.active
12         self.ws.append(['标题', '内容'])
13     def process_item(self, item, spider):
14         line = [item['title'], item['content']]
15         self.ws.append(line)
16         return item
17     def close_spider(self, spider):
18         self.wb.save('duozhi.xlsx')
19         self.wb.close()

```

main.py文件

```

1 from scrapy import cmdline
2 cmdline.execute(['scrapy', 'crawl', 'duozhi'])

```

其他文件里面不需要动