

Day08回顾

scrapy框架

■ 五大组件+工作流程+常用命令

```
1  【1】五大组件
2      1.1) 引擎 (Engine)
3      1.2) 爬虫程序 (Spider)
4      1.3) 调度器 (Scheduler)
5      1.4) 下载器 (Downloader)
6      1.5) 管道文件 (Pipeline)
7      1.6) 下载器中间件 (Downloader Middlewares)
8      1.7) 蜘蛛中间件 (Spider Middlewares)
9
10 【2】工作流程
11     2.1) Engine向Spider索要URL,交给Scheduler入队列
12     2.2) Scheduler处理后出队列,通过Downloader Middlewares交给Downloader去下载
13     2.3) Downloader得到响应后,通过Spider Middlewares交给Spider
14     2.4) Spider数据提取:
15         a) 数据交给Pipeline处理
16         b) 需要跟进URL,继续交给Scheduler入队列,依次循环
17
18 【3】常用命令
19     3.1) scrapy startproject 项目名
20     3.2) scrapy genspider 爬虫名 域名
21     3.3) scrapy crawl 爬虫名
```

完成scrapy项目完整流程

■ 完整流程

```
1  【1】crapy startproject Tencent
2  【2】cd Tencent
3  【3】scrapy genspider tencent tencent.com
4  【4】items.py(定义爬取数据结构)
5      import scrapy
6      class TencentItem(scrapy.Item):
7          job_name = scrapy.Field()
8
9  【5】tencent.py (写爬虫文件)
10     import scrapy
11     class TencentSpider(scrapy.Spider):
```

```

12         name = 'tencent'
13         allowed_domains = ['tencent.com']
14         start_urls = ['http://tencent.com/']
15         def parse(self, response):
16             xxx
17             yield item
18 【6】 pipelines.py(数据处理)
19     class TencentPipeline(object):
20         def process_item(self, item, spider):
21             return item
22 【7】 settings.py(全局配置)
23     LOG_LEVEL = ''
24     LOG_FILE = ''
25     FEED_EXPORT_ENCODING = ''
26 【8】 run.py
27     scrapy crawl tencent

```

我们必须记住

■ 熟练记住

```

1 【1】 响应对象response属性及方法
2     1.1) response.text : 获取响应内容 - 字符串
3     1.2) response.body : 获取bytes数据类型
4     1.3) response.xpath('')
5     1.4) response.xpath('').extract() : 提取文本内容,将列表中所有元素序列化为Unicode字符串
6     1.5) response.xpath('').extract_first() : 序列化提取列表中第1个文本内容
7     1.6) response.xpath('').get() : 提取列表中第1个文本内容(等同于extract_first())
8
9 【2】 settings.py中常用变量
10    2.1) 设置日志级别
11        LOG_LEVEL = ''
12    2.2) 保存到日志文件(不在终端输出)
13        LOG_FILE = ''
14    2.3) 设置数据导出编码(主要针对于json文件)
15        FEED_EXPORT_ENCODING = 'utf-8'
16    2.4) 设置User-Agent
17        USER_AGENT = ''
18    2.5) 设置最大并发数(默认为16)
19        CONCURRENT_REQUESTS = 32
20    2.6) 下载延迟时间(每隔多长时间请求一个网页)
21        DOWNLOAD_DELAY = 0.5
22    2.7) 请求头
23        DEFAULT_REQUEST_HEADERS = {}
24    2.8) 添加项目管道
25        ITEM_PIPELINES = {'项目名.pipelines.类名':200}
26    2.9) cookie(默认禁用,取消注释-True|False都为开启)
27        COOKIES_ENABLED = False
28    2.10) 非结构化数据存储路径
29        IMAGES_STORE = ''
30        FILES_STORE = ''
31    2.11) 添加下载器中间件
32        DOWNLOADER_MIDDLEWARES = {'项目名.middlewares.类名':200}

```

```
33
34 【3】 日志级别
35     DEBUG < INFO < WARNING < ERROR < CRITICAL
```

爬虫项目启动方式

■ 启动方式

```
1 【1】 方式一
2     1.1) 从爬虫文件(spider)的start_urls变量中遍历URL地址交给调度器入队列,
3     1.2) 把下载器返回的响应对象(response) 交给爬虫文件的parse(self,response)函数处理
4
5 【2】 方式二
6     重写start_requests()方法, 从此方法中获取URL, 交给指定的callback解析函数处理
7     2.1) 去掉start_urls变量
8     2.2) def start_requests(self):
9           # 生成要爬取的URL地址, 利用scrapy.Request()方法交给调度器
```

数据持久化存储

■ MySQL-MongoDB-Json-csv

```
1 *****存入MySQL、MongoDB*****
2
3 【1】 在setting.py中定义相关变量
4 【2】 pipelines.py中新建管道类, 并导入settings模块
5     def open_spider(self,spider):
6         # 爬虫开始执行1次,用于数据库连接
7
8     def process_item(self,item,spider):
9         # 用于处理抓取的item数据
10        return item
11
12    def close_spider(self,spider):
13        # 爬虫结束时执行1次,用于断开数据库连接
14
15 【3】 settings.py中添加此管道
16     ITEM_PIPELINES = {'':200}
17
18 【注意】 process_item() 函数中一定要 return item
19
20 *****存入JSON、CSV文件*****
21 scrapy crawl maoyan -o maoyan.csv
22 scrapy crawl maoyan -o maoyan.json
23 【注意】
24     存入json文件时候需要添加变量(settings.py) : FEED_EXPORT_ENCODING = 'utf-8'
```

多级页面抓取之爬虫文件

■ 多级页面攻略

```
1  【场景1】只抓取一级页面的情况
2  """
3  一级页面: 名称(name)、爱好(likes)
4  """
5  import scrapy
6  from ..items import OneItem
7
8  class OneSpider(scrapy.Spider):
9      name = 'one'
10     allowed_domains = ['www.one.com']
11     start_urls = ['http://www.one.com/']
12     def parse(self, response):
13         # 抓取1条数据,将item对象传给管道1次,可以只创建1次
14         item = OneItem()
15         dd_list = response.xpath('//dd')
16         for dd in dd_list:
17             item['name'] = dd.xpath('./a/text()').get()
18             item['likes'] = dd.xpath('./a/text()').get()
19
20         yield item
21
22  【场景2】二级页面数据抓取
23  """
24  一级页面: 名称(name)、详情页链接(url)-需要继续跟进
25  二级页面: 详情页内容(content)
26  """
27  import scrapy
28  from ..items import TwoItem
29
30  class TwoSpider(scrapy.Spider):
31      name = 'two'
32      allowed_domains = ['www.two.com']
33      start_urls = ['http://www.two.com/']
34      def parse(self, response):
35          """一级页面解析函数,提取 name 和 url(详情页链接,需要继续请求)"""
36          dd_list = response.xpath('//dd')
37          for dd in dd_list:
38              item = TwoItem()
39              item['name'] = dd.xpath('./text()').get()
40              item['url'] = dd.xpath('./@href').get()
41              # 生成1个需要跟进的URL地址,将此item对象交给调度器入队列
42              yield scrapy.Request(
43                  url=item['url'], meta={'item': item}, callback=self.parse_two_page)
44
45      def parse_two_page(self, response):
46          """二级页面解析函数,提取内容(content)"""
47          item = response.meta['item']
48          item['content'] = response.xpath('//content/text()').get()
49
50          # 所有字段提取完成,yield给管道文件
51          yield item
52
```

```

53
54 【场景3】三级页面抓取
55 """
56 一级页面：名称(one_name)、详情页链接(one_url)-需要继续跟进
57 二级页面：名称(two_name)、下载页链接(two_url)-需要继续跟进
58 三级页面：具体所需内容(content)
59 """
60 import scrapy
61 from ..items import ThreeItem
62
63 class ThreeSpider(scrapy.Spider):
64     name = 'three'
65     allowed_domains = ['www.three.com']
66     start_urls = ['http://www.three.com/']
67
68     def parse(self, response):
69         """一级页面解析函数 - one_name、one_url"""
70         dd_list = response.xpath('//dd')
71         for dd in dd_list:
72             # 此item需要交给调度器入队列了
73             item = ThreeItem()
74             item['one_name'] = dd.xpath('./text()').get()
75             item['one_url'] = dd.xpath('./@href').get()
76             # 交给调度器入队列
77             yield scrapy.Request(
78                 url=item['one_url'], meta={'one_item': item}, callback=self.parse_two)
79
80     def parse_two(self, response):
81         """二级页面解析函数 - two_name、two_url(需要跟进的链接有多个)"""
82         one_item = response.meta['item']
83         li_list = response.xpath('//li')
84         for li in li_list:
85             # 此处提取的链接需要继续跟进,所以要创建item对象
86             item = ThreeItem()
87             item['two_name'] = li.xpath('./text()').get()
88             item['two_url'] = li.xpath('./@href').get()
89             # 此时item对象中只有 two_name、two_url, 并没有one_name、one_url
90             item['one_name'] = one_item['one_name']
91             item['one_url'] = one_item['one_url']
92             # 交给调度器入队列
93             yield scrapy.Request(
94                 url=item['two_url'], meta={'two_item': item}, callback=self.parse_three)
95
96     def parse_three(self, response):
97         """三级页面解析函数 - content"""
98         item = response.meta['two_item']
99         item['content'] = response.xpath('//content/text()').get()
100
101         # 至此,1条完整的item数据提取完成,交给调度器入队列
102         yield item

```

Day09笔记

腾讯招聘职位信息抓取

1、创建项目+爬虫文件

```
1 scrapy startproject Tencent
2 cd Tencent
3 scrapy genspider tencent careers.tencent.com
4
5 # 一级页面(postId):
6 https://careers.tencent.com/tencentcareer/api/post/Query?
timestamp=1566266592644&countryId=&cityId=&bgIds=&productId=&categoryId=&parentCategoryId=&a
ttrId=&keyword={}&pageIndex={}&pageSize=10&language=zh-cn&area=cn
7
8 # 二级页面(名称+类别+职责+要求+地址+时间)
9 https://careers.tencent.com/tencentcareer/api/post/ByPostId?timestamp=1566266695175&postId=
{}&language=zh-cn
```

2、定义爬取的数据结构

```
1 import scrapy
2
3 class TencentItem(scrapy.Item):
4     # 名称+类别+职责+要求+地址+时间
5     job_name = scrapy.Field()
6     job_type = scrapy.Field()
7     job_duty = scrapy.Field()
8     job_require = scrapy.Field()
9     job_address = scrapy.Field()
10    job_time = scrapy.Field()
11    # 具体职位链接
12    job_url = scrapy.Field()
13    post_id = scrapy.Field()
```

3、爬虫文件

```
1 # -*- coding: utf-8 -*-
2 import scrapy
3 from urllib import parse
4 import requests
5 import json
6 from ..items import TencentItem
7
8
9 class TencentSpider(scrapy.Spider):
10     name = 'tencent'
11     allowed_domains = ['careers.tencent.com']
12     # 定义常用变量
```

```

13     one_url = 'https://careers.tencent.com/tencentcareer/api/post/Query?
timestamp=1566266592644&countryId=&cityId=&bgIds=&productId=&categoryId=&parentCategoryId=&
attrId=&keyword={}&pageIndex={}&pageSize=10&language=zh-cn&area=cn'
14     two_url = 'https://careers.tencent.com/tencentcareer/api/post/ByPostId?
timestamp=1566266695175&postId={}&language=zh-cn'
15     headers = {'User-Agent': 'Mozilla/5.0'}
16     keyword = input('请输入职位类别:')
17     keyword = parse.quote(keyword)
18
19     # 重写start_requests()方法
20     def start_requests(self):
21         total = self.get_total()
22         # 生成一级页面所有页的URL地址,交给调度器
23         for index in range(1,total+1):
24             url = self.one_url.format(self.keyword,index)
25             yield scrapy.Request(url=url,callback=self.parse_one_page)
26
27     # 获取总页数
28     def get_total(self):
29         url = self.one_url.format(self.keyword, 1)
30         html = requests.get(url=url, headers=self.headers).json()
31         count = html['Data']['Count']
32         total = count//10 if count%10==0 else count//10 + 1
33
34         return total
35
36     def parse_one_page(self, response):
37         html = json.loads(response.text)
38         for one in html['Data']['Posts']:
39             # 此处是不是有URL需要交给调度器去入队列了? - 创建item对象!
40             item = TencentItem()
41             item['post_id'] = one['PostId']
42             item['job_url'] = self.two_url.format(item['post_id'])
43             # 创建1个item对象,请将其交给调度器入队列
44             yield scrapy.Request(url=item['job_url'],meta=
{'item':item},callback=self.detail_page)
45
46     def detail_page(self,response):
47         """二级页面:详情页数据解析"""
48         item = response.meta['item']
49         # 将响应内容转为python数据类型
50         html = json.loads(response.text)
51         # 名称+类别+职责+要求+地址+时间
52         item['job_name'] = html['Data']['RecruitPostName']
53         item['job_type'] = html['Data']['CategoryName']
54         item['job_duty'] = html['Data']['Responsibility']
55         item['job_require'] = html['Data']['Requirement']
56         item['job_address'] = html['Data']['LocationName']
57         item['job_time'] = html['Data']['LastUpdateTime']
58
59         # 至此: 1条完整数据提取完成,没有继续送往调度器的请求了,交给管道文件
60         yield item

```

■ 4、提前建库建表 - MySQL

```

1 create database tencentdb charset utf8;
2 use tencentdb;
3 create table tencenttab(
4     job_name varchar(500),
5     job_type varchar(200),
6     job_duty varchar(5000),
7     job_require varchar(5000),
8     job_address varchar(100),
9     job_time varchar(100)
10 )charset=utf8;

```

■ 5、管道文件

```

1 class TencentPipeline(object):
2     def process_item(self, item, spider):
3         return item
4
5 import pymysql
6 from .settings import *
7
8 class TencentMysqlPipeline(object):
9     def open_spider(self, spider):
10        """爬虫项目启动时,连接数据库1次"""
11        self.db = pymysql.connect(MYSQL_HOST, MYSQL_USER, MYSQL_PWD, MYSQL_DB, charset=CHARSET)
12        self.cursor = self.db.cursor()
13
14    def process_item(self, item, spider):
15        ins='insert into tencenttab values(%s,%s,%s,%s,%s,%s)'
16        job_li = [
17            item['job_name'],
18            item['job_type'],
19            item['job_duty'],
20            item['job_require'],
21            item['job_address'],
22            item['job_time']
23        ]
24        self.cursor.execute(ins, job_li)
25        self.db.commit()
26
27        return item
28
29    def close_spider(self, spider):
30        """爬虫项目结束时,断开数据库1次"""
31        self.cursor.close()
32        self.db.close()

```

■ 6、settings.py

```

1 ROBOTS_TXT = False
2 DOWNLOAD_DELAY = 0.5
3 DEFAULT_REQUEST_HEADERS = {
4     'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
5     'Accept-Language': 'en',
6     'User-Agent': 'Mozilla/5.0',
7 }

```



```

8 ITEM_PIPELINES = {
9     'Tencent.pipelines.TencentPipeline': 300,
10    'Tencent.pipelines.TencentMysqlPipeline': 500,
11 }
12 # MySQL相关变量
13 MYSQL_HOST = 'localhost'
14 MYSQL_USER = 'root'
15 MYSQL_PWD = '123456'
16 MYSQL_DB = 'tencentdb'
17 CHARSET = 'utf8'

```

盗墓笔记小说抓取

目标

```

1 【1】URL地址 : http://www.daomubiji.com/
2 【2】要求 : 抓取目标网站中盗墓笔记所有章节的所有小说的具体内容, 保存到本地文件
3     ./data/novel/盗墓笔记1:七星鲁王宫/七星鲁王_第一章_血尸.txt
4     ./data/novel/盗墓笔记1:七星鲁王宫/七星鲁王_第二章_五十年后.txt

```

准备工作xpath

```

1 【1】一级页面 - 大章节标题、链接:
2     1.1) 基准xpath匹配a节点对象列表: '//li[contains(@id,"menu-item-20")]/a'
3     1.2) 大章节标题: './text()'
4     1.3) 大章节链接: './@href'
5
6 【2】二级页面 - 小章节标题、链接
7     2.1) 基准xpath匹配article节点对象列表: '//article'
8     2.2) 小章节标题: './a/text()'
9     2.3) 小章节链接: './a/@href'
10
11 【3】三级页面 - 小说内容
12     3.1) p节点列表: '//article[@class="article-content"]/p/text()'
13     3.2) 利用join()进行拼接: ' '.join(['p1','p2','p3',''])

```

项目实施

1、创建项目及爬虫文件

```

1 scrapy startproject Daomu
2 cd Daomu
3 scrapy genspider daomu www.daomubiji.com

```

2、定义要爬取的数据结构 - itemspy

```

1 import scrapy
2
3 class DaomuItem(scrapy.Item):
4     # 1. 一级页面标题+链接
5     parent_title = scrapy.Field()
6     parent_url = scrapy.Field()
7     # 2. 二级页面标题+链接
8     son_title = scrapy.Field()
9     son_url = scrapy.Field()
10    # 3. 目录
11    directory = scrapy.Field()
12    # 4. 小说内容
13    content = scrapy.Field()

```

■ 3、爬虫文件实现数据抓取 - daomu.py

```

1 # -*- coding: utf-8 -*-
2 import scrapy
3 from ..items import DaomuItem
4 import os
5
6 class DaomuSpider(scrapy.Spider):
7     name = 'daomu'
8     allowed_domains = ['www.daomubiji.com']
9     start_urls = ['http://www.daomubiji.com/']
10
11    def parse(self, response):
12        """一级页面解析函数"""
13        # 基准xpath
14        a_list = response.xpath('//li[contains(@id,"menu-item-20")]/a')
15        for a in a_list:
16            # 此处是不是有需要继续交给调度器的请求了? - 创建item对象!!!
17            item = DaomuItem()
18            item['parent_title'] = a.xpath('./text()').get()
19            item['parent_url'] = a.xpath('./@href').get()
20            item['directory'] = './novel/{}/'.format(item['parent_title'])
21            if not os.path.exists(item['directory']):
22                os.makedirs(item['directory'])
23
24            # 交给调度器入队列
25            yield scrapy.Request(url=item['parent_url'], meta={'meta_1': item},
callback=self.detail_page)
26
27    def detail_page(self, response):
28        """二级页面解析: 提取小章节名称、链接"""
29        meta_1 = response.meta['meta_1']
30        # 基准xpath, 获取所有章节的节点对象列表
31        article_list = response.xpath('//article')
32        for article in article_list:
33            # 此处是不是有继续交给调度器的请求了? - 创建item对象!!!
34            item = DaomuItem()
35            item['son_title'] = article.xpath('./a/text()').get()
36            item['son_url'] = article.xpath('./a/@href').get()
37            item['parent_title'] = meta_1['parent_title']
38            item['parent_url'] = meta_1['parent_url']

```

```

39         item['directory'] = meta_1['directory']
40
41         # 交给调度器入队列,创建1个,交1个
42         yield scrapy.Request(url=item['son_url'],meta=
{'meta_2':item},callback=self.get_content)
43
44     def get_content(self,response):
45         """三级页面解析函数: 获取小说内容"""
46         item = response.meta['meta_2']
47         # p_list: ['段落1','段落2','段落3']
48         p_list = response.xpath('//article[@class="article-content"]//p/text()).extract()
49         content = '\n'.join(p_list)
50         item['content'] = content
51
52         # 没有继续交给调度器的请求了,所以不用创建item对象,直接交给管道文件处理
53         yield item

```

■ 4、管道文件实现数据处理 - pipelines.py

```

1 class DaomuPipeline(object):
2     def process_item(self, item, spider):
3         filename = item['directory'] + item['son_title'].replace(' ','_') + '.txt'
4         with open(filename,'w') as f:
5             f.write(item['content'])
6
7         return item

```

■ 5、全局配置 - setting.py

```

1 ROBOTSTXT_OBEY = False
2 DOWNLOAD_DELAY = 0.5
3 DEFAULT_REQUEST_HEADERS = {
4     'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
5     'Accept-Language': 'en',
6     'User-Agent': 'Mozilla/5.0'
7 }
8 ITEM_PIPELINES = {
9     'Daomu.pipelines.DaomuPipeline': 300,
10 }

```

图片管道(360图片抓取案例)

■ 目标

```

1 【1】 URL地址
2     1.1) www.so.com -> 图片 -> 美女
3     1.2) 即: https://image.so.com/z?ch=beauty
4
5 【2】 图片保存路径
6     ./images/xxx.jpg

```

■ 抓取网络数据包

```
1  【1】通过分析, 该网站为Ajax动态加载
2  【2】F12抓包, 抓取到json地址 和 查询参数(QueryString)
3      2.1) url = 'https://image.so.com/zjl?ch=beauty&sn={}&listtype=new&temp=1'
4      2.2) 查询参数
5              ch: beauty
6              sn: 0 # 发现sn的值在变, 0 30 60 90 120 ... ...
7              listtype: new
8              temp: 1
```

项目实施

■ 1、创建爬虫项目和爬虫文件

```
1 scrapy startproject So
2 cd So
3 scrapy genspider so image.so.com
```

■ 2、定义要爬取的数据结构(items.py)

```
1 img_url = scrapy.Field()
2 img_title = scrapy.Field()
```

■ 3、爬虫文件实现图片链接+名字抓取

```
1 import scrapy
2 import json
3 from ..items import SoItem
4
5 class SoSpider(scrapy.Spider):
6     name = 'so'
7     allowed_domains = ['image.so.com']
8     # 重写start_requests()方法
9     url = 'https://image.so.com/zjl?ch=beauty&sn={}&listtype=new&temp=1'
10
11     def start_requests(self):
12         for sn in range(0, 91, 30):
13             full_url = self.url.format(sn)
14             # 扔给调度器入队列
15             yield scrapy.Request(url=full_url, callback=self.parse_image)
16
17     def parse_image(self, response):
18         html = json.loads(response.text)
19         item = SoItem()
20         for img_dict in html['list']:
21             item['img_url'] = img_dict['qimg_url']
22             item['img_title'] = img_dict['title']
23
24         yield item
```

■ 4、管道文件 (pipelines.py)

```
1 from scrapy.pipelines.images import ImagesPipeline
2 import scrapy
3
4 class SoPipeline(ImagesPipeline):
5     # 重写get_media_requests()方法
6     def get_media_requests(self, item, info):
7         yield scrapy.Request(url=item['img_url'], meta={'name': item['img_title']})
8
9     # 重写file_path()方法, 自定义文件名
10    def file_path(self, request, response=None, info=None):
11        img_link = request.url
12        # request.meta属性
13        filename = request.meta['name'] + '.' + img_link.split('.')[ -1]
14        return filename
```

■ 5、全局配置(settings.py)

```
1 ROBOTSTXT_OBEY = False
2 DOWNLOAD_DELAY = 0.1
3 DEFAULT_REQUEST_HEADERS = {
4     'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
5     'Accept-Language': 'en',
6     'User-Agent': 'Mozilla/5.0',
7 }
8 ITEM_PIPELINES = {
9     'So.pipelines.SoPipeline': 300,
10 }
11 IMAGES_STORE = 'D:/AID1910/spider_day09_code/So/images/'
```

■ 6、运行爬虫(run.py)

```
1 from scrapy import cmdline
2
3 cmdline.execute('scrapy crawl so'.split())
```

图片管道使用方法总结

```

1  【1】爬虫文件：将图片链接yield到管道
2  【2】管道文件：
3      from scrapy.pipelines.images import ImagesPipeline
4      class XxxPipeline(ImagesPipeline):
5          def get_media_requests(self,xxx):
6              pass
7
8          def file_path(self,xxx):
9              pass
10
11  【3】 settings.py中：
12      IMAGES_STORE = '绝对路径'

```

文件管道使用方法总结

```

1  【1】爬虫文件：将文件链接yield到管道
2  【2】管道文件：
3      from scrapy.pipelines.images import FilesPipeline
4      class XxxPipeline(FilesPipeline):
5          def get_media_requests(self,xxx):
6              pass
7
8          def file_path(self,xxx):
9              return filename
10
11  【3】 settings.py中：
12      FILES_STORE = '绝对路径'

```

scrapy - post请求

▪ 方法+参数

```

1  scrapy.FormRequest(
2      url=posturl,
3      formdata=formdata,
4      callback=self.parse
5  )

```

有道翻译案例实现

▪ 1、创建项目+爬虫文件

```
1 scrapy startproject Youdao
2 cd Youdao
3 scrapy genspider youdao fanyi.youdao.com
```

■ 2、items.py

```
1 result = scrapy.Field()
```

■ 3、youdao.py

```
1 # -*- coding: utf-8 -*-
2 import scrapy
3 import time
4 import random
5 from hashlib import md5
6 import json
7 from ..items import YoudaoItem
8
9 class YoudaoSpider(scrapy.Spider):
10     name = 'youdao'
11     allowed_domains = ['fanyi.youdao.com']
12     word = input('请输入要翻译的单词:')
13
14     def start_requests(self):
15         post_url = 'http://fanyi.youdao.com/translate_o?smartresult=dict&smartresult=rule'
16         salt, sign, ts = self.get_salt_sign_ts(self.word)
17         formdata = {
18             'i': self.word,
19             'from': 'AUTO',
20             'to': 'AUTO',
21             'smartresult': 'dict',
22             'client': 'fanyideskweb',
23             'salt': salt,
24             'sign': sign,
25             'ts': ts,
26             'bv': 'cf156b581152bd0b259b90070b1120e6',
27             'doctype': 'json',
28             'version': '2.1',
29             'keyfrom': 'fanyi.web',
30             'action': 'FY_BY_REALTIME'
31         }
32         # 发送post请求的方法
33         yield scrapy.FormRequest(url=post_url, formdata=formdata)
34
35     def get_salt_sign_ts(self, word):
36         # salt
37         salt = str(int(time.time() * 1000)) + str(random.randint(0, 9))
38         # sign
39         string = "fanyideskweb" + word + salt + "n%A-rKaT5fb[Gy?;N5@Tj"
40         s = md5()
41         s.update(string.encode())
42         sign = s.hexdigest()
43         # ts
44         ts = str(int(time.time() * 1000))
45         return salt, sign, ts
```

```

46
47     def parse(self, response):
48         item = YoudaoItem()
49         html = json.loads(response.text)
50         item['result'] = html['translateResult'][0][0]['tgt']
51
52         yield item

```

■ 4、pipelines.py

```

1 class YoudaoPipeline(object):
2     def process_item(self, item, spider):
3         print('翻译结果:', item['result'])
4         return item

```

■ 5、settings.py

```

1 ROBOTSTXT_OBEY = False
2 LOG_LEVEL = 'WARNING'
3 COOKIES_ENABLED = False
4 DEFAULT_REQUEST_HEADERS = {
5     "Cookie": "OUTFOX_SEARCH_USER_ID=970246104@10.169.0.83;
OUTFOX_SEARCH_USER_ID_NCOO=570559528.1224236;
_ntes_nnid=96bc13a2f5ce64962adfd6a278467214,1551873108952;
JSESSIONID=aaae9i7p1XPlKaJH_gkYw; td_cookie=18446744072941336803;
SESSION_FROM_COOKIE=unknown; __rl__test__cookies=1565689460872",
6     "Referer": "http://fanyi.youdao.com/",
7     "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML,
like Gecko) Chrome/76.0.3809.100 Safari/537.36",
8 }
9 ITEM_PIPELINES = {
10     'Youdao.pipelines.YoudaoPipeline': 300,
11 }

```

scrapy添加cookie的三种方式

```

1 【1】修改 settings.py 文件
2     1.1) COOKIES_ENABLED = False -> 取消注释,开启cookie,检查headers中的cookie
3     1.2) DEFAULT_REQUEST_HEADERS = {} 添加Cookie
4
5 【2】利用cookies参数
6     1.1) settings.py: COOKIES_ENABLED = True # 修改为TRUE后, 检查 Request()方法中cookies
7     1.2) def start_requests(self):
8         yield scrapy.Request(url=url, cookies={}, callback=xxx)
9
10 【3】DownloadMiddleware设置中间件
11     3.1) settings.py: COOKIES_ENABLED = TRUE # 找Request()方法中cookies参数
12     3.2) middlewares.py
13         def process_request(self, request, spider):
14             request.cookies={}

```


scrapy shell的使用

▪ 定义

- 1 【1】调试蜘蛛的工具
- 2 【2】交互式shell, 可在不运行spider的前提下,快速调试 scrapy 代码(主要测试xpath表达式)

▪ 基本使用

```
1 # scrapy shell URL地址
2 *1、 request.url      : 请求URL地址
3 *2、 request.headers  : 请求头(字典)
4 *3、 request.meta     : item数据传递, 定义代理(字典)
5
6 4、 response.text     : 字符串
7 5、 response.body     : bytes
8 6、 response.xpath('')
9 7、 response.status   : HTTP响应码
10
11 # 可用方法
12 shelp() : 帮助
13 fetch(request) : 从给定的请求中获取新的响应, 并更新所有相关对象
14 view(response) : 在本地Web浏览器中打开给定的响应以进行检查
```

▪ scrapy.Request()参数

```
1 1、 url
2 2、 callback
3 3、 headers
4 4、 meta : 传递数据,定义代理
5 5、 dont_filter : 是否忽略域组限制
6     默认False,检查allowed_domains['']
7 6、 cookies
```

设置中间件(随机User-Agent)

▪ 少量User-Agent切换

```
1 【1】方法一 : settings.py
2     1.1) USER_AGENT = ''
3     1.2) DEFAULT_REQUEST_HEADERS = {}
4
5 【2】方法二 : 爬虫文件
6     yield scrapy.Request(url,callback=函数名,headers={})
```

▪ 大量User-Agent切换 (middlewares.py设置中间件)

```
1 【1】获取User-Agent方式
2     1.1) 方法1 : 新建Useragents.py,存放大量User-Agent, random模块随机切换
3     1.2) 方法2 : 安装fake_useragent模块(sudo pip3 install fack_useragent)
```

```

4         from fake_useragent import UserAgent
5         agent = UserAgent().random
6
7     【2】middlewares.py新建中间件类
8     class RandomUseragentMiddleware(object):
9         def process_request(self, request, spider):
10             agent = UserAgent().random
11             request.headers['User-Agent'] = agent
12
13     【3】settings.py添加此下载器中间件
14     DOWNLOADER_MIDDLEWARES = {'': 优先级}

```

设置中间件(随机代理)

```

1 class RandomProxyDownloaderMiddleware(object):
2     def process_request(self, request, spider):
3         request.meta['proxy'] = xxx
4
5     def process_exception(self, request, exception, spider):
6         return request

```

■ 练习

1 有道翻译,将cookie以中间件的方式添加的scrapy项目中

今日作业

```

1 【1】URL地址: http://www.1ppt.com/xiazai/
2 【2】目标:
3     2.1) 爬取所有栏目分类下的,所有页的PPT
4     2.2) 数据存放: /home/tarena/ppts/工作总结PPT/xxx
5                   /home/tarena/ppts/个人简历PPT/xxx
6 【提示】: 使用FilesPipeline,并重写方法

```