

爬虫案例2(下载文件和图片)

利用scrapy提供的2个数据流处理管道

1.FilesPipeline 2.ImagesPipeline

一.下载matplotlib例子进行文件下载的学习

爬取网址: <https://matplotlib.org/examples/index.html#>

详情页面: https://matplotlib.org/examples/api/agg_oo.html

步骤1.利用scrapy shell 分析页面

注意然后可以 fetch 重新请求页面 在分析详情页面

思路: 分析列表页面,分析详情页面

步骤2:思路实现

创建项目 —— 配置文件处理管道 —— 指定下载目录 —— 构建实体类 —— 构建逻辑

1.创建项目

```
MichaelYun:PycharmProjects Yun$ scrapy startproject matplotlib
New Scrapy project 'matplotlib', using template directory '/usr/local/lib/python
3.7/site-packages/scrapy/templates/project', created in:
  /Users/Yun/PycharmProjects/matplotlib

You can start your first spider with:
  cd matplotlib
  scrapy genspider example example.com
MichaelYun:PycharmProjects Yun$ cd matplotlib/
MichaelYun:matplotlib Yun$ scrapy genspider example matplotlib.org
Created spider 'example' using template 'basic' in module:
  matplotlib.spiders.example
MichaelYun:matplotlib Yun$
```

2.配置文件的管理的处理 并 指定下载目录

```

settings.py
55 DOWNLOADER_MIDDLEWARES = {
56     # 'matplotlib.middlewares.MatplotlibDownloaderMiddleware': 543,
57     #}
58
59 # Enable or disable extensions
60 # See https://doc.scrapy.org/en/latest/topics/extensions.html
61 #EXTENSIONS = {
62     # 'scrapy.extensions.telnet.TelnetConsole': None,
63     #}
64
65 # Configure item pipelines
66 # See https://doc.scrapy.org/en/latest/topics/item-pipeline.html
67 ITEM_PIPELINES = {
68     'scrapy.pipelines.files.FilesPipeline': 1,
69 }
70 FILES_STORE = 'examples_src'
71

```

3.配置实体类

```

1 import scrapy
2
3
4 class MatplotlibItem(scrapy.Item):
5     file_urls = scrapy.Field()
6     files = scrapy.Field()

```

4.实现逻辑代码

```

1 # -*- coding: utf-8 -*-
2 import scrapy
3 from scrapy.linkextractors import LinkExtractor
4 from ..items import MatplotlibItem
5
6
7 class ExampleSpider(scrapy.Spider):
8     name = 'example'
9     allowed_domains = ['matplotlib.org']
10    start_urls = ['https://matplotlib.org/examples/index.html']
11
12    def parse(self, response):
13        le = LinkExtractor(restrict_css='div.toc-tree-wrapper.compound', deny='/index.html$')
14        #deny 正则表达式 排除这个地址不提取
15        for link in le.extract_links(response):
16            yield scrapy.Request(link.url, callback=self.parse_example)
17
18    def parse_example(self, response):
19        href = response.css('a.reference.external::attr(href)').extract_first()
20        url = response.urljoin(href)
21        example = MatplotlibItem()
22        example['file_urls'] = [url] # file_urls 必须是一个集合
23        return example

```



额外的:

改进文件额外的下载

```

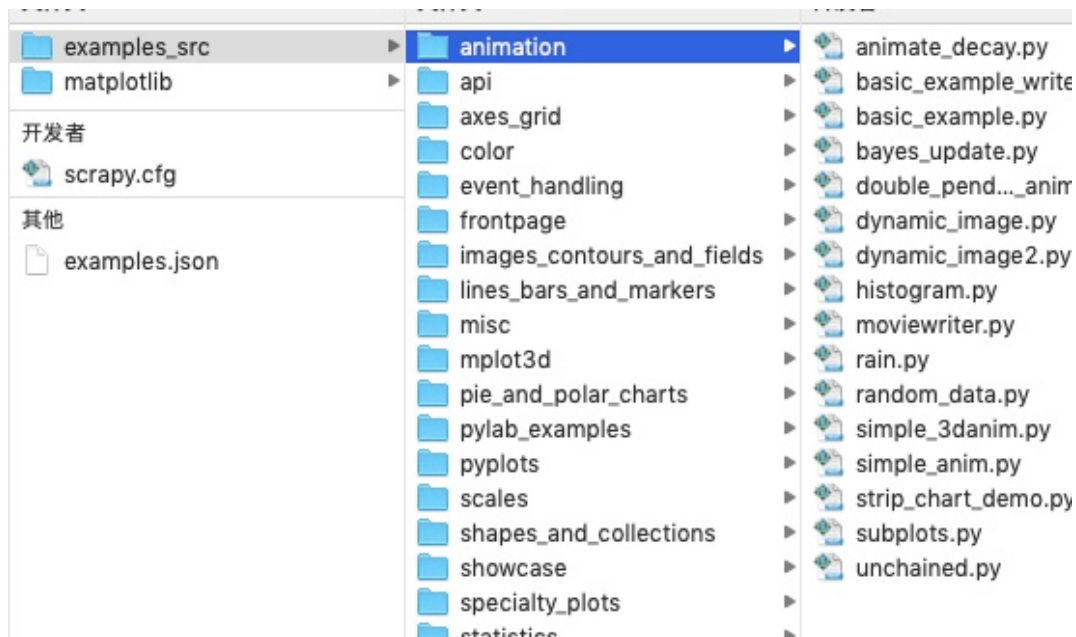
1 from scrapy.pipelines.files import FilesPipeline
2 from urllib.parse import urlparse
3 from os.path import basename,dirname,join
4
5 class MatplotlibPipeline(FilesPipeline):
6     def file_path(self, request, response=None, info=None):
7         path = urlparse(request.url).path
8         print(path)
9         return join(basename(dirname(path)), basename(path))

```

```

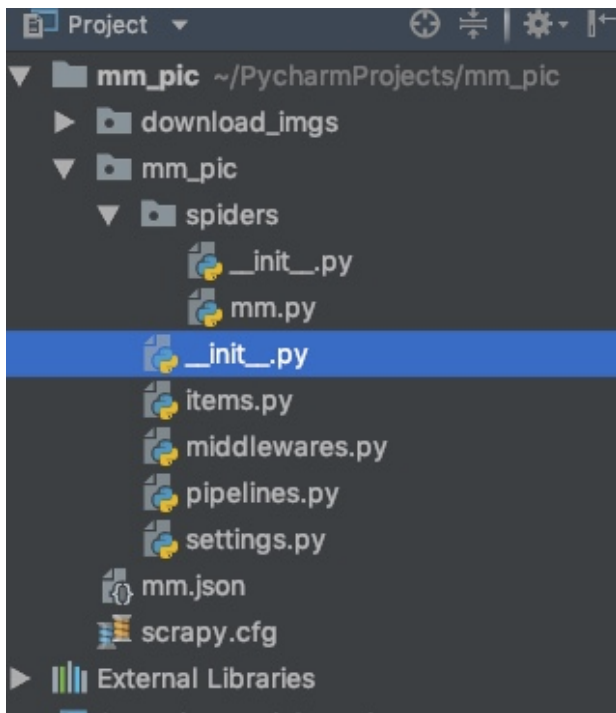
1 ITEM_PIPELINES = {
2     # 'scrapy.pipelines.files.FilesPipeline': 1,
3     'matplotlib.pipelines.MatplotlibPipeline':1,
4 }

```



-----分隔线(下载校花网图片)

步骤1.创建项目



步骤2 设置管道和位置

```
ITEM_PIPELINES = {  
    'mm_pic.pipelines.MmPicPipeline': 1,  
    # 'scrapy.pipelines.images.ImagesPipeline': 1  
}  
  
IMAGES_STORE = 'download_imgs'
```

步骤3 实现自定义管道

```
1 # -*- coding: utf-8 -*-  
2  
3 # Define your item pipelines here  
4 #  
5 # Don't forget to add your pipeline to the ITEM_PIPELINES setting  
6 # See: https://doc.scrapy.org/en/latest/topics/item-pipeline.html  
7  
8 from scrapy.pipelines.images import ImagesPipeline  
9 from scrapy import Request  
10 class MmPicPipeline(ImagesPipeline):  
11     def get_media_requests(self, item, info):  
12         # for image_url in item['image_urls']:  
13             # yield Request(meta={'mm_name': item['mm_name']})  
14             yield Request(item['image_urls'][0], meta={'mm': item['mm_name']})  
15  
16     def file_path(self, request, response=None, info=None):  
17         file = (request.meta['mm'])[0:2]  
18         filename = request.meta['mm']  
19         MM_FileName = u'full/{0}/{1}.jpg'.format(file, filename)  
20         return MM_FileName  
21  
22 # 顺序从上至下
```

步骤4.实现逻辑 处理路径-提取路径-返回对应的实体

```

1 # -*- coding: utf-8 -*-
2 import scrapy
3 import re
4 from ..items import MmPicItem
5 from scrapy.linkextractors import LinkExtractor
6
7
8 class MmSpider(scrapy.Spider):
9     name = 'mm'
10    allowed_domains = ['www.xiaohuar.com']
11    start_urls = ['http://www.xiaohuar.com/hua/']
12
13    def parse(self, response):
14        img_names = response.css('div.item_t img::attr(alt)').extract()
15        img_urls = response.css('div.item_t img::attr(src)').extract()
16        # 处理不带地址的链接
17        for index in range(len(img_urls)):
18            if re.match('/d',img_urls[index]):
19                img_urls[index] = response.urljoin(img_urls[index])
20        for index in range(len(img_names)):
21            mm = MmPicItem()
22            mm['mm_name'] = img_names[index]
23            mm['image_urls'] = [img_urls[index]]
24            yield mm
25        # 接着捕获下一页
26        le = LinkExtractor(restrict_css='div.page_num a:nth-last-child(2)')
27        links = le.extract_links(response)
28        if links[0].text == '下一页':
29            next_url = links[0].url
30            yield scrapy.Request(next_url, callback=self.parse)
31
32
33
34

```

步骤5.执行操作

scrapy crawl mm -o mm.json