

bash

```
1 scrapy startproject xxxPro  # 创建一个工程
2 cd xxxPro
3 scrapy genspider -t crawl xxx www.xxx.com
4 # run:
5 scrapy crawl sun
```

settings.py

```
1 USER_AGENT = 'Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36 (KHTML, like Gecko) \
2                                     Chrome/101.0.4951.64 Safari/537.36'
3 ROBOTSTXT_OBEY = False
4 LOG_LEVEL = 'ERROR'
```

items.py

```
1 import scrapy
2
3 class SunproItem(scrapy.Item):
4     title = scrapy.Field()
5     new_num = scrapy.Field()
6
7 class DetailItem(scrapy.Item):
8     new_id = scrapy.Field()
9     content = scrapy.Field()
```

pipelines.py

```
1 class SunproPipeline(object):
2     def process_item(self, item, spider):
3         if item.__class__.__name__ == 'DetailItem':
4             print(item['new_id'], item['content'])
5         else:
6             print(item['new_num'], item['title'])
```

main.py

```

1 from scrapy.linkextractors import LinkExtractor
2 from scrapy.spiders import CrawlSpider, Rule
3 from sunPro.items import SunproItem, DetailItem
4
5 class SunSpider(CrawlSpider):
6     name = 'sun'
7     # allowed_domains = ['www.xxx.com']
8     start_urls = ['http://www.xxx.com']
9     # 链接提取器：根据指定规则进行指定链接的提取
10    link = LinkExtractor(allow=r'type=4&page=\d+')
11    link_detail = LinkExtractor(allow=r'question/\d+/\d+\.html')
12
13    # 规则解析器
14    rules = (
15        Rule(link, callback='parse_item', follow=True)
16        Rule(link_detail, callback='parse_detail', follow=False)
17    )
18
19    def parse_item(self, response):
20        tr_list = response.xpath('//*[@id="morelist"]/div//table/tr')
21        for tr in tr_list:
22            new_num = tr.xpath('./td[1]/text()').extract_first()
23            new_title = tr.xpath('./td[2]/a[2]@title').extract_first()
24            item = SunproItem()
25            item['title'] = new_title
26            item['new_num'] = new_num
27
28            yield item
29
30    def parse_detail(self, response):
31        new_id = response.xpath('//span[2]/text()').extract_first()
32        new_content = response.xpath('//tr[1]/text()').extract()
33        new_content = ''.join(new_content)
34
35        item = DetailItem()
36        item['content'] = new_content
37        item['new_id'] = new_id
38
39        yield item

```