

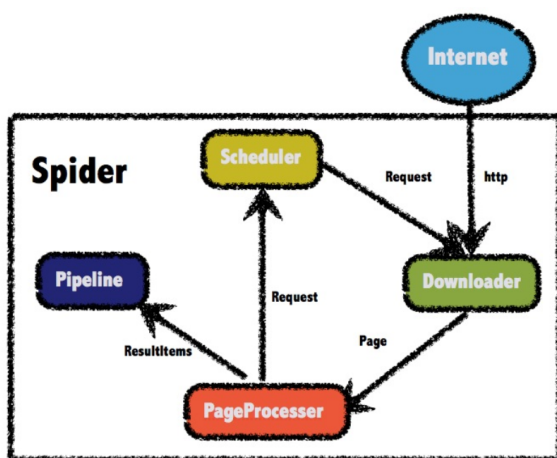
使用WebMagic爬取小说

一、环境搭建

- 开发工具： IDEA
- 爬虫工具： WebMagic 不使用Maven， 将下载的jar包在IDE的项目中导入

二、编写代码

0.理解WebMagic处理流程



- 只需编写自己的PageProcessor提取网页信息， Pipeline对信息进行持久化

1.MyPageProcessor类

- 继承自PageProcessor， 实现process函数定义如何处理页面信息
- 本次爬虫所爬取的页面分为列表页和文章页， 列表页需要抓取文章页url的地址并放入处理队列， 文章页需要抓取文本内容而无需额外抓取页面内url地址

- 使用Xpath与正则表达式提取所在网页Html内的所需信息

```
public void process(Page page) {  
    // 部分二：定义如何抽取页面信息，并保存下来  
  
    //区分列表页和文章页  
    if(page.getUrl().regex(s: "http://www.jianlaixiaoshuo.com/book/").match()){  
        //获取当前文章id  
        String id = page.getUrl().regex(s: "[1-9]\\d*").toString();  
        //更新最大id  
        if(Integer.parseInt(id) > maxPageId) maxPageId = Integer.parseInt(id);  
        //存储该页标题  
        page.putField(key: "title"+id, field: page.getHtml().xpath("//h1/text()").toString() + "\n");  
        //获取该页文本，分段，加换行并存储  
        List<String> text = page.getHtml().xpath("//div[@id='BookText']/p/text()").all();  
        String textString = "";  
        for(int i = 0; i < text.size(); i++){  
            textString = textString + text.get(i) + "\n";  
        }  
        page.putField(key: "text"+id, textString);  
    }else{  
        //将列表页的所有文章链接加入请求队列  
        page.addTargetRequests(page.getHtml().links().regex(s: "http://www.jianlaixiaoshuo.com/book/[1  
        maxPageId = 0;  
    }  
}
```

2.MyPipeline类

- 继承自Pipeline，实现process函数定义如何持久化爬取信息

- 使用文件输出流将爬取内容转化为txt文件中的文本内容

```
public void process(ResultItems items, Task tasks){
    Map<String,Object> map = items.getAll();
    try{
        File f = new File( pathname: "C:\\Users\\Frey\\Desktop\\MyNovel.txt");
        FileOutputStream fop = new FileOutputStream(f, append: true);
        String title;
        String text;
        for(int i = 1;i <= MyPageProcessor.maxPageId;i++) {
            title = "title" + String.valueOf(i);
            text = "text" + String.valueOf(i);
            if(map.get(title) != null){
                System.out.println(map.get(title).toString() + " ✓");
                fop.write(map.get(title).toString().getBytes());
                fop.write(map.get(text).toString().getBytes());
                fop.flush();
            }
        }
        //debug
        fop.close();
    }catch (IOException e) {
        System.out.print("Exception");
    }
}
```

3.测试与结果

- MyPageProcessor中可调整线程数与休眠时间，爬取时Console输出已经爬取到的文章页的

```
第二十四章 相赠
✓
第二十五章 离别
✓
第二十六章 好说话
✓
第二十七章 点睛
✓
第二十八章 财迷
✓
第二十九章 狐魅
✓
第三十章 暗室
✓
标题
```

- 在MyPipeline中可定义输出文件地址，默认为 "C:\\Users\\Frey\\Desktop\\MyNovel.txt" , 需要根据用户修改地址

● 结果文件截图

