

```
1 一、介绍
2 jsoup 是一款 Java 的 HTML 解析器，可直接解析某个 URL 地址、HTML 文本内容。它提供了一套非常省力的 API，可通过 DOM、CSS 以及类似
取出和操作数据。
3     1. 从一个URL，文件或字符串中解析HTML；
4     2. 使用DOM或CSS选择器来查找、取出数据；
5     3. 可操作HTML元素、属性、文本；
6 htmlunit是一款开源的Java页面分析工具，读取页面后，可以有效的使用htmlunit 分析页面上的内容。项目可以模拟浏览器运行，被誉为Java浏
界面的浏览器，运行速度也是非常迅速的。
7
8 二、依赖引用
9 <!-- jsoup -->
10 <dependency>
11     <groupId>org.jsoup</groupId>
12     <artifactId>jsoup</artifactId>
13     <version>1.11.3</version>
14 </dependency>
15 <!-- htmlunit -->
16 <dependency>
17     <groupId>net.sourceforge.htmlunit</groupId>
18     <artifactId>htmlunit</artifactId>
19     <version>2.34.1</version>
20 </dependency>
21
22 三、爬取图片
23 1. 下载工具类
24 public class JsoupUtils {
25     public static void downImages(String filePath, String fileName, String imgUrl) {
26         // 若指定文件夹没有，则先创建
27         File dir = new File(filePath);
28         if (!dir.exists()) {
29             dir.mkdirs();
30         }
31         // 写出的路径
32         File file = new File(filePath + File.separator + fileName);
33         try {
34             // 获取图片URL
35             URL url = new URL(imgUrl);
36             // 获得连接
37             URLConnection connection = url.openConnection();
38             // 设置10秒的相应时间
39             connection.setConnectTimeout(10 * 1000);
40             // 获得输入流
41             InputStream in = connection.getInputStream();
42             // 获得输出流
43             BufferedOutputStream out = new BufferedOutputStream(new FileOutputStream(file));
44             // 构建缓冲区
45             byte[] buf = new byte[1024];
46             int size;
47             // 写入到文件
48             while (-1 != (size = in.read(buf))) {
49                 out.write(buf, 0, size);
50             }
51             out.close();
52             in.close();
53         } catch (MalformedURLException e) {
```

```

54         e.printStackTrace();
55     } catch (IOException e) {
56         e.printStackTrace();
57     }
58 }
59
60 public static String formatNumber(Integer number){
61     return String.format("%s_%05d.jpg", "bd_img", number);
62 }
63
64 /**
65  * 获取页面文档字符串(等待异步JS执行)
66  */
67 public static Document getDocumentByHtmlUnit(String url) throws Exception {
68     WebClient webClient = new WebClient(BrowserVersion.CHROME);
69     WebClientOptions options = webClient.getOptions();
70     options.setUseInsecureSSL(true); //是否使用不安全的SSL
71     options.setThrowExceptionOnScriptError(false); //当JS执行出错的时候是否抛出异常
72     options.setThrowExceptionOnFailingStatusCode(false); //当HTTP的状态非200时是否抛出异常
73     options.setActiveXNative(false); //是否允许使用ActiveX
74     options.setCssEnabled(false); //是否启用CSS
75     options.setJavaScriptEnabled(true); //很重要，启用JS
76     options.setTimeout(30000); //设置"浏览器"的请求超时时间
77     options.setDoNotTrackEnabled(false); //不跟踪抓取
78     webClient.setAjaxController(new NicelyResynchronizingAjaxController()); //很重要，设置支持AJAX
79     webClient.setJavaScriptTimeout(30*1000); //设置JS执行的超时时间
80
81     HtmlPage htmlPage = webClient.getPage(url); //模拟浏览器打开一个目标网址
82     webClient.waitForBackgroundJavaScript(50); //该方法阻塞线程
83     String xml = htmlPage.asXml();
84     webClient.close();
85     return Jsoup.parse(xml);
86 }
87 }
88
89 2.测试
90 @Slf4j
91 public class JsoupTests extends BaseTest {
92
93     private final static String baseUri = "https://image.baidu.com/";
94     private final static String rootDir = "d:/img/";
95     private final static Map<String, Map<String, String>> groupMap = new HashMap<>(); //[[分组名, [图片名, 图片url]]
96     private final static AtomicInteger imgNo = new AtomicInteger(0);
97
98     @Test
99     public void testDownImages() throws Exception{
100         log.info("开始解析网页...");
101         Connection connect = Jsoup.connect(baseUri);
102         try {
103             //得到Document对象
104             Document document = connect.get();
105
106             //图片分组
107             Elements groups = document.getElementsByClass("bd-home-content-album-item");
108             for (Element group : groups) {
109                 Elements titles = group.getElementsByClass("bd-home-content-album-item-title");
110                 String groupTitle = titles.first().text();
111                 //获取分组详情页面
112                 String childUri = group.attr("abs:href");
113
114                 //Document childDoc = Jsoup.connect(childUri).get(); //获取静态页面

```

```
114 Document childDoc = JsoupUtils.getDocumentByHtmlUnit(childUri); //获取动态加载后的内容
115 Element imgDiv = childDoc.getElementById("imgList");
116 Elements imgList = imgDiv.getElementsByClass("albumsdetail-item-img");
117 Map<String, String> imgMap = new HashMap<>();
118 for (Element imgElement : imgList) {
119     String imgUrl = imgElement.attr("abs:src");
120     String imgName = JsoupUtils.formatNumber(imgNo.getAndIncrement());
121     imgMap.put(imgName, imgUrl);
122 }
123 groupMap.put(groupTitle, imgMap);
124 }
125
126 //下载图片
127 log.info("解析完成, 开始下载>>>");
128 for (Map.Entry<String, Map<String, String>> stringMapEntry : groupMap.entrySet()) {
129     String groupDir = stringMapEntry.getKey();
130     if (groupDir.contains(">")) continue;
131     log.info("分组下载: {}", groupDir);
132     stringMapEntry.getValue().forEach((fileName, fileUrl) ->
133         JsoupUtils.downloadImages(rootDir.concat(groupDir), fileName, fileUrl)
134     );
135 }
136 log.info("下载完成");
137 } catch (IOException e) {
138     e.printStackTrace();
139 }
140 }
141 }
```

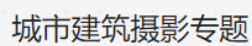
## 结果

```
>> ✓ Tests passed: 1 of 1 test - 13 sec 561 ms
ms 2022-05-10 14:52:50.776 ERROR 17556 --- [main] c.g.h.javascript.StrictErrorReporter : runtimeError: message=[An invalid or illegal selector was :
ms 2022-05-10 14:52:50.854 INFO 17556 --- [main] .g.h.NicelyResynchronizingAjaxController : Re-synchronized call to https://image.baidu.com/user/login
2022-05-10 14:52:50.901 WARN 17556 --- [main] c.g.htmlunit.IncorrectnessListenerImpl : Obsolete content type encountered: 'text/javascript'.
2022-05-10 14:52:51.089 WARN 17556 --- [Client@21f7c0be] c.g.htmlunit.IncorrectnessListenerImpl : Obsolete content type encountered: 'application/x-javascrip
2022-05-10 14:52:51.095 INFO 17556 --- [main] com.formssi.mall.redis.JsoupTests : 解析完成, 开始下载>>>
2022-05-10 14:52:51.095 INFO 17556 --- [main] com.formssi.mall.redis.JsoupTests : 分组下载: 渐变风格插画
2022-05-10 14:52:51.625 INFO 17556 --- [main] com.formssi.mall.redis.JsoupTests : 分组下载: 城市建筑摄影专题
2022-05-10 14:52:52.231 INFO 17556 --- [main] com.formssi.mall.redis.JsoupTests : 分组下载: 航拍地球系列
2022-05-10 14:52:52.703 INFO 17556 --- [main] com.formssi.mall.redis.JsoupTests : 分组下载: 宠物图片
2022-05-10 14:52:53.075 INFO 17556 --- [main] com.formssi.mall.redis.JsoupTests : 分组下载: 皮影
2022-05-10 14:52:53.751 INFO 17556 --- [main] com.formssi.mall.redis.JsoupTests : 下载完成
```

## 效果







共 791 张图片

img.albumsdetail-item-img 310.39 × 310



```
<script type="text/javascript">...</script>  
class="bd-albumsdetail-wrapper">  
  class="bd-albumsdetail-header"></div>  
  id="bd-albumsdetail-content" style="width: 1616px;">  
    iv class="albumsdetail-cover clearfix">...</div>  
    iv id="topWidget" class="albumsdetail-topwidget cleagearfix" style="width: 1616px;"></div>  
    iv id="noPage">...</div>  
    iv id="imgList">  
      <div class="albumsdetail-column" style="width: 310.4px;">  
        <a class="albumsdetail-item" href="/search/detail?tn=baiduiimagetail&word=%E5%9F%8E%E5%B8%82%E5%B8%BA%E7%AD_%2F%2f7.baidu.com%2Fit%2Fu%3D1595072465%C3644073269%26fm%3D193%26fx%3DGIF" target="_blank" index="0" width="310.4" style="width: 310.4px; height: 310px;">  
           == $0  
          <div class="albumsdetail-item-inner-border"></div>  
        </a>  
        <a class="albumsdetail-item" href="/search/detail?tn=baiduiimagetail&word=%E5%9F%8E%E5%B8%82%E5%B8%BA%E7%AD_%2F%2f7.baidu.com%2Fit%2Fu%3D33779234486%C2C1094031034%26fm%3D193%26fx%3DGIF" target="_blank" index="9" width="310.4" style="width: 310.4px; height: 206px;">...</a>  
        <a class="albumsdetail-item" href="/search/detail?tn=baiduiimagetail&word=%E5%9F%8E%E5%B8%82%E5%B8%BA%E7%AD_%2F%2f7.baidu.com%2Fit%2Fu%3D18887420%C2894941323%26fm%3D193%26fx%3DGIF" target="_blank" index="14" width="310.4" style="width: 310.4px; height: 206px;">...</a>  
        <a class="albumsdetail-item" href="/search/detail?tn=baiduiimagetail&word=%E5%9F%8E%E5%B8%82%E5%B8%BA%E7%AD_%2F%2f7.baidu.com%2Fit%2Fu%3D1285847167%C2C3193778276%26fm%3D193%26fx%3DGIF" target="_blank" index="19" width="310.4" style="width: 310.4px; height: 206px;">...</a>  
        <a class="albumsdetail-item" href="/search/detail?tn=baiduiimagetail&word=%E5%9F%8E%E5%B8%82%E5%B8%BA%E7%AD_%2F%2f7.baidu.com%2Fit%2Fu%3D3124693600%C2C356058981%26fm%3D193%26fx%3DGIF" target="_blank" index="24" width="310.4" style="width: 310.4px; height: 206px;">...</a>
```