

转载链接：<http://www.xdemo.org/jsoup-html-parse/>

Jsoup应该说是最简单快速的Html解析程序了，完善的API以及与JS类似的操作方式，为Java的Html解析带来极大的方便，结合多线程适合做一些网络数据的抓取，本文从一下几个方面介绍一下，篇幅有限，化繁为简。

下载Jsoup<http://jsoup.org/download>

查看官方提供的手册：<http://jsoup.org/cookbook/>

教程：<https://my.oschina.net/flashsword/blog?catalog=390084>

官方文档中文版：<http://www.open-open.com/jsoup/parsing-a-document.htm>

1. 获取一个Document，这是Jsoup最核心的一个对象

有三种途径来加载Document：字符串，URL地址，文件

```
1  /**
2   *
3   */
4  package org.xdemo.example.jsoupdemo.input;
5
6  import java.io.File;
7  import java.io.IOException;
8
9  import org.jsoup.Jsoup;
10 import org.jsoup.nodes.Document;
11 import org.jsoup.nodes.Element;
12 import org.jsoup.safety.Whitelist;
13
14 /**
15  * @作者 Goofy
16  * @邮件 252878950@qq.com
17  * @日期 2014-4-2上午10:54:53
18  * @描述
19  */
20 public class ParseDocument {
21
22     /**
23      * 将String转换成Document
24      * @return org.jsoup.nodes.Document
25      */
26     public static Document parseHtmlFromString(){
27         String html = "<html><head><title>标题</title></head>"
28             + "<body><p>段落</p></body></html>";
29         Document doc = Jsoup.parse(html);
30         return doc;
31     }
32
33     /**
34      * 注意：这是一个不安全的方法
35      * 将String转换成Html片段,注意防止跨站脚本攻击
36      * @return Element
37      */
38     public static Element parseHtmlFragmentFromStringNotSafe(){
39         String html = "<div><p>Lorem ipsum.</p>";
40         Document doc = Jsoup.parseBodyFragment(html);
41         Element body = doc.body();
42         return body;
43     }
44 }
```

```

44
45 /**
46  * 这是一个安全的方法
47  * 将String转换成Html片段,注意防止跨站脚本攻击
48  * @return Element
49  */
50 public static Element parseHtmlFragmentFromStringSafe(){
51     String html = "<div><p>Lorem ipsum.</p>";
52     //白名单列表定义了哪些元素和属性可以通过过滤器,其他的元素和属性一律移除
53     Whitelist wl=new Whitelist();
54     //比较松散的过滤,包括
55     // "a", "b", "blockquote", "br", "caption", "cite", "code", "col",
56     // "colgroup", "dd", "div", "dl", "dt", "em", "h1", "h2", "h3", "h4", "h5", "h6"
57     // "i", "img", "li", "ol", "p", "pre", "q", "small", "strike", "strong",
58     // "sub", "sup", "table", "tbody", "td", "tfoot", "th", "thead", "tr", "u",
59     // "ul"
60     Whitelist.relaxed();
61     //没有任何标签,只有文本
62     Whitelist.none();
63     //常规的过滤器
64     // "a", "b", "blockquote", "br", "cite", "code", "dd", "dl", "dt", "em",
65     // "i", "li", "ol", "p", "pre", "q", "small", "strike", "strong", "sub",
66     // "sup", "u", "ul"
67     Whitelist.basic();
68     //常规的过滤器,多了一个img标签
69     Whitelist.basicWithImages();
70     //文本类型的标签
71     // "b", "em", "i", "strong", "u"
72     Whitelist.simpleText();
73     //另外还可以自定义过滤规则,例如
74     wl.addTags("a");
75     //执行过滤
76     Jsoup.clean(html, wl);
77     Document doc = Jsoup.parseBodyFragment(html);
78     Element body = doc.body();
79     return body;
80 }
81
82 /**
83  * 从URL加载
84  * @return Document
85  */
86 public static Document parseDocumentFromUrl(){
87     Document doc = null;
88     try {
89         doc = Jsoup.connect("http://www.google.com/").get();
90         //获取标题
91         String title = doc.title();
92         System.out.println(title); //输出: Google
93         //data(key,value)是该URL要求的参数
94         //userAgent制定用户使用的代理类型
95         //cookie带上cookie,如cookie("JSESSIONID","FDE234242342342423432432")
96         //连接超时时间
97         //post或者get方法
98         doc = Jsoup.connect("http://www.xxxxx.com/")
99             .data("query", "Java")
100             .userAgent("Mozilla")
101             .cookie("auth", "token")
102             .timeout(3000)
103             .post();
104
105     } catch (IOException e) {
106         e.printStackTrace();
107     }
108     return doc;
109 }
110 /**
111  * 从文件加载
112  * @return Document
113  */
114 public static Document parseDocumentFromFile(){
115     File input = new File("/tmp/input.html");
116     Document doc=null;
117     try {
118         //从文件加载Document文档

```

```

119         doc = Jsoup.parse(input, "UTF-8");
120         System.out.println(doc.title());
121     } catch (IOException e) {
122         e.printStackTrace();
123     }
124     return doc;
125 }
126
127
128
129 }

```

2.选择器

```

1  package org.xdemo.example.jsoupdemo.extracter;
2
3  import java.util.regex.Pattern;
4
5  import org.jsoup.Jsoup;
6  import org.jsoup.nodes.Document;
7  import org.jsoup.nodes.Element;
8
9  /**
10   * @作者 Goofy
11   * @邮件 252878950@qq.com
12   * @日期 2014-4-2上午10:41:19
13   * @描述 选择器 操作示例
14   */
15  public class Selector {
16
17      public static void main(String[] args) {
18          Document doc;
19          try {
20              //获取文档
21              doc=Jsoup.connect("http://xxx.com/").get();
22
23              /*****获取单一元素*****/
24              //与JS类似的根据ID选择的选择器<div id="content"></div>
25              Element content = doc.getElementById("content");
26
27              /*****一下方法的返回值都是Elements集合*****/
28
29              //获取所有的a标签<a href="#"></a>
30              content.getElementsByTag("a");
31              //类选择器<div></div>
32              doc.getElementsByClass("divClass");
33              //获取Document的所有元素
34              doc.getAllElements();
35              //根据属性获取元素<a href="#"></a>
36              doc.getElementsByAttribute("href");
37              //根据属性前缀获取元素 <li data-name="Peter Liu" data-city="ShangHai" data-lang="CSharp"
38              doc.getElementsByAttributeStarting("data-");
39              //根据key-value选择如<a href="http://xdemo.org"></a>
40              doc.getElementsByAttributeValue("href", "http://xdemo.org");
41              //和上面的正好相反
42              doc.getElementsByAttributeValueNot("href", "http://xdemo.org");
43              //根据key-value,其中value可能是key对应属性的一个子字符串,选择如<a href="http://xdemo.org">
44              doc.getElementsByAttributeValueContaining("href", "xdemo");
45              //根据key-value,其中key对应值的结尾是value,选择如<a href="http://xdemo.org"></a>
46              doc.getElementsByAttributeValueEnding("href", "org");
47              //和上面的正好相反
48              doc.getElementsByAttributeValueStarting("href", "http://xdemo");
49              //正则匹配, value需要满足正则表达式, <a href="http://xdemo.org"></a>,如href的值含有汉字
50              doc.getElementsByAttributeValueMatching("href", Pattern.compile("[\u4e00-\u9fa5]"));
51              //同上
52              doc.getElementsByAttributeValueMatching("href", "[\u4e00-\u9fa5]");
53              //根据元素所在的z-index获取元素
54              doc.getElementsByIndexEquals(0);

```

```
55 //获取z-index大于x的元素
56 doc.getElementsByIndexGreaterThan(0);
57 //和上面的正好相反
58 doc.getElementsByIndexLessThan(10);
59
60 //遍历标签
61 for (Element link : content.getElementsByTag("a")) {
62     String linkHref = link.attr("href");
63     String linkText = link.text();
64 }
65
66 /*****一些其他常用的方法*****/
67 //获取网页标题
68 doc.title();
69 //获取页面的所有文本
70 doc.text();
71
72 //为元素添加一个css class
73 content.addClass("newClass");
74 //根据属性获取值
75 content.attr("id");
76 //获取所有子元素
77 content.children();
78 //获取元素内的所有文本
79 content.text();
80 //获取同级元素
81 content.siblingElements();
82
83
84 } catch (Exception e) {
85     e.printStackTrace();
86 }
87
88 }
89
90 }
```

3.最后说一点，就是安全问题，解析html的时候要防止跨站脚本攻击cross-site scripting (XSS)，作者也考虑到了这一点，所以真正使用时需要注意。

[转载请注明来源:<http://www.xdemo.org/jsoup-html-parse/>]