

Evan

Only let oneself become strong enough, good enough, can afford the life that you want to.

☰ 目录视图

☰ 摘要视图

🔒 订阅专栏

从创业到再就业，浅述对程序员职业生涯的看法 征文 | 你会为 AI 转型么？ 赠书：7月大咖新书机器学习/Android/python

使用Jsoup 抓取页面的数据

标签：jsoup html android java

2015-12-22 09:46 299人阅读 评论(0) 收藏 举报

☰ 分类： 移动开发 ( 38 ) ▾

需要使用的是jsoup-1.7.3. jar包 如果需要看文档我下载请借一步到官网：<http://jsoup.org/>

这里贴一下我用到的 Java工程的测试代码

```
[java]
01. <span style="font-size:14px;">package com.javen.Jsoup;
02.
03. import java.io.IOException;
04.
05. import org.jsoup.Jsoup;
06. import org.jsoup.nodes.Document;
07. import org.jsoup.nodes.Element;
08. import org.jsoup.select.Elements;
09.
10. public class JsoupTest {
11.     static String url="http://www.cnblogs.com/zyw-205520/archive/2012/12/20/2826402.html";
12.     /**
13.      * @param args
14.      * @throws Exception
15.      */
16.     public static void main(String[] args) throws Exception {
17.
18.         // TODO Auto-generated method stub
19.         BolgBody();
20.         //test();
21.         //Blog();
22.         /*
23.          * Document doc = Jsoup.connect("http://www.oschina.net/")
24.          * .data("query", "Java") // 请求参数 .userAgent("I ' m jsoup") // 设置
25.          * User-Agent .cookie("auth", "token") // 设置 cookie .timeout(3000) //
26.          * 设置连接超时时间 .post();
27.          */// 使用 POST 方法访问 URL
28.
29.         /*
30.          * // 从文件中加载 HTML 文档 File input = new File("D:/test.html"); Document doc
31.          * = Jsoup.parse(input,"UTF-8","http://www.oschina.net/");
32.          */
33.     }
34.
35.     /**
36.      * 获取指定HTML 文档指定的body
37.      * @throws IOException
38.      */
39.     private static void BolgBody() throws IOException {
40.         // 直接从字符串中输入 HTML 文档
41.         String html = "<html><head><title> 开源中国社区 </title></head>"
42.             + "<body><p> 这里是 jsoup 项目的相关文章 </p></body></html>";
43.         Document doc = Jsoup.parse(html);
44.         System.out.println(doc.body());
45.
46.     }
```

```

47. // 从 URL 直接加载 HTML 文档
48. Document doc2 = Jsoup.connect(url).get();
49. String title = doc2.body().toString();
50. System.out.println(title);
51. }
52.
53. /**
54.  * 获取博客上的文章标题和链接
55.  */
56. public static void article() {
57.     Document doc;
58.     try {
59.         doc = Jsoup.connect("http://www.cnblogs.com/zyw-205520/").get();
60.         Elements ListDiv = doc.getElementsByAttributeValue("class", "postTitle");
61.         for (Element element : ListDiv) {
62.             Elements links = element.getElementsByTag("a");
63.             for (Element link : links) {
64.                 String linkHref = link.attr("href");
65.                 String linkText = link.text().trim();
66.                 System.out.println(linkHref);
67.                 System.out.println(linkText);
68.             }
69.         }
70.     } catch (IOException e) {
71.         // TODO Auto-generated catch block
72.         e.printStackTrace();
73.     }
74. }
75.
76. /**
77.  * 获取指定博客文章的内容
78.  */
79. public static void Blog() {
80.     Document doc;
81.     try {
82.         doc = Jsoup.connect("http://www.cnblogs.com/zyw-205520/archive/2012/12/20/2826402.html").get();
83.         Elements ListDiv = doc.getElementsByAttributeValue("class", "postBody");
84.         for (Element element : ListDiv) {
85.             System.out.println(element.html());
86.         }
87.     } catch (IOException e) {
88.         // TODO Auto-generated catch block
89.         e.printStackTrace();
90.     }
91. }
92. }
93.
94. }</span>

```

下面来介绍android中使用Jsoup异步解析网页的数据 **请注意：这里很容易遇到一个乱码的问题**

1. 配置文件：AndroidManifest.xml中加 权限 <uses-permission android:name="android.permission.INTERNET"></uses-permission>
2. layout的布局文件

```

[html]
01. <span style="font-size:14px;"><LinearLayout xmlns:android="http://schemas.android.com/apk/res/android"
02.     xmlns:tools="http://schemas.android.com/tools"
03.     android:layout_width="match_parent"
04.     android:layout_height="match_parent"
05.     android:orientation="vertical" >
06.
07.     <WebView
08.         android:id="@+id/webView"
09.         android:layout_width="fill_parent"
10.         android:layout_height="200dp" />
11.
12.     <ScrollView
13.         android:layout_width="wrap_content"
14.         android:layout_height="wrap_content" >
15.
16.         <TextView
17.             android:id="@+id/textView"

```

```
18.         android:layout_width="wrap_content"
19.         android:layout_height="wrap_content"
20.         android:text="@string/hello_world" />
21.     </ScrollView>
22.
23. </LinearLayout></span>
```

### 3.主要异步加载数据的代码

```
[java]
01. <span style="font-size:14px;">package com.javen.aaa;
02.
03. import java.io.BufferedReader;
04. import java.io.IOException;
05. import java.io.InputStreamReader;
06. import java.net.URL;
07.
08. import org.jsoup.Jsoup;
09. import org.jsoup.nodes.Document;
10. import org.jsoup.nodes.Element;
11. import org.jsoup.select.Elements;
12.
13. import android.app.Activity;
14. import android.app.Dialog;
15. import android.app.ProgressDialog;
16. import android.os.AsyncTask;
17. import android.os.Bundle;
18. import android.util.Log;
19. import android.webkit.WebView;
20. import android.widget.TextView;
21.
22. public class MainActivity extends Activity {
23.     private WebView webView;
24.     private TextView textView;
25.     private static final int DIALOG_KEY = 0;
26.     @Override
27.     protected void onCreate(Bundle savedInstanceState) {
28.         super.onCreate(savedInstanceState);
29.         setContentView(R.layout.main);
30.         webView = (WebView) findViewById(R.id.webView);
31.         textView=(TextView) findViewById(R.id.textView);
32.         try {
33.             ProgressAsyncTask asyncTask=new ProgressAsyncTask(webView,textView);
34.             asyncTask.execute(10000);
35.         } catch (Exception e) {
36.             // TODO Auto-generated catch block
37.             e.printStackTrace();
38.         }
39.     }
40.
41.     public String test() {
42.         StringBuffer buffer=new StringBuffer();
43.         Document doc;
44.         try {
45.             doc = Jsoup.connect("http://www.cnblogs.com/zyw-205520/").get();
46.             Elements ListDiv = doc.getElementsByAttributeValue("class","postTitle");
47.             for (Element element :ListDiv) {
48.                 Elements links = element.getElementsByTag("a");
49.                 for (Element link : links) {
50.                     String linkHref = link.attr("href");
51.                     String linkText = link.text().trim();
52.                     buffer.append("linkHref=="+linkHref);
53.                     buffer.append("linkText=="+linkText);
54.
55.                     System.out.println(linkHref);
56.                     System.out.println(linkText);
57.                 }
58.             }
59.         } catch (IOException e) {
60.             // TODO Auto-generated catch block
61.             e.printStackTrace();
62.         }
63.         return buffer.toString();
64.
65.     }
```

```

66.
67. // 弹出"查看"对话框
68. @Override
69. protected Dialog onCreateDialog(int id) {
70.     switch (id) {
71.         case DIALOG_KEY: {
72.             ProgressDialog dialog = new ProgressDialog(this);
73.             dialog.setMessage("获取数据中 请稍候...");
74.             dialog.setIndeterminate(true);
75.             dialog.setCancelable(true);
76.             return dialog;
77.         }
78.     }
79.     return null;
80. }
81.
82. public static String readHtml(String myurl) {
83.     StringBuffer sb = new StringBuffer("");
84.     URL url;
85.     try {
86.         url = new URL(myurl);
87.         BufferedReader br = new BufferedReader(new InputStreamReader(url.openStream(), "gbk"));
88.         String s = "";
89.         while ((s = br.readLine()) != null) {
90.             sb.append(s + "\r\n");
91.         }
92.     } catch (Exception e) {
93.         e.printStackTrace();
94.     }
95.     return sb.toString();
96. }
97.
98. class ProgressAsyncTask extends AsyncTask<Integer, Integer, String> {
99.
100.     private WebView webView;
101.     private TextView textView;
102.     public ProgressAsyncTask(WebView webView, TextView textView) {
103.         super();
104.         this.webView=webView;
105.         this.textView=textView;
106.     }
107.
108.     /**
109.      * 这里的Integer参数对应AsyncTask中的第一个参数 这里的String返回值对应AsyncTask的第三个参数
110.      * 该方法并不运行在UI线程当中，主要用于异步操作，所有在该方法中不能对UI当中的空间进行设置和修改
111.      * 但是可以调用publish Progress方法触发onProgressUpdate对UI进行操作
112.      */
113.     @Override
114.     protected String doInBackground(Integer... params) {
115.         String str =null;
116.         Document doc = null;
117.         try {
118.             // String url ="http://www.cnblogs.com/zyw-205520/p/3355681.html";
119.             //
120.             // doc= Jsoup.parse(new URL(url).openStream(),"utf-8", url);
121.             // //doc = Jsoup.parse(readHtml(url));
122.             // //doc=Jsoup.connect(url).get();
123.             // str=doc.body().toString();
124.             doc = Jsoup.connect("http://www.cnblogs.com/zyw-205520/archive/2012/12/20/2826402.html").get();
125.             Elements ListDiv = doc.getElementsByAttributeValue("class", "postBody");
126.             for (Element element :ListDiv) {
127.                 str=element.html();
128.                 System.out.println(element.html());
129.             }
130.             Log.d("doInBackground", str.toString());
131.             System.out.println(str);
132.             //你可以试试GBK或UTF-8
133.         } catch (Exception e) {
134.             // TODO Auto-generated catch block
135.             e.printStackTrace();
136.         }
137.         return str.toString() ;
138.         //return test();
139.     }
140.
141.     /**
142.      * 这里的String参数对应AsyncTask中的第三个参数（也就是接收doInBackground的返回值）
143.      * 在doInBackground方法执行结束之后在运行，并且运行在UI线程当中 可以对UI空间进行设置
144.      */

```

```
145.         @Override
146.         protected void onPostExecute(String result) {
147.             webView.loadData(result, "text/html;charset=utf-8", null);
148.             textView.setText(result);
149.             removeDialog(DIALOG_KEY);
150.         }
151.
152.         // 该方法运行在UI线程当中,并且运行在UI线程当中 可以对UI空间进行设置
153.         @Override
154.         protected void onPreExecute() {
155.             showDialog(DIALOG_KEY);
156.         }
157.
158.         /**
159.          * 这里的Integer参数对应AsyncTask中的第二个参数
160.          * 在doInBackground方法当中, 每次调用publishProgress方法都会触发onProgressUpdate执行
161.          * onProgressUpdate是在UI线程中执行, 所有可以对UI空间进行操作
162.          */
163.         @Override
164.         protected void onProgressUpdate(Integer... values) {
165.
166.         }
167.     }
168.
169. }</span>
```

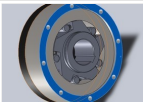
原文 : <http://www.cnblogs.com/zyw-205520/p/3421687.html>

顶 1      踩 0

- 上一篇    Fragment+ ViewPager实现仿微信点击和滑动切换界面
- 下一篇    android AsyncTask介绍 AsyncTask和Handler对比

相关文章推荐

- 使用Jsoup 抓取页面的数据
- jsoup抓取其他网站的页面代码
- 使用java的html解析器jsoup和jQuery实现一个白...
- android 使用Jsoup 抓取页面的数据
- HttpURLConnection 和HttpClient+Jsoup处理...
- Android 利用jsoup 抓取腾讯应用市场的软件APP...
- jSoup Cookbook-提取数据 7 使用选择器语法查...
- 在android中使用jsoup解析页面链接
- jsoup 默认抓取页面大小为1M
- 使用Jsoup抓取页面的数据



逆止器



python培训



网站制作



电磁加热设备



电动执行器



社保代理代缴



射频功率放大



断路器测试仪



网站建设

猜你在找

- |                     |                             |
|---------------------|-----------------------------|
| 机器学习之概率与统计推断        | 机器学习之数学基础                   |
| 机器学习之凸优化            | 机器学习之矩阵                     |
| 响应式布局全新探索           | 探究Linux的总线、设备、驱动模型          |
| 深度学习基础与TensorFlow实践 | 深度学习之神经网络原理与实战技巧            |
| 前端开发在线峰会            | TensorFlow实战进阶：手把手教你做图像识别应用 |

查看评论