

java实现爬虫爬网站图片 - javaxiaojian的专栏 - 博客频道

版权声明：本文为博主原创文章，未经博主允许不得转载。

第一步，实现 LinkQueue，对url进行过滤和存储的操作

```
1. import java.util.ArrayList;
2. import java.util.Collections;
3. import java.util.HashSet;
4. import java.util.List;
5. import java.util.Set;
6. public class LinkQueue {
7.     // 已访问的 url 集合
8.     private static Set<String> visitedUrl = Collections.synchronizedSet(new HashSet<String>
9.         ());
10.    // 未访问的url
11.    private static List<String> unVisitedUrl = Collections.synchronizedList(new ArrayList<String>
12.        >());
13.    // 未访问的URL出队列
14.    public static String unVisitedUrlDeQueue() {
15.        if (unVisitedUrl.size() > 0) {
16.            String url = unVisitedUrl.remove(0);
17.            visitedUrl.add(url);
18.            return url;
19.        }
20.        return null;
21.    }
22.    // 新的url添加进来的时候进行验证，保证只是添加一次
23.    public static void addUnvisitedUrl(String url) {
24.        if (url != null && !url.trim().equals("") && !visitedUrl.contains(url)
25.            && !unVisitedUrl.contains(url))
26.            unVisitedUrl.add(url);
27.    }
28.    // 判断未访问的URL队列中是否为空
29.    public static boolean unVisitedUrlsEmpty() {
30.        return unVisitedUrl.isEmpty();
31.    }
32.}
```

29. }

30. }

```
1. import java.util.HashSet;
2. import java.util.Set;
3. import org.htmlparser.Node;
4. import org.htmlparser.NodeFilter;
5. import org.htmlparser.Parser;
6. import org.htmlparser.filters.NodeClassFilter;
7. import org.htmlparser.filters.OrFilter;
8. import org.htmlparser.tags.LinkTag;
9. import org.htmlparser.util.NodeList;
10. import org.htmlparser.util.ParserException;
11. /**
12.  * 过滤http的url，获取可以符合规则的url
13.  * @author Administrator
14.  *
15.  */
16. public class ParserHttpUrl {
17.     // 获取一个网站上的链接,filter 用来过滤链接
18.     public static Set<String> extracLinks(String url, LinkFilter filter) {
19.         Set<String> links = new HashSet<String>();
20.         try {
21.             Parser parser = new Parser(url);
22.             // 过滤 <frame >标签的 filter，用来提取 frame 标签里的 src 属性所表示的链接
23.             NodeFilter frameFilter = new NodeFilter() {
24.                 public boolean accept(Node node) {
25.                     if (node.getText().startsWith("frame src=")) {
26.                         return true;
27.                     } else {
28.                         return false;
29.                     }
30.                 }
31.             };
32.             // OrFilter 来设置过滤 <a> 标签，和 <frame> 标签
33.             OrFilter linkFilter = new OrFilter(new NodeClassFilter(
34.                 LinkTag.class), frameFilter);
35.             // 得到所有经过过滤的标签
36.             NodeList list = parser.extractAllNodesThatMatch(linkFilter);
```

```

37.         for (int i = 0; i < list.size(); i++) {
38.             Node tag = list.elementAt(i);
39.             if (tag instanceof LinkTag) // <a> 标签
40.             {
41.                 LinkTag link = (LinkTag) tag;
42.                 String linkUrl = link.getLink(); // url
43.                 if (filter.accept(linkUrl))
44.                     links.add(linkUrl);
45.             } else // <frame> 标签
46.             {
47.                 // 提取 frame 里 src 属性的链接如 <frame src="test.html"/>
48.                 String frame = tag.getText();
49.                 int start = frame.indexOf("src=");
50.                 frame = frame.substring(start);
51.                 int end = frame.indexOf(" ");
52.                 if (end == -1)
53.                     end = frame.indexOf(">");
54.                 String frameUrl = frame.substring(5, end - 1);
55.                 if (filter.accept(frameUrl))
56.                     links.add(frameUrl);
57.             }
58.         }
59.     } catch (ParserException e) {
60.         e.printStackTrace();
61.     }
62.     return links;
63. }
64. }

```

```

1. import java.io.File;
2. import java.io.FileOutputStream;
3. import java.io.InputStream;
4. import java.net.URL;
5. import java.net.URLConnection;
6. import java.util.ArrayList;
7. import java.util.List;
8. import java.util.regex.Matcher;
9. import java.util.regex.Pattern;
10. /**

```

```
11. * java抓取网络图片
12. *
13. * @author swinglife
14. *
15. */
16. public class DownLoadPic {
17.     // 编码
18.     private static final String ECODING = "UTF-8";
19.     // 获取img标签正则
20.     private static final String IMGURL_REG = "<img.*src=(.*?)[^>]*?>";
21.     // 获取src路径的正则
22.     private static final String IMGSRG_REG = "http:\\\"?(.*?)(\\\"|>|\\s+)\"";
23.     public static void downloadPic(String url) {
24.         // 获得html文本内容
25.         String HTML = null;
26.         try {
27.             HTML = DownLoadPic.getHTML(url);
28.         } catch (Exception e) {
29.             e.printStackTrace();
30.         }
31.         if (null != HTML && !"".equals(HTML)) {
32.             // 获取图片标签
33.             List<String> imgUrl = DownLoadPic.getImageUrl(HTML);
34.             // 获取图片src地址
35.             List<String> imgSrc = DownLoadPic.getImageSrc(imgUrl);
36.             // 下载图片
37.             DownLoadPic.download(imgSrc);
38.         }
39.     }
40.     /**
41.      * 获取HTML内容
42.      *
43.      * @param url
44.      * @return
45.      * @throws Exception
46.      */
47.     private static String getHTML(String url) throws Exception {
48.         URL uri = new URL(url);
```

```
49.     URLConnection connection = uri.openConnection();
50.     InputStream in = connection.getInputStream();
51.     byte[] buf = new byte[1024];
52.     int length = 0;
53.     StringBuffer sb = new StringBuffer();
54.     while ((length = in.read(buf, 0, buf.length)) > 0) {
55.         sb.append(new String(buf, ECODING));
56.     }
57.     in.close();
58.     return sb.toString();
59. }
60. /**
61.  * 获取ImageUrl地址
62.  *
63.  * @param HTML
64.  * @return
65.  */
66. private static List<String> getImageUrl(String HTML) {
67.     Matcher matcher = Pattern.compile(IMGURL_REG).matcher(HTML);
68.     List<String> listImgUrl = new ArrayList<String>();
69.     while (matcher.find()) {
70.         listImgUrl.add(matcher.group());
71.     }
72.     return listImgUrl;
73. }
74. /**
75.  * 获取ImageSrc地址
76.  *
77.  * @param listImageUrl
78.  * @return
79.  */
80. private static List<String> getImageSrc(List<String> listImageUrl) {
81.     List<String> listImgSrc = new ArrayList<String>();
82.     for (String image : listImageUrl) {
83.         Matcher matcher = Pattern.compile(IMGSRC_REG).matcher(image);
84.         while (matcher.find()) {
85.             listImgSrc.add(matcher.group().substring(0,
86.                 matcher.group().length() - 1));
```

```

87.     }
88. }
89.     return listImgSrc;
90. }
91. /**
92.  * 下载图片
93.  *
94.  * @param listImgSrc
95.  */
96. private static void download(List<String> listImgSrc) {
97.     for (String url : listImgSrc) {
98.         try {
99.             String imageName = url.substring(url.lastIndexOf("/") + 1,
100.                url.length());
101.             URL uri = new URL(url);
102.             InputStream in = uri.openStream();
103.             FileOutputStream fo = new FileOutputStream(new File(imageName));
104.             byte[] buf = new byte[1024];
105.             int length = 0;
106.             while ((length = in.read(buf, 0, buf.length)) != -1) {
107.                 fo.write(buf, 0, length);
108.             }
109.             in.close();
110.             fo.close();
111.         } catch (Exception e) {
112.             e.printStackTrace();
113.         }
114.     }
115. }
116. }

```

```

1. public class Crawler {
2.     /**
3.     * 抓取过程
4.     *
5.     * @return
6.     * @param seeds
7.     */
8.     public void crawling(String url) { // 定义过滤器

```

```
9.     Filter filter = new Filter() {
10.         public boolean accept(String url) {
11.             //这里过滤规则随需要爬的网站的规则进行改变，推荐使用正则实现，本人是爬豆瓣
            网站
12.             if(url.indexOf("douban.com/group/topic") != -1 || url.indexOf("douban.com/group/haix
            iuzu/discussion?start") != -1 )
13.                 return true;
14.             else
15.                 return false;
16.         }
17.     };
18.     // 初始化 URL 队列
19.     LinkQueue.addUnvisitedUrl(url);
20.     // 循环条件，待抓取的链接不空
21.     while (!LinkQueue.unVisitedUrlsEmpty()) {
22.         // 队头URL出队列
23.         String visitUrl = (String) LinkQueue.unVisitedUrlDeQueue();
24.         if (visitUrl == null)
25.             continue;
26.         DownLoadPic.downloadPic(visitUrl);
27.         // 提取出下载网页中的 URL
28.         Set<String> links = ParserHttpUrl.extracLinks(visitUrl, filter);
29.         // 新的未访问的 URL 入队
30.         for (String link : links) {
31.             LinkQueue.addUnvisitedUrl(link);
32.         }
33.     }
34. }
35. // main 方法入口
36. public static void main(String[] args) {
37.     Crawler crawler = new Crawler();
38.     crawler.crawling("http://www.douban.com/group/haixiuzu/discussion?start=0");
39. }
40. }
```