



网络知识

爬虫技术

爬虫不要滥用，否则可能会触犯法律！

爬取代码追求以下几点：

- 准确性：可以爬取到，并且数据通过处理正好是自己需要的。
- 稳定性：可以全部爬到而不崩溃。
- 效率：爬取速度快。

跳过异常

如果使用java循环发起大量请求时，可以跳过部分请求失败的地址：

```
1 List<String> urlList = new ArrayList<>();
2 for(String url : urlList){
3     try{
4         http.get(url);
5     }catch(Exception e){
6         continue;
7     }
8 }
```

Java

当然也可以将异常的地址存下，待下次再重新获取：

```
1 List<String> urlList = new ArrayList<>();
2 List<String> errList = new ArrayList<>();
3 for(String url : urlList){
4     try{
5         http.get(url);
6         ...
7     }catch(Exception e){
8         errList.add(url);
9         continue;
10    }
11 }
```

Java

可以使用递归一次性解决：

```
1 public static void toGet(List<String> urlList){
2     // 无地址
3     if(urlList.size() == 0){
4         return urlList;
5     }
6     // 请求并处理
7     List<String> errList = new ArrayList<>();
8     for(String url : urlList){
9         try{
10            http.get(url);
11            ...
12        }catch(Exception e){
13            errList.add(url);
14        }
15    }
16    toGet(errList);
17 }
```

Java

```
14         continue;
15     }
16 }
17 // 递归调用
18 toGet(errList);
19 }
```

注意：此法只是不断重试失败地址，如果地址本身就无法访问，那么就需要进行相应改造。

HTML 解析

Jsoup

jsoup 可以通过 css 选择器来解析 dom 树。

引入

maven 依赖：

```
1 <dependency>
2   <groupId>org.jsoup</groupId>
3   <artifactId>jsoup</artifactId>
4   <version>1.11.3</version>
5 </dependency>
```

XML

document 解析

获取 dom 树：

```
1 // 从html文本加载dom树
2 Document document = Jsoup.parse("<html>...</html>");
3
4 // 从远程地址获取dom树
5 Document document = Jsoup.connect(url).get();
```

Java

查询节点集：

```
1 // 使用css选择器进行查询
2 Elements elements = document.select("#name .value");
```

Java

查询节点：

```
1 // 节点集直接取出
2 Element element = elements.get(0);
3
4 // 使用Dom语法直接获取节点
5 Element element = document.getElementById("user");
```

Java

获取节点值：

```
1 // 获取html文本
2 String html = element.html();
3
4 // 获取标签名
5 String tag = element.tagName();
6
7 // 获取属性值
8 String attr = element.attr("href");
9
10 // 获取节点文本
11 String text = element.text();
```

Java

修改节点值：

```
1 // 设置html文本
2 Element element = element1.html("<div>test</div>");
3
```

Java

```
4 // 设置属性值
5 Element element = element1.attr("key", "value");
6
7 // 设置标签文本值
8 Element element = element1.text("test");
```

操作节点:

```
1 element.append("<div>test</div>");
```

Java

使用示例

解析远程 dom 示例:

```
1 public static void jsonDom(String url) throws IOException {
2     //Jsoup.parse("");
3     Document document = Jsoup.connect(url).get();
4     // 获取全部
5     Elements elements = document.select(".name .value");
6     for (Element element : elements) {
7         Elements select = element.select("td");
8         String code = select.get(0).text();
9         String type = select.get(1).text();
10        String name = select.get(2).text();
11        System.out.println(code + " | " + type + " | " + name);
12    }
13 }
```

Java

HtmlParser