

明明是悟空

术到极致，几近于道！

博客园 首页 新随笔 联系 订阅 管理

Java判断文件编码格式

转自：http://blog.csdn.net/zhangzh332/article/details/6719025

一般情况下我们遇到的文件编码格式为GBK或者UTF-8。由于中文Windows默认的编码是GBK，所以一般只要判定UTF-8编码格式。对于UTF-8编码格式的文本文件，其前3个字节的值就是-17、-69、-65，所以，判定是否是UTF-8编码格式的代码片段如下：

```
Java代码
1. java.io.File f=new java.io.File("待判定的文本文件名");
2. try{
3.     java.io.InputStream ios=new java.io.FileInputStream(f);
4.     byte[] b=new byte[3];
5.     ios.read(b);
6.     ios.close();
7.     if(b[0]==-17&&b[1]==-69&&b[2]==-65)
8.         System.out.println(f.getName()+"编码为UTF-8");
9.     else System.out.println(f.getName()+"可能是GBK");
10. }catch(Exception e){
11.     e.printStackTrace();
12. }
```

上述代码只是简单判定了是否是UTF-8格式编码的文本文件，如果项目对要判定的文本文件编码不可控（比如用户上传的一些HTML、XML等文本），可以采用一个现成的开源项目：cpdetector，它所在的网址是：<http://cpdetector.sourceforge.net/>。它的类库很小，只有500K左右，利用该类库判定文本文件的代码如下（由于cpdetector的算法使用概率统计，所以结果并不是100%准确的，但是是迄今为止我见过的最准确的....@\_@）：

```
Java代码
1. /*-----
2.     detector是探测器，它把探测任务交给具体的探测实现类的实例完成。
3.     cpDetector内置了一些常用的探测实现类，这些探测实现类的实例可以通过add方法
4.     加进来，如ParsingDetector、 JChardetFacade、 ASCIIDetector、 UnicodeDetector。
5.     detector按照“谁最先返回非空的探测结果，就以该结果为准”的原则返回探测到的
6.     字符集编码。
7.     -----*/
8. cpdetector.io.CodepageDetectorProxy detector =
9. cpdetector.io.CodepageDetectorProxy.getInstance();
10. /*-----
11.     ParsingDetector可用于检查HTML、XML等文件或字符流的编码,构造方法中的参数用于
12.     指示是否显示探测过程的详细信息，为false不显示。
13.     -----*/
14. detector.add(new cpdetector.io.ParsingDetector(false));
15. /*-----
16.     JChardetFacade封装了由Mozilla组织提供的JChardet，它可以完成大多数文件的编码
17.     测定。所以，一般有了这个探测器就可满足大多数项目的要求，如果你还不放心，可以
18.     再多加几个探测器，比如下面的ASCIIDetector、UnicodeDetector等。
19.     -----*/
20. detector.add(cpdetector.io.JChardetFacade.getInstance());
21. //ASCIIDetector用于ASCII编码测定
22. detector.add(cpdetector.io.ASCIIDetector.getInstance());
23. //UnicodeDetector用于Unicode家族编码的测定
```

公告

昵称：明明是悟空  
园龄：6年6个月  
粉丝：92  
关注：4  
+加关注

<				2018年1:
日	一	二	三	
28	29	30	31	
4	5	6	7	
11	12	13	14	
18	19	20	21	
25	26	27	28	
2	3	4	5	

搜索

我的标签

- Linux(127)
- web开发(84)
- java(63)
- c/c++(59)
- android(45)
- Linux内核(38)
- chromium(29)
- webrtc(23)

```
24. detector.add(cpdetector.io.UnicodeDetector.getInstance());
25. java.nio.charset.Charset charset = null;
26. File f=new File("待测的文本文件名");
27. try {
28.     charset = detector.detectCodepage(f.toURL());
29. } catch (Exception ex) {ex.printStackTrace();}
30. if(charset!=null){
31.     System.out.println(f.getName()+"编码是: "+charset.name());
32. }else
33.     System.out.println(f.getName()+"未知");
```

上面代码中的detector不仅可以用于探测文件的编码，也可以探测任意输入的文本流的编码，方法是调用其重载形式：

Java代码

```
1. charset=detector.detectCodepage(待测的文本输入流,测量该流所需的读入字节数);
```

上面的字节数由程序员指定，字节数越多，判定越准确，当然时间也花得越长。要注意，字节数的指定不能超过文本流的最大长度。

判定文件编码的具体应用举例：

属性文件(.properties)是Java程序中的常用文本存储方式，象STRUTS框架就是利用属性文件存储程序中的字符串资源。它的内容如下所示：

Java代码

```
1. #注释语句
2. 属性名=属性值
```

读入属性文件的一般方法是：

Java代码

```
1. FileInputStream ios=new FileInputStream("属性文件名");
2. Properties prop=new Properties();
3. prop.load(ios);
4. ios.close();
```

利用java.io.Properties的load方法读入属性文件虽然方便，但如果属性文件中有中文，在读入之后就会出现乱码现象。发生这个原因是load方法使用字节流读入文本，在读入后需要将字节流编码成为字符串，而它使用的编码是"iso-8859-1",这个字符集是ASCII码字符集，不支持中文编码，所以这时需要使用显式的转码：

Java代码

```
1. String value=prop.getProperty("属性名");
2. String encValue=new String(value.getBytes("iso-8859-1"),"属性文件的实际编码");
3.
```

标签: web开发

好文要顶

关注我

收藏该文

明明是悟空

关注 - 4

粉丝 - 92

+加关注

« 上一篇：FileUpload之FileItem

» 下一篇：SAE Java相关问题小结

posted @ 2014-05-16 21:08 明明是悟空 阅读(5247) 评论(0) 编辑 收藏

H264(16)

数据库(16)

更多

随笔档案

2018年10月 (3)

2018年9月 (1)

2018年8月 (3)

2018年7月 (9)

2018年6月 (2)

2018年5月 (8)

2018年4月 (8)

2018年3月 (18)

2018年2月 (6)

2018年1月 (3)

2017年12月 (3)

2017年11月 (7)

2017年10月 (4)

2017年9月 (5)

2017年8月 (1)

2017年7月 (3)

2017年6月 (5)

2017年5月 (8)

2017年4月 (13)

2017年3月 (11)

2017年2月 (3)

2017年1月 (2)

刷新评论 刷新页面 返回顶部

2018/11/11		Java判断文件编码格式 - 明明是悟空 - 博客园	
<div>注册用户登录后才能发表评论，请 <a href="#">登录</a> 或 <a href="#">注册</a>，<a href="#">访问网站首页</a>。</div>		2016年12月 (2)	
<div><div>【推荐】超50万VC++源码: 大型组态工控、电力仿真CAD与GIS源码库！</div><div>【活动】11.1-11.11，3000元神券限量开抢，51CTO全场课程5折起，还送精美礼品！</div><div>【推荐】华为云11.11普惠季 血拼风暴 一促即发</div><div>【工具】SpreadJS纯前端表格控件，可嵌入应用开发的在线Excel</div><div>【腾讯云】拼团福利，AMD云服务器8元/月</div></div>		2016年11月 (10)	
<div></div>		2016年10月 (7)	
<div><div>相关博文：</div><div><div>· cpdetector获取文件编码</div><div>· Java如何获取文件编码格式</div><div>· 代码文件的编码不统一导致的坑</div><div>· xml、文件操作功能类</div><div>· Linux查看文件编码格式及文件编码转换&lt;转&gt;</div></div></div>		2016年9月 (27)	
<div></div>		2016年8月 (23)	
<div><div>最新新闻：</div><div><div>· 生物在量子进化？你可能正在利用量子现象生存</div><div>· 特斯拉大幅调整Model X和S价格和配置</div><div>· Drive.ai、Waymo率先商业化，智能驾驶加速进入冲击阶段</div><div>· 自由软件基金会认为商业性条款是非自由的</div><div>· 天猫双十一销售额破1682亿元：超过去年全天 提前8小时</div><div>» 更多新闻...</div></div></div>		2016年7月 (28)	
		2016年6月 (13)	
		2016年5月 (7)	
		2016年4月 (9)	
		2016年3月 (26)	
		2016年2月 (9)	
		2016年1月 (8)	
		2015年12月 (13)	
		2015年11月 (3)	
		2015年10月 (14)	
		2015年9月 (22)	
		2015年8月 (12)	
		2015年7月 (9)	
		2015年6月 (9)	
		2015年5月 (26)	
		2015年4月 (27)	
		2015年3月 (27)	
		2015年2月 (18)	
		2015年1月 (8)	
		2014年12月 (8)	
		2014年11月 (1)	

2014年10月 (13)
2014年9月 (15)
2014年8月 (37)
2014年7月 (2)
2014年6月 (3)
2014年5月 (26)
2014年4月 (17)
2014年3月 (11)
2014年2月 (7)
2014年1月 (6)
2013年11月 (14)
2013年10月 (8)
2013年9月 (12)
2013年8月 (20)

最新评论
1. Re:Jacob的使用出招  ***.ocx 是什么东西？ 到？？
2. Re:线程安全的单例  感谢大佬
3. Re:如何在java程序 或者shell脚本  大哥，如何在系统或者t 只能执行指定的几个命 行mkdir ，不让执行rm
4. Re:jsp放在web-inf  最近正好在学Servlet和 p访问的问题,弄了好久,

解决办法;  
博主的文章给我带来了帮助!!  
很有帮助,谢谢.

5. Re:二进制的计算 ( 计算补码存储数据 )  
  
@bamboo~1+ ( -1  
原码运算= ( 0001+11  
000=0...

阅读排行榜

- 1. 如何在java程序中调用shell脚本(39314)
- 2. linux mysql 操作命令
- 3. 1080P、720P、4C理论带宽(21157)
- 4. 修改Tomcat编码方式(39)
- 5. 实例讲解教你读懂路由

评论排行榜

- 1. JSP/Servlet的编码问题
- 2. struts2文件下载及inputName="inputName">inputName>的理解(2)
- 3. 二进制的计算 ( 计算补码存储数据 ) (2)
- 4. ( 三 ) WebRTC手记 (2)
- 5. 谷歌开源项目Chromium与项目构建 ( Win7+vs2010)

推荐排行榜

- 1. JAVA 的wait(), notify()同步机制(3)
- 2. jsp放在web-inf下的原因

3. C语言中结构体的位域  
(2)

4. JAVA 正则表达式、正则代码(2)

5. 表现层(jsp)、持久层  
务层 ( 逻辑层、service  
abean )、控制层 ( act

Copyright ©2018 明明是悟空