

正则表达式在编程题自动阅卷中的应用

佘石泉, 周肆清

(中南大学, 湖南 长沙 410083)

摘要:分析了编程题的自动阅卷的现状及存在的不足。为了让计算机能够更加合理和有效地进行编程题的阅卷,提出了一种用正则表达式来分析程序代码,模拟人工阅卷的方法。介绍了此方法的原理和主要功能,给出了方法实现的具体步骤,对方法的关键部分如 Java 处理正则表达式的各个类以及方法的具体实施等都进行了详细的解释并举例进行了说明。其核心思想是利用正则表达式来抽象标准答案,再利用正则表达式强大的文本匹配功能来进行程序代码的分析,把不变的文本转化为具有一般性的表达式,大幅度增加了匹配的灵活性,从而达到简化阅卷的目的。

关键词:正则表达式;抽象;文本匹配

中图分类号:G434

文献标识码:A

文章编号:1673-629X(2007)07-0244-03

Application of Regex in Auto-Checking Paper of Programs

SHE Shi-quan, ZHOU Si-qing

(Central South University, Changsha 410083, China)

Abstract: Analyzed the status quo and the shortages of auto-checking answer sheet of programs. In order to look over the examination papers by computer more reasonably and effectively, a method was given which used regular expressions to analyze the codes of program and simulate the artificial model. This paper also introduced the elements and chief function of the method, and the implementary process was given too; moreover, the pivotal parts of the method like the classes of Java that used to manage regular expressions and concretely actualization of the method were explained amply and examples were additional. The central idea is to changeing for complicated standard solution to a series of regular expressions, and utilize their powerful function in text matching to analyze the programming codes, thereby the aim of predigesting for looking over answer papers is achieved.

Key words: regular expressions; abstract; text matching

1 问题的提出

进入 21 世纪以来,随着计算机的普及,计算机自动阅卷技术的应用已经越来越广泛。客观题的自动阅卷技术已经发展成熟,但主观题乃至编程题的阅卷技术却停滞难前,一直在束缚着计算机自动阅卷技术的发展。目前对编程题阅卷普遍采用的办法是查看考生程序的运行结果,对则满分,错则零分。这种做法与人工阅卷是背道而驰的,很不公平,甚至还给了学生投机舞弊的机会。而有些阅卷系统中虽然加入了对考生结果代码进行关键词的比较机制,但只是进行简单的字词比较,况且存在如何选择关键词的问题。

基于以上所述,有人提出了一种对考生所编写的源程序进行分析检测的办法^[1]。也即对考生的源程序

和预先做好的一个尽可能全面的标准答案进行比较,根据比较的结果,得到一个称之为代码可信度的数据,把这个数据列入到评分模式之中。代码可信度越高,则得分越高。这种做法比之前者有了较大的改进,但该方法还只是对文本的一种简单比较,效果当然不是那么理想,所以在实现方法上还可以进一步去研究。

文中提出了一种基于正则表达式来分析代码可信度的方法,即利用正则表达式来进行文本的检测,此想法是首先把程序的代码按照功能分块,看考生的源程序是否含有类似的功能模块,得到代码可信度。

2 人工阅卷模式分析

众所周知,任何一种编程语言实现某种功能的方法都不唯一,就拿 C 语言来说,实现一个循环便有 while, for, do-while 三种方法,稍微复杂的程序其代码便可能有很多种,因此难以以一个简单的文本作为标准答案。在人工阅卷的时候,如果未对题目解题方法做明确规定的話,各种可能的答案都应是正确的,因此

收稿日期:2006-09-14

作者简介:佘石泉(1980-),男,湖南邵阳人,硕士研究生,研究方向为数据库技术、人工智能;周肆清,副教授,硕士生导师,研究方向为图像处理与数据库技术。

阅卷的教师要考生按照题意完成了题目的要求,就可以算是正确的,并不会计较考生以何种方法完成。在给分的时候,阅卷老师也是按照考生的代码对题目直接或间接要求的一些功能的实现情况来给分,而不仅仅依靠程序的运行结果。下面举例进行说明。

题目:读入一个4位正整数,然后把它逆序输出。

分析:这个题目的关键点是把这个4位正整数进行分解,然后把各个数位逆转过来,再重新组成一个新的正整数。考生如果能够实现这一段代码,应该要获得较高的分数。另外输入和输出以及运行结果也占有一定分值。假如这个题目共占8分,整数分解重组的代码块占5分,记做块A,输入、输出、运行结果各占1分,分别记做块B、C、D。阅卷教师即可按照得分点来进行评分了。比如,某考生只做了输入和输出,其C语言代码如下所示,则得2分,即获得B和C块的得分,D块的得分必须依赖于A块,否则所做出的运行结果即使正确也不能得分。

```
void main()
{
    int n;
    scanf("%d",&n);
    printf("%d",n);
}
```

当然,还要考虑考生代码存在错误书写或语法错误等情况,这也是无法避免的,但只要错误不多,阅卷老师一般都会不扣分或者少扣分。

自动阅卷应当以人工阅卷模式为参考,尽量模拟人工阅卷,只有这样,才可能做到公平合理。因此,文中提出一种利用正则表达式来对考生代码进行检测的方法。当然此方法只是整个阅卷系统中的一部分,它的最终结果只是要得到一个代码可信度,并不能单独来充当阅卷的工具。具体的步骤如下:

- 1) 按照题意,分析出程序的得分点;
 - 2) 制定每个得分点的正则表达式;
 - 3) 扫描考生程序,利用正则表达式与考生程序进行匹配;
 - 4) 根据匹配情况得到代码可信度;
- 此方法的关键步骤是2)和3),因此文中将重点分析2)和3)的实现。

3 Java 处理正则表达式简介

正则表达式的英文名是 Regular Expression,缩写为 Regex,其起源最早可以追溯到 1956 年,后来被应用到信息技术领域。正则表达式第一个实用应用程序是 Unix 中的 qed 编辑器,而现在,正则表达式引擎已

经被 Unix 中的许多工具所实现。另外,正则表达式在许多脚本语言中也有着非常重要的地位。它的主要的功能是处理文本和匹配模式,常用于在文本中进行查找和替换,或者检验文本是否符合某一模式的规则。比如可以用正则表达式查找某文本中的所有的字符串“hi”,并把查到的全部字符串替换成“hello”;也可以用正则表达式来检验某个输入字符串是否符合 Email 地址的格式^[2]。

下面简单介绍一下 Java 处理正则表达式的两个重要的类:Pattern 类及 Matcher 类。Pattern 类提供了一个 compile 方法,这是一个静态的方法,它接收一个正则表达式作为参数,返回此正则表达式的一个 Pattern 类对象;然后再利用该对象来调用 Matcher 方法。Matcher 方法接收一个字符序列作为参数,即应用中的待匹配的文本,返回一个 Matcher 类对象。

Matcher 类提供诸如 matches, find, lookingAt, group, replaceAll 等方法。其中的 matches 方法与 Pattern 类中的 matches 方法完成相同的任务,但它不接收任何参数,因为匹配模式(正则表达式)和匹配对象(待匹配的文本)都已经封装在了 Matcher 类的对象中^[3]。

这里,笔者用一个简单的例子来介绍一下 Java 中如何使用正则表达式进行文本匹配工作的。正则表达式为“- \ w + @ \ w + \ . [a - z A - Z] { 2, 3 } \$”,用来检测某个 Email 地址的书写格式。 \ w + 表示一个或多个词字符,词字符包括大小写字母、数字和下划线; \ . 匹配字面意义符号“.”,而非正则表达式中“.”; [a - z A - Z] { 2, 3 } 表示 2 个或者 3 个字母的组合。在 Java 中,有些字符必须经过转义,而在正则表达式中也有转义,因此正则表达式“ \ w”在 Java 中会被写成“ \ \ w”。如下代码将产生一次成功的匹配^[4]:

```
Pattern p = Pattern.compile("- \ w + @ \ w + \ . [ a - z A - Z ] { 2, 3 } $");
Matcher m = p.matcher("HelloRegex@126.com");
m.matches();
```

4 方法具体实现

再来分析前面所举的例子,得分点不再赘述,正则表达式则可以这样来构造:

1) 对于代码块 B 和 C,要检测考生是否做了正确的输入和输出,正则表达式如下。

```
输入: "scanf ( ( \ "%d \ s * \ ", &[ a - z A - Z ] \ w * \ )";
输出: "printf ( ( \ "%d \ s * \ ", [ a - z A - Z ] \ w * \ )";
```

2) 代码块 A 的检测:对于这一代码块,主要检测

