基于规则的大规模试卷文本语块识别方法的研究*

郭凯红^{1, 2}, 李文立¹

(1.大连理工大学 管理学院, 辽宁 大连 116024, 2 辽宁大学 信息学院, 沈阳 110036)

摘 要:提出了一种基于规则的试卷文本语块识别方法,有效解决了试题库中大规模试题数据的初始化问题。通过定义文本语块识别规则,构建自动机识别模型,在理论上描述了试卷文本的识别过程。实验表明,该模型具有良好的性能,在此基础上,实现了一个原型系统,通过具体的应用实例验证了该方法的可行性和有效性。

关键词: 规则: 语块: 试卷文本: 识别模型

中图分类号: TP391.43 文献标志码: A 文章编号: 1001-3695(2009)04-1391-03

Study of massive paper texts chunking based on rules

GUO Ka÷hong^{1, 2}, LIW en-li¹

(1. School of M anagoment, Dalian University of Technology, Dalian Liaoning 116024, China; 2 College of Information, Liaoning University, Shenyang 110036, China)

Abstract To solve the initiating of massive examination questions in database efficiently, proposed a paper texts chunking method based on rules. Defining recognition rules of paper texts and constructing automata recognition model, described the recognition processing of paper texts theoretically. Experiment results show that this model has better performance. By these works, implemented a prototype system, and verified its feasibility and effectiveness by a practical application.

Key words rules, chunk, examination paper texts, recognition model

0 引言

浅层句法分析(shallow parsing), 也称做部分句法分析 (partial parsing)或语块分析 (chunk parsing) [1], 是近年来自然 语言处理领域出现的一个新的语言处理策略。它是与完全句 法分析相对的, 不要求得到完全的句法分析树, 只要求识别其中的某些结构相对简单的成分, 如非递归的名词短语、动词短语等。这些识别出来的结构通常被称做语块 (chunk)。 概括起来, 句法分析的方法基本上可以分成两类 [2]: 基于统计的方法和基于规则的方法。统计方法主要来自概率统计和信息论, 而规则方法则是根据人工书写的或 (半)自动获取的语法规则标注出短语的边界和短语的类型。

浅层句法分析的主要任务是语块的识别和分析^[3],这就使句法分析的任务在某种程度上得到简化,同时也有利于句法分析技术在大规模真实文本处理系统中迅速得到应用,目前已取得大量成果^[24~11]。文献 [12]从组成形式和上下文语境两个方面来识别汉语文本中的特殊符号串;文献 [13]提出了基于规则的数据收集策略,将数据收集过程分成共性和特性两种情况,将提炼出的若干规则形成不同的规则集以便于开发和定制。这些研究成果主要集中在"语块库建立"的问题上。本文面向具体应用,研究由大规模试卷文本实现试题库高效初始化的方法,即通过语块识别程序,根据预先定义的识别规则,应用规则方法对集中的大规模试卷文本进行识别,逐级分析有限状态层叠^[14],获取试卷文本结构及各试题的属性,将各属性值实

时写入数据库,自动完成试题库的初始化工作。本文方法改变 了传统的人工录入试题库的操作模式,有效解决了大规模试题 库的初始化问题。

1 文本结构分析

教育考试的普通试卷(笔试或上机考试)主要有选择、填空、改错、问答等有限几种类型题,每种类型题又由若干道具体的试题组成。任何一道试题总伴有题目、答案、分值等几种相关数据项,这些实际上就是每道试题的具体信息;将若干条这样的试题按类型分组,就构成一组类型题;可以将多种类型题组合在一起,从而构成一套试卷。由此得出,每道试题及其相关属性(答案、分值等)是构成试卷的基本元素,试题库实际存储的是这些属性值,而不是整套试卷;每道试题通过其他相关属性标志它所属的试题类型及试卷套数。

根据以上分析, 现提取出试题必要的数据项, 即试卷序号、试题类型序号及类型说明、试题序号及题目、试题答案、试题分值。这实际上给出了试题的数据结构, 识别系统将按这种格式实现对试题数据的分析操作。

2 识别规则

根据前述的试题结构, 考虑到试卷文本的可读写性, 定义标记符号为"@""[]""()""#""%"。其中: "@"是试卷开始符, 符号后必须紧跟数字, 代表试卷的序号; "[]"是试题类型开始符, 符号"["与"]"之间必须是数字, 代表试题类型的序

收稿日期: 2008-06-21; **修回日期**: 2008-08-30 **基金项目**: 国家自然科学基金资助项目(70572099); 辽宁省自然科学基金资助项目(1050349)

作者简介: 郭凯红 (1973-), 界, 河南镇平人, 博士研究生, 主要研究方向为信息管理、系统工程 (guokh@ 126 cm); 李文立 (1969-), 界, 河南平 顶山人, 副教授, 博导, 博士, 主要研究方向为电子商务, 信息管理、系统工程, © 1994-2012 China Academic Journal Electronic Publishing House. All rights reserved. http://www.cnki.net

号; "()"是题目开始符, 符号"("与")"之间必须是数字, 代表 试题的序号: "#"是试题答案开始符,符号后紧跟的内容是当 前试题的答案: "%"是试题分值开始符. 符号后必须紧跟数 字, 代表当前试题的分值。总体要求是, 一套试卷必须有试卷 序号、试题类型序号及类型说明、试题序号及题目内容、试题答 案内容、试题分值等属性值: 暂不考虑试卷文本中空格和回车 换行符的识别问题, 其处理方法将在实现中灵活处理。

上述标记符号被称做保留符号,如同高级语言中的保留字 一样, 由语块识别程序当做标记控制符号使用。如果文本中的 试题内容含有上述保留符号,则将该符号连续双写,如试题中 含有内容 "# include (stdin h)", 在试卷文本中应写成 "## include (std in h)", 语块识别程序在分析这段文本时, 自动地将 符号串"#"解析成输出符号"#", 而不是标记控制符号。满足 本识别规则的试卷文本如下所示:

@ 1

[1]选择题 1

(1)下面哪种操作系统是非图形化的?

A. Linux

B. DOS

 $C\!.\,W$ in dow s

D. Macin to sh

#B

% 2

(2)下面哪种软件用于专业图像编辑?

A. Word

B. WPS

C. Photoshop

D. Borland C++

C

% 2

[4]简答题

(1)什么是计算机病毒?

#关于计算机病毒目前没有一个公认的定义, 我国公安部 计算机安全检察司对病毒的定义是: 计算机病毒是指编制或者 在计算机程序中插入的破坏计算机功能或者毁坏数据. 影响计 算机使用,并能自我复制的一组计算机指令或者程序代码。

% 5

.

3 自动机识别模型

根据前述定义的识别规则,即通过增加标记定义文本语块 的边界, 现建立文本语块的自动机识别模型, 描述试卷文本的 识别过程。首先构造机器 M 的输入字母表和输出字母表。由 于试卷文本是建立在 ASC II码上的字符串, 定义输入字母表为 全体可视 ASCII码集, 用 Σ 表示, 即 $S = \Sigma$ 。定义输出字母表 时,要求一个输入串/试卷文本)的响应在逻辑上具有可识别 的特征, 也就是说, 建立在输出字母表上的字符串能够抽象地 表示清楚每套试卷及每道试题的逻辑结构。这里用 0表示试 卷开始, 当连续输出两个以上 0时, 表示输入串 (试卷文本)被 M 拒绝; 1表示文本识别成功,输入串 (试卷文本)被 M 接受; * 抽象地表示试题内容; f(x, T) = (x, t) = (x, t) 表示识别规则中的标 记符号,即标记符号依照原样输出。所以定义输出字母表为 $R = \{0, 1, *, [,], (,), #, \% \}_{\circ}$

1)初态 A M 从状态 A 出发, 这是试卷文本识别过程的开 始, 在状态A下输入试卷文本的第一个符号。

2)陷阱状态 B 陷阱状态用做在其他状态下发现输入串 不可能是该自动机所识别的内容时所进入的状态,在此状态 下, 自动机读完输入串中剩余的字符。如果输入串 / 试卷文 本)的第一个符号是非@,此文本必被M所拒绝,此时M进入 状态 B, 输出 Q 一旦 M 进入状态 B, 无论输入什么符号, 输出 都为 0 并且不能再转向其他状态。

3)强制接受试卷序号状态 C M 进入状态 C 表示强制要 求试卷必须有数字序号。如果输入串(试卷文本)的第一个符 号是@,M 进入状态 C,输出 0,表示试卷开始:继续输入的符 号如果是数字 $1 \sim 9$ 则 M 进入状态 C', 输出*; 否则 M 进入状 态 B, 输出 Q

4)接受试卷序号状态 $C^{'}$ M 进入状态 $C^{'}$, 表示输入的试 卷序号暂时可以被接受。继续输入的符号如果是数字 0~9 则 M 仍处于状态 C', 输出*;继续输入的符号如果是 I, 则 M进入状态 D, 输出 f; 否则, M 进入状态 B, 输出 Q

5)强制接受试题类型序号状态 D M 进入状态 D, 表示强 制要求试题类型必须有数字序号。继续输入的符号如果是数 字 $1 \sim 9 \cup M$ 进入状态 D', 输出*; 否则 M 进入状态 B, 输出 Q

6)接受试题类型序号状态 D' M 进入状态 D',表示输入 的试题类型序号暂时可以被接受。继续输入的符号如果是数 字 $0 \sim 9$ 则 M 仍处于状态 D', 输出*:继续输入的符号如果 是 /, 则 M 进入状态 D'', 输出 /; 否则 M 进入状态 B, 输出 C0,

(7)强制接受试题类型说明状态 $(D^{''})$ (M) 进入状态 $(D^{''})$ 表示 强制要求试题类型必须有文本说明。继续输入的符号如果是 非 (M, M, H) 进入状态 (M, H) 进入状态 (M, H) 强输出*; 否则 (M, H) 进入状态 (M, H) 强流 (M, H)

8)接受试题类型说明状态 D⑤ M 进入状态 D⑤表示输入 的试题类型文本说明暂时可以被接受。继续输入的符号如果 是非(, M, M, M, M) 仍处于状态(D) 输出*;继续输入的符号如果是 (, 则 M 进入状态 E, 输出 (。)

9)强制接受试题序号状态 E M 进入状态 E 表示强制要 求试题必须有数字序号。继续输入的符号如果是数字 1~9 则 M 进入状态 E', 输出*; 否则 M 进入状态 B, 输出 Q

10)接受试题序号状态 $E^{'}$ M 进入状态 $E^{'}$, 表示输入的试 题序号暂时可以被接受。继续输入的符号如果是数字 0~9 则 M 仍处于状态 E', 输出*;继续输入的符号如果是),则 M进入状态 E'', 输出); 否则 M 进入状态 B, 输出 Q

11)强制接受题目状态 $E^{''}$ M 进入状态 $E^{''}$ 表示强制要求 题目必须有文本内容。继续输入的符号如果是非#,则 M 进入 状态 E⑤输出*: 否则M 进入状态 B. 输出 D0.

12)接受题目状态 E extstyle extstyle M 进入状态 E extstyle extstyl内容暂时可以被接受。继续输入的符号如果是非#,则 M 仍处 于状态 E⑤输出*:继续输入的符号如果是#.则 M 进入状态 F, 输出#。

13)强制接受试题答案状态 F M 进入状态 F, 表示强制 要求试题答案必须有文本内容。继续输入的符号如果是非%, 则 M 进入状态 F', 输出*; 否则 M 进入状态 B, 输出 Q

(14)接受试题答案 状态 F' M 进入状态 F', 表示输入的试 题答案暂时可以被接受。继续输入的符号如果是非%,则 M 仍处于状态 F' 输出*;继续输入的符号如果是%,则M 进入

下面考察 M 至少需要多少个状态。 Tunda Symma Academic Journal Electronic Publish

状态 G, 输出%。

15)强制接受试题分值状态 G M 进入状态 G, 表示强制要求试题必须有分值。继续输入的符号如果是数字 $1 \sim 9$ 则M 进入状态 G', 输出 1, 否则 M 进入状态 B, 输出 0,

16)接受试卷文本状态 G' M 进入状态 G', 表示试卷文本中的一道试题被识别成功, 输入串 (试卷文本)被 M 接受。继续输入的符号如果是数字 $0 \sim 9$, 则 M 仍处于状态 G', 输出 1; 继续输入的符号如果是@,则 M 进入状态 C, 开始下一套试卷的识别, 输出 0 继续输入的符号如果是 [f],则 M 进入状态 D, 开始下一组试题类型的识别, 输出 [f] 继续输入的符号如果是 [f]则 M 进入状态 [f] 开始下一道试题的识别, 输出 [f] 否则, 输入的符号如果不为空的话, [f] 进入状态 [f] [f]

至此,得到机器 M 的状态集: $Q = \{A, B, C, C', D, D', D', D \in E, E', E', E \in F, F', G, G'\}$ 。所以, $M = \{Q, \Sigma, \{Q, 1, *, [,], (,), #, \%\}$, f, g, A)。其中,状态转换函数 f为

$$f(A, @) = C \qquad \qquad f(A, \overline{@}) = B$$

$$f(C, 1...9) = C' \qquad \qquad f(C, \overline{1...9}) = B$$

$$f(C', 0...9) = C' \qquad \qquad f(C', \overline{1}) = D$$

$$f(C', \overline{1} \land \overline{0...9}) = B \qquad \qquad f(D, 1...9) = D'$$

$$f(D, \overline{1...9}) = B \qquad \qquad f(D', \overline{0...9}) = D'$$

$$f(D', \overline{1}) = D'' \qquad \qquad f(D', \overline{1} \land \overline{0...9}) = B$$

$$f(D'', \overline{1}) = D \oplus \qquad \qquad f(D'', \overline{1} \land \overline{0...9}) = B$$

$$f(D, \overline{1}) = D \oplus \qquad \qquad f(D, \overline{1}) = B$$

$$f(D, \overline{1}) = D \oplus \qquad \qquad f(D, \overline{1}) = B$$

$$f(D, \overline{1}) = D \oplus \qquad \qquad f(D, \overline{1}) = B$$

$$f(D, \overline{1}) = D \oplus \qquad \qquad f(D, \overline{1}) = B$$

$$f(E, \overline{1...9}) = E' \qquad \qquad f(E, \overline{1...9}) = B$$

$$f(E', \overline{1}) = D \oplus \qquad \qquad f(E', \overline{1}) = E \oplus$$

$$f(E'', \overline{1}) = D \oplus \qquad \qquad f(E, \overline{1}) = F'$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = F'$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = F'$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = F'$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = F'$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = F'$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = D \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E, \overline{1}) = B \qquad \qquad f(E, \overline{1}) = B$$

$$f(E,$$

输出函数 g为

$$\begin{split} g(A, @) &= 0 & g(A, \overline{@}) &= 0 \\ g(C, 1 \dots 9) &= * & g(C, \overline{1 \dots 9}) &= 0 \\ g(C', 0 \dots 9) &= * & g(C', f) &= f \\ g(C', \overline{f} \land \overline{0 \dots 9}) &= 0 & g(D, 1 \dots 9) &= * \\ g(D, \overline{1 \dots 9}) &= 0 & g(D', \overline{f} \land \overline{0 \dots 9}) &= * \\ g(D', \overline{f}) &= f & g(D', \overline{f} \land \overline{0 \dots 9}) &= 0 \\ g(D', \overline{f}) &= * & g(D', \overline{f} \land \overline{0 \dots 9}) &= 0 \\ g(D, \overline{f}) &= * & g(D, \overline{f}) &= 0 \\ g(D, \overline{f}) &= * & g(D, \overline{f}) &= 0 \\ g(D, \overline{f}) &= * & g(D, \overline{f}) &= 0 \\ g(D, \overline{f}) &= * & g(D, \overline{f}) &= 0 \\ g(D, \overline{f}) &= * & g(D, \overline{f}) &= 0 \\ g(E, \overline{f}) &= * & g(E, \overline{f}) &= 0 \\ g(E, \overline{f}) &= * & g(E, \overline{f}) &= * \\ g(E, \overline{f}) &= * & g(E, \overline{f}) &= * \\ g(E, \overline{f}) &= * & g(F, \overline{f}) &= * \\ g(F, \overline{f}) &= * & g(F, \overline{f}) &= * \\ g(F, \overline{f}) &= * & g(F, \overline{f}) &= f \\ g(G, \overline{f}) &= 0 & g(G, \overline{f}) &= f \\ g(G, \overline{f}) &= 0 & g(G, \overline{f}) &= f \\ g(G, \overline{f}) &= 0 & g(G, \overline{f}) &= f \\ g(G, \overline{f}) &= f & g(B, \overline{f}) &= 0 \\ g(G, \overline{f}) &= f & g(B, \overline{f}) &= f \\ g(G, \overline{f}) &= f & g(B, \overline{f}) &= f \\ g(G, \overline{f}) &= f & g(B, \overline{f}) &= f \\ g(G, \overline{f}) &= f & g(B, \overline{f}) &= f \\ g(G,$$

M 的状态如图 1所示。

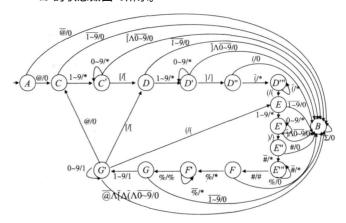


图1 试卷文本自动机识别模型状态图

4 实验分析及改进

应用本模型对 23 套模拟试卷, 每套 5, 6种类型题 (选择、填空、判断、改错、简答、计算、程序设计等)共计 416 道试题进行语块识别, 以考察机器 M 识别的准确率及激励与响应的关系。对于满足识别规则要求的输入串 (试卷文本), M 识别的准确率达 100%。输入串 (试卷文本)经 M 识别后得到其响应的逻辑结构式为

$$\underbrace{0*[\underbrace{*}]^{**}(*)^{**}\#^{**}\underbrace{0}^{*}[(*)^{**}\#^{**}\%1...[\underbrace{*}]^{*}(*)^{*}\#^{*}\%11}_{\text{item}}\underbrace{0*[\underbrace{*}]^{*}(*)^{*}\#^{*}\%1...\%1...}_{\text{paper1}}\underbrace{1...\%1...}_{\text{paper2}}$$

分析上式可以看出,从符号 [或 (开始至随后出现的第一个 1之间的部分或两个邻近 1之间的部分,是试卷中的一道试题 (item)的抽象表示,且试题包括了试卷序号、试题类型序号及类型说明、试题序号及题目、试题答案、试题分值等重要内容;从符号 0开始至下一个 0之间的非空部分或者说两个邻近0之间的非空部分,是一套试卷 (paper)的抽象表示,且试卷包括了若干道试题。在试卷的抽象表示中,注意到第二个 0的前一个符号一定是 1, 这说明输入串 (试卷文本)经 M 识别后成功地到达了状态 G', 试卷包括了至少一道满足要求的试题。如果至少两个连续的 0出现在逻辑结构式中,说明输入串 (试卷文本)在识别过程中被 M 拒绝,进入了陷阱状态 B,并且从最早出现的两个连续 0的位置开始,其后也全部是连续的 Q

本模型采用逐个读取输入串(试卷文本)字符的方式,即首先读取一个输入串字符,判明其类型并给出相应的处理,循环反复,直至读取的字符为空,表示输入串(试卷文本)结束,因此系统识别效率不高。注意到试卷文本的可读性特点,即试卷序号、试题类型序号及类型说明、试题序号及题目、试题答案、试题分值等各属性值均以新一行开始,且各标记控制符号总是位于每一属性值之首,因此在编码实现时,为了提高效率,可以采用逐行读取输入串(试卷文本)的方式,即每次读取试卷文本的一行数据,同时判断这行数据的开头几个字符是标记控制符号还是普通文本数据符号,给出相应的处理;循环反复,直至输入串(试卷文本)结束。这种实现方式突破了自动机模型在状态转换过程中只能处理单一字符的束缚,识别程序执行速度快、效率高。

5 应用实例

根据前述自动机识别模型及实现策略,在 Borland C++

5 系统实现架构分析

数据库和模型库的分析是从系统实现的逻辑结构来分析的。系统的实现架构指的是以何种物理组成方式实现系统,各数据库位于何处、模型库位于何处、以何种网络形式进行连接等。对于具有一定地域覆盖度的系统而言,通常流行的架构方式有浏览器 /服务器 (B/S),客户机 /服务器 (C/S)等。如果 GLD_DSS采用的是 B/S 这类瘦客户端的逻辑结构,那么所有的信息必须位于调度中心服务器上,这要求连接各工作站的网络的带宽足够大。其优点是系统易于维护和管理,各用户只通过浏览器进入系统,用户端无须安装任何软件。对于小型粮食企业,由于其通常只有一个粮食仓储中心,建成高速局域网也较为简单,适合采用这种方式。对于大型粮食物流企业,由于其地域覆盖较广,业务量也较大,采用 C/S结构较为合理。

6 结束语

本文从逻辑和实现角度对基于 GIS的粮食物流配送决策支持系统 (GLD_DSS)进行了总体分析。结合 GIS的信息需求,对数据库系统进行了较为细致的分析;结合了 GLD_DSS的主要决策问题,对模型库的功能进行了分析。此外,还考虑了系统的实现架构。粮食物流配送决策支持系统是一个较为复杂的系统,本文所作的分析只是宏观性的,在具体的实现过程中还要更进一步地进行细化。由于此系统涉及面宽,实际情况比较复杂,与粮食企业的经济效益有较大的联系,需要不断进行更加深入的研究和开发。

参考文献:

[1] 孙吉贵,白洪涛,于海鸿,等. 粮食调拨决策支持系统的设计与实

(上接第 1393页) 统作为湖北省教育考试院组织实施的"高等教育自学上机考试系统(forWindows)"之核心子系统,负责后台数据库大规模试题数据的初始化工作。该软件在上机考试系统中成功应用,并在试卷文本语块识别过程中表现出良好的性能。

6 结束语

本文给出了一种基于规则的大规模试卷文本语块识别方法及识别模型,并以该模型为基础实现了一个原型系统,较好地解决了数据库中大规模试题数据的高效初始化问题。实验表明,该识别模型具有良好的性能,对于试卷在未来可能发生的各种变化,如题型改变或题量增减,模型依然有较强的适应性。限于篇幅,本文未给出模型实现的源码。由于系统使用之前必须预先了解试题库的数据结构及类型,对于未知结构的关系表,目前尚未考虑涉及动态获取并绑定表属性等问题,在一定程度上限制了本系统的推广及应用范围。这将是下一步研究工作的重点。

参考文献:

- [1] ABNEY S P. Parsing by chunks [M] //BERW ICK R, ABNEY S, TENNY C, et al. Principle-based parsing Dordercht K lawer A cadem ic Publishers, 1991: 257-278.
- [2] SANG E F T K, BUCHHOLZ S Introduction to the CoNLL-2000 shared task: chunking[C] //Proc of the 2nd Workshop on Learning Language in Logic Morristown Association for Computational Linguistics, 2000: 127-132.

- 现[]]. 吉林大学学报:信息科学版, 2005, 23(1): 81-85.
- [2] 麂应荣. 粮食物流系统优化研究[D]. 长春: 吉林大学, 2007.
- [3] HARPER P R, SHAHANIA K A decision support system for the care of H IV and A IDS patients in India[J]. European Journal of Operational Research, 2003, 147(1): 187-197.
- [4] BERGEY PK, RAGSDALETC, HOSKOTEM. A decision support system for the electrical power districting problem [J]. Decision Support Systems, 2003, 36(1): 1-17.
- [5] HAASTRUP P, MAN EZZO V. A decision support system for urban waste management [J]. European Journal of Operational Research 1998 109(2): 330-341.
- [6] PALK, PALMER O. A decision-support system for business acquisitions [J]. Decision Support Systems, 2000, 27(4): 411-429.
- [7] WATTAU C J AKOKA J. Logistics information system auditing using expert system technology [J]. Expert Systems with Applications, 1996, 11(4): 463-473.
- [8] ALD N N, STAHRE F. Electronic commerce, marketing channels and logistics platforms a wholesaler perspective [J]. European Journal of Operational Research, 2003, 144(2): 270-279.
- [9] 唐孝飞,孙壮志,胡思继.物流配送决策支持系统的分析[J].北方交通大学学报,2002 26(5): 92-97
- [10] 陈述彭, 鲁学军, 周成虎. 地理信息系统导论 [M]. 北京: 科学出版 社, 2000.
- [11] ZHANG Q in-wen, ZHEN Tong ZHU Yu-hua, et al. A hybrid in telligent algorithm for the vehicle routing with time windows [C] //Proc of International Conference on Intelligent Computing Berlin Springer-Verlag, 2008, 47-54.
- [12] 甄彤, 张秋闻, 马志. 基于改 进蚁群算法的粮食物流调度研究[J]. 河南工业大学学报:自然科学版, 2008, 29(3): 62-65.
- [13] MOYN HAN G P, SRAJ STER NG J U P, et al Decision support system for strategic logistics planning [J]. Computer in Industry, 1995, 26(1): 75-84.

(2): 74-83

- [4] ARGAMON S. DAGAN I KRYMOLOWSKIY. A memory-based approach to learning shallow natural language patterns [C] //Proc of the 36th Annua Meeting of the Association for Computational Linguistics. Morristown Association for Computational Linguistics, 1998: 67-73
- [5] ZHANG Tong DAMERAU F, JOHNSON D. Text chunking based on a generalization of winnow [J]. Journal of Machine Learning Research, 2002, 2 615-637.
- [6] 周强, 孙茂松. 汉语句子的组成分析体系 [J]. 计算机学报, 1999, 22(11): 1158-1165.
- [7] 梁颖红,赵铁军.基于关联度评价的中心词扩展的英语文本语块识别[J]. 计算机研究与发展, 2006 43(1): 153-158.
- [8] 梁颖红, 赵铁军, 于浩. 基于改进 K-均值聚类的汉语语块识别 [J]. 哈尔滨工业大学学报, 2007, 39(7): 1106-1109.
- [9] 魏玮, 杜金华, 徐拔. 基于分层语块分析的统计翻译研究 [J]. 中文信息学报, 2007, 21(5): 87-90.
- [10] 秦玉平, 王秀坤, 艾青, 等. 多主题文本分类的实现算法 [J]. 计算机工程, 2008, 34(2): 190-192.
- [11]潘大志,成琥,黄青松.基于规则、串频统计和上下文关系的现代 汉语分词系统的实现[J].内蒙古师范大学学报:自然科学版, 2008 37(1):71-74
- [12] 李宏乔, 樊孝忠. 汉语文本中特殊符号串的自动识别技术 [J]. 计算机工程, 2004, 30(12): 114-115
- [13] 陈永府, 杨小献, 黄正东. 基于规则的数据收集研究[J]. 计算机工程与设计, 2007, 28(1): 158-161.
- [14] ABNEY S. Partial parsing via finite-state cascades[J]. Natural Lan-

[3] 孙宏林, 俞士汶 浅层句法分析方法概述 [J]. 当低语言学, 2000 2 guage Engineering, 1996, 2(4): 337-344/www.cnki.net