

最大匹配算法研究

万 廷

(武汉理工大学 计算机科学与技术学院, 湖北 武汉 430063)

摘要: 最大匹配算法是中文分词中最常用的方法,但其有着过分依赖于词典的弊端。对最大匹配算法进行了深入探讨与研究,使用 n -gram 技术更新词典解决其弊端,从而提高分词效果。最后通过双向匹配算法与 n -gram 相结合的实验验证了该方案的可行性,并对该方案进行了总结。

关键词: 中文分词;最大匹配; n -gram;词频;双向匹配

中图分类号: TP391.1

文献标识码: A

文章编号: 1674-7720(2012)08-0062-02

The research of maximum match algorithm

Wan Ting

(College of Computer Science and Technology, Wuhan University of Technology, Wuhan 430063, China)

Abstract: As the most common algorithm in Chinese segment, maximum match algorithm also has a problem that excessive depend on the dictionary. After the intensive study of maximum match algorithm, this paper presents a way to update the dictionary by using the technology of n -gram to solve its problem. Thus improve the quality of maximum match algorithm. Finally do a test that use bilateral match and n -gram technology together to prove the way is feasible and the way is summarized.

Key words: Chinese segment; maximum match; n -gram; word frequency; bilateral match

作为计算机信息处理中最基础、最关键的技术,中文分词一直是人们研究的热点。中文分词就是将连续的汉字序列按照一定的规律分割成一个个单独的词的过程^[1]。在英文句子中,单词之间是以空格作为自然分界符的,所以英文分词比较简单;而中文以字为基本单位,将一序列字串联在一起形成句子,从而表达意思,中文的句和段能通过明显的分界符来划分,但是词没有一个形式上的分界符,所以中文分词比英文分词相对困难许多。中文分词方法总结起来大致可分为三大类:基于词典直接匹配的分词方法、基于规则和理解的分词方法和基于统计模型的分词方法^[2]。本文主要讨论基于词典匹配算法中的最大匹配算法,针对其过分依赖词典这一弊端进行了探讨并提出了对策。

1 最大匹配算法

最大匹配算法是最常用也是最基本的字符串匹配算法之一。它能够保证切分出来的词长度最大,同时易于实现^[3]。最大匹配算法包括正向最大匹配算法、逆向最大匹配算法和双向最大匹配算法。

1.1 正向最大匹配算法

正向最大匹配算法流程^[4]如图 1 所示。

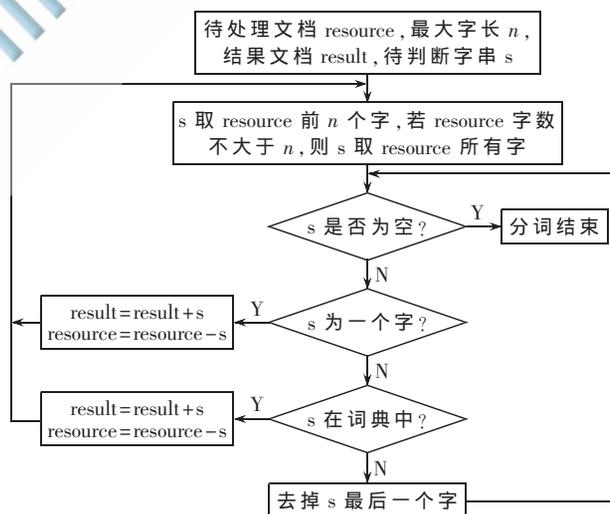


图 1 正向最大匹配算法流程图

以“中华人民共和国简称中国”为例,设定取词长度 n 为 8,待匹配字符串为 s ,按照上述步骤处理过程为:

(1) s 为“中华人民共和国简”,查找词典进行匹配操作,发现没有该词;

(2) s 去掉最后一个字,变为“中华人民共和国”,查

技术与方法 Technique and Method

- 词算法[J]. 科技创新导报, 2009(9): 248.
- [4] 赵源. 基于最大匹配的中文分词改进算法研究[J]. 科技信息, 2010(35): 487, 496.
- [5] 王瑞雷, 栾静, 潘晓花, 等. 一种改进的中文分词正向最大匹配算法[J]. 计算机应用与软件, 2011, 28(3): 195-197.
- [6] 吴胜远. 一种汉语分词方法[J]. 计算机研究与发展, 1996, 33(4): 306-311.
- [7] 李文, 洪亲, 滕忠坚, 等. 基于 n-gram 的字符串分割技术的算法实现[J]. 计算机与现代化, 2010(9): 85-87.
- [8] 张磊, 张代远. 中文分词算法解析[J]. 电脑知识与技术, 2009, 5(1): 192-193.

(收稿日期: 2011-12-19)

作者简介:

万蕊, 男, 1988年生, 硕士研究生, 主要研究方向: 软件工程方法与技术。

