

数据库、信号与信息处理

改进的 χ^2 统计文本特征选择方法

肖婷, 唐雁

西南大学 计算机与信息科学学院, 重庆 400715

收稿日期 2008-11-20 修回日期 2009-2-13 网络版发布日期 2009-5-8 接受日期

摘要 特征选择是当今研究领域的一个热点, 尤其是文本分类领域中的热点。针对 χ^2 统计方法的两个缺陷: 降低了低频词的权重和提高了很少在指定类中出现但普遍存在于其他类的特征在该类中的权重, 对 χ^2 统计方法进行改进, 并通过做模拟和对比实验, 对比改进前后的方法对文本分类的影响。在模拟和对比实验中, 改进后方法的分类效果要好于传统的方法。

关键词 [文本分类](#) [特征选择](#) [\$\chi^2\$ 统计](#)

分类号

Improved χ^2 statistics method for text feature selection

XIAO Ting, TANG Yan

School of Computer & Information Science, Southwest University, Chongqing 400715, China

Abstract

Feature selection is a hot topic in current search field, especially in the field of text categorization. In this paper, χ^2 statistical method has two defects. One is reducing the weight of the low-frequency words. The other is increasing the weight of the characteristics in the designated class. The characteristics little appear in designated class but other classes. Through simulation and comparison experiment, the result is better than before.

Key words [text categorization](#) [feature selection](#) [\$\chi^2\$ statistics](#)

DOI: 10.3778/j.issn.1002-8331.2009.14.041

通讯作者 肖婷 xiaott2006@163.com

扩展功能

本文信息

▶ [Supporting info](#)

▶ [PDF\(376KB\)](#)

▶ [\[HTML全文\]\(0KB\)](#)

▶ [参考文献](#)

服务与反馈

▶ [把本文推荐给朋友](#)

▶ [加入我的书架](#)

▶ [加入引用管理器](#)

▶ [复制索引](#)

▶ [Email Alert](#)

▶ [文章反馈](#)

▶ [浏览反馈信息](#)

相关信息

▶ [本刊中 包含“文本分类”的相关文章](#)

▶ [本文作者相关文章](#)

· [肖婷](#)

· [唐雁](#)