

知识组织与知识管理

基于熵的新闻网页抽取方法的研究

朱红灿¹;龙朝阳^{1,2}

湘潭大学管理学院¹

收稿日期 2007-2-5 修回日期 2007-2-28 网络版发布日期 2007-4-30 接受日期

摘要 为了减少或根除新闻网站中大量非主题信息的干扰,提出一种新闻网页抽取方法,采用基于熵的计算和DOM树的知识,从新闻网页中抽取主题文档和相关链接。

Abstract In this paper,an approach for news article extraction from Web page is proposed and this approach applies information theory to DOM tree. Experiment on several news Web sites shows that it is practical.

关键词 [熵](#) [信息抽取](#) [信息块](#) [DOM](#)

Key words Entropy; Information extraction; Informative block; DOM

分类号 [TP181](#)

DOI:

通讯作者:

朱红灿 zhuhongcan@xtu.edu.cn

作者个人主页: 朱红灿 龙朝阳

扩展功能

本文信息

▶ [Supporting info](#)

▶ [PDF](#)(OKB)

▶ [\[HTML全文\]](#)(OKB)

▶ [参考文献\[PDF\]](#)

▶ [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](#)

▶ [引用本文](#)

▶ [Email Alert](#)

▶ [文章反馈](#)

▶ [浏览反馈信息](#)

相关信息

▶ [本刊中 包含“熵”的 相关文章](#)

▶ 本文作者相关文章

• [朱红灿](#)

• [龙朝阳](#)

•