

人工智能及识别技术

含有语义特征的网页新闻自动抽取

施 洋, 张 奇, 黄萱菁

(复旦大学计算机科学技术学院, 上海 200433)

收稿日期 修回日期 网络版发布日期 接受日期

摘要 通过分析新闻网页的语义特征以及网页之间存在的通用性质, 提出一种含有语义特征的网页新闻自动抽取方法, 包括利用语义分类器识别新闻网页中的种子信息以及页面中的局部信息来完成抽取。在分类器中加入语义特征可以使F1值达到94.2%。在语义分类器与局部特征结合的情况下, F1值可以达到96.9%。实验结果证明, 该方法能有效提高网页信息抽取算法的精度, 降低机器学习所需要的标注成本。

关键词 [网络信息抽取; 语义特征; 局部特征](#)

分类号 [TP393](#)

DOI:

通讯作者:

作者个人主页:

施 洋; 张 奇; 黄萱菁

扩展功能

本文信息

▶ [Supporting info](#)

▶ [PDF](#) (407KB)

▶ [\[HTML全文\]](#) (0KB)

▶ [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“网络信息抽取; 语义特征; 局部特征”的 相关文章](#)

▶ [本文作者相关文章](#)