

人工智能

一种基于HITS的主题敏感爬行方法

蒋宗礼¹;徐学可¹;李帅²

北京工业大学计算机学院¹

清华大学电子工程系²

收稿日期 2007-10-25 修回日期 2007-11-10 网络版发布日期 2008-4-28 接受日期

摘要 基于主题的信息采集是信息检索领域内一个新兴且实用的方法,通过将下载页面限定在特定的主题领域,来提高搜索引擎的效率和提供信息的质量。其思想是在爬行过程中按预先定义好的主题有选择地收集相关网页,避免下载主题不相关的网页,其目标是更准确地找到对用户有用的信息。探讨了主题爬虫的一些关键问题,通过改进主题模型、链接分类模型的学习方法及链接分析方法来提高下载网页的主题相关度及质量。在此基础上设计并实现了一个主题爬虫系统,该系统利用主题敏感HITS来计算网页优先级。实验表明效果良好。

关键词 [主题爬虫](#) [超链接引导的主题搜索](#) [主题模型](#)

分类号

DOI:

对应的英文版文章: [A7105855](#)

通讯作者:

蒋宗礼 jiangzli@bjut.edu.cn

作者个人主页: 蒋宗礼 徐学可 李帅

扩展功能

本文信息

- ▶ [Supporting info](#)
- ▶ [PDF \(839KB\)](#)
- ▶ [\[HTML全文\]\(0KB\)](#)
- ▶ [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“主题爬虫”的 相关文章](#)
- ▶ 本文作者相关文章

- [蒋宗礼](#)
- [徐学可](#)
- [李帅](#)