

数据库与软件技术

不平衡数据集学习中基于初分类的过抽样算法

韩慧<sup>1</sup>; 王路<sup>2</sup>; 温明<sup>2</sup>; 王文渊<sup>2,2</sup>

清华大学自动化系<sup>1</sup>

收稿日期 2006-3-1 修回日期 网络版发布日期 2006-8-1 接受日期

**摘要** 为了有效地提高不平衡数据集中少数类的分类性能,提出了基于初分类的过抽样算法。首先,对测试集进行初分类,以尽可能多地保留多数类的有用信息;其次,对于被初分类预测为少数类的样本进行再次分类,以有效地提高少数类的分类性能。使用美国加州大学欧文分校的数据集将基于初分类的过抽样算法与合成少数类过抽样算法、欠抽样方法进行了实验比较。结果表明,基于初分类的过抽样算法的少数类与多数类的分类性能都优于其他两种算法。

**关键词** [不平衡数据集](#) [过抽样](#) [欠抽样](#)

分类号

**DOI:**

对应的英文版文章: [6020959](#)

通讯作者:

韩慧 [hanh01@mails.tsinghua.edu.cn](mailto:hanh01@mails.tsinghua.edu.cn)

作者个人主页: 韩慧 王路 温明 王文渊

## 扩展功能

本文信息

▶ [Supporting info](#)

▶ [PDF \(753KB\)](#)

▶ [\[HTML全文\]\(0KB\)](#)

▶ [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“不平衡数据集”的  
相关文章](#)

▶ 本文作者相关文章

· [韩慧](#)

· [王路](#)

· [温明](#)

· [王文渊](#)