Accepted Manuscript

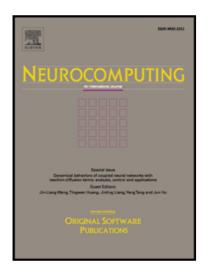
The Gradual Resampling Ensemble for mining imbalanced data streams with concept drift

Siqi Ren, Bo Liao, Wen Zhu, Zeng Li, Wei Liu, Keqin Li

 PII:
 S0925-2312(18)30096-1

 DOI:
 10.1016/j.neucom.2018.01.063

 Reference:
 NEUCOM 19265



To appear in: Neurocomputing

Received date:9 March 2017Revised date:10 November 2017Accepted date:24 January 2018

Please cite this article as: Siqi Ren, Bo Liao, Wen Zhu, Zeng Li, Wei Liu, Keqin Li, The Gradual Resampling Ensemble for mining imbalanced data streams with concept drift, *Neurocomputing* (2018), doi: 10.1016/j.neucom.2018.01.063

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

The Gradual Resampling Ensemble for mining imbalanced data streams with concept drift

Siqi Ren^a, Bo Liao^{a,*}, Wen Zhu^a, Zeng Li^b, Wei Liu^a, Keqin Li^{a,c}

^aCollege of Information Science and Engineering, Hunan University, Changsha 410082, Hunan, China

⁵ ^bSchool of Computer Science and Technology, University of Science and Technology of China, Hefei 230027, Anhui, China

^cDepartment of Computer Science, State University of New York, New Paltz, New York 12561, USA

8 Abstract

1

2

3

4

7

Knowledge extraction from data streams has received increasing interest in recent years. Howg ever, most of the existing studies assume that the class distribution of data streams is relatively 10 balanced. The reaction of concept drifts is more difficult if a data stream is class imbalanced. 11 Current oversampling methods generally selectively absorb the previously received minority ex-12 amples into the current minority set by evaluating similarities of past minority examples and the 13 current minority set. However, the similarity evaluation is easily affected by data difficulty fac-14 tors. Meanwhile, these oversampling techniques have ignored the majority class distribution, thus 15 risking class overlapping. 16 To overcome these issues, we propose an ensemble classifier called Gradual Resampling En-17 semble (GRE). GRE could handle data streams which exhibit concept drifts and class imbalance. 18 On the one hand, a selectively resampling method, where drifting data can be avoidable, is ap-19 plied to select a part of previous minority examples for amplifying the current minority set. The 20 disjuncts can be discovered by the DBSCAN clustering, and thus the influences of small disjuncts 21 and outliers on the similarity evaluation can be avoidable. Only those minority examples with low 22 probability of overlapping with the current majority set can be selected for resampling the current 23 minority set. On the other hand, previous component classifiers are updated using latest instances. 24 Thus, the ensemble could quickly adapt to a new condition, regardless types of concept drifts. 25 Through the gradual oversampling of previous chunks using the current minority events, the class 26 distribution of past chunks can be balanced. Favorable results in comparison to other algorithms 27 suggest that GRE can maintain good performance on minority class, without sacrificing majority 28

²⁹ class performance.

30 Keywords:

³¹ Concept drift, Data stream mining, Ensemble classifier, Class imbalance .

*Corresponding author

Preprint submitted to Neurocomputing

Email addresses: siqirenzl@163.com (Siqi Ren), dragonbw@163.com (Bo Liao), syzhuwen@163.com (Wen Zhu), lizeng@mail.ustc.edu.cn (Zeng Li), lw2001184@163.com (Wei Liu), lik@newpaltz.edu (Keqin Li)

Download English Version:

https://daneshyari.com/en/article/6864414

Download Persian Version:

https://daneshyari.com/article/6864414

Daneshyari.com