



Data and text mining

METAbolomics data Balancing with Over-sampling Algorithms (META-BOA): an online resource for addressing class imbalance

Emily Hashimoto-Roth^{1,2}, Anuradha Surendra³, Mathieu Lavallée-Adam¹,
Steffany A.L. Bennett^{1,2,4,*} and Miroslava Čuperlović-Culf^{1,2,3,*}

¹Department of Biochemistry, Microbiology and Immunology, Ottawa Institute of Systems Biology, Ottawa, ON, Canada, ²Neural Regeneration Laboratory and India Taylor Lipidomic Research Platform, Ottawa, ON K1H 8M5, Canada, ³Digital Technologies Research Centre, National Research Council of Canada, Ottawa, ON K1A 0R6, Canada and ⁴Department of Chemistry and Biomolecular Sciences, Centre for Catalysis Research and Innovation, University of Ottawa, Ottawa, ON K1N 6N5, Canada

*To whom correspondence should be addressed.

Associate Editor: Jonathan Wren

Received on April 5, 2022; revised on September 4, 2022; editorial decision on September 20, 2022

Abstract

Motivation: Class imbalance, or unequal sample sizes between classes, is an increasing concern in machine learning for metabolomic and lipidomic data mining, which can result in overfitting for the over-represented class. Numerous methods have been developed for handling class imbalance, but they are not readily accessible to users with limited computational experience. Moreover, there is no resource that enables users to easily evaluate the effect of different over-sampling algorithms.

Results: METAbolomics data Balancing with Over-sampling Algorithms (META-BOA) is a web-based application that enables users to select between four different methods for class balancing, followed by data visualization and classification of the sample to observe the augmentation effects. META-BOA outputs a newly balanced dataset, generating additional samples in the minority class, according to the user's choice of Synthetic Minority Over-sampling Technique (SMOTE), Borderline-SMOTE (BSMOTE), Adaptive Synthetic (ADASYN) or Random Over-Sampling Examples (ROSE). To present the effect of over-sampling on the data META-BOA further displays both principal component analysis and t-distributed stochastic neighbor embedding visualization of data pre- and post-over-sampling. Random forest classification is utilized to compare sample classification in both the original and balanced datasets, enabling users to select the most appropriate method for their further analyses.

Availability and implementation: META-BOA is available at <https://complimet.ca/meta-boa>.

Contact: miroslava.cuperlovic-culf@nrc-cnrc.gc.ca or steffany.bennett@gmail.com.

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

Metabolomics and lipidomics are emerging fields with unprecedented potential to elucidate novel pathogenic metabolic processes, producing diagnostic and treatment-relevant knowledge. Both 'omics' disciplines generate high-dimensional data. Yet classes of samples are often not equally balanced within datasets, presenting a notable challenge for machine learning strategies. While some machine learning models, such as support vector machines, are relatively robust to imbalanced datasets, most models are subject to overfitting, favoring the over-represented class. These challenges can bias overall analyses, contrary to the intent of an unbiased

machine learning approach. Class imbalance can be addressed by under-sampling the majority class or over-sampling the minority class using a variety of computational algorithms [reviewed in Kovács (2019) and Sharma *et al.* (2022)]. However, to our knowledge, there is no tool that enables researchers to rapidly perform and compare different over-sampling methods. METAbolomics data Balancing with Over-sampling Algorithms (META-BOA) is a user-friendly web-based application that generates balanced datasets, for both binary and multiclass datasets. The effect of each method on the dataset may also be evaluated, through visual representations of unsupervised and supervised analyses, comprising of

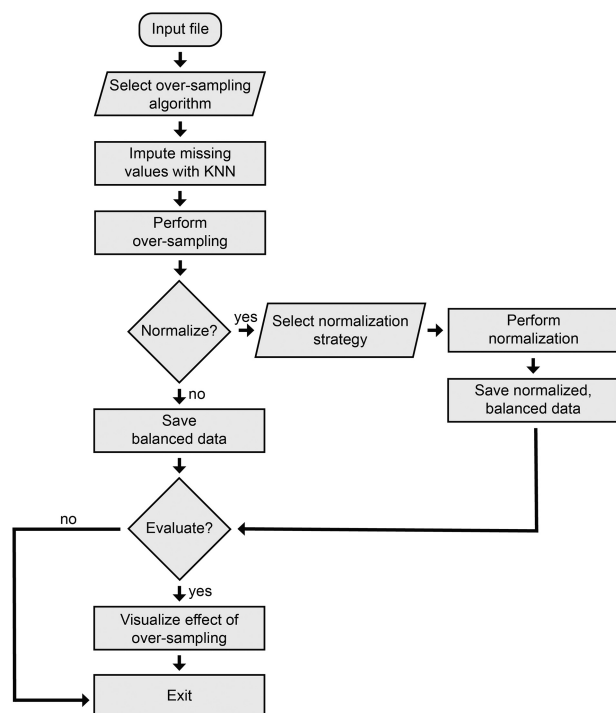


Fig. 1. Outline of META-BOA workflow. Output is class-balanced dataset with or without normalization. It is recommended that the user impute data prior to submitting to META-BOA. Remaining missing data are imputed using the k -nearest neighbors method of the *Scikit-learn* package. Example of META-BOA results is provided in [Supplementary Materials](#)

principal component analysis (PCA) and t-distributed stochastic neighbor embedding (t-SNE) visualizations and Random Forest classification.

2 Implementation

META-BOA is an application, accessible through any web browser that, through over-sampling, generates balanced sample groups for both binary and multiclass datasets. META-BOA provides users with the option to select between four over-sampling methods: Synthetic Minority Over-sampling Technique (SMOTE), Borderline-SMOTE (BSMOTE), Adaptive Synthetic (ADASYN) and Random Over-Sampling Examples (ROSE). These algorithms were chosen to maximize user alternatives: SMOTE randomly generates new synthetic samples within the known minority class samples without replication (Chawla et al., 2002). BSMOTE generates synthetic samples on the borderline between majority and minority instances (Han et al., 2005). ADASYN creates more samples in the neighborhood of minority samples that are in the vicinity of a larger number of the majority class cases (He et al., 2008). ROSE is a bootstrap-based approach that creates synthetic samples in the neighborhood of minority class features (Lunardon et al., 2014).

The META-BOA workflow is presented in Figure 1. The platform was coded in Python 3.7 and deployed with an R Shiny graphical user interface. To run META-BOA, the user selects a single. CSV file containing their imbalanced dataset. The user then specifies:

1. Choice of over-sampling algorithm (default = SMOTE);
2. Normalization applied to data (default = no normalization); and
3. Dimensionality reduction and visualization (default = yes).

Users can choose from four over-sampling algorithms: (i) SMOTE, (ii) BSMOTE, (iii) ADASYN and (iv) ROSE. Each option

has been implemented using the *imblearn* package (version 0.8.0) (Lemaitre et al., 2017). Data visualization and Random Forest classification provide a means to evaluate the effect of the selected over-sampling algorithm. META-BOA performs both PCA and t-SNE dimensionality reduction, evaluating linear and non-linear reduction, respectively. The Random Forest classification model is trained and tested with a 90/10 random split and 250 trees, implementing a 10-fold cross-validation scheme, wherein only the training data is over-sampled. These evaluations are intended to provide a preliminary analysis of the newly balanced dataset, while showing the effect of over-sampling on the visualization and classification results. The over-sampled dataset is provided for further analysis externally. Dimensionality reduction and classification are implemented using the *Scikit-learn* package (version 0.24.0) (Pedregosa et al., 2011). Result files are made available for download immediately upon analysis completion, as a zipped file which includes a .CSV datafile and .SVG image files.

3 Conclusion

META-BOA is an open-access web-based application for balancing sample classes in metabolomics and lipidomics datasets, or any other type of data, using over-sampling algorithms. This tool will help facilitate future analyses of metabolomic and lipidomic data. We encourage users to implement META-BOA in their data pre-processing pipelines and test the different over-sampling algorithms available to evaluate their effects on their given datasets.

Funding

This work was supported in part by RGPIN-2019-06796 to SALB from the Natural Sciences and Engineering Research Council of Canada (NSERC), as well as an NSERC CREATE Matrix Metabolomics Training grant to S.A.L.B. and M.L.-A., an operating grant AI-4D-102-3 to S.A.L.B. and M.C.-C. from the National Research Council AI for Design Challenge Program, and an NSERC Discovery Grant to M.L.-A. This research was enabled by support provided by Compute Ontario and Compute Canada with Resource Allocation to M.L.-A. E.H.-R. received an NSERC CREATE Matrix Metabolomics Scholarship.

Conflict of Interest: none declared.

References

- Chawla, N.V. et al. (2002) SMOTE: synthetic minority over-sampling technique. *J. Artif. Int. Res.*, **16**, 321–357.
- Han, H. et al. (2005) Borderline-SMOTE: a new over-sampling method in imbalanced data sets learning. In: Huang, D.-S. et al. (eds.) *Advances in Intelligent Computing*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 878–887.
- He, H. et al. (2008) ADASYN: adaptive synthetic sampling approach for imbalanced learning. In: *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*. pp. 1322–1328.
- Kovács, G. (2019) An empirical comparison and evaluation of minority over-sampling techniques on a large number of imbalanced datasets. *Appl. Soft Comput.*, **83**, 105662.
- Lemaitre, G. et al. (2017) Imbalanced-Learn: a python toolbox to tackle the curse of imbalanced datasets in machine learning. *J. Mach. Learn. Res.*, **18**, 559–563.
- Lunardon, N. et al. (2014) ROSE: a package for binary imbalanced learning. *R J.*, **6**, 79.
- Pedregosa, F. et al. (2011) Scikit-learn: machine learning in python. *J. Mach. Learn. Res.*, **12**, 2825–2830.
- Sharma, S. et al. (2022) A review of the oversampling techniques in class imbalance problem. In: Khanna, A. (ed.) *International Conference on Innovative Computing and Communications*. Springer Singapore, Singapore, pp. 459–472.