

Imbalanced-learn: A Python Toolbox to Tackle the Curse of Imbalanced Datasets in Machine Learning

Guillaume Lemaitre, Fernando Nogueira, Christos Aridas

► **To cite this version:**

Guillaume Lemaitre, Fernando Nogueira, Christos Aridas. Imbalanced-learn: A Python Toolbox to Tackle the Curse of Imbalanced Datasets in Machine Learning. *Journal of Machine Learning Research*, *Journal of Machine Learning Research*, 2017, 18, pp.1 - 5. hal-01516244

HAL Id: hal-01516244

<https://hal.inria.fr/hal-01516244>

Submitted on 29 Apr 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Imbalanced-learn: A Python Toolbox to Tackle the Curse of Imbalanced Datasets in Machine Learning

Guillaume Lemaître

*Parietal team, Inria, CEA, Université Paris-Saclay
1 Rue Honoré d'Estienne d'Orves, 91120 Palaiseau, France*

GUILLAUME.LEMAITRE@INRIA.FR

Fernando Nogueira

*ShoppeAI
488 Wellington Street West, Suite 304, Toronto, Ontario M5V 1E3, Canada*

FMFNOGUEIRA@GMAIL.COM

Christos K. Aridas

*Computational Intelligence Laboratory
Department of Mathematics
University of Patras
GR-26504 Patras, Greece*

CHAR@UPATRAS.GR

Editor: Geoff Holmes

Abstract

`imbalanced-learn` is an open-source python toolbox aiming at providing a wide range of methods to cope with the problem of imbalanced dataset frequently encountered in machine learning and pattern recognition. The implemented state-of-the-art methods can be categorized into 4 groups: (i) under-sampling, (ii) over-sampling, (iii) combination of over- and under-sampling, and (iv) ensemble learning methods. The proposed toolbox depends only on `numpy`, `scipy`, and `scikit-learn` and is distributed under MIT license. Furthermore, it is fully compatible with `scikit-learn` and is part of the `scikit-learn-contrib` supported project. Documentation, unit tests as well as integration tests are provided to ease usage and contribution. Source code, binaries, and documentation can be downloaded from <https://github.com/scikit-learn-contrib/imbalanced-learn>.

Keywords: Imbalanced Dataset, Over-Sampling, Under-Sampling, Ensemble Learning, Machine Learning, Python.

1. Introduction

Real world datasets commonly show the particularity to have a number of samples of a given class under-represented compared to other classes. This imbalance gives rise to the “class imbalance” problem (Prati et al., 2009) (or “curse of imbalanced datasets”) which is the problem of learning a concept from the class that has a small number of samples.

The class imbalance problem has been encountered in multiple areas such as telecommunication managements, bioinformatics, fraud detection, and medical diagnosis, and has been considered one of the top 10 problems in data mining and pattern recognition (Yang and Wu, 2006; Rastgoo et al., 2016). Imbalanced data substantially compromises the learning process, since most of the standard machine learning algorithms expect balanced class distribution or an equal misclassification cost (He and Garcia, 2009). For this reason, several

approaches have been specifically proposed to handle such datasets. Some of these methods have been implemented mainly in R language (Torgo, 2010; Kuhn, 2015; Dal Pozzolo et al., 2013). Up to our knowledge, there is no python toolbox allowing such processing while cutting edge machine learning toolboxes are available (Pedregosa et al., 2011; Sonnenburg et al., 2010).

In this paper, we present the `imbalanced-learn` API, a *python toolbox to tackle the curse of imbalanced datasets in machine learning*. The following sections present the project vision, a snapshot of the API, an overview of the implemented methods, and finally, we conclude this work by including future functionalities for the `imbalanced-learn` API.

2. Project management

Quality assurance In order to ensure code quality, a set of unit tests is provided leading to a coverage of 99 % for the release 0.2 of the toolbox. Furthermore, the code consistency is ensured by following `PEP8` standards and each new contribution is automatically checked through `landscape`, which provides metrics related to code quality.

Continuous integration To allow both the user and the developer to either use or contribute to this toolbox, Travis CI is used to easily integrate new code and ensure back-compatibility.

Community-based development All the development is performed in a collaborative manner. Tools such as `git`, `GitHub`, and `gitter` are used to ease collaborative programming, issue tracking, code integration, and idea discussions.

Documentation A consistent API documentation is provided using `sphinx` and `numpydoc`. An additional installation guide and examples are also provided and centralized on `GitHub`¹.

Project relevance At the edition time, the repository is visited no less than 2,000 times per week, attracting about 300 unique visitors per week. Additionally, the toolbox is supported by `scikit-learn` through the `scikit-learn-contrib` projects.

3. Implementation design

```

1 from sklearn.datasets import make_classification
2 from sklearn.decomposition import PCA
3 from imblearn.over_sampling import SMOTE
4
5 # Generate the dataset
6 X, y = make_classification(n_classes=2, weights=[0.1, 0.9],
7                           n_features=20, n_samples=5000)
8
9 # Apply the SMOTE over-sampling
10 sm = SMOTE(ratio='auto', kind='regular')
11 X_resampled, y_resampled = sm.fit_sample(X, y)

```

Listing 1: Code snippet to over-sample a dataset using SMOTE.

The implementation relies on `numpy`, `scipy`, and `scikit-learn`. Each sampler class implements three main methods inspired from the `scikit-learn` API: (i) `fit` computes several statistics which are later needed to resample the data into a balanced set; (ii)

1. <https://github.com/scikit-learn-contrib/imbalanced-learn>

Method	Over-sampling		Under-sampling	
	Binary	Mutli-class	Binary	Multiclass
ADASYN (He et al., 2008)	✓	✗	✗	✗
SMOTE (Chawla et al., 2002; Han et al., 2005; Nguyen et al., 2011)	✓	✗	✗	✗
ROS	✓	✓	✗	✗
CC	✗	✗	✓	✓
CNN (Hart, 1968)	✗	✗	✓	✓
ENN (Wilson, 1972)	✗	✗	✓	✓
RENN	✗	✗	✓	✓
AKNN	✗	✗	✓	✓
NM (Mani and Zhang, 2003)	✗	✗	✓	✓
NCL (Laurikkala, 2001)	✗	✗	✓	✓
OSS (Kubat et al., 1997)	✗	✗	✓	✓
RUS	✗	✗	✓	✓
IHT (Smith et al., 2014)	✗	✗	✓	✗
TL (Tomek, 1976)	✗	✗	✓	✗
BC (Liu et al., 2009)	✗	✗	✓	✗
EE (Liu et al., 2009)	✗	✗	✓	✓
SMOTE + ENN (Batista et al., 2003)	✓	✗	✓	✗
SMOTE + TL (Batista et al., 2003)	✓	✗	✓	✗

sample performs the sampling and returns the data with the desired balancing ratio; and (iii) `fit_sample` is equivalent to calling the method `fit` followed by the method `sample`. A class `Pipeline` is inherited from the `scikit-learn` toolbox to automatically combine `samplers`, `transformers`, and `estimators`. Additionally, we provide some specific state-of-the-art metrics to evaluate classification performance.

4. Implemented methods

The `imbalanced-learn` toolbox provides four different strategies to tackle the problem of imbalanced dataset: (i) under-sampling, (ii) over-sampling, (iii) a combination of both, and (iv) ensemble learning. The following subsections give an overview of the techniques implemented.

4.1 Notation and background

Let χ be an imbalanced dataset with χ_{min} and χ_{maj} being the subset of samples belonging to the minority and majority class, respectively. The balancing ratio of the dataset χ is defined as:

$$r_\chi = \frac{|\chi_{min}|}{|\chi_{maj}|}, \quad (1)$$

where $|\cdot|$ denotes the cardinality of a set. The balancing process is equivalent to resample χ into a new dataset χ_{res} such that $r_\chi > r_{\chi_{res}}$.

Under-sampling Under-sampling refers to the process of reducing the number of samples in χ_{maj} . The implemented methods can be categorized into 2 groups: (i) fixed under-sampling and (ii) cleaning under-sampling. *Fixed under-sampling* refer to the methods which perform under-sampling to obtain the appropriate balancing ratio $r_{\chi_{res}}$. Contrary to the previous methods, *cleaning under-sampling* do not allow to reach specifically the balancing ratio $r_{\chi_{res}}$, but rather clean the feature space based on some empirical criteria.

Over-sampling Contrary to under-sampling, data balancing can be performed by over-sampling such that new samples are generated in χ_{min} to reach the balancing ratio $r_{\chi_{res}}$.

Combination of over- and under-sampling Over-sampling can lead to over-fitting which can be avoided by applying cleaning under-sampling methods (Prati et al., 2009).

Ensemble learning Under-sampling methods imply that samples of the majority class are lost during the balancing procedure. Ensemble methods offer an alternative to use most of the samples. In fact, an ensemble of balanced sets is created and used to later train any classifier.

5. Future plans and conclusion

In this paper, we shortly presented the foundations of the `imbalanced-learn` toolbox vision and API. As avenues for future works, additional methods based on prototype/instance selection, generation, and reduction will be added as well as additional user guides.

References

- G. E. Batista, A. L. Bazzan, and M. C. Monard. Balancing training data for automated annotation of keywords: a case study. In *WOB*, pages 10–18, 2003.
- N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. SMOTE: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, pages 321–357, 2002.
- A. Dal Pozzolo, O. Caelen, S. Waterschoot, and G. Bontempi. Racing for unbalanced methods selection. In *International Conference on Intelligent Data Engineering and Automated Learning*, pages 24–31. Springer, 2013.
- H. Han, W.-Y. Wang, and B.-H. Mao. Borderline-smote: a new over-sampling method in imbalanced data sets learning. In *International Conference on Intelligent Computing*, pages 878–887. Springer, 2005.
- P. Hart. The condensed nearest neighbor rule. *Information Theory, IEEE Transactions on*, 14(3):515–516, May 1968.
- H. He and E. Garcia. Learning from imbalanced data. *Knowledge and Data Engineering, IEEE Transactions on*, 21(9):1263–1284, 2009.
- Haibo He, Yang Bai, Eduardo A Garcia, and Shutao Li. Adasyn: Adaptive synthetic sampling approach for imbalanced learning. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, pages 1322–1328. IEEE, 2008.
- M. Kubat, S. Matwin, et al. Addressing the curse of imbalanced training sets: one-sided selection. In *International Conference in Machine Learning*, volume 97, pages 179–186. Nashville, USA, 1997.

- M. Kuhn. Caret: classification and regression training. *Astrophysics Source Code Library*, 1:05003, 2015.
- J. Laurikkala. *Improving identification of difficult small classes by balancing class distribution*. Springer, 2001.
- X.-Y. Liu, J. Wu, and Z.-H. Zhou. Exploratory undersampling for class-imbalance learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 39(2):539–550, 2009.
- I. Mani and I. Zhang. knn approach to unbalanced data distributions: a case study involving information extraction. In *Proceedings of Workshop on Learning from Imbalanced Datasets*, 2003.
- H. M. Nguyen, E. W. Cooper, and K. Kamei. Borderline over-sampling for imbalanced data classification. *International Journal of Knowledge Engineering and Soft Data Paradigms*, 3(1):4–21, 2011.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12(Oct):2825–2830, 2011.
- R. C. Prati, G. E. Batista, and M. C. Monard. Data mining with imbalanced class distributions: concepts and methods. In *Indian International Conference Artificial Intelligence*, pages 359–376, 2009.
- M. Rastgoo, G. Lemaitre, J. Massich, O. Morel, F. Marzani, R. Garcia, and F. Meriaudeau. Tackling the problem of data imbalancing for melanoma classification. In *Bioimaging*, 2016.
- M. R. Smith, T. Martinez, and C. Giraud-Carrier. An instance level analysis of data complexity. *Machine learning*, 95(2):225–256, 2014.
- S. C. Sonnenburg, S. Henschel, C. Widmer, J. Behr, A. Zien, F. de Bona, A. Binder, C. Gehl, V. Franc, et al. The SHOGUN machine learning toolbox. *Journal of Machine Learning Research*, 11(Jun):1799–1802, 2010.
- I. Tomek. Two modifications of CNN. *Systems, Man, and Cybernetics, IEEE Transactions on*, 6:769–772, 1976.
- L. Torgo. *Data mining with R: learning with case studies*. Chapman & Hall/CRC, 2010.
- D. L. Wilson. Asymptotic properties of nearest neighbor rules using edited data. *Systems, Man and Cybernetics, IEEE Transactions on*, (3):408–421, 1972.
- Q. Yang and X. Wu. 10 challenging problems in data mining research. *International Journal of Information Technology & Decision Making*, 5(04):597–604, 2006.