

# geometric-smote: A package for flexible and efficient over-sampling

Georgios Douzas, Fernando Bacao\*

*NOVA Information Management School, Universidade Nova de Lisboa*

---

## Abstract

Classification of imbalanced datasets is a challenging task for standard algorithms. G-SMOTE is a state-of-the-art oversampling algorithm that has been shown to outperform the popular SMOTE oversampler and its variations. In this paper, a Python implementation of G-SMOTE is provided that integrates with the Scikit-Learn ecosystem. Therefore, machine learning researchers and practitioners can benefit from its use in a straightforward manner.

*Keywords:* Machine learning, Imbalanced learning problem, Oversampling

---

---

\*Postal Address: NOVA Information Management School, Campus de Campolide, 1070-312 Lisboa, Portugal, Telephone: +351 21 382 8610

*Email addresses:* `gdouzas@novaims.unl.pt` (Georgios Douzas),  
`bacao@novaims.unl.pt` (Fernando Bacao)

Code metadata	
Current code version	v0.1.2
Permanent link to code/repository used for this code version	<a href="https://github.com/AlgoWit/geometric-smote">https://github.com/AlgoWit/geometric-smote</a>
Legal Code License	MIT
Code versioning system used	git
Software code languages, tools, and services used	Python, Travis CI, AppVeyor, Read the Docs, Codecov, CircleCI, zenodo, Anaconda Cloud
Compilation requirements, operating environments & dependencies	Linux, Mac OS, Windows
If available Link to developer documentation/manual	<a href="https://geometric-smote.readthedocs.io/">https://geometric-smote.readthedocs.io/</a>
Support email for questions	georgios.douzas@gmail.com

Table 1: Code metadata

## 1. Motivation and significance

### 1.1. Introduction

The imbalanced learning problem is defined as a machine learning classification task using datasets with binary or multi-class targets where one of the classes, called the majority class, outnumbers significantly the remaining classes, called the minority class(es) [1]. Learning from imbalanced data is a frequent non-trivial problem for academic researchers and industry practitioners. The imbalance learning problem can be found in multiple domains such as chemical and biochemical engineering, financial management, information technology, security, business, agriculture or emergency management [2].

The Imbalance Ratio (IR) is defined as the ratio between the number of samples of the majority class and each of the minority classes. IR values between 100 and 100.000 have been observed [3], [4]. Standard machine learning classification algorithms induce a bias towards the majority class during training. This results in low performance when metrics suitable for imbalanced data are used for the classifier’s evaluation. Figure 1 shows an example of imbalanced data in two dimensions and the resulting decision boundary of a typical classifier when they are used as a training set.

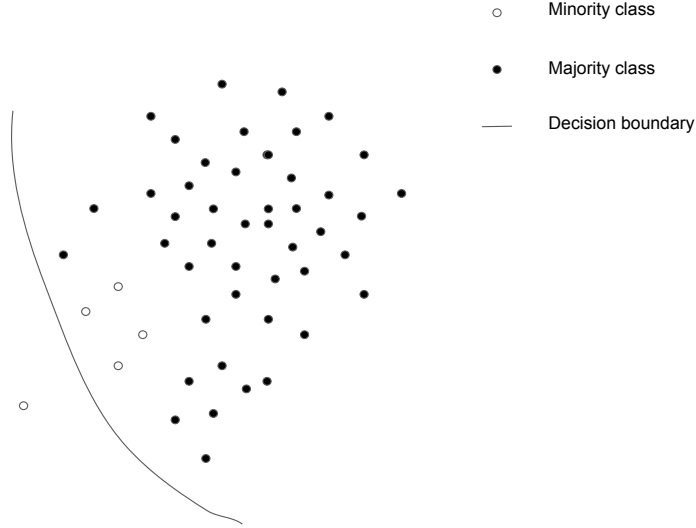


Figure 1: Imbalanced data in two dimensions. The decision boundary of a classifier shows a bias towards the majority class,

### 20 1.2. Oversampling algorithms

21 Various approaches have been proposed to deal with the imbalanced learn-  
 22 ing problem [5]. One general approach is the modification at the data level by  
 23 oversampling the minority class(es). Synthetic Minority Oversampling Tech-  
 24 nique (SMOTE) [3] is the most popular oversampling algorithm. It generates  
 25 synthetic instances along a line segment that joins minority class instances.  
 26 Although SMOTE has been shown to be effective for oversampling imbal-  
 27 anced data, it also has some weaknesses [6]. In order to improve the quality  
 28 of the generated data, many variations of SMOTE have been proposed. Nev-  
 29 ertheless, all of these variations use the same data generation mechanism, i.e.  
 30 linear interpolation between minority class samples.

31 A Python implementation of SMOTE and a couple of its variations is  
 32 available in the Imbalanced-Learn [7] library. Imbalanced-Learn is fully com-  
 33 patible with Scikit-Learn [8], one of the most popular machine learning li-  
 34 braries.

### 35 1.3. Geometric SMOTE

36 Geometric SMOTE (G-SMOTE) [9] uses a different approach compared  
 37 to existing SMOTE’s variations. More specifically, G-SMOTE oversampling  
 38 algorithm substitutes the data generation mechanism of SMOTE by defin-  
 39 ing a flexible geometric region around each minority class instance and gen-  
 40 erating synthetic instances inside the boundaries of this region. The al-

41 gorithm requires the selection of the hyperparameters `truncation_factor`,  
 42 `deformation_factor`, `selection_strategy` and `k_neighbors`. The first three  
 43 of them, called geometric hyperparameters, control the shape of the geomet-  
 44 ric region while the later adjusts its size. Figure 2 presents a visual compar-  
 45 ison between the data generation mechanisms of SMOTE and G-SMOTE.

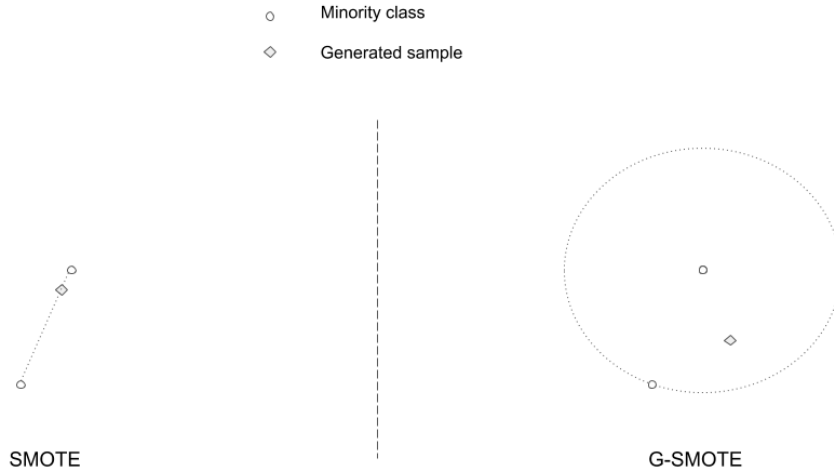


Figure 2: Comparison between the data generation mechanisms of SMOTE and G-SMOTE. SMOTE uses linear interpolation, while G-SMOTE defines a circle as the permissible data generation area.

46 G-SMOTE algorithm has been shown to outperform SMOTE across 69  
 47 imbalanced datasets for various classifiers and evaluation metrics. This pa-  
 48 per, presents a Python implementation of G-SMOTE. In section 2, the soft-  
 49 ware description is given while section 3 provides a demonstrative example  
 50 of its functionalities.

## 51 2. Software description

52 The `geometric-smote` project is written in Python 3.7. It contains an  
 53 object-oriented implementation of the G-SMOTE algorithm as well as an  
 54 extensive online documentation. The implementation provides an API that  
 55 is compatible with Imbalanced-Learn and Scikit-Learn libraries, therefore it  
 56 makes full use of various features that support standard machine learning  
 57 functionalities.

## 58 2.1. Software Architecture

59 The `geometric-smote` project contains the Python package `gsmote`. The  
60 main module of `gsmote` is called `geometric-smote`. It contains the class  
61 `GeometricSMOTE` that implements the G-SMOTE algorithm. The initializa-  
62 tion of an `GeometricSMOTE` instance includes G-SMOTE’s hyperparameters  
63 that control the generation of synthetic data. Additionally, `GeometricSMOTE`  
64 inherits from the `BaseOverSampler` class of Imbalanced-Learn library. There-  
65 fore, an instance of `GeometricSMOTE` class provides the `fit` and `fit_resample`  
66 methods, the two main methods for resampling as explained in subsection 2.2.  
67 This is achieved by implementing the `_fit_resample` abstract method of the  
68 parent class `BaseOverSampler`. More specifically, the function `_make_geometric_sample`  
69 implements the data generation mechanism of G-SMOTE as shortly de-  
70 scribed in section 1.3. This function is called in the `_make_geometric_samples`  
71 method of the `GeometricSMOTE` class in order to generate the appropriate  
72 number of synthetic data for a particular minority class. Finally, the method  
73 `_make_geometric_samples` is called in `_fit_resample` method to generate  
74 synthetic data for all minority classes. Figure 3 provides a visual represen-  
75 tation of the above classes and functions hierarchy.

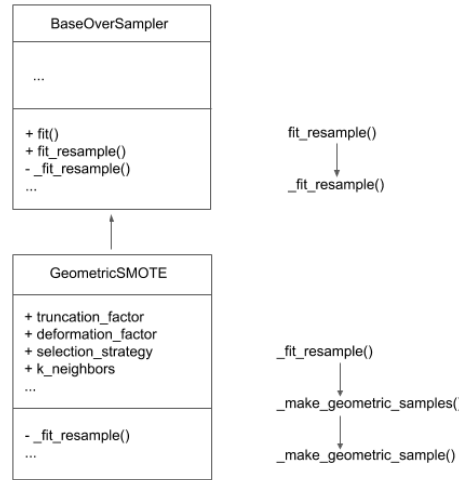


Figure 3: UML class diagrams and callgraphs of main classes and methods.

## 76 2.2. Software Functionalities

77 As it was mentioned in subsection 2.1, the class `GeometricSMOTE` repre-  
78 sents the G-SMOTE oversampler. The initializer of `GeometricSMOTE` includes  
79 the following G-SMOTE’s hyperparameters: `truncation_factor`, `deformation_factor`,

80 `selection_strategy` and `k_neighbors` as explained in subsection 1.3. Once  
 81 the `GeometricSMOTE` object is initialized with a specific parametrization, it  
 82 can be used to resample the imbalanced data represented by the input ma-  
 83 trix `X` and the target labels `y`. Following the Scikit-Learn API, both `X`, `y` are  
 84 array-like objects of appropriate shape.

85 Resampling is achieved by using the two main methods of `fit` and `fit_resample`  
 86 the `GeometricSMOTE` object. More specifically, both of them take as input  
 87 parameters the `X` and `y`. The first method computes various statistics which  
 88 are used to resample `X` while the second method does the same but addition-  
 89 ally returns a resampled version of `X` and `y`.

90 The `geometric-smote` project has been designed to integrate with the  
 91 Imbalanced-Learn toolbox and Scikit-Learn ecosystem. Therefore the `GeometricSMOTE`  
 92 object can be used in a machine learning pipeline, through Imbalanced-  
 93 Learn’s class `Pipeline`, that automatically combines `samplers`, `transformers`  
 94 and `estimators`. The next section provides examples of the above function-  
 95 alities.

## 96 3. Illustrative Examples

### 97 3.1. Basic example

98 An example of resampling multi-class imbalanced data using the `fit_resample`  
 99 method is presented in Listing 1. Initially, a 3-class imbalanced dataset is  
 100 generated. Next, `GeometricSMOTE` object is initialized with default values for  
 101 the hyperparameters, i.e. `truncation_factor = 1.0`, `deformation_factor =`  
 102 `0.0`, `selection_strategy = combined`. Finally, the object’s `fit_resample`  
 103 method is used to resample the data. Printing the class distribution before  
 104 and after resampling confirms that the resampled data are perfectly bal-  
 105 anced. The resampled data `X_res`, `y_res` can be used as training data for  
 106 any classifier in the place of `X`, `y`.

Listing 1: Resampling of imbalanced data using the `fit_resample` method.

```
107 # Import classes and functions.
108 from collections import Counter
109 from gsmote import GeometricSMOTE
110 from sklearn.datasets import make_classification
111
112 # Generate an imbalanced 3-class dataset.
113 X, y = make_classification(
114     random_state=23,
115     n_classes=3,
116     n_informative=5,
```

```

117     n_samples=500,
118     weights=[0.8, 0.15, 0.05]
119 )
120
121 # Create a GeometricSMOTE object with default hyperparameters.
122 gsmote = GeometricSMOTE(random_state=10)
123
124 # Resample the imbalanced dataset.
125 X_res, y_res = gsmote.fit_resample(X, y)
126
127 # Print number of samples per class for initial and resampled data.
128 init_count = list(Counter(y).values())
129 resampled_count = list(Counter(y_res).values())
130
131 print(f'Initial class distribution: {init_count}.')
132 # Initial class distribution: [400, 75, 25].
133
134 print(f'Resampled class distribution: {resampled_count}.')
135 # Resampled class distribution: [400, 400, 400].

```

### 136 3.2. Machine learning pipeline

137 The `GeometricSMOTE` object can be used as a part of a machine learning  
138 pipeline. Listing 2 presents a pipeline composed by a G-SMOTE oversampler,  
139 a PCA tranformation and a decision tree classifier. The pipeline is trained  
140 on imbalanced binary-class data and evaluated on a hold-out set. The user  
141 applies the process in a small number of steps while the internal details of  
142 the calculations are hidden.

Listing 2: Training and evaluation of a machine learning pipeline that contains the `GeometricSMOTE` object.

```

143 # Import classes and functions.
144 from gsmote import GeometricSMOTE
145 from sklearn.datasets import make_classification
146 from sklearn.decomposition import PCA
147 from sklearn.tree import DecisionTreeClassifier
148 from sklearn.model_selection import train_test_split
149 from sklearn.metrics import f1_score
150 from imblearn.pipeline import make_pipeline
151
152 # Generate an imbalanced binary-class dataset.
153 X, y = make_classification(

```

```

154         random_state=23,
155         n_classes=2,
156         n_samples=500,
157         weights=[0.8, 0.2]
158     )
159
160     # Split the data to training and hold-out sets.
161     X_train, X_test, y_train, y_test = train_test_split(X, y, random_state
162
163     # Create the pipeline's objects with default hyperparameters.
164     gsmote = GeometricSMOTE(random_state=11)
165     pca = PCA()
166     clf = DecisionTreeClassifier(random_state=3)
167
168     # Create the pipeline.
169     pip = make_pipeline(gsmote, pca, clf)
170
171     # Fit the pipeline to the training set.
172     pip.fit(X_train, y_train)
173
174     # Evaluate the pipeline on the hold-out set using the F-score.
175     test_score = f1_score(y_test, pip.predict(X_test))
176
177     print(f'F-score on hold-out set: {test_score}.')
178     # F-score on hold-out set: 0.7.

```

#### 179 4. Impact and conclusions

180 Classification of imbalanced datasets is a challenging task for standard  
181 machine learning algorithms. G-SMOTE, as an enhancement of the SMOTE  
182 data generation mechanism, provides a flexible and effective way for resam-  
183 pling the imbalanced data that has been shown to outperform SMOTE and  
184 its variations. Machine learning researchers and industry practitioners can  
185 benefit from using G-SMOTE in their work since the imbalanced learning  
186 problem is an active research area as well as a common feature of real-world  
187 applications.

188 The `geometric-smote` project provides a Python implementation of the  
189 state-of-the-art oversampling algorithm G-SMOTE. The main advantage of  
190 this implementation is that it is built on top of Scikit-Learn's ecosystem.  
191 Therefore, using the oversampler in typical machine learning workflows is  
192 an effortless task. Also, the public API of the main class `GeometricSMOTE`



193 is identical to the one implemented in Imbalanced-Learn for all oversam-  
194 plers. Consequently, users of Imbalanced-Learn and Scikit-Learn, that apply  
195 oversampling on imbalanced data, can integrate the `gsmote` package in their  
196 existing work in a straightforward manner.

## 197 5. Conflict of Interest

198 We wish to confirm that there are no known conflicts of interest associated  
199 with this publication and there has been no significant financial support for  
200 this work that could have influenced its outcome.

## 201 References

- 202 [1] N. V. Chawla, A. Lazarevic, L. Hall, K. Boyer, SMOTEBoost: improving  
203 prediction of the minority class in boosting, *Principles of Knowledge Dis-*  
204 *covery in Databases, PKDD-2003* (2003) 107–119doi:10.1007/b13634.  
205 URL [http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.](http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.80.1499)  
206 [1.80.1499](http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.80.1499)
- 207 [2] G. Haixiang, L. Yijing, J. Shang, G. Mingyun, H. Yuanyue, G. Bing,  
208 Learning from class-imbalanced data: Review of methods and appli-  
209 cations, *Expert Systems with Applications* 73 (2017) 220–239. doi:  
210 10.1016/j.eswa.2016.12.035.  
211 URL <https://doi.org/10.1016/j.eswa.2016.12.035>
- 212 [3] N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, SMOTE:  
213 Synthetic minority over-sampling technique, *Journal of Artificial Intel-*  
214 *ligence Research* 16 (2002) 321–357. arXiv:1106.1813, doi:10.1613/  
215 jair.953.
- 216 [4] S. Barua, M. M. Islam, X. Yao, K. Murase, MWMOTE - Majority  
217 weighted minority oversampling technique for imbalanced data set learn-  
218 ing, *IEEE Transactions on Knowledge and Data Engineering* 26 (2) (2014)  
219 405–425. doi:10.1109/TKDE.2012.232.
- 220 [5] A. Fernández, V. López, M. Galar, M. J. del Jesus, F. Herrera, Analysing  
221 the classification of imbalanced data-sets with multiple classes: Bina-  
222 rization techniques and ad-hoc approaches, *Knowledge-Based Systems*  
223 42 (2013) 97–110. doi:http://dx.doi.org/10.1016/j.knosys.2013.  
224 01.018.  
225 URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0950705113000300)  
226 [S0950705113000300](http://www.sciencedirect.com/science/article/pii/S0950705113000300)

- 227 [6] H. He, E. A. Garcia, Learning from Imbalanced Data, IEEE Transactions  
228 on Knowledge and Data Engineering 21 (9) (2009) 1263–1284. `arXiv:`  
229 `arXiv:1011.1669v3`, `doi:10.1109/TKDE.2008.239`.
- 230 [7] G. Lemaitre, F. Nogueira, C. K. Aridas, Imbalanced-learn: A Python  
231 Toolbox to Tackle the Curse of Imbalanced Datasets in Machine Learning,  
232 Journal of Machine Learning Research 18 (2016) 1–5. `arXiv:1609.06570`,  
233 `doi:http://www.jmlr.org/papers/volume18/16-365/16-365.pdf`.  
234 URL `http://arxiv.org/abs/1609.06570`
- 235 [8] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion,  
236 O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vander-  
237 plas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay,  
238 Scikit-learn: Machine learning in Python, Vol. 12, 2011. `arXiv:arXiv:`  
239 `1201.0490v2`, `doi:10.1007/s13398-014-0173-7.2`.  
240 URL `http://dl.acm.org/citation.cfm?id=2078195`
- 241 [9] G. Douzas, F. Bacao, Geometric SMOTE a geometrically enhanced  
242 drop-in replacement for SMOTE, Information Sciences 501 (2019)  
243 118–135. `doi:10.1016/J.INS.2019.06.007`.  
244 URL `https://www.sciencedirect.com/science/article/pii/`  
245 `S0020025519305353?via={%}3Dihub`