

---

# **multi-imbalance**

*Release 0.0.4*

**Damian Horna, Jacek Grycza, Hanna Klimczak, Kamil Pluciński**

**Jun 22, 2022**



# DOCSTRING

<b>1</b>	<b>Installation:</b>	<b>3</b>
<b>2</b>	<b>Example of code:</b>	<b>5</b>
2.1	multi_imbalance . . . . .	5
2.1.1	multi_imbalance package . . . . .	5
2.1.1.1	Subpackages . . . . .	5
2.1.1.2	Module contents . . . . .	19
2.2	License . . . . .	19
2.3	Contact . . . . .	19
	<b>Python Module Index</b>	<b>21</b>
	<b>Index</b>	<b>23</b>



Multi-class imbalance is a common problem occurring in real-world supervised classifications tasks. While there has already been some research on the specialized methods aiming to tackle that challenging problem, most of them still lack coherent Python implementation that is simple, intuitive and easy to use. multi-imbalance is a python package tackling the problem of multi-class imbalanced datasets in machine learning.



**INSTALLATION:**

```
pip install multi-imbalance
```





## EXAMPLE OF CODE:

```
from multi_imbalance.resampling.mdo import MDO

# Mahalanbois Distance Oversampling
mdo = MDO(k=9, k1_frac=0, seed=0)

# read the data
X_train, y_train, X_test, y_test = ...

# preprocess
X_train_resampled, y_train_resampled = mdo.fit_transform(np.copy(X_train), np.copy(y_
→train))

# train the classifier on preprocessed data
clf_tree = DecisionTreeClassifier(random_state=0)
clf_tree.fit(X_train_resampled, y_train_resampled)

# make predictions
y_pred = clf_tree.predict(X_test)
```

## 2.1 multi\_imbalance

### 2.1.1 multi\_imbalance package

#### 2.1.1.1 Subpackages

**multi\_imbalance.datasets package**

**Subpackages**

**multi\_imbalance.datasets.tests package**

**Submodules**

**multi\_imbalance.datasets.tests.test\_data\_loader module**

Test the datasets loader.

```
multi_imbalance.datasets.tests.test_data_loader.test_load_datasets()
```

## Module contents

### Module contents

`multi_imbalance.datasets.load_datasets` (*data\_home*='./../data')

Load the benchmark datasets.

**Parameters** `data_home` – Default catalogue in which the data is stored in .tar.gz format.

#### Returns

OrderedDict of Bunch object. Each Bunch object referred as dataset have the following attributes:

- **dataset.data** : ndarray, shape (n\_samples, n\_features)
- **dataset.target** : ndarray, shape (n\_samples, )
- **dataset.DESCR** : string Description of the each dataset.

### multi\_imbalance.ensemble package

#### Subpackages

#### multi\_imbalance.ensemble.tests package

#### Submodules

#### multi\_imbalance.ensemble.tests.test\_ecoc module

`multi_imbalance.ensemble.tests.test_ecoc.test_dense_and_sparse_with_not_cached_matrices` (*enc*

`multi_imbalance.ensemble.tests.test_ecoc.test_ecoc_with_sklearn_pipeline` (*encoding\_strategy*,  
*over-*  
*sam-*  
*pling*)

`multi_imbalance.ensemble.tests.test_ecoc.test_encoding` (*encoding\_strategy*, *over-*  
*sampling*)

`multi_imbalance.ensemble.tests.test_ecoc.test_hamming_distance` ()

`multi_imbalance.ensemble.tests.test_ecoc.test_no_oversampling` ()

`multi_imbalance.ensemble.tests.test_ecoc.test_own_classifier_without_predict_and_fit` ()

`multi_imbalance.ensemble.tests.test_ecoc.test_own_preprocessing_without_fit_transform` ()

`multi_imbalance.ensemble.tests.test_ecoc.test_predefined_classifiers_and_weighting_without`

`multi_imbalance.ensemble.tests.test_ecoc.test_random_oversampling` ()

`multi_imbalance.ensemble.tests.test_ecoc.test_unknown_classifier` ()

`multi_imbalance.ensemble.tests.test_ecoc.test_unknown_preprocessing` ()

`multi_imbalance.ensemble.tests.test_ecoc.test_with_own_classifier` ()

`multi_imbalance.ensemble.tests.test_ecoc.test_with_own_preprocessing` ()

**multi\_imbalance.ensemble.tests.test\_mrbbagging module**

**class** multi\_imbalance.ensemble.tests.test\_mrbbagging.**TestMRBBagging** (*methodName='runTest'*)  
 Bases: unittest.case.TestCase

Create an instance of the class that will use the named test method when executed. Raises a ValueError if the instance does not have a method with the specified name.

**test\_group\_data** ()

**test\_group\_data\_with\_none** ()

**test\_api** ()

**test\_api\_multiple\_trees** ()

**test\_api\_with\_feature\_selection** ()

**test\_api\_with\_feature\_selection\_sqrt\_features** ()

**test\_api\_with\_random\_feature\_selection** ()

**test\_fit\_with\_invalid\_classifier** ()

**test\_fit\_with\_invalid\_labels** ()

**test\_with\_invalid\_k** ()

**multi\_imbalance.ensemble.tests.test\_ovo module**

multi\_imbalance.ensemble.tests.test\_ovo.**test\_binary\_classifiers** (*classifier*)

multi\_imbalance.ensemble.tests.test\_ovo.**test\_ecoc\_with\_sklearn\_pipeline** (*preprocessing\_btwn,*  
*clas-*  
*si-*  
*fier,*  
*pre-*  
*pro-*  
*cess-*  
*ing*)

multi\_imbalance.ensemble.tests.test\_ovo.**test\_fit\_predict** ()

multi\_imbalance.ensemble.tests.test\_ovo.**test\_max\_voting** ()

multi\_imbalance.ensemble.tests.test\_ovo.**test\_own\_preprocessing\_without\_fit\_resample** ()

multi\_imbalance.ensemble.tests.test\_ovo.**test\_predefined\_classifiers\_and\_preprocessings\_witl**

multi\_imbalance.ensemble.tests.test\_ovo.**test\_unknown\_preprocessing** ()

multi\_imbalance.ensemble.tests.test\_ovo.**test\_unknown\_preprocessing\_between\_strategy\_raises**

multi\_imbalance.ensemble.tests.test\_ovo.**test\_with\_own\_classifier** ()

`multi_imbalance.ensemble.tests.test_ovo.test_with_own_preprocessing()`

### multi\_imbalance.ensemble.tests.test\_soupbagging module

`multi_imbalance.ensemble.tests.test_soupbagging.test_default_classifier()`

`multi_imbalance.ensemble.tests.test_soupbagging.test_exception()`

`multi_imbalance.ensemble.tests.test_soupbagging.test_fit_classifier_classifier()`

`multi_imbalance.ensemble.tests.test_soupbagging.test_soupbagging()`

## Module contents

### Submodules

#### multi\_imbalance.ensemble.ecoc module

```
class multi_imbalance.ensemble.ecoc.ECOC (binary_classifier='KNN', preprocessing='SOUP', encoding='OVO', n_neighbors=3, weights=None)
```

Bases: `sklearn.ensemble._bagging.BaggingClassifier`

ECOC (Error Correcting Output Codes) is ensemble method for multi-class classification problems. Each class is encoded with unique binary or ternary code (where 0 means that class is excluded from training set of binary classifier). Then in the learning phase each binary classifier is learned. In the decoding phase the class which is closest to test instance in the sense of Hamming distance is chosen.

#### Parameters

- **binary\_classifier** – binary classifier used by the algorithm. Possible classifiers:
  - **'tree'**: Decision Tree Classifier,
  - **'NB'**: Naive Bayes Classifier,
  - **'KNN'**: K-Nearest Neighbors
  - **'ClassifierMixin'**: An instance of a class that implements ClassifierMixin
- **preprocessing** – method for oversampling between aggregated classes in each dichotomy. Possible methods:
  - **None**: no oversampling applied,
  - **'globalCS'**: random oversampling - randomly chosen instances of minority classes are duplicated
  - **'SMOTE'**: Synthetic Minority Oversampling Technique
  - **'SOUP'**: Similarity Oversampling Undersampling Preprocessing
  - **'TransformerMixin'**: An instance of a class that implements TransformerMixin
- **encoding** – algorithm for encoding classes. Possible encodings:
  - **'dense'**:  $\lceil 10 \log_2(\text{num\_of\_classes}) \rceil$  dichotomies, -1 and 1 with probability 0.5 each
  - **'sparse'**:  $\lceil 10 \log_2(\text{num\_of\_classes}) \rceil$  dichotomies, 0 with probability 0.5, -1 and 1 with probability 0.25 each

- **'OVO'** : 'one vs one' -  $n(n-1)/2$  dichotomies, where  $n$  is number of classes, one for each pair of classes. Each column has one 1 and one -1 for classes included in particular pair, 0s for remaining classes.
- **'OVA'** : 'one vs all' - number of dichotomies is equal to number of classes. Each column has one 1 and -1 for all remaining rows
- **'complete'** [ $2^{n-1}-1$  dichotomies, reference] T. G. Dietterich and G. Bakiri. Solving multiclass learning problems via error-correcting output codes. Journal of Artificial Intelligence Research, 2:263–286, 1995.
- **n\_neighbors** –
- **weights** – strategy for dichotomies weighting. Possible values:
  - **None** : no weighting applied
  - **'acc'** : accuracy-based weights
  - **'avg\_tpr\_min'** : weights based on average true positive rates of dichotomies

**fit** ( $X, y, minority\_classes=None$ )

#### Parameters

- **X** – two dimensional numpy array (number of samples x number of features) with float numbers
- **y** – one dimensional numpy array with labels for rows in X
- **minority\_classes** – list of classes considered to be minority classes

**Returns** self: object

**predict** ( $X$ )

**Parameters** **X** – two dimensional numpy array (number of samples x number of features) with float numbers

**Returns** numpy array, shape = [number of samples]. Predicted target values for X.

## multi\_imbalance.ensemble.mrbagging module

```
class multi_imbalance.ensemble.mrbagging.MRBBagging(k, learning_algorithm, un-
dersampling=True, feature
selection=False,
random_fs=False,
half_features=True, ran-
dom_state=None)
```

Bases: sklearn.ensemble.\_bagging.BaggingClassifier

Multi-class Roughly Balanced Bagging (MRBBagging) is a generalization of MRBBagging for adapting to multiple minority classes.

Reference: M. Lango, J. Stefanowski: Multi-class and feature selection extensions of RoughlyBalanced Bagging for imbalanced data. J. Intell Inf Syst (2018) 50: 97

#### Parameters

- **k** – number of classifiers (multiplied by 3 when choosing feature selection)
- **learning\_algorithm** – classifier to be used

- **undersampling** – (optional) boolean value to determine if undersampling or oversampling should be performed
- **feature\_selection** – (optional) boolean value to determine if feature selection should be performed
- **random\_fs** – (optional) boolean value to determine if feature selection should be all random (if False, chi<sup>2</sup>, F test and random feature selection are performed)
- **half\_features** – (optional) boolean value to determine if the number of features to be selected should be 50% (if False, it is set to the square root of the base number of features)
- **random\_state** – (optional) the seed of the pseudo random number generator

**fit** (*x*, *y*, **\*\*kwargs**)

Build a MRBBagging ensemble of estimators from the training data.

**Parameters**

- **x** – Two dimensional numpy array (number of samples x number of features) with float numbers.
- **y** – One dimensional numpy array with labels for rows in X.

**Returns** self (object)

**predict** (*data*)

Predict classes for examples in data.

**Parameters data** – Two dimensional numpy array (number of samples x number of features) with float numbers.

**multi\_imbalance.ensemble.ovo module**

**class** multi\_imbalance.ensemble.ovo.OVO (*binary\_classifier='tree', n\_neighbors=3, preprocessing='SOUP', preprocessing\_between='all'*)

Bases: sklearn.ensemble.\_bagging.BaggingClassifier

OVO (One vs One) is an ensemble method that makes predictions for multi-class problems. OVO decomposes problem into  $m(m-1)/2$  binary problems, where *m* is number of classes. Each of binary classifiers distinguishes between two classes. In the learning phase each classifier is learned only with instances from particular two classes. In prediction phase each classifier decides between these two classes. Results are aggregated and final output is derived depending on chosen aggregation model.

**Parameters**

- **binary\_classifier** – binary classifier. Possible classifiers:
  - **'tree'**: Decision Tree Classifier,
  - **'KNN'**: K-Nearest Neighbors
  - **'NB'**: Naive Bayes
  - **'ClassifierMixin'**: An instance of a class that implements ClassifierMixin
- **n\_neighbors** – number of nearest neighbors in KNN, works only if `binary_classifier=='KNN'`
- **preprocessing** – method for preprocessing of pairs of classes in the learning phase of ensemble. Possible values:
  - **None**: no preprocessing applied

- **'globalCS'**: oversampling with globalCS algorithm
- **'SMOTE'**: oversampling with SMOTE algorithm
- **'SOUP'**: oversampling and undersampling with SOUP algorithm
- **'TransformerMixin'**: An instance of a class that implements TransformerMixin
- **preprocessing\_between** – types of classes between which resampling should be applied. Possible values:
  - **'all'**: oversampling between each pair of classes
  - **'maj-min'**: oversampling only between majority and minority classes

**fit** (*X*, *y*, *minority\_classes=None*)

**Parameters**

- **X** – two dimensional numpy array (number of samples x number of features) with float numbers
- **y** – one dimensional numpy array with labels for rows in X
- **minority\_classes** – list of classes considered to be minority

**Returns** self: object

**predict** (*X*)

**Parameters** **X** – two dimensional numpy array (number of samples x number of features) with float numbers

**Returns** numpy array, shape = [number of samples]. Predicted target values for X.

**should\_perform\_oversampling** (*first\_class*, *second\_class*)

## multi\_imbalance.ensemble.soup\_bagging module

**class** multi\_imbalance.ensemble.soup\_bagging.**SOUPBagging** (*classifier=None*,  
*maj\_int\_min=None*,  
*n\_classifiers=5*)

Bases: sklearn.ensemble.\_bagging.BaggingClassifier

Version of Bagging that applies SOUP in each classifier

Reference: Lango, M., and Stefanowski, J. SOUP-Bagging: a new approach for multi-class imbalanced data classification. PP-RAI '19: Polskie Porozumienie na Rzecz Sztucznej Inteligencji (2019).

**Parameters**

- **classifier** – Instance of classifier
- **maj\_int\_min** – dict {'maj': majority class labels, 'min': minority class labels}
- **n\_classifiers** – number of classifiers

**fit** (*X*, *y*, *\*\*kwargs*)

**Parameters**

- **X** – array-like, sparse matrix of shape = [n\_samples, n\_features] The training input samples.
- **y** – array-like, shape = [n\_samples]. The target values (class labels).
- **\*\*kwargs** – dict (optional)

**Returns** self object

**static fit\_classifier** (*args*)

**predict** (*X*, *strategy*: *str* = 'average')

Predict class for *X*. The predicted class of an input sample is computed as the class with the highest sum of predicted probability.

**Parameters**

- **X** – {array-like, sparse matrix} of shape = [*n\_samples*, *n\_features*]. The training input samples.
- **strategy** –
  - 'average': takes max from average values in prediction
  - 'optimistic': takes always best value of probability
  - 'pessimistic': takes always the worst value of probability
  - 'mixed': for minority classes takes optimistic strategy, and pessimistic for others. It requires *maj\_int\_min*

**Returns** array of shape = [*n\_samples*]. The predicted classes.

**predict\_proba** (*X*)

Predict class probabilities for *X*.

**Parameters** **X** – {array-like, sparse matrix} of shape = [*n\_samples*, *n\_features*]. The training input samples.

**Returns** array of shape = [*n\_classifiers*, *n\_samples*, *n\_classes*]. The class probabilities of the input samples.

`multi_imbalance.ensemble.soup_bagging.fit_clf` (*args*)

## Module contents

`multi_imbalance.resampling` package

Subpackages

`multi_imbalance.resampling.tests` package

Submodules

`multi_imbalance.resampling.tests.test_globalcs` module

`multi_imbalance.resampling.tests.test_globalcs.calc_duplicates_quantities` (*X*,  
*y*,  
*X\_oversampled*)

`multi_imbalance.resampling.tests.test_globalcs.get_goal_quantity` (*y*)

`multi_imbalance.resampling.tests.test_globalcs.global_cs_mock` ()

`multi_imbalance.resampling.tests.test_globalcs.test_output_equal_replication` (*X*,  
*y*,  
*global\_cs\_mock*)



---

```
multi_imbalance.resampling.tests.test_globalcs.test_output_length_validate(X,
                                                                              y,
                                                                              global_cs_mock)
```

### multi\_imbalance.resampling.tests.test\_mdo module

```
multi_imbalance.resampling.tests.test_mdo.mdo_mock()
```

```
multi_imbalance.resampling.tests.test_mdo.test_choose_samples(X, y,
                                                                sc_minor_expected,
                                                                weights_expected,
                                                                mdo_mock)
```

```
multi_imbalance.resampling.tests.test_mdo.test_choose_samples_when_correct(mdo_mock)
```

```
multi_imbalance.resampling.tests.test_mdo.test_choose_samples_when_zero_samples_expected(m
```

```
multi_imbalance.resampling.tests.test_mdo.test_mdo_api(mdo_mock)
```

```
multi_imbalance.resampling.tests.test_mdo.test_zero_variance(mdo_mock)
```

### multi\_imbalance.resampling.tests.test\_soup module

```
multi_imbalance.resampling.tests.test_soup.soup_mock()
```

```
multi_imbalance.resampling.tests.test_soup.test_calculating_safe_levels_for_class(X,
                                                                              y,
                                                                              zero_safe_le
                                                                              one_safe_lev
                                                                              first_sample
                                                                              soup_mock)
```

```
multi_imbalance.resampling.tests.test_soup.test_calculating_safe_levels_for_sample(X,
                                                                              y,
                                                                              zero_safe_
                                                                              one_safe_l
                                                                              first_sampl
                                                                              soup_mock)
```

```
multi_imbalance.resampling.tests.test_soup.test_oversample(X, y, class_name, ex-
                                                                pected_undersampling,
                                                                ex-
                                                                pected_oversampling,
                                                                soup_mock)
```

```
multi_imbalance.resampling.tests.test_soup.test_undersample(X, y, class_name, ex-
                                                                pected_undersampling,
                                                                ex-
                                                                pected_oversampling,
                                                                soup_mock)
```

### multi\_imbalance.resampling.tests.test\_spider module

```
multi_imbalance.resampling.tests.test_spider.test_estimate_cost_matrix()
multi_imbalance.resampling.tests.test_spider.test_fit_resample()
multi_imbalance.resampling.tests.test_spider.test_intersect()
multi_imbalance.resampling.tests.test_spider.test_knn()
multi_imbalance.resampling.tests.test_spider.test_min_cost_classes()
multi_imbalance.resampling.tests.test_spider.test_setdiff()
multi_imbalance.resampling.tests.test_spider.test_union()
```

### multi\_imbalance.resampling.tests.test\_static\_smote module

```
multi_imbalance.resampling.tests.test_static_smote.test_static_smote()
```

## Module contents

### Submodules

#### multi\_imbalance.resampling.global\_cs module

```
class multi_imbalance.resampling.global_cs.GlobalCS(shuffle: bool = True)
    Bases: imblearn.base.BaseSampler
```

Global CS is an algorithm that equalizes number of samples in each class. It duplicates all samples equally for each class to achieve majority class size

#### multi\_imbalance.resampling.mdo module

```
class multi_imbalance.resampling.mdo.MDO(k=5, k1_frac=0.4, seed=0, prop=1,
                                         maj_int_min=None)
    Bases: imblearn.base.BaseSampler
```

Mahalanbois Distance Oversampling is an algorithm that oversamples all classes to a quantity of the major class. Samples for oversampling are chosen based on their k neighbours and new samples are created in random place but with the same Mahalanbois distance from the centre of class to chosen sample.

#### Parameters

- **k** – Number of neighbours considered during the neighbourhood analysis
- **k1\_frac** – Ratio of the number of neighbours in the sample class to all neighbours in the neighbourhood. If the ratio is greater, the example will not be considered noise
- **seed** –
- **prop** – Oversampling ratio, if equal to one the class size after resampling will be equal to the size of the largest class
- **maj\_int\_min** – dict {'maj': majority class labels, 'min': minority class labels}

```
calculate_same_class_neighbour_quantities(S_minor, S_minor_label)
```

## multi\_imbalance.resampling.soup module

**class** multi\_imbalance.resampling.soup.**SOUP** (*k: int = 7, shuffle=False, maj\_int\_min=None*)  
 Bases: imblearn.base.BaseSampler

Similarity Oversampling and Undersampling Preprocessing (SOUP) is an algorithm that equalizes number of samples in each class. It also takes care of the similarity between classes, which means that it removes samples from majority class, that are close to samples from the other class and duplicate samples from the minority classes, which are in the safest area in space

### Parameters

- **k** – number of neighbors
- **shuffle** – bool - output will be shuffled
- **maj\_int\_min** – dict {‘maj’: majority class labels, ‘min’: minority class labels}

## multi\_imbalance.resampling.spider module

**class** multi\_imbalance.resampling.spider.**SPIDER3** (*k, maj\_int\_min=None, cost=None*)  
 Bases: imblearn.base.BaseSampler

SPIDER3 algorithm implementation for selective preprocessing of multi-class imbalanced data sets.

Reference: Wojciechowski, S., Wilk, S., Stefanowski, J.: An Algorithm for Selective Preprocessing of Multi-class Imbalanced Data. Proceedings of the 10th International Conference on Computer Recognition Systems CORES 2017

### Parameters

- **k** – Number of nearest neighbors considered while resampling.
- **maj\_int\_min** – Dict that contains lists of majority, intermediate and minority classes labels.
- **cost** – The cost matrix. An element  $c[i, j]$  of this matrix represents the cost associated with misclassifying an example from class  $i$  as class one from class  $j$ .

**amplify** (*int\_min\_class*)

**clean** (*int\_min\_class*)

**relabel** (*int\_min\_class*)

## multi\_imbalance.resampling.static\_smote module

**class** multi\_imbalance.resampling.static\_smote.**StaticSMOTE**  
 Bases: imblearn.base.BaseSampler

Static SMOTE implementation:

Reference: Fernández-Navarro, F., Hervás-Martínez, C., Gutiérrez, P.A.: A dynamic over-sampling procedure based on sensitivity for multi-class problems. Pattern Recognit. 44, 1821–1833 (2011)

## Module contents

### multi\_imbalance.utils package

#### Submodules

#### multi\_imbalance.utils.array\_util module

`multi_imbalance.utils.array_util.contains` (*dataset, example*)

Returns if dataset contains the example. :param dataset: :param example: :return: True or False depending on whether dataset contains the example.

`multi_imbalance.utils.array_util.index_of` (*arr, example*)

**Returns** Index of learning example in arr.

`multi_imbalance.utils.array_util.intersect` (*arr1, arr2*)

Performs the intersection operation over two numpy arrays (not removing duplicates).

#### Parameters

- **arr1** – Numpy array number 1.
- **arr2** – Numpy array number 2.

**Returns** The intersection of arr1 and arr2.

`multi_imbalance.utils.array_util.setdiff` (*arr1, arr2*)

Performs the difference over two numpy arrays.

#### Parameters

- **arr1** – Numpy array number 1.
- **arr2** – Numpy array number 2.

**Returns** Result of the difference of arr1 and arr2.

`multi_imbalance.utils.array_util.union` (*arr1, arr2*)

Performs the union over two numpy arrays (not removing duplicates, as it's how the algorithm SPIDER3 actually works).

#### Parameters

- **arr1** – Numpy array number 1.
- **arr2** – Numpy array number 2.

**Returns** The union of arr1 and arr2.

#### multi\_imbalance.utils.data module

`multi_imbalance.utils.data.construct_flat_2pc_df` (*X, y*) → `pan-`  
`das.core.frame.DataFrame`

This function takes two dimensional X and one dimensional y arrays, concatenates and returns them as data frame

#### Parameters

- **X** – two dimensional numpy array
- **y** – one dimensional numpy array with labels

**Returns** Data frame with 3 columns x1 x2 and y and with number of rows equal to number of rows in X

```
multi_imbalance.utils.data.construct_maj_int_min(y: numpy.ndarray, strategy='median') → collections.OrderedDict
```

This function creates dictionary with information which classes are minority or majority

#### Parameters

- **y** – One dimensional numpy array that contains class labels
- **strategy** – The principle according to which the division into minority and majority classes will be determined:
  - **'median'**: A class whose size is equal to the median of the class sizes will be considered “intermediate”
  - **'average'**: The average class size will be calculated, all classes that are smaller will be considered as minority and the rest will be considered majority

**Returns** dictionary with keys ‘maj’, ‘int’, ‘min’. The value for each key is a list containing the class labels belonging to the given group

```
multi_imbalance.utils.data.get_project_root() → pathlib.Path
```

Returns project root folder.

```
multi_imbalance.utils.data.load_arff_dataset(path: str, one_hot_encode: bool = True, return_non_cat_length: bool = False)
```

Load and return the dataset saved in arff type file

#### Parameters

- **path** (*str*) – location of dataset file
- **one\_hot\_encode** (*bool*) – flag, if true encodes categorical variables using OneHotEncoder
- **return\_non\_cat\_length** (*bool*) – flag, if true returns the number of non categorical variables

#### Returns

- ndarray X - dimensional numpy array where non categorical variables are stored in first columns followed by categorical variables
- ndarray y - one dimensional numpy array with the classification target
- bool non\_cat\_length - number of non categorical variables (only if return\_non\_cat\_length=True)

```
multi_imbalance.utils.data.load_datasets_arff(return_non_cat_length=False, dataset_paths=None)
```

### multi\_imbalance.utils.metrics module

multi\_imbalance.utils.metrics.**gmean\_score**(*y\_test*, *y\_pred*, *correction*: float = 0.001) → float

Calculate geometric mean score

#### Parameters

- **y\_test** – numpy array with labels
- **y\_pred** – numpy array with predicted labels
- **correction** – value that replaces 0 during multiplication to avoid zeroing the result

**Returns** geometric\_mean\_score: float

### multi\_imbalance.utils.min\_int\_maj module

### multi\_imbalance.utils.plot module

multi\_imbalance.utils.plot.**plot\_cardinality\_and\_2d\_data**(*X*, *y*, *dataset\_name*="") → None

Plots cardinality of classes from *y* as well as scatter plot of *X* transformed to two dimensions using PCA

#### Parameters

- **X** (*ndarray*) – two dimensional numpy array
- **y** (*ndarray*) – one dimensional numpy array
- **dataset\_name** (*str*) – title of chart

multi\_imbalance.utils.plot.**plot\_visual\_comparision\_datasets**(*X1*, *y1*, *X2*, *y2*, *dataset\_name1*="", *dataset\_name2*="") → None

Plots comparision of *X1* *y1* and *X2* *y2* using `plot_cardinality_and_2d_data`, which plots cardinality of classes from *y* as well as scatter plot of *X* transformed to two dimensions using PCA

#### Parameters

- **X1** (*ndarray*) – two dimensional numpy array with data from dataset1
- **y1** (*ndarray*) – one dimensional numpy array with target classes from dataset1
- **X2** (*ndarray*) – two dimensional numpy array with data from dataset2
- **y2** (*ndarray*) – one dimensional numpy array with target classes from dataset1
- **dataset\_name1** (*str*) – first dataset chart title
- **dataset\_name2** (*str*) – second dataset chart title

## Module contents

### 2.1.1.2 Module contents

## 2.2 License

The MIT License (MIT)

Copyright (c) 2019 Damian Horna, Kamil Pluciński, Hanna Klimczak, Jacek Grycza

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the “Software”), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED “AS IS”, WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

## 2.3 Contact

Question? Please contact [test@gmail.com](mailto:test@gmail.com)





## PYTHON MODULE INDEX

### m

- multi\_imbalance, 19
- multi\_imbalance.datasets, 6
- multi\_imbalance.datasets.tests, 6
- multi\_imbalance.datasets.tests.test\_data\_loader, 5
- multi\_imbalance.ensemble, 12
- multi\_imbalance.ensemble.ecoc, 8
- multi\_imbalance.ensemble.mrbbagging, 9
- multi\_imbalance.ensemble.ovo, 10
- multi\_imbalance.ensemble.soup\_bagging, 11
- multi\_imbalance.ensemble.tests, 8
- multi\_imbalance.ensemble.tests.test\_ecoc, 6
- multi\_imbalance.ensemble.tests.test\_mrbbagging, 7
- multi\_imbalance.ensemble.tests.test\_ovo, 7
- multi\_imbalance.ensemble.tests.test\_soupbagging, 8
- multi\_imbalance.resampling, 16
- multi\_imbalance.resampling.global\_cs, 14
- multi\_imbalance.resampling.mdo, 14
- multi\_imbalance.resampling.soup, 15
- multi\_imbalance.resampling.spider, 15
- multi\_imbalance.resampling.static\_smote, 15
- multi\_imbalance.resampling.tests, 14
- multi\_imbalance.resampling.tests.test\_globalcs, 12
- multi\_imbalance.resampling.tests.test\_mdo, 13
- multi\_imbalance.resampling.tests.test\_soup, 13
- multi\_imbalance.resampling.tests.test\_spider, 14
- multi\_imbalance.resampling.tests.test\_static\_smote, 14
- multi\_imbalance.utils, 19
- multi\_imbalance.utils.array\_util, 16
- multi\_imbalance.utils.data, 16
- multi\_imbalance.utils.metrics, 18
- multi\_imbalance.utils.min\_int\_maj, 18
- multi\_imbalance.utils.plot, 18



## INDEX

### A

`amplify()` (*multi\_imbalance.resampling.spider.SPIDER3* method), 15

### C

`calc_duplicates_quantities()` (*in module multi\_imbalance.resampling.tests.test\_globalcs*), 12

`calculate_same_class_neighbour_quantities()` (*multi\_imbalance.resampling.mdo.MDO* method), 14

`clean()` (*multi\_imbalance.resampling.spider.SPIDER3* method), 15

`construct_flat_2pc_df()` (*in module multi\_imbalance.utils.data*), 16

`construct_maj_int_min()` (*in module multi\_imbalance.utils.data*), 17

`contains()` (*in module multi\_imbalance.utils.array\_util*), 16

### E

`ECOC` (*class in multi\_imbalance.ensemble.ecoc*), 8

### F

`fit()` (*multi\_imbalance.ensemble.ecoc.ECOC* method), 9

`fit()` (*multi\_imbalance.ensemble.mrbbagging.MRBBagging* method), 10

`fit()` (*multi\_imbalance.ensemble.ovo.OVO* method), 11

`fit()` (*multi\_imbalance.ensemble.soup\_bagging.SOUPBagging* method), 11

`fit_classifier()` (*multi\_imbalance.ensemble.soup\_bagging.SOUPBagging* static method), 12

`fit_clf()` (*in module multi\_imbalance.ensemble.soup\_bagging*), 12

### G

`get_goal_quantity()` (*in module multi\_imbalance.resampling.tests.test\_globalcs*), 12

`get_project_root()` (*in module multi\_imbalance.utils.data*), 17

`global_cs_mock()` (*in module multi\_imbalance.resampling.tests.test\_globalcs*), 12

`GlobalCS` (*class in multi\_imbalance.resampling.global\_cs*), 14

`gmean_score()` (*in module multi\_imbalance.utils.metrics*), 18

### I

`index_of()` (*in module multi\_imbalance.utils.array\_util*), 16

`intersect()` (*in module multi\_imbalance.utils.array\_util*), 16

### L

`load_arff_dataset()` (*in module multi\_imbalance.utils.data*), 17

`load_datasets()` (*in module multi\_imbalance.datasets*), 6

`load_datasets_arff()` (*in module multi\_imbalance.utils.data*), 17

### M

`MDO` (*class in multi\_imbalance.resampling.mdo*), 14

`mdo_mock()` (*in module multi\_imbalance.resampling.tests.test\_mdo*), 13

module

`multi_imbalance`, 19

`multi_imbalance.datasets`, 6

`multi_imbalance.datasets.tests`, 6

`multi_imbalance.datasets.tests.test_data_loader`, 5

`multi_imbalance.ensemble`, 12

`multi_imbalance.ensemble.ecoc`, 8

`multi_imbalance.ensemble.mrbbagging`, 9

`multi_imbalance.ensemble.ovo`, 10

`multi_imbalance.ensemble.soup_bagging`, 11



P

plot\_cardinality\_and\_2d\_data() (in module *multi\_imbalance.utils.plot*), 18  
 plot\_visual\_comparision\_datasets() (in module *multi\_imbalance.utils.plot*), 18  
 predict() (*multi\_imbalance.ensemble.ecoc.ECOC* method), 9  
 predict() (*multi\_imbalance.ensemble.mrbbagging.MRBBagging* method), 10  
 predict() (*multi\_imbalance.ensemble.ovo.OVO* method), 11  
 predict() (*multi\_imbalance.ensemble.soup\_bagging.SOUPBagging* method), 12  
 predict\_proba() (*multi\_imbalance.ensemble.soup\_bagging.SOUPBagging* method), 12

R

relabel() (*multi\_imbalance.resampling.spider.SPIDER3* method), 15

S

setdiff() (in module *multi\_imbalance.utils.array\_util*), 16  
 should\_perform\_oversampling() (*multi\_imbalance.ensemble.ovo.OVO* method), 11  
 SOUP (class in *multi\_imbalance.resampling.soup*), 15  
 soup\_mock() (in module *multi\_imbalance.resampling.tests.test\_soup*), 13  
 SOUPBagging (class in *multi\_imbalance.ensemble.soup\_bagging*), 11  
 SPIDER3 (class in *multi\_imbalance.resampling.spider*), 15  
 StaticSMOTE (class in *multi\_imbalance.resampling.static\_smote*), 15

T

test\_\_group\_data() (*multi\_imbalance.ensemble.tests.test\_mrbbagging.TestMRBBagging* method), 7  
 test\_\_group\_data\_with\_none() (*multi\_imbalance.ensemble.tests.test\_mrbbagging.TestMRBBagging* method), 7  
 test\_api() (*multi\_imbalance.ensemble.tests.test\_mrbbagging.TestMRBBagging* method), 7  
 test\_api\_multiple\_trees() (*multi\_imbalance.ensemble.tests.test\_mrbbagging.TestMRBBagging* method), 7  
 test\_api\_with\_feature\_selection() (*multi\_imbalance.ensemble.tests.test\_mrbbagging.TestMRBBagging* method), 7

test\_api\_with\_feature\_selection\_sqrt\_features() (*multi\_imbalance.ensemble.tests.test\_mrbbagging.TestMRBBagging* method), 7  
 test\_api\_with\_random\_feature\_selection() (*multi\_imbalance.ensemble.tests.test\_mrbbagging.TestMRBBagging* method), 7  
 test\_binary\_classifiers() (in module *multi\_imbalance.ensemble.tests.test\_ovo*), 7  
 test\_calculating\_safe\_levels\_for\_class() (in module *multi\_imbalance.resampling.tests.test\_soup*), 13  
 test\_calculating\_safe\_levels\_for\_sample() (in module *multi\_imbalance.resampling.tests.test\_soup*), 13  
 test\_choose\_samples() (in module *multi\_imbalance.resampling.tests.test\_mdo*), 13  
 test\_choose\_samples\_when\_correct() (in module *multi\_imbalance.resampling.tests.test\_mdo*), 13  
 test\_choose\_samples\_when\_zero\_samples\_expected() (in module *multi\_imbalance.resampling.tests.test\_mdo*), 13  
 test\_default\_classifier() (in module *multi\_imbalance.ensemble.tests.test\_soupbagging*), 8  
 test\_dense\_and\_sparse\_with\_not\_cached\_matrices() (in module *multi\_imbalance.ensemble.tests.test\_ecoc*), 6  
 test\_ecoc\_with\_sklearn\_pipeline() (in module *multi\_imbalance.ensemble.tests.test\_ecoc*), 6  
 test\_ecoc\_with\_sklearn\_pipeline() (in module *multi\_imbalance.ensemble.tests.test\_ovo*), 7  
 test\_encoding() (in module *multi\_imbalance.ensemble.tests.test\_ecoc*), 6  
 test\_estimate\_cost\_matrix() (in module *multi\_imbalance.resampling.tests.test\_spider*), 14  
 test\_exception() (in module *multi\_imbalance.ensemble.tests.test\_soupbagging*), 8  
 test\_fit\_classifier\_classifier() (in module *multi\_imbalance.ensemble.tests.test\_soupbagging*), 8  
 test\_fit\_predict() (in module *multi\_imbalance.ensemble.tests.test\_ovo*), 7  
 test\_fit\_resample() (in module *multi\_imbalance.resampling.tests.test\_spider*), 14

```

test_fit_with_invalid_classifier()      test_random_oversampling() (in module
    (multi_imbalance.ensemble.tests.test_mrbbagging.TestMRBBAgging), 6
    method), 7
test_fit_with_invalid_labels()          test_setdiff() (in module
    (multi_imbalance.ensemble.tests.test_mrbbagging.TestMRBBAgging
    method), 7
    multi_imbalance.resampling.tests.test_spider),
test_hamming_distance() (in module
    multi_imbalance.ensemble.tests.test_ecoc),
    6
test_intersect() (in module
    multi_imbalance.resampling.tests.test_spider),
    14
test_knn() (in module
    multi_imbalance.resampling.tests.test_spider),
    14
test_load_datasets() (in module
    multi_imbalance.datasets.tests.test_data_loader),
    5
test_max_voting() (in module
    multi_imbalance.ensemble.tests.test_ovo),
    7
test_mdo_api() (in module
    multi_imbalance.resampling.tests.test_mdo),
    13
test_min_cost_classes() (in module
    multi_imbalance.resampling.tests.test_spider),
    14
test_no_oversampling() (in module
    multi_imbalance.ensemble.tests.test_ecoc),
    6
test_output_equal_replication() (in mod- test_with_invalid_k()
    (in module multi_imbalance.resampling.tests.test_globalcs),
    12
    multi_imbalance.ensemble.tests.test_mrbbagging.TestMRBBAgging
    method), 7
test_output_length_validate() (in module test_with_own_classifier() (in module
    multi_imbalance.resampling.tests.test_globalcs),
    12
    multi_imbalance.ensemble.tests.test_ecoc), 6
test_oversample() (in module
    multi_imbalance.resampling.tests.test_soup),
    13
    multi_imbalance.ensemble.tests.test_ovo), 7
test_own_classifier_without_predict_and_fit() TestMRBBAgging (class in
    (in module multi_imbalance.ensemble.tests.test_ecoc),
    6
    multi_imbalance.ensemble.tests.test_mrbbagging),
    7
test_own_preprocessing_without_fit_resample() (in module multi_imbalance.ensemble.tests.test_ovo),
    7
test_own_preprocessing_without_fit_transform()
    (in module multi_imbalance.ensemble.tests.test_ecoc),
    6
test_predefined_classifiers_and_preprocessings_without_errors()
    (in module multi_imbalance.ensemble.tests.test_ovo),
    7
test_predefined_classifiers_and_weighting_without_exceptions()
    (in module multi_imbalance.ensemble.tests.test_ecoc),
    6
test_soubagging() (in module
    multi_imbalance.ensemble.tests.test_soupbagging),
    8
test_static_smote() (in module
    multi_imbalance.resampling.tests.test_static_smote),
    14
test_undersample() (in module
    multi_imbalance.resampling.tests.test_soup),
    13
test_union() (in module
    multi_imbalance.resampling.tests.test_spider),
    14
test_unknown_classifier() (in module
    multi_imbalance.ensemble.tests.test_ecoc),
    6
test_unknown_preprocessing() (in module
    multi_imbalance.ensemble.tests.test_ecoc), 6
test_unknown_preprocessing() (in module
    multi_imbalance.ensemble.tests.test_ovo), 7
test_unknown_preprocessing_between_strategy_raises_
    (in module multi_imbalance.ensemble.tests.test_ovo),
    7
test_zero_variance() (in module
    multi_imbalance.resampling.tests.test_mdo),
    13
union() (in module multi_imbalance.utils.array_util),
    16

```