



HAL
open science

metric-learn: Metric Learning Algorithms in Python

William de Vazelhes, Cj Carey, Yuan Tang, Nathalie Vauquier, Aurélien Bellet

► **To cite this version:**

William de Vazelhes, Cj Carey, Yuan Tang, Nathalie Vauquier, Aurélien Bellet. metric-learn: Metric Learning Algorithms in Python. 2019. hal-02376986

HAL Id: hal-02376986

<https://inria.hal.science/hal-02376986v1>

Preprint submitted on 22 Nov 2019

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.

metric-learn: Metric Learning Algorithms in Python

William de Vazelhes

INRIA, France

WILLIAM.DE-VAZELHES@INRIA.FR

CJ Carey

Google LLC, United States

PERIMOSOCORDIAE@GMAIL.COM

Yuan Tang

Ant Financial, United States

TERRYTANGYUAN@GMAIL.COM

Nathalie Vauquier

INRIA, France

NATHALIE.VAUQUIER@INRIA.FR

Aurélien Bellet

INRIA, France

AURELIEN.BELLET@INRIA.FR

Editor:

Abstract

`metric-learn` is an open source Python package implementing supervised and weakly-supervised distance metric learning algorithms. As part of `scikit-learn-contrib`, it provides a unified interface compatible with `scikit-learn` which allows to easily perform cross-validation, model selection, and pipelining with other machine learning estimators. `metric-learn` is thoroughly tested and available on PyPi under the MIT licence.

Keywords: Machine Learning, Python, Metric Learning, Scikit-learn

1. Introduction

Many approaches in machine learning require a measure of distance between data points. Traditionally, practitioners would choose a standard distance metric (Euclidean, City-Block, Cosine, etc.) using a priori knowledge of the domain. However, it is often difficult to design metrics that are well-suited to the particular data and task of interest. Distance metric learning, or simply metric learning (Bellet et al., 2015), aims at automatically constructing task-specific distance metrics from data. A key advantage of metric learning is that it can be applied beyond the standard supervised learning setting (data points associated with labels), in situations where only weaker forms of supervision are available (e.g., pairs of points that should be similar/dissimilar). The learned distance metric can be used to perform retrieval tasks such as finding elements (images, documents) of a database that are semantically closest to a query element. It can also be plugged into other machine learning algorithms, for instance to improve the accuracy of nearest neighbors models (for classification, regression, anomaly detection...) or to bias the clusters found by clustering algorithms towards the intended semantics. Finally, metric learning can be used to perform dimensionality reduction. These use-cases highlight the importance of integrating metric learning with the rest of the machine learning pipeline and tools.



Figure 1: Different types of supervision for metric learning illustrated on face image data taken from the Labeled Faces in the Wild dataset (Huang et al., 2012).

`metric-learn` is an open source package for metric learning in Python, which implements many popular metric-learning algorithms with different levels of supervision through a unified interface. Its API is compatible with `scikit-learn` (Pedregosa et al., 2011), a prominent machine learning library in Python. This allows for streamlined model selection, evaluation, and pipelining with other estimators.

2. Background on Metric Learning

Metric learning is generally formulated as an optimization problem where one seeks to find the parameters of a distance function that minimize some objective function over the input data. All algorithms currently implemented in `metric-learn` learn so-called Mahalanobis distances. Given a real-valued parameter matrix L of shape $(\mathbf{n_components}, \mathbf{n_features})$ where $\mathbf{n_features}$ is the number of features describing the data, the associated Mahalanobis distance between two points x and x' is defined as $D_L(x, x') = \sqrt{(Lx - Lx')^\top (Lx - Lx')}$. This is equivalent to Euclidean distance after linear transformation of the feature space defined by L . Thus, if L is the identity matrix, standard Euclidean distance is recovered. Mahalanobis distance metric learning can thus be seen as learning a new embedding space, with potentially reduced dimension $\mathbf{n_components}$. Note that D_L can also be written as $D_L(x, x') = \sqrt{(x - x')^\top M (x - x')}$, where we refer to $M = L^\top L$ as the Mahalanobis matrix.

Metric learning algorithms can be categorized according to the form of data supervision they require to learn a metric. `metric-learn` currently implements algorithms that fall into the following categories. *Supervised learners* learn from a dataset with one label per training example, aiming to bring together points from the same class while spreading points from different classes. For instance, data points could be face images and the class could be the identity of the person (see Figure 1a). *Pair learners* require a set of pairs of points, with each pair labeled to indicate whether the two points are similar or not. These methods aim to learn a metric that brings pairs of similar points closer together and pushes pairs of dissimilar points further away from each other. Such supervision is often simpler to collect than class labels in applications when there are many labels. For instance, a human annotator can often quickly decide whether two face images correspond to the same person (Figure 1b) while matching a face to its identity among many possible people may be difficult. Finally, *quadruplet learners* consider 4-tuples of points and aim to learn a metric that brings the two first points of each quadruplet closer than the two last points. This can be used to learn a metric space in which closer points are more similar with respect to an attribute of interest, which may be continuous and difficult to annotate accurately

(e.g., the age of a person on an image, see Figure 1c). Quadruplet supervision is also used in problems with a class hierarchy.

3. Overview of the Package

The current release of `metric-learn` (v.0.5.0) can be installed from the Python Package Index (PyPI), for Python 2.7 and 3.5 or later. The source code is available on GitHub at <http://github.com/scikit-learn-contrib/metric-learn> and is free to use, provided under the MIT license. `metric-learn` depends on core libraries from the SciPy ecosystem: `numpy`, `scipy`, and `scikit-learn`. Detailed documentation (including installation guidelines, the description of the algorithms and the API, as well as examples) is available at <http://contrib.scikit-learn.org/metric-learn>. The development is collaborative and open to all contributors through the usual GitHub workflow of issues and pull requests. Community interest for the package has been demonstrated by its recent inclusion in the `scikit-learn-contrib` organization which hosts high-quality `scikit-learn-compatible` projects,¹ and by its more than 740 stars and 170 forks on GitHub at the time of writing. The quality of the code is ensured by a thorough test coverage (96% as of July 2019). Every new contribution is automatically checked by a continuous integration platform to enforce sufficient test coverage as well as syntax formatting with `flake8`.

Currently, `metric-learn` implements 9 popular metric learning algorithms. Supervised learners include Neighborhood Components Analysis (NCA, Goldberger et al., 2004), Large Margin Nearest Neighbors (LMNN, Weinberger and Saul, 2009), Relative Components Analysis (RCA, Shental et al., 2002),² Local Fisher Discriminant Analysis (LFDA, Sugiyama, 2007) and Metric Learning for Kernel Regression (MLKR, Weinberger and Tesauro, 2007). The latter is designed for regression problems with continuous labels. Pair learners include Mahalanobis Metric for Clustering (MMC, Xing et al., 2002), Information Theoretic Metric Learning (ITML, Davis et al., 2007) and Sparse High-Dimensional Metric Learning (SDML, Qi et al., 2009). The package implements one quadruplet learner: Metric Learning from Relative Comparisons by Minimizing Squared Residual (LSML, Liu et al., 2012). Detailed descriptions of these algorithms can be found in the package documentation.

4. Software Architecture and API

`metric-learn` provides a unified interface to all metric learning algorithms. It is designed to be fully compatible with the functionality of `scikit-learn`. All metric learners inherit from an abstract `BaseMetricLearner` class, which itself inherits from `scikit-learn`'s `BaseEstimator`. All classes inheriting from `BaseMetricLearner` should implement two methods: `get_metric` (returning a function that computes the distance, which can be plugged into `scikit-learn` estimators like `KMeansClustering`) and `score_pairs` (returning the distances between a set of pairs of points passed as a 3D array). Mahalanobis distance learning algorithms also inherit from a `MahalanobisMixin` interface, which has an attribute `components_` corresponding to the transformation matrix L of the Mahalanobis distance. `MahalanobisMixin` implements `get_metric` and `score_pairs` accordingly as well

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

2. RCA takes as input slightly weaker supervision in the form of *chunklets* (groups of points of same class).

as a few additional methods. In particular, `transform` allows to transform data using `components_`, and `get_mahalanobis_matrix` returns the Mahalanobis matrix $M = L^T L$.

Supervised metric learners inherit from `scikit-learn`'s base class `TransformerMixin`, the same base class used by `sklearn.LinearDiscriminantAnalysis` and others. As such, they are compatible for pipelining with other estimators via `sklearn.pipeline.Pipeline`.

Weakly supervised algorithms (pair and quadruplet learners) `fit` and `predict` on a set of tuples passed as a 3-dimensional array. Tuples can be pairs or quadruplets depending on the algorithm. Pair learners take as input an array-like `pairs` of shape `(n_pairs, 2, n_features)`, as well as an array-like `y_pairs` of shape `(n_pairs,)` giving labels (similar or dissimilar) for each pair. In order to `predict` the labels of new pairs, one needs to set a threshold on the distance value. This threshold can be set manually or automatically calibrated (at fit time or afterwards on a validation set) to optimize a given score such as accuracy or F1-score using the method `calibrate_threshold`. Quadruplet learners work on array-like of shape `(n_quadruplets, 4, n_features)`, where for each quadruplet the two first elements are the ones we want to be closer than the two last ones. They can naturally `predict` whether a new quadruplet is in the right order by comparing the two pairwise distances. These design choices enable use of `scikit-learn`'s scoring functions out of the box, as well as the standard routines for model selection, including `GridSearchCV`. To illustrate, the following code snippet computes cross validation scores for ITML (with default parameters) on pairs from Labeled Faces in the Wild (Huang et al., 2012).

```
>>> from sklearn.datasets import fetch_lfw_pairs
>>> from sklearn.model_selection import cross_validate, train_test_split
>>> from sklearn.decomposition import PCA
>>> from metric_learn import ITML
>>> pairs, y_pairs = [fetch_lfw_pairs()[key] for key in ['pairs', 'target']]
>>> pairs = PCA(n_components=25).fit_transform(pairs.reshape(4400, -1)).reshape(-1, 2, 25)
>>> pairs, _, y_pairs, _ = train_test_split(pairs, 2*y_pairs-1)
>>> cross_validate(ITML(), pairs, y_pairs, scoring='roc_auc', return_train_score=True)
```

5. Future Work

`metric-learn` is under active development. We list here some promising directions to further improve the package. To scale to large datasets, we would like to implement stochastic solvers (SGD and its variants), forming batches of tuples on the fly to avoid loading all data in memory at once. We also plan to incorporate new algorithms that provide added value to the package, in particular some that learn from triplet supervision (Schultz and Joachims, 2003), can deal with multi-label (Liu and Tsang, 2015) and high-dimensional problems (Liu and Bellet, 2019), or learn other forms of metrics (e.g., nonlinear ones, bilinear similarities, and multiple local metrics, see Bellet et al., 2015).

Acknowledgments

We are thankful to Inria for funding 2 years of development. We also thank `scikit-learn` developers from the Inria Parietal team (in particular Gaël Varoquaux, Alexandre Gramfort and Olivier Grisel) for fruitful discussions on the design of the API and funding to attend SciPy 2019, as well as `scikit-learn-contrib` reviewers for their valuable feedback.

References

- Aurélien Bellet, Amaury Habrard, and Marc Sebban. *Metric Learning*. Morgan & Claypool Publishers, 2015.
- Jason V. Davis, Brian Kulis, Prateek Jain, Suvrit Sra, and Inderjit S. Dhillon. Information-Theoretic Metric Learning. In *ICML*, 2007.
- Jacob Goldberger, Sam Roweis, Geoff Hinton, and Ruslan Salakhutdinov. Neighbourhood Components Analysis. In *NIPS*, 2004.
- Gary B. Huang, Marwan Mattar, Honglak Lee, and Erik Learned-Miller. Learning to Align from Scratch. In *NIPS*, 2012.
- E. Y. Liu, Z. Guo, X. Zhang, V. Jojic, and W. Wang. Metric Learning from Relative Comparisons by Minimizing Squared Residual. In *ICDM*, 2012.
- Kuan Liu and Aurélien Bellet. Escaping the Curse of Dimensionality in Similarity Learning: Efficient Frank-Wolfe Algorithm and Generalization Bounds. *Neurocomputing*, 333:185–199, 2019.
- Weiwei Liu and Ivor W. Tsang. Large Margin Metric Learning for Multi-Label Prediction. In *AAAI*, 2015.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Guo-Jun Qi, Jinhui Tang, Zheng-Jun Zha, Tat-Seng Chua, and Hong-Jiang Zhang. An Efficient Sparse Metric Learning in High-dimensional Space via L1-penalized Log-determinant Regularization. In *ICML*, 2009.
- Matthew Schultz and Thorsten Joachims. Learning a Distance Metric from Relative Comparisons. In *NIPS*, 2003.
- Noam Shental, Tomer Hertz, Daphna Weinshall, and Misha Pavel. Adjustment Learning and Relevant Component Analysis. In *ECCV*, 2002.
- Masashi Sugiyama. Dimensionality Reduction of Multimodal Labeled Data by Local Fisher Discriminant Analysis. *Journal of Machine Learning Research*, 8:1027–1061, 2007.
- Kilian Q. Weinberger and Lawrence K. Saul. Distance Metric Learning for Large Margin Nearest Neighbor Classification. *Journal of Machine Learning Research*, 10:207–244, 2009.
- Kilian Q. Weinberger and Gerald Tesauro. Metric Learning for Kernel Regression. In *AISTATS*, 2007.
- Eric P. Xing, Andrew Y. Ng, Michael I. Jordan, and Stuart J. Russell. Distance Metric Learning with Application to Clustering with Side-Information. In *NIPS*, 2002.