

# A Tutorial on Distance Metric Learning: Mathematical Foundations, Algorithms and Experiments

JUAN LUIS SUÁREZ, University of Granada  
 SALVADOR GARCÍA, University of Granada  
 FRANCISCO HERRERA, University of Granada

Distance metric learning is a branch of machine learning that aims to learn distances from the data. Distance metric learning can be useful to improve similarity learning algorithms, and also has applications in dimensionality reduction. This paper describes the distance metric learning problem and analyzes its main mathematical foundations. In addition, it also discusses some of the most popular distance metric learning techniques used in classification, showing their goals and the required information to understand and use them. Furthermore, some experiments to evaluate the performance of the different algorithms are also provided. Finally, this paper discusses several possibilities of future work in this topic.

Categories and Subject Descriptors: I.2.6 [**Artificial Intelligence**]: Learning; H.2.8 [**Database management**]: Database Applications—*Data mining*; G.1.3 [**Numerical Analysis**]: Numerical Linear Algebra; G.1.6 [**Numerical Analysis**]: Optimization

General Terms: Algorithms, Experimentation, Performance, Theory

Additional Key Words and Phrases: Distance Metric Learning, Classification, Mahalanobis Distance, Dimensionality, Similarity

## 1. INTRODUCTION

The use of distances in machine learning has been present since its inception. Distances provide a similarity measure between the data, so that close data will be considered similar, while remote data will be considered dissimilar. One of the most popular examples of this similarity learning is the well-known nearest neighbors rule for classification, where a new sample is labeled with the majority class within its nearest neighbors in the training set. This classifier was presented by Cover and Hart [1967], even though this idea had already been mentioned in earlier publications [Sebestyen 1962; Nilsson 1965].

Algorithms in the style of the nearest neighbors classifier are among the main motivators of distance metric learning. These kind of algorithms have usually used a standard distance, like the euclidean distance, to measure the data similarity. However, a standard distance may ignore some important properties available in our dataset, so that the learning results could be non optimal. The search for a distance that brings similar data as close as possible, while moving non similar data away, can significantly increase the quality of these algorithms.

Distance metric learning has applications beyond the improvement of such algorithms. We will see, for example, that learning distances is closely related to learning linear mappings, which in turn is closely related with dimensionality reduction [Cunningham and Ghahramani 2015].

Although techniques such as principal component analysis or linear discriminant analysis, which are considered distance metric learning techniques, have been popular in the statistical field since the middle of the 20th century, it is not until the beginning of the 21st century that distance metric learning is properly spoken of, and perhaps the

---

Our work has been supported by the research project TIN2017-89517-P and by a research scholarship (FPU18/05989), given to the author Juan Luis Suárez by the Spanish Ministry of Science, Innovation and Universities.

Author's addresses: J. L. Suárez, S. García and F. Herrera are with the Department of Computer Science and Artificial Intelligence, University of Granada, Spain.

algorithm from Xing et al. [2003] is responsible for drawing attention to this concept for the first time.

During the first decade of the 21st century some of the most well-known distance metric learning algorithms were developed, and are still popular today. The most relevant of these algorithms will be studied throughout this tutorial. Over recent years distance metric learning remains active, both in the search for new proposals for innovative distance metric learning algorithms, and in the refinement of techniques already employed over the past decade. Some of these techniques will also be shown. Nowadays, distance metric learning is still used in many real applications [Nguyen and De Baets 2019; Liang et al. 2019].

In this paper we make a theoretical study of supervised distance metric learning, in which we show the mathematical foundations of distance metric learning and its algorithms. Furthermore, we analyze several distance metric learning algorithms for classification, from the problems and the objective functions they try to optimize, to the methods that lead to the solution of these problems.

Regarding the theoretical background of distance metric learning, we have studied three mathematical fields closely related with this topic. The first one is convex analysis [Rockafellar 2015; Boyd and Vandenberghe 2004]. Convex analysis is present in many distance metric learning algorithms, since they try to optimize convex functions over convex sets. Some interesting properties about convex sets, as well as how to deal with constrained convex problems, will be shown in this study. We will also see how the use of matrices is a fundamental part of modeling our problem. Matrix analysis [Horn and Johnson 1990] will therefore be the second field. The third field is information theory [Cover and Thomas 2006], which is also used in some of the algorithms we will show. In addition, the theoretical approach on machine learning that we will follow is the one provided by Shalev-Shwartz and Ben-David [2014].

As explained before, our work focuses on supervised distance metric learning techniques. A large amount of algorithms have been proposed over the years. These algorithms were developed with different purposes and based on different ideas, so that we can classify them in different groups. In this way, we can find algorithms whose main goal is dimensionality reduction [Fisher 1936; Wang and Zhang 2007], algorithms specifically oriented to improve distance based classifiers, such as the nearest neighbors classifier [Weinberger and Saul 2009; Goldberger et al. 2005], or the nearest centroid classification [Mensink et al. 2012], and a few techniques are also based on information theory [Davis et al. 2007; Nguyen et al. 2017; Globerson and Roweis 2006]. Some of these algorithms also allow kernel versions [Torresani and Lee 2007; Mika et al. 1999; Wang and Zhang 2007; Nguyen et al. 2017], that allow for the extension of distance metric learning to highly dimensional spaces.

To complete this study, we carry out several experiments involving all the algorithms analyzed throughout this work, executed over 34 datasets. For that, we define different settings to explore their performance and capabilities when considering maximum dimension, centroid-based methods, different kernels and dimensionality reduction. Bayesian statistical tests are used to assess the significant differences among algorithms [Benavoli et al. 2017].

Several surveys on distance metric learning have been proposed. Among the well-known surveys we can find the work of Yang and Jin [2006], Kulis et al. [2013], Bellet et al. [2013] and Moutafis et al. [2017]. In our paper, we want to differ from these previous publications by focusing on a deeper analysis of the main concepts of distance metric learning, trying to get to its basic ideas, as well as providing an experimental framework with the most popular metric learning algorithms. Besides, we will discuss some opportunities for future work in this topic.

Our paper is organized as follows. Section 2 introduces the distance metric problem, explains the family of distances we will work with and shows several examples and applications. Section 3 discusses all the distance metric learning algorithms chosen for this tutorial. Section 4 describes the experiments done to evaluate the performance of the algorithms and shows the obtained results. Finally, Sections 5 and 6 conclude the paper by summarizing the work done and indicating possible future avenues of research in this area, respectively. Due to space limitations, an electronic supplement is provided, in which Appendix A presents the mathematical foundations of distance metric learning, structured in the three blocks discussed previously, and Appendix B provides the detailed explanation of the algorithms in Section 3.

## 2. DISTANCE METRIC LEARNING

In this section we will introduce the distance metric learning problem. To begin with, we will remember the concept of distance, with special emphasis on those distances that will allow us to model our problem. Next, we will describe the distance metric learning problem, and we will finish by showing some of its applications.

### 2.1. Mahalanobis Distances

We will start by reviewing the concept of distance and some of its properties.

*Definition 2.1.* Let  $X$  be a non empty set. A *distance* or *metric* over  $X$  is a map  $d: X \times X \rightarrow \mathbb{R}$  that satisfies the following properties:

- (1) Coincidence:  $d(x, y) = 0 \iff x = y$ , for every  $x, y \in X$ .
- (2) Symmetry:  $d(x, y) = d(y, x)$ , for every  $x, y \in X$ .
- (3) Triangle inequality:  $d(x, z) \leq d(x, y) + d(y, z)$ , for every  $x, y, z \in X$ .

The ordered pair  $(X, d)$  is called a *metric space*.

The coincidence property stated above will not be important for us. That is why we will also consider mappings known as *pseudodistances*, which demand only that  $d(x, x) = 0$ , instead of the coincidence property. In fact, pseudodistances are very related with dimensionality reduction, which is an important application of distance metric learning. From now on, when we talk about of distances, we will be considering proper distances as well as pseudodistances.

*Remark 2.2.* As an immediate consequence of the definition, we have the following additional properties about distances:

- (4) Non negativity:  $d(x, y) \geq 0$  for every  $x, y \in X$ .
- (5) Reverse triangle inequality:  $|d(x, y) - d(y, z)| \leq d(x, z)$  for every  $x, y, z \in X$ .
- (6) Generalized triangle inequality:  $d(x_1, x_n) \leq \sum_{i=1}^{n-1} d(x_i, x_{i+1})$  for  $x_1, \dots, x_n \in X$ .

When we work in the  $d$ -dimensional euclidean space we find a family of distances very useful in the computing field. These distances are parameterized by positive semidefinite matrices and are known as *Mahalanobis distances*. In what follows, we will refer to  $\mathcal{M}_{d' \times d}(\mathbb{R})$  (resp.  $\mathcal{M}_d(\mathbb{R})$ ) as the set of matrices of dimension  $d' \times d$  (resp. square matrices of dimension  $d$ ), and to  $S_d(\mathbb{R})_0^+$  as the set of positive semidefinite matrices of dimension  $d$ .

*Definition 2.3.* Let  $d \in \mathbb{N}$  and  $M \in S_d(\mathbb{R})_0^+$ . The *Mahalanobis distance* corresponding to the matrix  $M$  is the map  $d_M: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  given by

$$d_M(x, y) = \sqrt{(x - y)^T M (x - y)}, \quad x, y \in \mathbb{R}^d.$$

Mahalanobis distances come from the (semi-)dot products in  $\mathbb{R}^d$  defined by the positive semidefinite matrix  $M$ . When  $M$  is full-rank, Mahalanobis distances are proper distances. Otherwise, they are pseudodistances. Observe that the euclidean usual distance is a particular example of a Mahalanobis distance, when  $M$  is the identity matrix  $I$ . Mahalanobis distances have additional properties specific to distances over normed spaces.

- (7) Homogeneous:  $d(ax, ay) = |a|d(x, y)$ , for  $a \in \mathbb{R}$ , and  $x, y \in \mathbb{R}^d$ .
- (8) Translation invariance:  $d(x, y) = d(x + z, y + z)$ , for  $x, y, z \in \mathbb{R}^d$ .

Sometimes the term “Mahalanobis distance” is used to describe the squared distances of the form  $d_M^2(x, y) = (x - y)^T M (x - y)$ . In the area of computing, it is much more efficient to work with  $d_M^2$  rather than with  $d_M$ , as this avoids the calculation of square roots. Although  $d_M^2$  is not really a distance, it keeps the most useful properties of  $d_M$  from the distance metric learning perspective, as we will see, such as the greater or lesser closeness between different pairs of points. That is why the use of the term “Mahalanobis distance” for both  $d_M$  and  $d_M^2$  is quite widespread.

To conclude this part, we return to the issue of dimensionality reduction that we mentioned when introducing the concept of pseudodistance. When we work with a pseudodistance  $\sigma$  over a set  $X$ , it is possible to define an equivalence relationship given by  $x \sim y$  if and only if  $\sigma(x, y) = 0$ , for each  $x, y \in X$ . With this relationship we can consider the quotient space  $X/\sim$ , and the map  $\hat{\sigma}: X/\sim \times X/\sim \rightarrow \mathbb{R}$  given by  $\hat{\sigma}([x], [y]) = \sigma(x, y)$ , for each  $[x], [y] \in X/\sim$ . This map is well defined and is a distance over the quotient space. When  $\sigma$  is a Mahalanobis distance over  $\mathbb{R}^d$ , with rank  $d' < d$  (we define the rank of a Mahalanobis distance as the rank of the associated positive semidefinite matrix), then the previous quotient space becomes a vector space isomorphic to  $\mathbb{R}^{d'}$ , and the distance  $\hat{\sigma}$  is a full-rank Mahalanobis distance over  $\mathbb{R}^{d'}$ . That is why, when we have a Mahalanobis pseudodistance on  $\mathbb{R}^d$ , we can view this as a proper Mahalanobis distance over a lower dimensional space, hence we have obtained a dimensionality reduction.

## 2.2. Problem Description

One of the most important components in many human cognitive processes is the ability to detect similarities between different objects. This ability has been taken to the field of machine learning by designing algorithms that learn from a dataset according to the similarities between those data.

To measure the similarity between data, it is necessary to introduce a distance, which allows us to establish a measure whereby it is possible to determine when a pair of samples is more similar than another pair of samples. However, there is an infinite number of distances we can work with, and not all of them will adapt properly to our data. Therefore, the choice of an adequate distance is a crucial element in this type of algorithm. The search for an appropriate distance is the task that is carried out in distance metric learning.

*Distance metric learning* is a machine learning discipline with the purpose of learning distances from a dataset. In its most general version, a dataset  $\mathcal{X} = \{x_1, \dots, x_N\}$  is available, on which certain similarity measures between different pairs or triplets of data are collected. These similarities are determined by the sets

$$\begin{aligned} S &= \{(x_i, x_j) \in \mathcal{X} \times \mathcal{X} : x_i \text{ and } x_j \text{ are similar.}\}, \\ D &= \{(x_i, x_j) \in \mathcal{X} \times \mathcal{X} : x_i \text{ and } x_j \text{ are not similar.}\}, \\ R &= \{(x_i, x_j, x_l) \in \mathcal{X} \times \mathcal{X} \times \mathcal{X} : x_i \text{ is more similar to } x_j \text{ than to } x_l.\}. \end{aligned}$$

With these data and similarity constraints, the problem to be solved consists in finding, after establishing a family of distances  $\mathcal{D}$ , those distances that best adapt to the criteria specified by the similarity constraints. To do this, a certain loss function  $\ell$  is set, and the distances to seek will be those that solve the optimization problem

$$\min_{d \in \mathcal{D}} \ell(d, S, D, R).$$

When we focus on supervised learning, in addition to dataset  $\mathcal{X}$  we have a list of labels  $y_1, \dots, y_N$  corresponding to each sample in  $\mathcal{X}$ . The general formulation of the distance metric learning problem is easily adapted to this new situation, just by considering the sets  $S$  and  $D$  as

$$\begin{aligned} S &= \{(x_i, x_j) \in \mathcal{X} \times \mathcal{X} : y_i = y_j\}, \\ D &= \{(x_i, x_j) \in \mathcal{X} \times \mathcal{X} : y_i \neq y_j\}. \end{aligned}$$

In addition, the set  $R$  may be also available by defining triplets  $(x_i, x_j, x_l)$  where in general  $y_i = y_j \neq y_l$ , and also verifying certain conditions on the distance between  $x_i$  and  $x_j$ , as opposed to the distance between  $x_i$  and  $x_l$ . This is the case, for example, for impostors in the LMNN algorithm (see Section B.2.1 and Weinberger and Saul [2009]). In any case, labels have all the necessary information in the field of supervised distance metric learning. From now on we will focus on this kind of problem.

Furthermore, focusing on the nature of the dataset, practically all of the distance metric learning theory is developed for numerical data, due in part to the richness of the distances available to this kind of sets, and their ability to be parameterized computationally, and in part to the fact that nominal data can be converted to binary numerical variables, or ordinal variables, with an appropriate preprocessing. For this reason, from now on, we will focus on supervised learning problems with numerical datasets.

We will suppose then that  $\mathcal{X} \subset \mathbb{R}^d$ . As we saw in the previous section, for finite-dimensional vector spaces we have the family of Mahalanobis distances,  $\mathcal{D} = \{d_M : M \in S_d(\mathbb{R})_0^+\}$ . With this family, we have at our disposal all the distances associated with dot products in  $\mathbb{R}^d$  (and in lower dimensions). In addition, this family is determined by the set of positive semidefinite matrices, and therefore, we can use these matrices, which we will call *metric matrices*, to parameterize distances. In this way, the general problem adapted to supervised learning with Mahalanobis distances can be rewritten as

$$\min_{M \in S_d(\mathbb{R})_0^+} \ell(d_M, (x_1, y_1), \dots, (x_N, y_N)).$$

However this is not the only way to parameterize this type of problem. We know, from Theorem A.10, that if  $M \in S_d(\mathbb{R})_0^+$ , then there exists a matrix  $L \in \mathcal{M}_d(\mathbb{R})$  so that  $M = L^T L$ , and this matrix is unique except for an isometry. So then we get

$$d_M^2(x, y) = (x - y)^T M (x - y) = (x - y)^T L^T L (x - y) = (L(x - y))^T (L(x - y)) = \|L(x - y)\|_2^2.$$

Therefore, we can also parameterize Mahalanobis distances through any matrix, although in this case the interpretation is different. When we learn distances through positive semidefinite matrices we are learning a new metric over  $\mathbb{R}^d$ . When we learn distances with the previous  $L$  matrices, we are learning a linear map (given by  $x \mapsto Lx$ ) that transforms the data in the space, and the corresponding distance is the usual euclidean distance after projecting the data onto the new space through the linear map. Both approaches are equivalent thanks to Theorem A.10.

In relation to dimensionality, it is important to note that, when the learned metric  $M$  is not full-rank, we are actually learning a distance over a space of lower dimension (as

we mentioned in the previous section), which allows us to reduce the dimensionality of our dataset. The same occurs when we learn linear maps that are not full-rank. We can extend this case and opt to learn directly linear maps defined by  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , with  $d' < d$ . In this way, we ensure that data are directly projected into a space of dimension no greater than  $d'$ .

Both learning the metric matrix  $M$  and learning the linear transformation  $L$ , are useful approaches to model distance metric learning problems, each one with its advantages and disadvantages. For example, parameterizations via  $M$  usually lead to convex optimization problems. In contrast, convexity in problems parameterized by  $L$  is not so easy to achieve. On the other hand, parameterizations through  $L$  make it possible to learn projections directly onto lower dimensional spaces, while dimensional constraints for problems parameterized by  $M$  are not so easy to achieve. Let us examine these differences with simple examples.

*Example 2.4.* Many of the functions we will want to optimize will depend on the squared distance defined by the metric  $M$  or by the linear transformation  $L$ , that is, either they will have terms of the form  $\|v\|_M^2 = v^T M v$ , or of the form  $\|v\|_L^2 = \|Lv\|_2^2$ . Both the maps  $M \mapsto \|v\|_M^2$  and  $L \mapsto \|v\|_L^2$  are convex (the first is actually affine). However, if we want to subtract terms in this way, we lose convexity in  $L$ , because the mapping  $L \mapsto -\|v\|_L^2$  is no longer convex. In contrast, the mapping  $M \mapsto -\|v\|_M^2$  is still affine and, therefore, convex.

*Example 2.5.* Rank constraints are not convex, and therefore we may not dispose of a projection onto the set corresponding to those constraints, unless we learn the mapping (parameterized by  $L$ ) directly to the space with the desired dimension, as explained before. For example, if we consider the set  $C = \{M \in S_2(\mathbb{R})_0^+ : r(A) \leq 1\}$ , we get  $A = \begin{pmatrix} 2 & 0 \\ 0 & 0 \end{pmatrix} \in C$  and  $B = \begin{pmatrix} 0 & 0 \\ 0 & 2 \end{pmatrix} \in C$ . However,  $(1 - \lambda)A + \lambda B = I \notin C$ , for  $\lambda = 1/2$ .

### 2.3. Some Applications

This section describes some of the main applications of distance metric learning, illustrated with several examples.

— **Improve the performance of distance-based classifiers.** This is one of the main purposes of distance metric learning. Through such learning, a distance that fits well with the dataset and the classifier can be found, improving the performance of the classifier [Weinberger and Saul 2009; Goldberger et al. 2005]. An example is shown in Figure 1.

— **Dimensionality reduction.** As we have already commented, learning a low-rank metric implies a dimensionality reduction on the dataset we are working with. This dimensionality reduction provides numerous advantages, such as a reduction in the computational cost, both in space and time, of the algorithms that will be used later, or the removal of the possible noise introduced when picking up the data. In addition, some distance-based classifiers are exposed to a problem called *curse of dimensionality* (see, for example, Shalev-Shwartz and Ben-David [2014], sec. 19.2.2). By reducing the dimension of the dataset, this problem also becomes less serious. Finally, if deemed necessary, projections onto dimension 1, 2 and 3 would allow us to obtain visual representations of our data, as shown in Figure 2. In general, many real-world problems arise with a high dimensionality, and need a dimensionality reduction to be handled properly [Van Der Maaten et al. 2009].

— **Axes selection and data rearrangement.** Closely related to dimensionality reduction, this application is a result of algorithms that learn transformations which allow the coordinate axes to be moved (or selected according to the dimension), so that

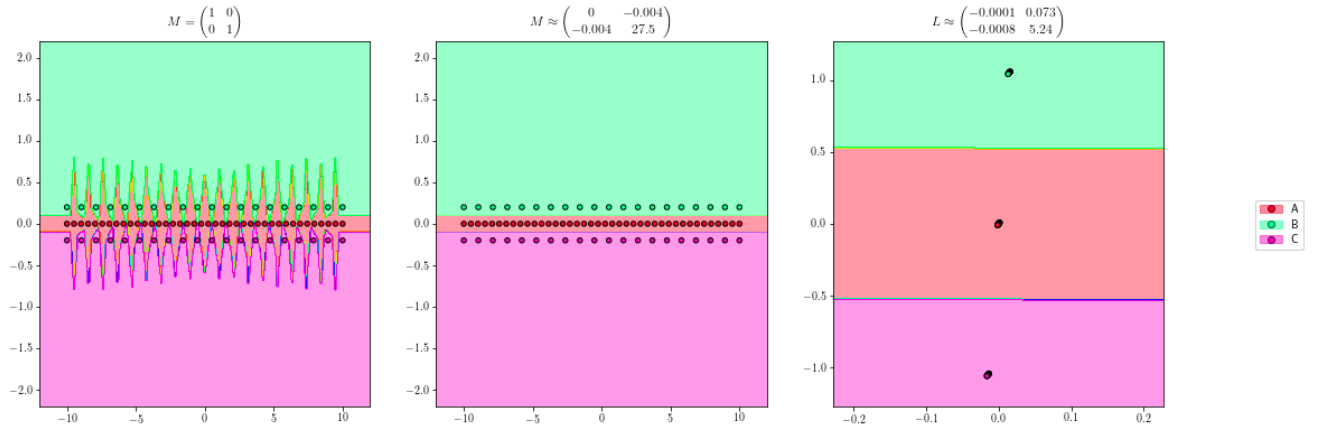


Fig. 1. Suppose we have a dataset in the plane, where data can belong to three different classes, whose regions are defined by parallel lines. Suppose that we want to classify new samples using the one nearest neighbor classifier. If we use euclidean distance, we would obtain the classification regions shown in the image on the left, because there is a greater separation between each sample in class B and class C than there is between the regions. However, if we learn an adequate distance and try to classify with the nearest neighbor classifier again, we obtain much more effective classification regions, as shown in the center image. Finally, as we have seen, learning a metric is equivalent to learning a linear map and to use euclidean distance in the transformed space. This is shown in the right figure. We can also observe that data are being projected, except for precision errors, onto a line, thus we are also reducing the dimensionality of the dataset.

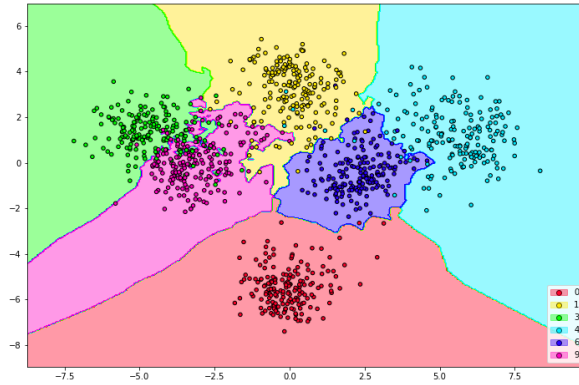


Fig. 2. 'Digits' dataset consists of 1797 examples. Each of them consists of a vector with 64 attributes, representing intensity values on an 8x8 image. The examples belong to 10 different classes, each of them representing the numbers from 0 to 9. By learning an appropriate transformation we are able to project most classes on the plane, so that we perceive clearly differentiated regions associated with each of the classes.

in the new coordinate system the vectors concentrate certain measures of information on their first components [Kokiopoulou et al. 2011]. An example is shown in Figure 3.

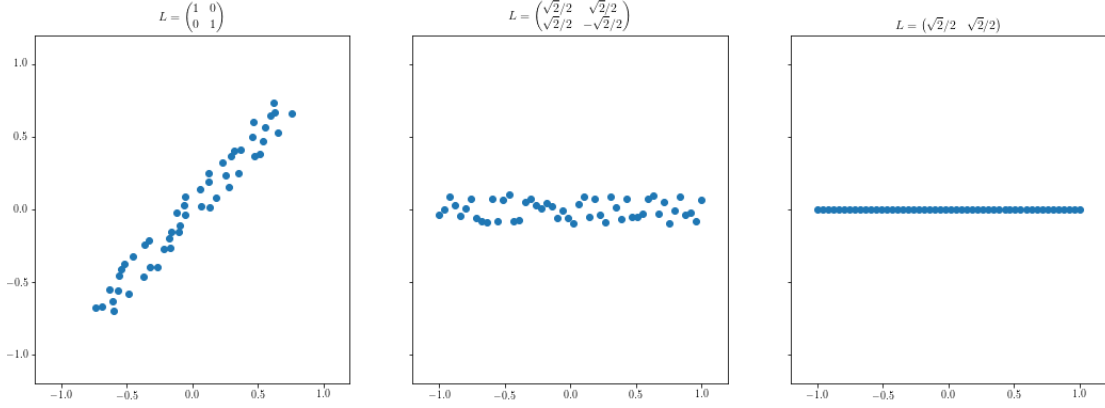


Fig. 3. The dataset in the left figure seems to concentrate most of its information on the diagonal line that joins the lower left and upper right corners. By learning an appropriate transformation, we can get that direction to fall on the horizontal axis, as shown in the center image. In this way, the first coordinate of the vectors in this new basis concentrates a large part of the variability of the vector. In addition, it seems reasonable to think that the values introduced by the vertical coordinate can be due to noise, so we can even keep only the first component, as shown in the right image.

— **Improve the performance of clustering algorithms.** Many of the clustering algorithms use a distance to measure the closeness between data, and thus establish the clusters so that data in the same cluster are considered close for that distance. Sometimes, although we do not know the ideal groupings of the data or the number of clusters to establish, we can know that certain pairs of points must be in the same cluster and that other specific pairs must be in different clusters [Xing et al. 2003]. This happens in numerous problems, for example, when clustering web documents [Aggarwal et al. 2012]. These documents have a lot of additional information, such as links between documents, which can be included as similarity constraints.

— **Semi-supervised learning.** Semi-supervised learning is a learning model in which there is one set of labeled data and another set (generally much larger) of unlabeled data. Both datasets are intended to learn a model that allows new data to be labeled. Semi-supervised learning arises from the fact that in many situations collecting unlabeled data is relatively straightforward, but assigning labels can require a supervisor to assign them manually, which may not be feasible. In contrast, when a lot of unlabeled data is used along with a small amount of labeled data, it is possible to considerably improve learning outcomes, as exemplified in Figure 4. Many of these techniques consist of constructing a graph with weighted edges from the data, where the value of the edges depends on the distances between the data. From this graph we try to infer the labels of the whole dataset, using different propagation algorithms [Zhu and Ghahramani 2002]. In the construction of the graph, the choice of a suitable distance is important, thus distance metric learning comes into play [Dhillon et al. 2010].



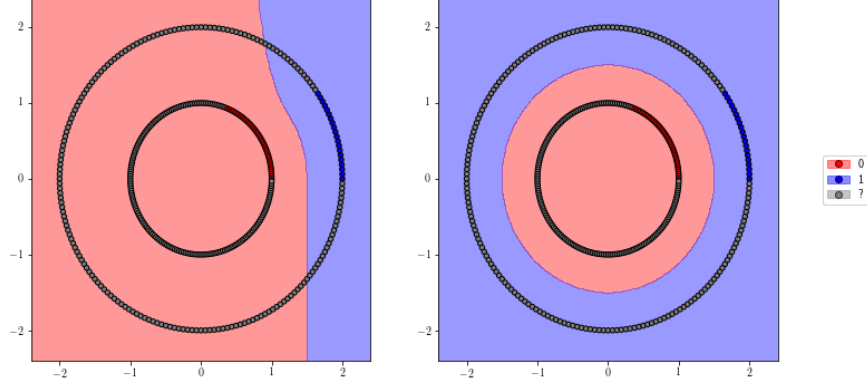


Fig. 4. Learning with only supervised information (left) versus learning with all unsupervised information (right).

From the applications we have seen, we can conclude that distance metric learning can be viewed as a preprocessing step for many distance-based learning algorithms. The algorithms analyzed in our work focus on the first three applications of the above enumeration.

### 3. ALGORITHMS FOR DISTANCE METRIC LEARNING

This section introduces some of the most popular techniques currently being used in supervised distance metric learning. Due to space issues, the section shows a brief description for each of the algorithms, while a detailed description can be found in Appendix B, where the problems the algorithms try to optimize are analyzed, together with their mathematical foundations and the techniques used to solve them.

Table I shows the algorithms studied throughout this work, including name, references and a short description. These algorithms will be empirically analyzed in the next section. This study is not intended to be exhaustive and therefore only some of the most popular algorithms have been selected for the theoretical study and the subsequent experimental analysis.

We will now provide a brief introduction to these algorithms. According to the main purpose of each algorithm, we can group them into different categories: dimensionality reduction techniques (Section 3.1), algorithms oriented to improve the nearest neighbors classifiers (Section 3.2), algorithms oriented to improve the nearest centroid classifiers (Section 3.3), or algorithms based on information theory (Section 3.4). These categories are not necessarily exclusive, but we have considered each of the algorithms in the category associated with their dominant purpose. We also introduce, in Section 3.5 several algorithms with less specific goals, and finally, in Section 3.6, the kernel versions for some of the algorithms studied.

For each of the techniques we will show the problem they try to solve. For more details of each algorithm, the reader can refer to its corresponding section of the appendix, as shown in Table I.

In what follows, we will assume a training set  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , with corresponding labels  $y_1, \dots, y_N$ .

Table I. Description of the distance metric learning algorithms analyzed in this study.

Name	References	App. Section	Description
PCA	Jolliffe [2002]	B.1.1	A dimensionality reduction technique that obtains directions maximizing variance. Although not supervised, it is important to allow dimensionality reduction in other algorithms that are not able on their own.
LDA	Fisher [1936]	B.1.2	A dimensionality reduction technique that obtains direction maximizing a ratio involving <i>between-class</i> variances and <i>within-class</i> variances.
ANMM	Wang and Zhang [2007]	B.1.3	A dimensionality reduction technique that aims at optimizing an average neighborhood margin between same-class neighbors and different-class neighbors, for each of the samples.
LMNN	Weinberger and Saul [2009]	B.2.1	An algorithm aimed at improving the accuracy of the $k$ -neighbors classifier. It tries to optimize a two-term error function that penalizes, on the one hand, large distance between each sample and its <i>target neighbors</i> , and on the other hand, small distances between different-class samples.
NCA	Goldberger et al. [2005]	B.2.2	An algorithm aimed at improving the accuracy of the $k$ -neighbors classifier, by optimizing the expected <i>leave-one-out</i> accuracy for the nearest neighbor classification.
NCMML	Mensink et al. [2012]	B.3.1	An algorithm aimed at improving the accuracy of the <i>nearest class mean</i> classifier, by optimizing a log-likelihood for the labeled data in the training set.
NCMC	Mensink et al. [2012]	B.3.2	A generalization of the NCMML algorithm aimed at improving nearest centroids classifiers that allow multiple centroids per class.
ITML	Davis et al. [2007]	B.4.1	An information theory based technique that aims at minimizing the Kullback-Leibler divergence with respect to an initial gaussian distribution, but while keeping certain similarity constraints between data.
DMLMJ	Nguyen et al. [2017]	B.4.2	An information theory based technique that aims at maximizing the Jeffrey divergence between two distributions, associated to similar and dissimilar points, respectively.
MCML	Globerson and Roweis [2006]	B.4.3	An information theory based technique that tries to collapse same-class points in a single point, as far as possible from the other classes collapsing points.
LSI	Xing et al. [2003]	B.5.1	A distance metric learning algorithm that globally minimizes the distances between same-class points, while fulfilling minimum-distance constraints for different-class points.
DML-eig	Ying and Li [2012]	B.5.2	A distance metric learning algorithm similar to LSI that offers a different resolution method based on eigenvalue optimization.
LDML	Guillaumin et al. [2009]	B.5.3	A probabilistic approach for distance metric learning based on the logistic function.
KLMNN	Weinberger and Saul [2009]; Torresani and Lee [2007]	B.6.1	The kernel version of LMNN.
KANMM	Wang and Zhang [2007]	B.6.2	The kernel version of ANMM.
KDMLMJ	Nguyen et al. [2017]	B.6.3	The kernel version of DMLMJ.
KDA	Mika et al. [1999]	B.6.4	The kernel version of LDA.

### 3.1. Dimensionality Reduction Techniques

Dimensionality reduction techniques try to learn a distance by searching a linear transformation from the dataset space to a lower dimensional euclidean space. We will describe the algorithms PCA [Jolliffe 2002], LDA [Fisher 1936] and ANMM [Wang and Zhang 2007].

**3.1.1. PCA.** PCA (*principal component analysis*) [Jolliffe 2002] is one of the most popular dimensionality reduction techniques in unsupervised distance metric learning. Although PCA is an unsupervised learning algorithm, it is necessary to talk about

it in our work, firstly because of its great relevance, and more particularly, because when a supervised distance metric learning algorithm does not allow a dimensionality reduction, PCA can be first applied to the data in order to be able to use the algorithm later in the lower dimensional space.

The purpose of PCA is to learn a linear transformation from the original space  $\mathbb{R}^d$  to a lower dimensional space  $\mathbb{R}^{d'}$  for which the loss when recomposing the data in the original space is minimized. This has been proven to be equivalent to iteratively finding orthogonal directions for which the projected variance of the dataset is maximized. The linear transformation is then the projection onto these directions. The optimization problem can be formulated as

$$\max_{\substack{L \in \mathcal{M}_{d' \times d}(\mathbb{R}) \\ LL^T = I}} \text{tr}(L\Sigma L^T),$$

where  $\Sigma$  is, except for a constant, the covariance matrix of  $\mathcal{X}$ , and  $\text{tr}$  is the trace operator. The solution to this problem can be obtained by taking as the rows of  $L$  the eigenvectors of  $\Sigma$  associated with its largest eigenvalues.

**3.1.2. LDA.** LDA (*linear discriminant analysis*) [Fisher 1936] is a classical distance metric learning technique with the purpose of learning a projection matrix that maximizes the separation between classes in the projected space using *within-class* and *between-class* variances. It follows a scheme similar to the one proposed by PCA, but in this case it takes into account the supervised information provided by the labels.

The optimization problem of LDA is formulated as

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \text{tr}((LS_w L^T)^{-1}(LS_b L^T)),$$

where  $S_b$  and  $S_w$  are, respectively, the between-class and within-class *scatter matrices*, which are defined as

$$S_b = \sum_{c \in \mathcal{C}} N_c (\mu_c - \mu)(\mu_c - \mu)^T,$$

$$S_w = \sum_{c \in \mathcal{C}} \sum_{i \in \mathcal{C}_c} (x_i - \mu_c)(x_i - \mu_c)^T,$$

where  $\mathcal{C}$  is the set of all the labels,  $\mathcal{C}_c$  is the set of indices  $i$  for which  $y_i = c \in \mathcal{C}$ ,  $N_c$  is the number of samples in  $\mathcal{X}$  with class  $c$ ,  $\mu_c$  is the mean of the training samples in class  $c$  and  $\mu$  is the mean of the whole training set. The solution to this problem can be found by taking as the rows of  $L$  the eigenvectors of  $S_w^{-1}S_b$  associated with its largest eigenvalues.

**3.1.3. ANMM.** ANMM (*average neighborhood margin maximization*) [Wang and Zhang 2007] is another distance metric learning technique specifically oriented to dimensionality reduction that tries to solve some of the limitations of PCA and LDA.

The objective of ANMM is to learn a linear transformation  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , with  $d' \leq d$  that maximizes an *average neighborhood margin* defined, for each sample, by the difference between its average distance to its nearest neighbors of different class and the average distance to its nearest neighbors of same class.

If we consider the set  $\mathcal{N}_i^o$  of the  $\xi$  samples in  $\mathcal{X}$  nearest to  $x_i$  and with the same class as  $x_i$ , and the set  $\mathcal{N}_i^e$  of the  $\zeta$  samples in  $\mathcal{X}$  nearest to  $x_i$  and with class different to  $x_i$ , we can express the average neighborhood margin, for the distance defined by  $L$ , as

$$\gamma^L = \text{tr}(L(S - C)L^T),$$

where  $S$  and  $C$  are, respectively, the *scatterness* and *compactness* matrices, defined as

$$S = \sum_i \sum_{k: x_k \in \mathcal{N}_i^e} \frac{(x_i - x_k)(x_i - x_k)^T}{|\mathcal{N}_i^e|}$$

$$C = \sum_i \sum_{j: x_j \in \mathcal{N}_i^o} \frac{(x_i - x_j)(x_i - x_j)^T}{|\mathcal{N}_i^o|}.$$

If we impose the scaling restriction  $LL^T = I$  (scaling would increase the average neighborhood margin indefinitely), the average neighborhood margin can be maximized by taking as the rows of  $L$  the eigenvectors of  $S - C$  associated with its largest eigenvalues.

### 3.2. Algorithms to Improve Nearest Neighbors Classifiers

One of the main applications of distance metric learning is to improve other distance based learning algorithms. Since the nearest neighbors classifier is one of the most popular distance based classifiers many distance metric learning algorithms are designed in order to improve this classifier, as is the case with LMNN [Weinberger and Saul 2009] and NCA [Goldberger et al. 2005].

**3.2.1. LMNN.** LMNN (*large margin nearest neighbors*) [Weinberger and Saul 2009] is a distance metric learning algorithm aimed specifically at improving the accuracy of the  $k$ -nearest neighbors classifier.

LMNN tries to bring each sample as close as possible to its *target neighbors*, which are  $k$  pre-selected same-class samples requested to become the nearest neighbors of the sample, while trying to prevent samples from other classes from invading a margin defined by those target neighbors. This setup allows the algorithm to locally separate the classes in an optimal way for  $k$ -neighbors classification.

Assuming the sets of target neighbors are chosen (usually they are taken as the nearest neighbors for Euclidean distance), the error function that LMNN minimizes is a two-term function. The first term is the target neighbors pulling term, given by

$$\varepsilon_{pull}(M) = \sum_{i=1}^N \sum_{j \rightsquigarrow i} d_M(x_i, x_j)^2,$$

where  $d_M$  is the Mahalanobis distance corresponding to  $M \in S_d(\mathbb{R})_0^+$  and  $j \rightsquigarrow i$  iff  $x_j$  is a target neighbor of  $x_i$ . The second term is the *impostors* pushing term, given by

$$\varepsilon_{push}(M) = \sum_{i=1}^N \sum_{j \rightsquigarrow i} \sum_{l=1}^N (1 - y_{il}) [1 + d_M(x_i, x_j)^2 - d_M(x_i, x_l)^2]_+,$$

where  $y_{il} = 1$  if  $y_i = y_l$  and 0 otherwise, and  $[\cdot]_+$  is defined as  $[z]_+ = \max\{z, 0\}$ . Finally, the objective function is given by

$$\varepsilon(M) = (1 - \mu)\varepsilon_{pull}(M) + \mu\varepsilon_{push}(M), \quad \mu \in ]0, 1[.$$

This function can be optimized using semidefinite programming. It is possible to optimize this function in terms of  $L$ , using gradient methods, as well. By optimizing in terms of  $M$  we gain convexity in the problem, while by optimizing in terms of  $L$  we can use the algorithm to force a dimensionality reduction.

**3.2.2. NCA.** NCA (*neighborhood components analysis*) [Goldberger et al. 2005] is another distance metric learning algorithm aimed specifically at improving the accuracy of the nearest neighbors classifiers. Its aim is to learn a linear transformation with the goal of minimizing the leave-one-out error expected by the nearest neighbor classification.

To do this, we define the probability that a sample  $x_i \in \mathcal{X}$  has  $x_j \in \mathcal{X}$  as its nearest neighbor for the distance defined by  $L \in \mathcal{M}_d(\mathbb{R})$ ,  $p_{ij}^L$ , as the softmax

$$p_{ij}^L = \frac{\exp(-\|Lx_i - Lx_j\|^2)}{\sum_{k \neq i} \exp(-\|Lx_i - Lx_k\|^2)} \quad (j \neq i), \quad p_{ii}^L = 0.$$

The expected number of correctly classified samples according to this probability is obtained as

$$f(L) = \sum_{i=1}^N \sum_{j \in C_i} p_{ij}^L,$$

where  $C_i$  is the set of indices  $j$  so that  $y_j = y_i$ . The function  $f$  can be maximized using gradient methods, and the distance resulting from this optimization is the one that minimizes the expected leave-one-out error, and therefore, the one that NCA learns.

### 3.3. Algorithms to Improve Nearest Centroids Classifiers

Apart from the nearest neighbors classifiers, other distance-based classifiers of interest are the so-called nearest centroid classifiers. These classifiers obtain a set of centroids for each class and classify a new sample by considering the nearest centroids to the sample. There are also distance metric learning algorithms designed for these classifiers, as is the case for NCML and NCMC [Mensink et al. 2012].

**3.3.1. NCML.** NCML (*nearest class mean metric learning*) [Mensink et al. 2012] is a distance metric learning algorithm specifically designed to improve the nearest class mean (NCM) classifier. To do this, it uses a probabilistic approach similar to that used by NCA to improve the accuracy of the nearest neighbors classifier.

In this case, we define the probability that a sample  $x_i \in \mathcal{X}$  will be labeled with the class  $c$ , according to the nearest class mean criterion, for the distance defined by  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , as

$$p_L(c|x) = \frac{\exp(-\frac{1}{2}\|L(x - \mu_c)\|^2)}{\sum_{c' \in \mathcal{C}} \exp(-\frac{1}{2}\|L(x - \mu_{c'})\|^2)},$$

where  $\mathcal{C}$  is the set of all the classes and  $\mu_c$  is the mean of the training samples with class  $c$ . The objective function that NCML tries to maximize is the log-likelihood for the labeled data in the training set, according to the probability defined above, that is,

$$\mathcal{L}(L) = \frac{1}{N} \sum_{i=1}^N \log p_L(y_i|x_i).$$

This function can be optimized using gradient methods.

**3.3.2. NCMC.** NCMC (*nearest class with multiple centroids*) is the generalization of the nearest class mean classifier. In this classifier, a set with an arbitrary number of centroids is calculated for each class, using a clustering algorithm. Then, a new sample is classified by assigning the label of its nearest centroid.

An immediate generalization of NCMML allows us to learn a distance oriented to improve NCMC. This distance metric learning algorithm is also referred to as NCMC. In this case, instead of the class means, we have a set of centroids  $\{m_{c_j}\}_{j=1}^{k_c}$ , for each class  $c \in \mathcal{C}$ . The generalized probability that a sample  $x_i \in \mathcal{X}$  will be labeled with the class  $c$  now is given by  $p_L(c|x) = \sum_{j=1}^{k_c} p_L(m_{c_j}|x)$ , where  $p_L(m_{c_j}|x)$  are the probabilities that  $m_{c_j}$  is the closest centroid to  $x$ , and is given by

$$p_L(m_{c_j}|x) = \frac{\exp\left(-\frac{1}{2}\|L(x - m_{c_j})\|^2\right)}{\sum_{c \in \mathcal{C}} \sum_{i=1}^{k_c} \exp\left(-\frac{1}{2}\|L(x - m_{c_i})\|^2\right)}.$$

Again, NCMC maximizes the log-likelihood function  $\mathcal{L}(L) = \frac{1}{N} \sum_{i=1}^N p_L(y_i|x_i)$  using gradient methods.

### 3.4. Information Theory Based Algorithms

Several distance metric learning algorithms rely on information theory to learn their corresponding distances. The information theory concepts used in the algorithms we will introduce below are described in Appendix A.3. The working scheme of these algorithms is similar. First, they establish different probability distributions on the data, and then they try to bring these distributions closer or further away using divergences. The information theory based algorithms we will study are ITML [Davis et al. 2007], DMLMJ [Nguyen et al. 2017] and MCML [Globerson and Roweis 2006].

**3.4.1. ITML.** ITML (*information theoretic metric learning*) [Davis et al. 2007] is a distance metric learning technique whose objective is to find a distance metric as close as possible to an initial pre-defined distance, on which similarity and dissimilarity constraints for same-class and different-class samples are satisfied. This approach tries to preserve the properties of the original distance while adapting it to our dataset thanks to the restrictions it adds.

We will denote the positive definite matrix associated with the initial distance as  $M_0$ . Given any positive definite matrix  $M \in S_d(\mathbb{R})^+$  and a fixed mean vector  $\mu$ , we can construct a normal distribution  $p(x|M)$  with mean  $\mu$  and covariance  $M$ . What ITML tries to minimize is the Kullback-Leibler divergence between  $p(x|M_0)$  and  $p(x|M)$ , subject to several similarity constraints on the data, that is

$$\begin{aligned} \min_{M \in S_d(\mathbb{R})^+} \quad & \text{KL}(p(x|M_0) \| p(x|M)) \\ \text{s.t.:} \quad & d_M(x_i, x_j) \leq u, \quad (i, j) \in S \\ & d_M(x_i, x_j) \geq l, \quad (i, j) \in D, \end{aligned}$$

where  $S$  and  $D$  are sets of pairs of indices on the elements of  $\mathcal{X}$  that represent the samples considered similar and not similar, respectively (normally, same-labeled pairs and different-labeled pairs), and  $u$  and  $l$  are, respectively, upper and lower bounds for the similarity and dissimilarity constraints. This problem can be optimized using gradient methods combined with iterated projections in order to fulfill the constraints.

**3.4.2. DMLMJ.** DMLMJ (*distance metric learning through the maximization of the Jeffrey divergence*) [Nguyen et al. 2017] is another distance metric learning technique based on information theory that tries to separate, with respect to the Jeffrey divergence, two probability distributions, the first associated with similar points while the second is associated with dissimilar points.

DMLMJ defines two difference spaces: a *k-positive difference space* that contains the differences between each sample in the dataset and its *k*-nearest neighbors from the same class, and a *k-negative difference space* that contains the differences between each sample and its *k*-nearest neighbors from different classes. Over these spaces, for a distance determined by a linear transformation  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , two gaussian distributions  $P_L$  and  $Q_L$  with equal mean are assumed. Then, the problem that DMLMJ optimizes is

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} f(L) = \text{JF}(P_L \| Q_L) = \text{KL}(P_L \| Q_L) + \text{KL}(Q_L \| P_L).$$

This problem can be transformed into a trace optimization problem similar to those of PCA and LDA, and can also be solved by taking eigenvectors from the covariance matrices involved in the problem.

**3.4.3. MCML.** MCML (*maximally collapsing metric learning*) [Globerson and Roweis 2006] is another distance metric learning technique based on information theory. The key idea of this algorithm is the fact that we would obtain an ideal class separation if we could project all the samples from the same class on a same point, far enough away from the points on which the rest of the classes would be projected.

In order to try to achieve this situation, MCML defines a probability that a sample  $x_j$  will be classified with the same label as  $x_i$ , with the distance given by a positive semidefinite matrix  $M \in S_d(\mathbb{R})_0^+$ ,  $p^M(j|i)$ , as the softmax

$$p^M(j|i) = \frac{\exp(-d_M(x_i, x_j)^2)}{\sum_{k \neq i} \exp(-d_M(x_i, x_k)^2)}.$$

Then, it also defines a probability  $p_0(j|i)$  for the ideal situation in which all the same-class samples collapse into the same point, far enough away from the collapsing points of the other classes, given by

$$p_0(j|i) \propto \begin{cases} 1, & y_i = y_j \\ 0, & y_i \neq y_j \end{cases}.$$

MCML tries to bring  $p^M(\cdot|i)$  as close to the ideal  $p^0(\cdot|i)$  as possible, for each  $i$ , using the Kullback-Leibler divergence between them. Therefore, the optimization problem is formulated as

$$\min_{M \in S_d(\mathbb{R})_0^+} f(M) = \sum_{i=1}^N \text{KL} [p_0(\cdot|i) \| p^M(\cdot|i)].$$

This function can be minimized using semidefinite programming.

### 3.5. Other Distance Metric Learning Techniques

In this section we will study some different proposals for distance metric learning techniques. The algorithms we will analyze are LSI [Xing et al. 2003], DML-eig [Ying and Li 2012] and LDML [Guillaumin et al. 2009].

**3.5.1. LSI.** LSI (*learning with side information*) [Xing et al. 2003], also sometimes referred to as MMC (*Mahalanobis metric for clustering*) is possibly one of the first algorithms that has helped make the concept of distance metric learning more well known. This algorithm is a global approach that tries to bring data of the same class closer together while keeping data from different classes far enough apart.

Assuming that the sets  $S$  and  $D$  represent, respectively, pairs of samples that should be considered similar or dissimilar (i.e. samples that belong to the same class or to different classes, respectively), LSI looks for a positive semidefinite matrix  $M \in S_d(\mathbb{R})_0^+$  that optimizes the following problem:

$$\begin{aligned} \min_M \quad & \sum_{(x_i, x_j) \in S} d_M(x_i, x_j)^2 \\ \text{s.t.:} \quad & \sum_{(x_i, x_j) \in D} d_M(x_i, x_j) \geq 1 \\ & M \in S_d(\mathbb{R})_0^+. \end{aligned}$$

This problem can be optimized using gradient descent together with iterated projections in order to fulfill the constraints.

**3.5.2. DML-eig.** DML-eig (*distance metric learning with eigenvalue optimization*) [Ying and Li 2012] is a distance metric learning algorithm inspired by the LSI algorithm of the previous section, proposing a very similar optimization problem but offering a completely different resolution method, based on eigenvalue optimization.

We consider again the two sets of pairs,  $S$  and  $D$ , of samples considered similar and dissimilar. DML-eig proposes an optimization problem slightly different from that proposed by the LSI algorithm, given by

$$\begin{aligned} \max_M \quad & \min_{(x_i, x_j) \in D} d_M(x_i, x_j)^2 \\ \text{s.t.:} \quad & \sum_{(x_i, x_j) \in S} d_M(x_i, x_j)^2 \leq 1 \\ & M \in S_d(\mathbb{R})_0^+. \end{aligned}$$

This problem can be transformed into a minimization problem for the largest eigenvalue of a symmetric matrix. This is a well-known problem and there are some iterative methods that allow this minimum to be reached [Overton 1988].

**3.5.3. LDML.** LDML (*logistic discriminant metric learning*) [Guillaumin et al. 2009] is a distance metric learning algorithm in which the optimization model makes use of the logistic function.

Recall that the *logistic* or *sigmoid* function is the map  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  given by

$$\sigma(x) = \frac{1}{1 + e^{-x}}.$$

In LDML, the logistic function is used to define a probability, which will assign the greater probability the smaller the distance between points. Given a positive semidefinite matrix  $M \in S_d(\mathbb{R})_0^+$ , this probability is expressed as

$$p_{ij,M} = \sigma(b - \|x_i - x_j\|_M^2),$$

where  $b$  is a positive threshold value that will determine the maximum value achievable by the logistic function, and that can be estimated by cross validation. What LDML tries to maximize is the log-likelihood given by

$$\mathcal{L}(M) = \sum_{i,j=1}^N y_{ij} \log p_{ij,M} + (1 - y_{ij}) \log(1 - p_{ij,M}),$$



where  $y_{ij}$  is a binary variable that takes the value 1 if  $y_i = y_j$  and 0 otherwise. This function can be optimized using semidefinite programming.

### 3.6. Kernel Distance Metric Learning

Kernel methods constitute a paradigm within machine learning that is very useful in many of the problems addressed in this discipline. They usually arise in problems where the learning algorithm capability is reduced, typically due to the shape of the dataset. A classic learning algorithm where the kernel trick is very useful is the *support vector machines* classifier [Burges 1998]. An example for this case is given in Figure 5.

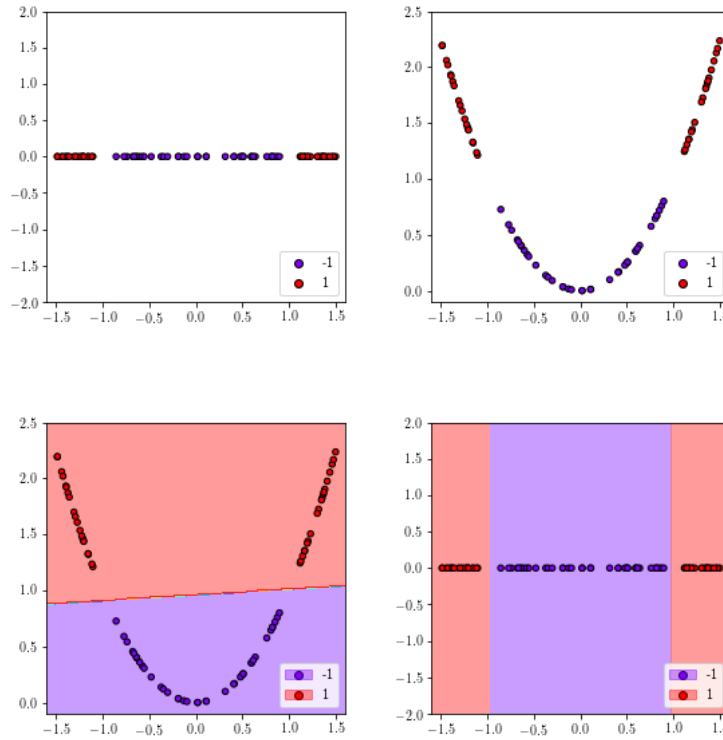


Fig. 5. Support vector machines and kernel trick. This binary classifier looks for the hyperplane that best separates both classes. Therefore, it is highly limited when the dataset is not separable by hyperplanes, as in the dataset of the upper-left image. A solution consists of sending the data into a higher dimensional space, where data can be separated by hyperplanes, and apply there the algorithm, as it is shown in the remaining images. The kernel trick allows us to execute the algorithm only in terms of the dot products of the samples in the new space, which makes it possible to work on very high dimensional spaces, or even infinite dimensional spaces. The existence of a *representer theorem* for support vector machines also allows the solution to be rewritten in terms of a vector with the size of the number of samples.

In distance metric learning, the usefulness of kernel learning is due to the limitations given by the Mahalanobis distances. Although learned metrics can later be used with non-linear classifiers, such as the nearest neighbors classifier, the metrics themselves are determined by linear transformations, which, in turn, are determined by the image of a basis in the departure space, which results in the fact that we only have the freedom to choose the image of as many data as the dimension has the space, mapping the rest of the vectors by linearity. When the amount of data is much larger than the space dimension this can become a limitation.

The kernel approach for distance metric learning follows a similar scheme to that of support vector machines. If we work with a dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , the idea is to send the data to a higher dimensional space, through a mapping  $\phi: \mathbb{R}^d \rightarrow \mathcal{F}$ , where  $\mathcal{F}$  is a Hilbert space called the *feature space*, and then to learn in the feature space using a distance metric learning algorithm. The way we will learn a distance in the feature space will be via a continuous linear transformation  $L: \mathcal{F} \rightarrow \mathbb{R}^{d'}$ , where  $d' \leq d$  (observe that  $L$  is not necessarily a matrix, since  $\mathcal{F}$  is not necessarily finite dimensional), which we will also denote  $L \in \mathcal{L}(\mathcal{F}, \mathbb{R}^{d'})$ .

As occurs with support vector machines, a great inconvenience arises when sending the data to the feature space, and that is that the problem dimension can highly increase, and therefore the application of the algorithms can be very expensive computationally. In addition, if we want to work in infinite dimensional feature spaces, it is impossible to deal with the data in this case, unless we turn to the kernel trick.

We define the *kernel function* as the mapping  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  given by  $K(x, x') = \langle \phi(x), \phi(x') \rangle$ . The success of kernel functions is due to the fact that many learning algorithms only need to know the dot products between the elements in the training set to be able to work. This will happen in the distance metric learning algorithms we will study later. We can observe, as an example, that the calculation of euclidean distances, which is essential in many distance metric learning algorithms, can be made using only the kernel function. Indeed, for  $x, x' \in \mathbb{R}^d$ , we have

$$\begin{aligned} \|\phi(x) - \phi(x')\|^2 &= \langle \phi(x) - \phi(x'), \phi(x) - \phi(x') \rangle \\ &= \langle \phi(x), \phi(x) \rangle - 2\langle \phi(x), \phi(x') \rangle + \langle \phi(x'), \phi(x') \rangle \\ &= K(x, x) + K(x', x') - 2K(x, x'). \end{aligned} \quad (1)$$

The next common problem for all the kernel-based distance metric learning algorithms is how to deal with the learned transformation. Since we are trying to learn a map  $L \in \mathcal{L}(\mathcal{F}, \mathbb{R}^{d'})$ , we may not be able to write it as a matrix, and when we can, this matrix may have dimensions that are too large. However, as  $L$  is continuous and linear, using the Riesz representation theorem, we can rewrite  $L$  as a vector of dot products by fixed vectors, that is,  $L = (\langle \cdot, w_1 \rangle, \dots, \langle \cdot, w_{d'} \rangle)$ , where  $w_1, \dots, w_{d'} \in \mathcal{F}$ . Furthermore, for the algorithms we will study, several *representer theorems* are known [Chatpatanasiri et al. 2010; Hofmann et al. 2008; Mika et al. 1999; Nguyen et al. 2017; Schölkopf et al. 1998]. These theorems allow the vectors  $w_i$  to be expressed as a linear combination of the samples in the feature space, that is, for each  $i \in \{1, \dots, d'\}$ , there is a vector  $\alpha^i = (\alpha_1^i, \dots, \alpha_N^i) \in \mathbb{R}^N$  so that  $w_i = \sum_{j=1}^N \alpha_j^i \phi(x_j)$ . Consequently, we can see that

$$L\phi(x) = A \begin{pmatrix} K(x_1, x) \\ \vdots \\ K(x_N, x) \end{pmatrix}, \quad (2)$$

where  $A \in \mathcal{M}_{d' \times N}(\mathbb{R})$  is given by  $A_{ij} = \alpha_j^i$ .

Thanks to these theorems, we can address the problem computationally as long as we are able to calculate the coefficients of matrix  $A$ . When transforming a new sample it will be enough to construct the previous column matrix evaluating the kernel function between the sample and each element in the training set, and then multiplying  $A$  by this matrix. On a final note, when training it is useful to view the kernel map as a matrix  $K \in S_N(\mathbb{R})$ , where  $K_{ij} = K(x_i, x_j)$ . A similar (in this case not necessarily square) matrix can be constructed when testing, with all the dot products between the train and test samples. Choosing the appropriate column of this matrix, we will be able to transform the corresponding test sample using Eq. 2.

Each distance metric learning technique that supports the use of kernels will use different tools for its performance, each one based on the original algorithms. Appendix B.6 will describe the kernelizations of some of the algorithms already introduced, namely, LMNN, ANMM, DMLMJ and LDA.

#### 4. EXPERIMENTAL FRAMEWORK AND RESULTS

With the algorithms introduced in the previous section, several experiments have been carried out. This section describes these experiments and shows the results.

##### 4.1. Description of the Experiments

For the distance metric learning algorithms studied, a collection of experiments has been developed, consisting of the following procedures.

- (1) Evaluation of all the algorithms capable of learning at maximum dimension, applied to the  $k$ -nearest neighbors classification, for different values of  $k$ .
- (2) Evaluation of the algorithms aimed at improving nearest centroid classifiers, applied to the corresponding centroid-based classifiers.
- (3) Evaluation of kernel-based algorithms, experimenting with different kernels, applied to the nearest neighbors classification.
- (4) Evaluation of algorithms capable of reducing dimensionality, for different dimensions, applied to the nearest neighbors classification.

When in experiment 1 we talk about “capable of learning at maximum dimension” we are excluding those dimensionality reduction algorithms that only learn a change of axes, as is the case of PCA and ANMM, which at maximum dimension learn a transformation whose associated distance is still the euclidean. LDA is kept, assuming that it will always take the maximum dimension that it is able to, according to the number of classes of the problem. The algorithms oriented to centroid-based classifiers are also excluded from experiment 1, together with those based on kernels, which will be analyzed in the experiments 2 and 3, respectively.

The stated experiments indicate that the magnitude with which we will measure the performance of the algorithms is the result of the  $k$ -neighbors classification, except in the case of the algorithms based on centroids, which will use their corresponding classifier. These classifiers will be evaluated by a 10-fold cross validation. The results obtained from the predictions on the training set will also be included, in order to evaluate possible overfitting.

To evaluate the algorithms we will use the implementations available in the Python library pyDML. The algorithms will be executed using their default parameters, which can be found in the pyDML’s documentation<sup>1</sup>. These default parameters have been set with standard values. The following exceptions to the default parameters have been made:

<sup>1</sup><https://pydml.readthedocs.io/>

- The LSI algorithm will have the parameter `supervised = True`, as it will be used for supervised learning.
- In the dimensionality reduction experiment (4), the algorithms will have the dimension number parameter set according to the dimension being evaluated.
- LMNN and KLMNN will have their parameter `k` equal to the number of neighbors being considered in the nearest neighbors classification.
- LMNN will be executed with stochastic gradient descent, instead of semidefinite programming, in dimensionality reduction experiments, thus learning a linear transformation instead of a metric.
- ANMM and KANMM will have their parameters `n_friends` and `n_enemies` equal to the number of neighbors being considered in the nearest neighbors classification.
- DMLMJ and KDMLMJ will have their parameter `n_neighbors` equal to the number of neighbors being considered in the nearest neighbors classification.
- NCMC will have its parameter `centroids_num` equal to the parameter `centroids_num` being considered in its corresponding classifier, `NCMC_Classifier`.

As for the datasets used in the experiments, up to 34 datasets have been collected, all of them available in KEEL<sup>2</sup>. All these datasets are numeric and without missing values, being oriented to standard classification problems. In addition, although some of the distance metric learning algorithms scale well with the number of samples, others cannot deal with datasets that are too large, so it was decided to select, for sets with a high number of samples, a subset of size that all algorithms can deal with, keeping the class distribution the same. Thus, the characteristics of the datasets are described in Table II. All datasets have been *min-max* normalized to the interval  $[0, 1]$ , feature to feature, prior to the execution of the experiments.

Finally, we describe the details of the experiments 1, 2, 3 and 4:

- (1) Algorithms will be evaluated with the classifiers 3-NN, 5-NN and 7-NN.
- (2) NCMLL will be evaluated with the Scikit-Learn NCM classifier, while NCMC will be evaluated with its associated classifier, available in pyDML, for two different values: 2 centroids per class and 3 centroids per class.
- (3) Algorithms will be evaluated with 3-NN classifier, using the following kernels: linear (Linear), grade-2 (Poly-2) and grade-3 (Poly-3) polynomials, gaussian (RBF) and laplacian (Laplacian). For the comparison, the kernel version of PCA<sup>3</sup> will be also included. In this case, only the smallest datasets will be considered, so that they can be applicable to the algorithms that scale the worst with the dimension (recall that the kernel trick forces algorithms to work in dimensions of the order of the number of samples).
- (4) Algorithms will be evaluated with the classifiers 3-NN, 5-NN and 7-NN. The dimensions to use will be: 1, 2, 3, 5, 10, 20, 30, 40, 50, the maximum dimension of the dataset, and the number of classes of the dataset minus 1. In this case, the following high-dimensionality datasets are selected: `sonar`, `movement_libras` and `spambase`. The algorithms to be evaluated in this experiment will be: PCA, LDA, ANMM, DMLMJ, LMNN and NCA.

## 4.2. Results

This section shows the results of the cross-validation for the different experiments. We will show in this text only the results of the 3-NN classifier, for those experiments that use nearest neighbors classifiers. The results obtained for the remaining  $k$ -NN used in

<sup>2</sup>KEEL, *knowledge extraction based on evolutionary learning* [Triguero et al. 2017]: <http://www.keel.es/>.

<sup>3</sup>It is implemented in Scikit-Learn: <http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.KernelPCA.html>. Its theoretical details can be found in Schölkopf et al. [1998].

Table II. Datasets used in the experiments.

Dataset	Number of samples	Number of features	Number of classes
appendicitis	106	7	2
balance	625	4	3
bupa	345	6	2
cleveland	297	13	5
glass	214	9	7
hepatitis	80	19	2
ionosphere	351	33	2
iris	150	4	3
monk-2	432	6	2
newthyroid	215	5	3
sonar	208	60	2
wine	176	13	3
movement_libras	360	90	15
pima	768	8	2
vehicle	846	18	4
vowel	990	13	11
wdbc	569	30	2
wisconsin	683	9	2
banana (20 %)	1,060	2	2
digits	1,797	64	10
letter (10 %)	2,010	16	26
magic (10 %)	1,903	10	2
optdigits	1,127	64	10
page-blocks (20 %)	1,089	10	4
phoneme (20 %)	1,081	5	2
ring (20 %)	1,480	20	2
satimage (20 %)	1,289	36	7
segment (20 %)	462	19	7
spambase (10 %)	460	57	2
texture (20 %)	1,100	40	11
thyroid (20 %)	1,440	21	3
titanic	2,201	3	2
twonorm (20 %)	1,481	20	2
winequality-red	1,599	11	11

the experiments are available on the pyDML-Stats<sup>4</sup> website, where the results of all these experiments have been stored. The scripts used to do the experiments can also be found in this website. To the results of the experiments 1, 2 and 3 we have added the average score obtained, and the average ranking. This ranking has been made by assigning integer values between 1 and  $m$ , where  $m$  is the number of algorithms being compared in each experiment (adding half fractions in case of a tie), according to the position of the algorithms over each dataset, 1 being the best algorithm, and  $m$  the worst one. The content of the different tables elaborated is described below.

- Table III shows the cross-validation results obtained for experiment 1, using the 3-NN score as evaluation measure. Some cells do not show results because the algorithm did not converge.
- Table IV shows the results of experiment 2. The evaluation measures were the NCM and NCMC classifiers with 2 and 3 centroids per class. For each classifier, the euclidean distance (Euclidean + CLF) and the distance learning algorithm associated with the classifier (NCML / NCMC (2 ctrd) / NCMC (3 ctrd)) have been evaluated.

<sup>4</sup>Source code: <https://github.com/jlsuarezdiaz/pyDML-Stats>. The current website is located at <https://jlsuarezdiaz.github.io/software/pyDML/stats/>

- Table V shows the cross-validation results obtained on the training set for the kernel-based algorithms using the 3-NN classifier. Table VI shows the corresponding results on the test set.
- Table VII shows the cross-validation results for the experiment 4 in dataset sonar, using the classifier 3-NN. On the left are the results for the training set, and on the right, the results for the test set. Each row shows the results for the different dimensions evaluated. Tables VIII and IX show the corresponding dimensionality results over the datasets `movement_libras` and `sonar`, respectively.

Table III. Results of cross-validation with 3-NN.

	Euclidean		LDA		ITML		DMMLJ		NCA		LMNN		LSI		DML-eig		MCMC		LDML	
	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test
appendicitis	.8428	.8339	.8428	.8522	.8533	.8604	.8490	.8256	<b>.8700</b>	.8504	.8407	.8422	.8659	<b>.8630</b>	.8585	.8622	.8502	.8513	.8669	.8422
balance	.8049	.8082	.8885	.8982	.8986	.8943	.8286	.8191	<b>.9592</b>	<b>.9684</b>	.8202	.8175	.9182	.9280	.8947	.8945	.8816	.8737	.8874	.8895
bupa	.6231	.6546	.6338	.6465	.6466	.6281	.6660	<b>.6776</b>	<b>.6943</b>	.6994	.6099	.6342	.6363	.6284	.5993	.6120	.5716	.5742	.5826	.5854
cleaveland	.5570	.5468	.5694	.5502	.5488	.5523	.5636	.5636	<b>.6804</b>	.5536	.5780	.5806	.5518	.5722	.5896	.5829	.5978	.5785	.5784	<b>.5972</b>
glass	.6759	.7015	.6235	.6231	.6453	.6549	<b>.7092</b>	.7041	.7065	.6917	.6780	<b>.7067</b>	.6495	.6235	.6407	.6263	.6319	.5850	.6242	.6063
hepatitis	.8236	.8325	.9402	.8609	.9002	.8815	.8894	.8894	<b>.9569</b>	.8325	.8514	.8418	.9139	.9130	.9125	<b>.9176</b>	.9250	.8829	.9458	.8547
ionosphere	.8569	.8550	.8834	.8394	.8771	.8862	.8752	.8752	<b>.9534</b>	<b>.9084</b>	.9281	.8859	.8898	.8768	.8904	.8741	.9053	.8630	.8907	.8512
iris	.9533	.9533	.9681	.9533	.9703	.9733	.9585	.9666	.9755	.9666	.9481	.9400	.9703	<b>.9800</b>	.9585	.9600	.9688	.9466	<b>.9807</b>	.9600
monk-2	.9578	.9655	.9451	.9561	.9223	.9352	.9709	.9724	<b>.1000</b>	<b>.1000</b>	.9812	.9816	<b>.1000</b>	<b>.1000</b>	.9878	.9909	.9665	.9676	.9382	.9495
newthyroid	.9421	.9538	.9596	.9586	.9452	.9398	.9436	.9448	<b>.9700</b>	.9722	.9658	<b>.9725</b>	.9591	.9634	.9602	.9629	.9565	.9582	.9509	.9675
sonar	.8317	.8570	.9011	.7782	.8435	.8120	.9097	.8361	.9823	.8703	<b>.9641</b>	<b>.8742</b>	.8531	.8506	.8547	.7975	.8755	.8563	.8766	.7886
wine	.9606	.9606	.9968	<b>.9888</b>	.9900	.9773	.9812	.9662	.9956	.9882	.9956	.9832	.9837	.9662	.9975	.9767	<b>.9975</b>	.9832	.9956	<b>.9888</b>
movement.libras	.7972	.8139	<b>.8685</b>	.6642	.8038	.7992	.8460	<b>.8649</b>	.8516	.8319	.9065	.8020	.7351	.7440	.7970	.7872	.8063	.8073	.7256	.7360
pima	.7372	.7396	.7259	<b>.7525</b>	.7148	.7149	.7366	.7422	<b>.7841</b>	.7370	.7290	.7278	.7206	.7395	.7174	.7266	.7173	.7239	.7285	.7240
vehicle	.7077	.7125	.7698	<b>.7623</b>	.7625	.7516	.7643	.7551	<b>.8186</b>	.7550	.6855	.6757	.6590	.6666	.6506	.6501	.7398	.7369	.7186	.7170
vowel	.9699	.9787	.9680	.9777	.9423	.9535	.9751	<b>.9808</b>	<b>.9799</b>	<b>.9808</b>	.9893	.9777	.9436	.9474	.6719	.6757	.8558	.8737	.8885	.9090
wtde	.9679	<b>.9716</b>	.9732	.9654	.9714	.9664	.9669	.9648	<b>.9751</b>	.9700	.9638	.9630	.9705	.9682	.9546	.9507	.9714	.9648	.9476	.9438
wisconsin	.9694	.9678	.9663	.9677	.9609	.9590	.9695	.9678	<b>.9723</b>	.9648	.9692	.9663	.9684	<b>.9722</b>	.9673	.9707	.9585	.9546	.9650	.9663
banana	.8543	.8555	.6504	.6469	.8536	.8556	.8550	.8565	<b>.8583</b>	<b>.8583</b>	<b>.8574</b>	.8563	.8535	.8517	.6718	.6878	.6252	.6102	.6268	.6319
digits	.9878	.9866	.9769	.9683	.9798	.9728	.9869	.9834	.9980	<b>.9894</b>	<b>.9993</b>	.9860	.9264	.9102	.8269	.8168	.9734	.9688	.9797	.9816
letter	.7174	.7208	.7955	.7967	.7161	.7195	.8163	.8204	<b>.8565</b>	<b>.8610</b>	.7048	.7162	.5396	.5496	.3191	.3214	.7600	.7534	.6217	.6372
magic	.8070	.8050	.7436	.7361	.8069	.8061	.8161	.8071	<b>.8396</b>	<b>.8145</b>	.7979	.7945	.7946	.7924	.7508	.7525	.7738	.7766	.7077	.6951
optdigits	.9495	.9495	.9777	.9671	.9512	.9731	.9669	.9770	.9761	.9956	<b>.9986</b>	<b>.9840</b>	.9398	.9306	.8164	.8022	.9761	.9591	.9596	.9591
page-blocks	.7957	.7992	.7321	.7243	.7853	.7770	.7960	.9504	.9637	.9577	.9459	.9439	-	-	.9515	.9523	.9613	.9642	.9438	.9404
phoneme	.6410	.6432	.7289	.7101	.7290	.7352	.6440	.6453	<b>.8002</b>	.7936	.7928	.7946	.7642	.7668	.7361	.7483	.7654	.7632	.7320	.7112
ring	.8585	.8564	.8541	.8387	.8495	.8341	.8670	<b>.8643</b>	<b>.8764</b>	.8511	.8565	.8558	.8490	.8465	.8153	.8130	.8246	.8171	.5427	.5501
satimage	.8970	.9020	.9353	<b>.9370</b>	.9357	.9265	.9071	.9081	<b>.9451</b>	.9187	.9076	.8928	.8898	.8853	.9095	.9068	.9367	.9319	.8818	.8710
segment	.8500	.8654	.9215	.8871	.8801	.8754	.8664	.8635	.8525	<b>.9391</b>	.9154	.9215	.9070	.9210	.9111	.9077	.9046	.9176	.9047	.9229
spambase	.9500	.9618	<b>.9983</b>	<b>.9981</b>	.9801	.9764	.9854	.9854	.9843	.9843	.9800	.9218	.9333	.9400	.8979	.9009	.9740	.9745	.8658	.8718
texture	.9313	.9319	.9375	.9450	.9397	.9402	.9355	.9361	.9459	.9395	.9320	.9319	.9357	.9354	.9458	.9485	.9377	.9320	<b>.9587</b>	<b>.9583</b>
thyroid	.7607	.7583	<b>.7727</b>	<b>.7804</b>	.7682	.7609	.7612	.7581	.5709	.6764	.6018	.6964	-	-	.7107	.7331	.7150	.7253	.7108	.7341
titanic	.9609	.9595	.9778	.9750	.9685	.9669	.9612	.9561	<b>.9817</b>	.9790	.9778	.9756	.9776	.9770	.9782	<b>.9810</b>	.9708	.9730	.9789	.9804
twonorm	.5808	<b>.5865</b>	.5657	.5733	.5754	.5828	.5828	.5860	<b>.6022</b>	.5766	.5647	.5772	.5656	.5809	.5281	.5292	.5675	.5611	.5376	.5471
winequality-red	.6558	.5661	.5147	.5426	.5808	.5382	.4926	.4544	<b>.1705</b>	<b>.3661</b>	.5220	.5279	.6411	.5382	.6691	.6220	.5735	.6279	.6794	.7161
AVG SCORE	.8383	.8425	.8515	.8363	.8500	.8470	.8560	.8526	<b>.8886</b>	<b>.8634</b>	.8490	.8492	.8410	.8405	.8059	.8041	.8438	.8359	.8125	.8062

Table IV. Results of the experiments with NCM and NCMC.

	Euclidean + NCM		NCMML		Euclidean + NCM (2 ctrd)		NCMC (2 ctrd)		Euclidean + NCM (3 ctrd)		NCMC (3 ctrd)	
	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test
appendicitis	.8365	<b>.8640</b>	<b>.8512</b>	.8440	.6479	.6031	.6248	.6246	.8102	.7800	.7557	.7266
balance	.7496	.7475	.6865	.6864	.7004	.6721	.8280	<b>.8225</b>	.7160	.6654	<b>.8321</b>	.8222
bupa	.5996	.6004	<b>.6628</b>	<b>.6407</b>	.6270	.6058	.6247	.6260	.6589	.5903	.6409	.5826
cleveland	.5742	.5377	.6296	<b>.5516</b>	.5727	.5050	.6280	.5093	.6203	.4871	<b>.6435</b>	.5135
glass	.5218	.4862	.6221	.5290	.6479	.5849	.5877	.5372	.6427	.5208	<b>.7461</b>	<b>.6421</b>
hepatitis	.8500	.8293	.9832	<b>.8688</b>	.8792	.8436	.9513	.8373	.8986	.7946	<b>.9833</b>	.8672
ionosphere	.7445	.7378	.9313	.8797	.9186	<b>.8889</b>	.9018	.8764	.9062	.8750	<b>.9477</b>	.8886
iris	.9318	.9133	<b>.9814</b>	.9600	.9696	<b>.9666</b>	.9711	.9600	.9696	.9466	.9800	.9600
monk-2	.8078	.8104	.7950	.7923	.8071	.8082	.8112	.8038	.8127	.7895	<b>.8457</b>	<b>.8237</b>
newthyroid	.9359	.9352	<b>.9798</b>	.9634	.9498	.9448	.9741	<b>.9725</b>	.9695	.9681	.9757	.9629
sonar	.7265	.7017	.9257	.7641	.7120	.6929	.9219	.7839	.8360	.7437	<b>.9412</b>	<b>.7982</b>
wine	.9687	.9495	<b>1.000</b>	.9663	.9825	.9659	.9918	<b>.9826</b>	.9762	.9432	.9912	.9704
movement_libras	.6358	.5946	.8176	.7575	.7777	.6764	.8803	.7852	.8717	.7749	<b>.9420</b>	<b>.8373</b>
pima	.7337	.7279	<b>.7717</b>	<b>.7604</b>	.7565	.7382	.6720	.6641	.7521	.7461	.7322	.7253
vehicle	.4545	.4491	<b>.7998</b>	<b>.7797</b>	.6066	.5824	.7429	.7219	.6537	.6179	.7558	.7244
vowel	.5367	.5070	.6689	.6383	.6087	.5717	.7272	.6868	.5732	.5333	<b>.7690</b>	<b>.7292</b>
wdbc	.9388	.9367	.9794	.9649	.9548	.9385	.9796	<b>.9736</b>	.9703	.9665	<b>.9810</b>	.9701
wisconsin	.9648	.9648	<b>.9681</b>	<b>.9662</b>	.9586	.9502	.9655	.9603	.9515	.9486	.9541	.9487
banana	.5769	.5737	.5571	.5558	.6417	.6359	.5846	.5859	.7516	.7573	<b>.7828</b>	<b>.7766</b>
digits	.9066	.8981	.8047	.8083	.9521	.9415	.8664	.8586	<b>.9723</b>	<b>.9644</b>	<b>.8893</b>	.8554
letter	.5707	.5341	.7114	.6868	.5895	.5282	.7251	.6902	.6601	.5714	<b>.7825</b>	<b>.7301</b>
magic	.7712	.7693	.7790	.7751	.7786	.7745	<b>.7947</b>	<b>.7861</b>	.7600	.7509	.7820	.7751
optdigits	.9173	.9104	.8018	.7978	.9479	.9352	.8542	.8185	<b>.9675</b>	<b>.9609</b>	.8923	.8641
page-blocks	.8133	.8165	<b>.9636</b>	<b>.9576</b>	.8219	.8221	.9090	.9071	.8680	.8696	.8966	.8953
phoneme	.7417	.7399	.7587	.7538	<b>.7683</b>	<b>.7667</b>	.7084	.7159	.7207	.7149	.7091	.7048
ring	.7780	.7723	.7819	<b>.7784</b>	.8002	.7601	.7197	.7012	<b>.8160</b>	.7750	.6844	.6636
satimage	.7868	.7844	<b>.8478</b>	<b>.8255</b>	.8052	.7921	.8342	.8123	.8244	.7844	.8362	.8007
segment	.8460	.8367	<b>.9437</b>	<b>.9037</b>	.8596	.8452	.9203	.8976	.8581	.8030	.9242	.8874
spambase	.8874	.8827	<b>.9584</b>	.9154	.8816	.8763	.9400	<b>.9241</b>	.8985	.8893	.9335	.9112
texture	.7445	.7372	<b>.9912</b>	<b>.9781</b>	.8586	.8500	.9694	.9581	.9079	.8900	.9759	.9654
thyroid	.4532	.4394	<b>.8163</b>	<b>.8082</b>	.5724	.5558	.6875	.6922	.5959	.5630	.7469	.7358
titanic	.7540	.7459	<b>.7825</b>	<b>.7854</b>	.6746	.6516	.5624	.6550	.5630	.6824	.7279	.7350
twonorm	.9807	<b>.9824</b>	.9854	.9797	.9799	.9723	.9847	.9790	.9787	.9743	<b>.9855</b>	.9777
winequality-red	.3519	.3371	<b>.4535</b>	<b>.4359</b>	.4067	.3838	.4031	.3870	.3940	.3582	.4024	.3738
AVG RANKING	4.941	4.441	<b>2.382</b>	<b>2.500</b>	4.117	4.000	3.411	2.911	3.764	4.235	<b>2.382</b>	2.911
AVG SCORE	.7468	.7369	.8233	.7958	.7770	.7538	.8014	.7793	.7978	.7647	<b>.8344</b>	<b>.7984</b>



[illegible][illegible]

Table VII. Results of dimensionality reduction experiments on sonar with 3-NN (train - test)

	PCA	LDA	ANMM	DMLMJ	NCA	LMNN		PCA	LDA	ANMM	DMLMJ	NCA	LMNN
1	.5016	.9011	.6965	.7826	.9214	.7237	1	.5619	.7782	.6770	.7256	.8073	.6640
2	.5891	-	.7670	.8050	.9807	.8782	2	.6293	-	.7541	.7113	.8077	.7593
3	.7729	-	.8359	.8333	.9770	.9513	3	.7641	-	.8395	.7741	.8265	.8077
5	.8215	-	.8904	.9033	.9759	.9914	5	.8075	-	.8263	.8182	.8220	.8408
10	.8600	-	.8958	.9652	.9764	.9994	10	.8699	-	.8751	.8651	.8270	.8654
20	.8541	-	.8872	.9583	.9668	<b>1.000</b>	20	.8601	-	.8749	<b>.8844</b>	.8699	.8703
30	.8456	-	.8627	.9508	.9706	<b>1.000</b>	30	.8610	-	.8649	.8749	.8653	.8754
40	.8365	-	.8424	.9460	.9839	<b>1.000</b>	40	.8465	-	.8610	.8697	.8792	.8613
50	.8312	-	.8370	.9263	.9850	<b>1.000</b>	50	.8565	-	.8515	.8654	.8558	.8706
Max. Dimension	.8317	-	.8317	.9097	.9823	<b>1.000</b>	Max. Dimension	.8370	-	.8370	.8361	.8703	.8706
N. Classes - 1	.5016	.9011	.6965	.7826	.9214	.7237	N. Classes - 1	.5619	.7782	.6770	.7256	.8073	.6640

Table VIII. Results of dimensionality reduction experiments on movement\_libras with 3-NN (train - test)

	PCA	LDA	ANMM	DMLMJ	NCA	LMNN		PCA	LDA	ANMM	DMLMJ	NCA	LMNN
1	.1938	.3339	.2414	.2720	.3360	.2547	1	.1747	.3169	.2675	.2694	.2606	.2673
2	.2813	.5362	.4597	.4720	.6638	.5416	2	.2574	.4553	.4800	.4476	.6181	.5251
3	.5232	.6143	.6435	.6684	.7195	.6900	3	.5483	.4978	.6680	.6684	.6920	.6499
5	.6873	.7211	.7473	.7918	.8188	.8156	5	.7177	.5938	.7763	.7774	.7655	.8012
10	.7831	.8661	.8053	.8857	.8383	.8485	10	.8007	.7001	.8119	.8711	.8017	.8220
20	.7978	-	.7972	.8705	.8442	.8490	20	.8139	-	.8106	<b>.8829</b>	.8143	.8333
30	.7981	-	.7978	.8652	.8438	.8514	30	.8139	-	.8139	.8696	.8191	.8133
40	.7972	-	.7975	.8594	.8469	.8526	40	.8139	-	.8139	.8605	.8323	.8233
50	.7972	-	.7972	.8538	.8431	.8498	50	.8139	-	.8139	.8627	.8310	.8255
Max. Dimension	.7972	-	.7972	.8460	.8516	.8490	Max. Dimension	.8139	-	.8139	.8649	.8319	.8133
N. Classes - 1	.7932	.8685	.8061	<b>.8901</b>	.8398	.8438	N. Classes - 1	.8185	.6642	.8137	.8811	.8274	.8211

Table IX. Results of dimensionality reduction experiments on spambase with 3-NN (train - test)

	PCA	LDA	ANMM	DMLMJ	NCA	LMNN		PCA	LDA	ANMM	DMLMJ	NCA	LMNN
1	.8369	.9215	.8567	.6995	<b>.9420</b>	.9340	1	.8106	.8871	.8587	.6588	.9044	.8872
2	.8316	-	.8869	.7724	<b>.9420</b>	.9386	2	.8261	-	.8850	.7173	<b>.9197</b>	.8958
3	.8487	-	.8973	.8886	.9388	.9335	3	.8543	-	.9090	.8807	.9152	.9068
5	.8784	-	.9079	.9009	.9415	.9335	5	.8782	-	.9049	.8700	.9111	.9069
10	.8681	-	.9222	.9195	.9400	.9318	10	.8826	-	.9198	.9044	.9153	.9113
20	.8700	-	.9067	.9217	.9400	.9297	20	.8695	-	.9048	.8937	.9155	.9005
30	.8586	-	.8787	.8867	.9403	.9328	30	.8502	-	.8675	.8851	.9154	.9027
40	.8572	-	.8654	.8727	.9369	.9318	40	.8547	-	.8567	.8611	.9111	.9005
50	.8536	-	.8560	.8596	.9374	.9299	50	.8655	-	.8633	.8569	.9133	.9070
Max. Dimension	.8500	-	.8500	.8635	.9391	.9285	Max. Dimension	.8654	-	.8654	.8525	.9154	.9092
N. Classes - 1	.8369	.9215	.8567	.6995	<b>.9420</b>	.9340	N. Classes - 1	.8106	.8871	.8587	.6588	.9044	.8872

#### 4.3. Analysis of Results

4.3.1. *In-depth analysis.* Below we will describe the main details observed in the algorithms for the different experiments carried out.

- **NCA.** On the results obtained in the first experiment, we can clearly see that NCA is the one that has obtained the best results. This is partly due to the fact that the algorithms have been evaluated with nearest neighbors classifiers, and NCA was specifically designed to improve this classifier. NCA got the first place in most of the validations over the training set, showing its ability to fit to the data, but it has also obtained clear victories in many of the datasets over the test set, thus also demonstrating a great capacity for generalization.

- **LMNN and DMLMJ.** We can also see that DMLMJ and LMNN algorithms stand out, although at a considerable distance from NCA. These algorithms were also oriented to nearest neighbors classification, which justifies these good results. About LMNN we also conjecture that it has a slow convergence with the projected gradient method, and it could have achieved better results with a greater number of iterations. In fact, in the analysis of dimensionality reduction experiments we will see a much better performance of LMNN with the stochastic gradient descent method.
- **LSI.** LSI is another algorithm capable of obtaining very good results on certain datasets, but it is penalized by many others, where it is not able to optimize enough, even not converging in several datasets.
- **ITML and MCML.** ITML and MCML are two algorithms that, despite getting the best results in a very small number of cases, they get decent results in most datasets, resulting in quite a stable performance. ITML does not learn too much from the characteristics of the training set, but is able to generalize what has been learned in a quite effective way, being possibly the algorithm that loses the least accuracy over the test set, with respect to the training set. On the other hand, MCML has more learning capacity, even showing a slight overfitting, as its results are worse than those of many algorithms on the test set.
- **LDA.** Another algorithm in which we can see overfitting, perhaps more clearly, is LDA. This algorithm is capable of getting very good results on the training set, surpassing most of the algorithms, but it gets noticeably worse when evaluated on the test dataset. Recall that LDA is able to learn only a maximum dimension equal to the number of classes of the dataset minus one. This may be causing the loss of important information on many datasets by the projection it learns.
- **DML-eig and LDML.** Finally, although DML-eig and LDML are able to get better results than euclidean distance on the training sets, on several datasets they have obtained quite low quality results. On many of the test datasets, they are surpassed by the euclidean distance.
- **NCMML and NCMC.** If we analyze the results of the centroid-based classifiers, we can easily observe that in the vast majority of cases the classifier has worked much better after learning the distance with its associated learning algorithm, than using the euclidean distance. It can also be observed that the results are subject to great variability, depending on the number of centroids chosen. This shows that the choice of an adequate number of centroids, that adapts well to the disposition of the different classes, is fundamental to achieve a successful learning with these algorithms.
- **Kernel algorithms.** Focusing now on the kernel-based algorithms, it is interesting to note how KLMNN with laplacian kernel is able to adjust as much as possible to the data, getting a 100 % success rate on most of the datasets. This success rate is not transferred, in general, to the test data, showing that this algorithm overfits with laplacian kernel. We can also observe that the best results are distributed in a varied way among the different evaluated options. The choice of a suitable kernel that fits well with the disposition of the data is decisive for the performance of kernel-based algorithms.
- **Dimensionality reduction experiments.** To conclude our analysis, dimensionality experiments allow us to observe that the best results are not always obtained when considering the maximum dimension. This may be due to the fact that the algorithms are able to denoise the data, ensuring that the classifier used later does not overfit. We also see that we cannot reduce the dimension as much as we want, because at some point we start losing information, which happens in many cases with LDA, which is its great limitation. In general, we can observe that all algorithms improve their results by reducing dimensionality until a certain value,

although the best results are provided by LMNN, DMLMJ and NCA. The results obtained by LMNN open the possibility of using this algorithm with stochastic gradient descent, instead of the semidefinite programming algorithm used in the first experiment, since the results it provides are quite good. Although these algorithms have obtained better results, the use of ANMM and LDA (as long as the dimension allows it) is important for the estimation of an adequate dimension, since they are much more efficient than the first ones. As for PCA, it gets the worst results in low dimensions, probably due to not considering the information of the labels.

**4.3.2. Global analysis.** To complete the verbal analysis carried out, we have developed a series of Bayesian statistical tests to assess the extent to which the performance of the different algorithms analyzed outperforms the others. To do this, we have elaborated several pairwise Bayesian sign tests [Benavoli et al. 2017]. In these tests, we will consider the differences between the obtained scores of two algorithms, assuming that their prior distribution is a Dirichlet Process [Benavoli et al. 2014], defined by a prior strength  $s = 1$  and a prior pseudo-observation  $z_0 = 0$ . After considering the score observations obtained for each dataset, we obtain a posterior distribution which gives us the probabilities that one algorithm outperforms the other. We also introduce a rope (region of practically equivalent) region, in which we consider the algorithms to have an equivalent performance. We have designated the rope region as the one where the score differences are in the interval  $[-0.01, 0.01]$ . In summary, from the posterior distribution we obtain three probabilities: the probability that the first algorithm outperforms the second, the probability that the second algorithm outperforms the first one, and the probability that both algorithms are equivalent. These probabilities can be visualized in a simplex plot for a sample of the posterior distribution, in which a greater tendency of the points towards one of the regions will represent a greater probability.

To do the Bayesian sign tests, we have used the R package `rNPBST` [Carrasco et al. 2017]. In Figure 6 we pairwise compare some of the algorithms that seem to have a better performance in experiment 1 with 3-NN (NCA, DMLMJ and LMNN) with the results of the 3-NN classifier for euclidean distance. In the comparison between euclidean distance and NCA, we can clearly see that the points are concentrated close to the [NCA, rope] segment. This shows us that euclidean distance is unlikely to outperform NCA, and there is also a high probability for NCA to outperform euclidean distance, since a big concentration of points is in the NCA region. We obtain similar conclusions for DMLMJ against euclidean distance, although in this case, despite the fact that euclidean distance is still unlikely to win, there is a greater concentration of points in the rope region. In the comparison between LMNN and euclidean distance, we see a more centered concentration of points, but slightly weighted towards the LMNN region. In the comparisons between the distance metric learning algorithms we observe the points weighted to the [NCA, rope] segment, concluding the difficulty of outperforming NCA, and between DMLMJ and LMNN we can see a pretty level playing field, but slightly biased to the DMLMJ algorithm.

The outperforming of euclidean distance is even more clear in the results of experiment 2. For these algorithms, we can clearly observe the points concentrated in the region corresponding to the nearest centroid metric learning algorithm, as shown in Figure 7. We have elaborated more pairwise Bayesian sign tests for the rest of algorithms in experiment 1. The results of these tests can be found on the `pyDML-Stats` website<sup>5</sup>, a web page where the results of all these experiments have been stored.

<sup>5</sup>Source code: <https://github.com/jlsuarezdiaz/pyDML-Stats>. The current website is located at <https://jlsuarezdiaz.github.io/software/pyDML/stats/>

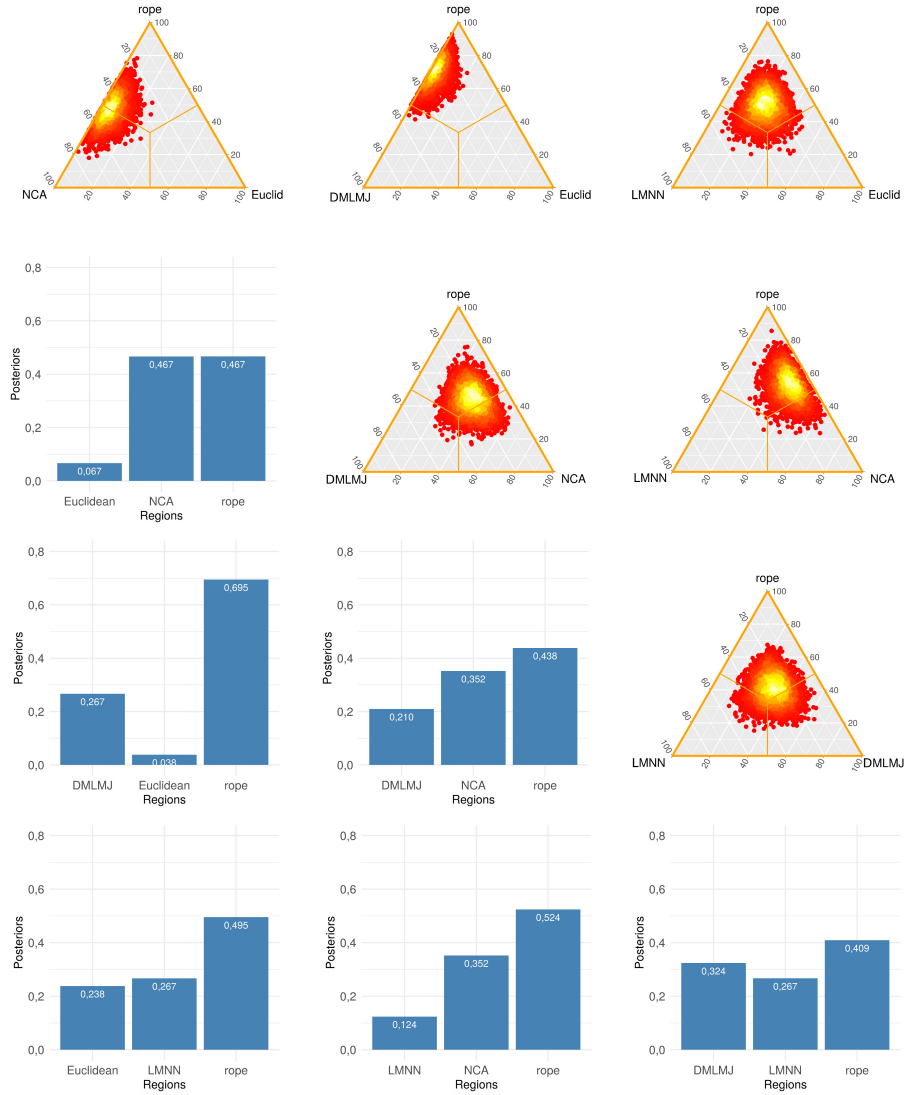


Fig. 6. Bayesian sign results for NCA, DMLMJ, LMNN and euclidean distance with 3-NN.

## 5. CONCLUSIONS

In this tutorial we have studied the concept of distance metric learning, showing what it is, what its applications are, how to design its algorithms, and the theoretical foundations of this discipline. We have also studied some of the most popular algorithms in this field, also with their theoretical foundations, and explaining different resolution techniques.

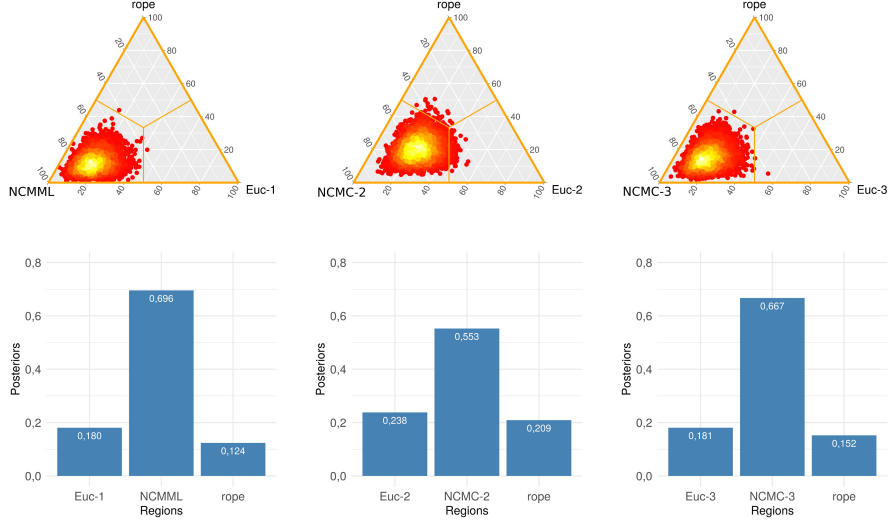


Fig. 7. Bayesian sign test results for the comparison between scores of nearest centroid classifiers with their corresponding distance metric learning algorithm against the same classifier with euclidean distance. The results are shown for nearest class mean classifier (left), nearest class with 2 centroids (center) and nearest class with 3 centroids (right).

In order to understand the theoretical foundations of distance metric learning and its algorithms, it has been necessary to delve into three different mathematical theories: convex analysis, matrix analysis and information theory. Convex analysis has made it possible to present many of the optimization problems studied in the algorithms, along with some methods for solving them. Matrix analysis has provided many useful tools to understand this discipline, from how to parameterize Mahalanobis distances, to the optimization with eigenvectors, going through the most basic algorithm of semidefinite programming. Finally, information theory has motivated several of the algorithms we have studied.

In addition, several experiments have been developed that have allowed to evaluate the performance of the algorithms analyzed in this study. The results of these experiments have allowed us to observe how algorithms such as LMNN, DMLMJ, and specially NCA can considerably improve the nearest neighbors classification, and how centroid-based distance learning algorithms also improve their corresponding classifiers. We have also seen the wide variety of possibilities offered by kernel-based algorithms, and the advantages that an appropriate reduction of the dimensionality of the datasets can offer.

## 6. FUTURE WORK

The work carried out in this paper has paved the way for future research in the study of distance metric learning. Some of these possibilities are:

- **Other approaches for the concept of distance.** Most of the current distance metric learning theory focus on Mahalanobis distances. However, some articles open a door to learning about other possible distances, such as local Mahalanobis distances, that lead to a multi-metric learning [Weinberger and Saul 2009]. Another approach

that has gained popularity in recent years is based on the use of neural networks to learn distances, which is being referred to as *deep metric learning* [Yi et al. 2014; Zhe et al. 2019; Cakir et al. 2019]. By developing new approaches, we will have a greater variety of distances to learn, and thus have a greater chance of success.

- **Kernelization of existing algorithms.** The kernelization of distance metric learning algorithms can be extended to other algorithms besides those presented. The search for a suitable parametrization and a representer theorem that allows the kernel trick to be applied is another possible task to carry out.

- **Other optimization mechanisms.** The algorithms studied optimize their objective functions by applying gradient descent method, regardless of whether the objective function is convex. The use of other optimization techniques, such as meta-heuristics, can be useful to improve those algorithms that do not have convex objective functions. The evolutionary approach to distance metric learning has been explored recently in several problems [Kalintha et al. 2017; Ali et al. 2018].

- **High dimensionality datasets.** Distance metric learning is of great interest in many real problems in high dimensionality, such as face recognition, where it is very useful to be able to measure the similarity between different images [Moutafis et al. 2017]. When we work with datasets of even greater dimensionality, the treatment of distances can become too expensive, since it would be necessary to store matrices of very large dimensions. In these situations it may be of interest to combine distance metric learning with feature selection techniques prepared for very high dimensional data [Tan et al. 2014].

- **Big Data solutions.** The problem of learning when the amount of data we have is huge and heterogeneous is one of the challenges of machine learning nowadays [Wu et al. 2014]. In the case of the distance metric learning algorithms, although many of them, specially those based on gradient descent, are quite slow and do not scale well with the number of samples, they can be largely parallelized in both matrix computations and gradient descent batches. In this way, distance metric learning can be extended to handle Big Data by developing specialized algorithms and integrating them with frameworks such as Spark [Meng et al. 2016] and Cloud Computing architectures [Hashem et al. 2015].

- **Singular problems.** In this paper, we have focused on distance metric learning for usual problems, like standard classification and dimensionality reduction, and we have also mentioned its applications for clustering and semi-supervised learning. However, distance metric learning can be useful in a wide variety of learning tasks [Charte et al. 2019], and can be carried out either by designing new algorithms or by adapting known algorithms from standard problems to these tasks. In recent years, several distance metric learning proposals have been made in problems like regression [Nguyen et al. 2016], multi-dimensional classification [Ma and Chen 2018], ordinal classification [Nguyen et al. 2018], multi-output learning [Liu et al. 2019] and even transfer learning [Luo et al. 2019].

- **Hybridization with other learning techniques.** Over the years, some distance-based algorithms, or some of their ideas or foundations, have been combined with other algorithms in order to improve their learning capabilities in certain problems. For example, the concept of nearest-neighbors has been combined with classifiers such as Naive-Bayes, obtaining a Naive-Bayes classifier whose feature distributions are determined by the nearest neighbors of each class [Yang and Tian 2012]; with neural networks, to find the best neural network architecture [Wang et al. 2017]; with random forests, by exploiting the relationship between voting points and potential nearest neighbors [Lin and Jeon 2006]; with deep learning, to provide interpretability and robustness to deep neural networks [Papernot and McDaniel 2018]; with ensemble methods, like bootstrap [Steele 2009; Hamamoto et al. 1997]; with support vector

machines, training them locally in neighborhoods [Zhang et al. 2006]; or with rule-learning algorithms, obtaining the so-called *nested generalized exemplar* algorithms [Wettschereck and Dietterich 1995]. The distances used in these combinations of algorithms can condition their performance, so designing appropriate distance learning algorithms for each of these tasks can be helpful for achieving good results. Following this topic, another option is to hybridize directly distance metric learning with other techniques, like ensemble learning [Mu et al. 2013].

## REFERENCES

- Charu C Aggarwal, Yuchen Zhao, and S Yu Philip. 2012. On text clustering with side information. In *Data Engineering (ICDE), 2012 IEEE 28th International Conference on*. IEEE, 894–904.
- Bassel Ali, Wasin Kalintha, Koichi Moriyama, Masayuki Numao, and Ken-ichi Fukui. 2018. Reinforcement learning for evolutionary distance metric learning systems improvement. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion*. ACM, 155–156.
- Aurélien Bellet, Amaury Habrard, and Marc Sebban. 2013. A survey on metric learning for feature vectors and structured data. *arXiv preprint arXiv:1306.6709* (2013).
- Alessio Benavoli, Giorgio Corani, Janez Demšar, and Marco Zaffalon. 2017. Time for a change: a tutorial for comparing multiple classifiers through Bayesian analysis. *The Journal of Machine Learning Research* 18, 1 (2017), 2653–2688.
- Alessio Benavoli, Giorgio Corani, Francesca Mangili, Marco Zaffalon, and Fabrizio Ruggeri. 2014. A Bayesian Wilcoxon signed-rank test based on the Dirichlet process. In *International conference on machine learning*. 1026–1034.
- Stephen Boyd and Jon Dattorro. 2003. *Alternating projections*. Technical Report. Stanford University. Lecture notes of EE392o, Autumn Quarter.
- Stephen Boyd and Lieven Vandenbergh. 2004. *Convex optimization*. Cambridge university press.
- Lev M Bregman. 1967. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR computational mathematics and mathematical physics* 7, 3 (1967), 200–217.
- Christopher JC Burges. 1998. A tutorial on support vector machines for pattern recognition. *Data mining and knowledge discovery* 2, 2 (1998), 121–167.
- Fatih Cakir, Kun He, Xide Xia, Brian Kulis, and Stan Sclaroff. 2019. Deep Metric Learning to Rank. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1861–1870.
- Tadeusz Caliński and Jerzy Harabasz. 1974. A dendrite method for cluster analysis. *Communications in Statistics-theory and Methods* 3, 1 (1974), 1–27.
- Jacinto Carrasco, Salvador García, María del Mar Rueda, and Francisco Herrera. 2017. rNPBST: An R Package Covering Non-parametric and Bayesian Statistical Tests. In *International Conference on Hybrid Artificial Intelligence Systems*. Springer, 281–292.
- David Charte, Francisco Charte, Salvador García, and Francisco Herrera. 2019. A snapshot on nonstandard supervised learning problems: taxonomy, relationships, problem transformations and algorithm adaptations. *Progress in Artificial Intelligence* in press (2019).
- Rathachai Chatpatanasiri, Teesid Korsrilabutr, Pasakorn Tangchanachaianan, and Boonserm Kijssirikul. 2010. A new kernelization framework for Mahalanobis distance learning algorithms. *Neurocomputing* 73, 10-12 (2010), 1570–1579.
- Thomas Cover and Peter Hart. 1967. Nearest neighbor pattern classification. *IEEE transactions on information theory* 13, 1 (1967), 21–27.
- Thomas M. Cover and Joy A. Thomas. 2006. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience.
- John P Cunningham and Zoubin Ghahramani. 2015. Linear dimensionality reduction: Survey, insights, and generalizations. *The Journal of Machine Learning Research* 16, 1 (2015), 2859–2900.
- Bernard Dacorogna. 2007. *Direct methods in the calculus of variations*. Vol. 78. Springer Science & Business Media. 67–148 pages.
- Jason V Davis and Inderjit S Dhillon. 2007. Differential entropic clustering of multivariate gaussians. In *Advances in Neural Information Processing Systems*. 337–344.



- Jason V Davis, Brian Kulis, Prateek Jain, Suvrit Sra, and Inderjit S Dhillon. 2007. Information-theoretic metric learning. In *Proceedings of the 24th international conference on Machine learning*. ACM, 209–216.
- Paramveer S Dhillon, Partha Pratim Talukdar, and Koby Crammer. 2010. Inference-driven metric learning for graph construction. In *4th North East Student Colloquium on Artificial Intelligence*.
- Ronald A Fisher. 1936. The use of multiple measurements in taxonomic problems. *Annals of eugenics* 7, 2 (1936), 179–188.
- Amir Globerson and Sam T Roweis. 2006. Metric learning by collapsing classes. In *Advances in neural information processing systems*. 451–458.
- Jacob Goldberger, Geoffrey E Hinton, Sam T Roweis, and Ruslan R Salakhutdinov. 2005. Neighbourhood components analysis. In *Advances in neural information processing systems*. 513–520.
- Matthieu Guillaumin, Jakob Verbeek, and Cordelia Schmid. 2009. Is that you? Metric learning approaches for face identification. In *Computer Vision, 2009 IEEE 12th international conference on*. 498–505.
- Yoshihiko Hamamoto, Shunji Uchimura, and Shingo Tomita. 1997. A bootstrap technique for nearest neighbor classifier design. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19, 1 (1997), 73–79.
- Ibrahim Abaker Targio Hashem, Ibrar Yaqoob, Nor Badrul Anuar, Salimah Mokhtar, Abdullah Gani, and Samee Ullah Khan. 2015. The rise of big data on cloud computing: Review and open research issues. *Information systems* 47 (2015), 98–115.
- Nicholas J Higham. 1988. Computing a nearest symmetric positive semidefinite matrix. *Linear algebra and its applications* 103 (1988), 103–118.
- Thomas Hofmann, Bernhard Schölkopf, and Alexander J Smola. 2008. Kernel methods in machine learning. *The annals of statistics* (2008), 1171–1220.
- Roger A Horn and Charles R Johnson. 1990. *Matrix analysis*. Cambridge university press.
- I.T. Jolliffe. 2002. *Principal Component Analysis*. Springer.
- Wasin Kalintha, Satoshi Ono, Masayuki Numao, and Ken-ichi Fukui. 2017. Kernelized Evolutionary Distance Metric Learning for Semi-Supervised Clustering.. In *AAAI*. 4945–4946.
- Effrosini Kokiopoulou, Jie Chen, and Yousef Saad. 2011. Trace optimization and eigenproblems in dimension reduction methods. *Numerical Linear Algebra with Applications* 18, 3 (2011), 565–602.
- Brian Kulis and others. 2013. Metric learning: A survey. *Foundations and Trends® in Machine Learning* 5, 4 (2013), 287–364.
- John A Lee and Michel Verleysen. 2007. *Nonlinear dimensionality reduction*. Springer Science & Business Media.
- J. Liang, Q. Hu, C. Dang, and W. Zuo. 2019. Weighted Graph Embedding-Based Metric Learning for Kinship Verification. *IEEE Transactions on Image Processing* 28, 3 (2019), 1149–1162.
- Yi Lin and Yongho Jeon. 2006. Random forests and adaptive nearest neighbors. *J. Amer. Statist. Assoc.* 101, 474 (2006), 578–590.
- W. Liu, D. Xu, I.W. Tsang, and W. Zhang. 2019. Metric Learning for Multi-Output Tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 2 (2019), 408–422.
- Y. Luo, Y. Wen, T. Liu, and D. Tao. 2019. Transferring Knowledge Fragments for Learning Distance Metric from a Heterogeneous Domain. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 4 (2019), 1013–1026.
- Zhongchen Ma and Songcan Chen. 2018. Multi-dimensional classification via a metric approach. *Neurocomputing* 275 (2018), 1121–1131.
- Xiangrui Meng, Joseph Bradley, Burak Yavuz, Evan Sparks, Shivaram Venkataraman, Davies Liu, Jeremy Freeman, DB Tsai, Manish Amde, Sean Owen, and others. 2016. Mllib: Machine learning in apache spark. *The Journal of Machine Learning Research* 17, 1 (2016), 1235–1241.
- Thomas Mensink, Jakob Verbeek, Florent Perronnin, and Gabriela Csurka. 2012. Metric learning for large scale image classification: Generalizing to new classes at near-zero cost. In *Computer Vision–ECCV 2012*. Springer, 488–501.
- Sebastian Mika, Gunnar Ratsch, Jason Weston, Bernhard Scholkopf, and Klaus-Robert Mullers. 1999. Fisher discriminant analysis with kernels. In *Neural networks for signal processing IX, 1999. Proceedings of the 1999 IEEE signal processing society workshop*. Ieee, 41–48.
- Panagiotis Moutafis, Mengjun Leng, and Ioannis A Kakadiaris. 2017. An overview and empirical comparison of distance metric learning methods. *IEEE transactions on cybernetics* 47, 3 (2017), 612–625.
- Yang Mu, Wei Ding, and Dacheng Tao. 2013. Local discriminative distance metrics ensemble learning. *Pattern Recognition* 46, 8 (2013), 2337–2349.

- B. Nguyen and B. De Baets. 2019. Kernel distance metric learning using pairwise constraints for person re-identification. *IEEE Transactions on Image Processing* 28, 2 (2019), 589–600.
- Bac Nguyen, Carlos Morell, and Bernard De Baets. 2016. Large-scale distance metric learning for k-nearest neighbors regression. *Neurocomputing* 214 (2016), 805–814.
- Bac Nguyen, Carlos Morell, and Bernard De Baets. 2017. Supervised distance metric learning through maximization of the Jeffrey divergence. *Pattern Recognition* 64 (2017), 215–225.
- Bac Nguyen, Carlos Morell, and Bernard De Baets. 2018. Distance metric learning for ordinal classification based on triplet constraints. *Knowledge-Based Systems* 142 (2018), 17–28.
- Nils J Nilsson. 1965. *Learning machines: foundations of trainable pattern-classifying systems*. McGraw-Hill.
- Michael L Overton. 1988. On minimizing the maximum eigenvalue of a symmetric matrix. *SIAM J. Matrix Anal. Appl.* 9, 2 (1988), 256–268.
- Nicolas Papernot and Patrick McDaniel. 2018. Deep k-nearest neighbors: Towards confident, interpretable and robust deep learning. *arXiv preprint arXiv:1803.04765* (2018).
- Ralph Tyrell Rockafellar. 2015. *Convex analysis*. Princeton university press.
- Walter Rudin. 1987. *Real and complex analysis*. Tata McGraw-Hill Education.
- Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. 1998. Nonlinear component analysis as a kernel eigenvalue problem. *Neural computation* 10, 5 (1998), 1299–1319.
- George S Sebestyen. 1962. *Decision-making processes in pattern recognition*. Macmillan.
- Shai Shalev-Shwartz and Shai Ben-David. 2014. *Understanding machine learning: From theory to algorithms*. Cambridge university press.
- Brian M Steele. 2009. Exact bootstrap k-nearest neighbor learners. *Machine Learning* 74, 3 (2009), 235–255.
- Mingkui Tan, Ivor W Tsang, and Li Wang. 2014. Towards ultrahigh dimensional feature selection for big data. *The Journal of Machine Learning Research* 15, 1 (2014), 1371–1429.
- Lorenzo Torresani and Kuang-chih Lee. 2007. Large margin component analysis. In *Advances in neural information processing systems*. 1385–1392.
- Isaac Triguero, Sergio González, Jose M Moyano, Salvador García, Jesús Alcalá-Fdez, Julián Luengo, Alberto Fernández, Maria José del Jesús, Luciano Sánchez, and Francisco Herrera. 2017. KEEL 3.0: an open source software for multi-stage analysis in data mining. *International Journal of Computational Intelligence Systems* 10, 1 (2017), 1238–1249.
- Laurens Van Der Maaten, Eric Postma, and Jaap Van den Herik. 2009. Dimensionality reduction: a comparative. *J Mach Learn Res* 10 (2009), 66–71.
- Fei Wang and Changshui Zhang. 2007. Feature extraction by maximizing the average neighborhood margin. In *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*. 1–8.
- Lin Wang, Bo Yang, Yuehui Chen, Xiaoqian Zhang, and Jeff Orchard. 2017. Improving neural-network classifiers using nearest neighbor partitioning. *IEEE transactions on neural networks and learning systems* 28, 10 (2017), 2255–2267.
- Kilian Q Weinberger and Lawrence K Saul. 2009. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research* 10, Feb (2009), 207–244.
- Dietrich Wettschereck and Thomas G Dietterich. 1995. An experimental comparison of the nearest-neighbor and nearest-hyperrectangle algorithms. *Machine learning* 19, 1 (1995), 5–27.
- Xindong Wu, Xingquan Zhu, Gong-Qing Wu, and Wei Ding. 2014. Data mining with big data. *IEEE transactions on knowledge and data engineering* 26, 1 (2014), 97–107.
- Eric P Xing, Michael I Jordan, Stuart J Russell, and Andrew Y Ng. 2003. Distance metric learning with application to clustering with side-information. In *Advances in neural information processing systems*. 521–528.
- Rui Xu and Donald Wunsch. 2005. Survey of clustering algorithms. *IEEE Transactions on neural networks* 16, 3 (2005), 645–678.
- Liu Yang and Rong Jin. 2006. Distance metric learning: A comprehensive survey. *Michigan State University* 2, 2 (2006), 4.
- Xiaodong Yang and Ying Li Tian. 2012. Eigenjoints-based action recognition using naive-bayes-nearest-neighbor. In *Computer vision and pattern recognition workshops (CVPRW), 2012 IEEE computer society conference on*. IEEE, 14–19.
- Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li. 2014. Deep metric learning for person re-identification. In *2014 22nd International Conference on Pattern Recognition*. IEEE, 34–39.
- Yiming Ying and Peng Li. 2012. Distance metric learning with eigenvalue optimization. *Journal of Machine Learning Research* 13, Jan (2012), 1–26.

- Hao Zhang, Alexander C Berg, Michael Maire, and Jitendra Malik. 2006. SVM-KNN: Discriminative nearest neighbor classification for visual category recognition. In *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, Vol. 2. IEEE, 2126–2136.
- Xuefei Zhe, Shifeng Chen, and Hong Yan. 2019. Directional statistics-based deep metric learning for image classification and retrieval. *Pattern Recognition* 93 (2019), 113–123.
- Xiaojin Zhu and Zoubin Ghahramani. 2002. *Learning from labeled and unlabeled data with label propagation*. Technical Report. Carnegie Mellon University.

## Appendix to: A Tutorial on Distance Metric Learning: Mathematical Foundations, Algorithms and Experiments

JUAN LUIS SUÁREZ, University of Granada  
SALVADOR GARCÍA, University of Granada  
FRANCISCO HERRERA, University of Granada

---

### A. MATHEMATICAL BACKGROUND

In this appendix we will study three mathematical blocks that make up the foundations of distance metric learning: convex analysis, matrix analysis and information theory.

#### A.1. Convex Analysis

Convex analysis is a fundamental field of study for many optimization problems. This field studies the convex sets, functions and problems. Convex functions have very useful properties in optimization tasks, and allow tools to be built to solve numerous types of convex optimization problems.

We will highlight some results of convex analysis in our work. First, we will show some important geometric properties of convex sets, such as the convex projection theorem, and then we will analyze some optimization methods that will be used later.

We start with the geometry of convex sets. We will work in the euclidean  $d$ -dimensional space,  $\mathbb{R}^d$ , where we note the dot product as  $\langle \cdot, \cdot \rangle$ .

*A.1.1. Convex Set Results.* Recall that convex sets are those for which any segment between two points in the set remains within the set, that is, a set  $K \subset \mathbb{R}^d$  is convex iff  $[x, y] = \{(1 - \lambda)x + \lambda y : \lambda \in [0, 1]\} \subset K$ , for every  $x, y \in K$ . An important result from convex sets states that, at every point on the border of a closed set, we can setup a hyperplane so that the convex set and the hyperplane intersect only at the boundary of the set, and the whole set lies on one side of the hyperplane. Furthermore, this property characterizes the closed convex sets with non empty interiors. This result is known as the supporting hyperplane theorem and we discuss it below.

*Definition A.1.*

Let  $T: \mathbb{R}^d \rightarrow \mathbb{R}$  be a linear map,  $\alpha \in \mathbb{R}$  and  $P = \{x \in \mathbb{R}^d : T(x) = \alpha\}$  be an hyperplane. Associated with  $P$ , we define  $P^+ = \{x \in \mathbb{R}^d : T(x) \geq \alpha\}$  and  $P^- = \{x \in \mathbb{R}^d : T(x) \leq \alpha\}$ .

We say that  $P$  is a *supporting hyperplane* for the set  $K \subset \mathbb{R}^d$  if  $P \cap \bar{K} \neq \emptyset$ , and either  $K \subset P^+$  or  $K \subset P^-$ . We refer to *supporting half-space* as the half-space that contains  $K$ , between  $P^+$  and  $P^-$ .

**THEOREM A.2 (SUPPORTING HYPERPLANE THEOREM).**

- (1) If  $K \subset \mathbb{R}^d$  is a closed convex set, then for each  $x_0 \in \partial K$  there is a supporting hyperplane  $P$  for  $K$  so that  $x_0 \in P$ .
- (2) Every proper closed convex set in  $\mathbb{R}^d$  is the intersection of all its supporting half-spaces.
- (3) Let  $K \subset \mathbb{R}^d$  be a closed set with non empty interior. Then,  $K$  is convex if and only if for every  $x \in \partial K$  there is a supporting hyperplane  $P$  for  $K$  with  $x \in P$ .

Proof of this result can be found in Dacorogna [2007] (chap. 2, theorem 2.7). We will use this theorem in the following results. The following property is fundamental to be

---

able to make sense of the optimization tools shown in this paper. We will see that, given a closed convex set and a point in  $\mathbb{R}^d$ , we can find a nearest point to the given point in the convex set, and it is unique, that is, there is a projection for the given point onto the convex set. In other words, projections onto convex sets are well defined. We prove this result below. We will see that projections will help us to deal with constrained convex problems.

**THEOREM A.3 (CONVEX PROJECTION).** *Let  $K \subset \mathbb{R}^d$  be a non empty closed convex set. Then, for every  $x \in \mathbb{R}^d$  there is a single point  $x_0 \in K$  with  $d(x, K) = d(x, x_0)$ , where we have defined the distance to the set  $K$  by*

$$d(x, K) = \inf\{d(x, y) : y \in K\}.$$

*The point  $x_0$  is called the projection of  $x$  onto  $K$  and it is usually denoted by  $P_K(x)$ . The function  $P_K: \mathbb{R}^d \rightarrow K$  given by the mapping  $x \mapsto P_K(x)$  is therefore well defined and it is called the projection onto  $K$ . In addition, for each  $x \in \mathbb{R}^d \setminus K$ , the half-space  $\{y \in \mathbb{R}^d : \langle x - P_K(x), y - P_K(x) \rangle \leq 0\}$  is a supporting half-space for  $K$  in  $P_K(x)$ .*

**PROOF.** First, we will prove the existence of a point in  $K$  in which the distance to  $K$  is achieved. In fact, this is true for every closed and not necessarily convex set. Let  $x \in \mathbb{R}^d$ . As  $K$  is closed, we can choose  $R > 0$  so that  $K \cap \bar{B}(x, R)$  is a compact and non empty set. We consider the distance to  $x$  in this set, that is, we define the map  $d_x: K \cap \bar{B}(x, R) \rightarrow \mathbb{R}_0^+$  by  $d_x(y) = d(x, y) = \|x - y\|$ .  $d_x$  is continuous and it is defined over a compact set, so it attains a minimum at a point  $x_0 \in K \cap \bar{B}(x, R)$ .

If we now take  $y \in K \cap \bar{B}(x, R)$ , we get  $d(x, y) = d_x(y) \geq d_x(x_0) = d(x, x_0)$ . On the other hand, if we take  $y \in K \setminus \bar{B}(x, R)$ , we get  $d(x, y) > R \geq d(x, x_0)$ . We have obtained that  $d(x, y) \geq d(x, x_0)$  for every  $y \in K$ , and therefore  $d(x, K) \geq d(x, x_0)$ . The remaining inequality is clear, since  $x_0 \in K$ , that is,  $x_0$  is the point we were looking for.

We will see now the uniqueness of the point found. Suppose that  $x_1, x_2 \in K$  verify that  $d(x, x_1) = d(x, K) = d(x, x_2)$ . We define  $x_0$  as the half point in the segment  $[x_1, x_2]$ . We have that  $x_0 \in K$ , since  $K$  is convex. Let us note that

$$\langle x_1 - x_2, x - x_0 \rangle = \langle x_1 - x_2, x - \frac{1}{2}(x_1 + x_2) \rangle = \frac{1}{2} \langle x_1 - x_2, 2x - x_1 - x_2 \rangle.$$

If we substitute  $x_1 - x_2 = (x - x_2) - (x - x_1)$  and  $2x - x_1 - x_2 = (x - x_2) + (x - x_1)$ , we obtain

$$\begin{aligned} \langle x_1 - x_2, x - x_0 \rangle &= \frac{1}{2} \langle (x - x_2) - (x - x_1), (x - x_2) + (x - x_1) \rangle \\ &= \frac{1}{2} (\|x - x_2\|^2 - \|x - x_1\|^2) \\ &= \frac{1}{2} (d(x, K)^2 - d(x, K)^2) = 0. \end{aligned}$$

Therefore, the vectors  $x_1 - x_2$  and  $x - x_0$  are orthogonal, and consequently so are  $x - x_0$  and  $x_0 - x_2 = (x_1 - x_2)/2$ . Applying Pythagorean theorem we have

$$d(x, K)^2 = \|x - x_2\|^2 = \|x - x_0\|^2 + \|x_0 - x_2\|^2 \geq \|x - x_0\|^2 \geq d(x, K)^2,$$

that is, the equality holds in the previous inequality. In particular, we obtain that  $\|x_0 - x_2\|^2 = 0$ , and then  $x_0 = x_2$ . Since  $x_0$  was the half point of  $[x_1, x_2]$  we conclude that  $x_1 = x_2$ , proving the uniqueness.

Finally we will prove the last assertion in the theorem. Let  $x \in \mathbb{R}^d \setminus K$  and suppose that there exists  $y \in K$  with  $\langle x - P_K(x), y - P_K(x) \rangle > 0$ . Since  $K$  is convex, the segment  $[y, P_K(x)]$  is contained in  $K$ , and therefore we have  $y_t = P_K(x) + t(y - P_K(x)) \in K$ , for

every  $t \in [0, 1]$ . We define the map  $f: [0, 1] \rightarrow \mathbb{R}$  by

$$\begin{aligned} f(t) &= \|y_t - x\|^2 = \|P_K(x) - x + t(y - P_K(x))\|^2 \\ &= \|P_K(x) - x\|^2 + 2t\langle P_K(x) - x, y - P_K(x) \rangle + t^2\|y - P_K(x)\|^2. \end{aligned}$$

$f$  is a polynomial in  $t$ , so it is differentiable, and

$$f'(0) = 2\langle P_K(x) - x, y - P_K(x) \rangle = -2\langle x - P_K(x), y - P_K(x) \rangle < 0.$$

Last expression implies that  $f$  is strictly decreasing in a neighborhood of 0, that is, there exists  $\varepsilon > 0$  so that  $\|y_t - x\|^2 < \|y_0 - x\|^2 = \|P_K(x) - x\|^2$ , for  $0 < t < \varepsilon$ , which results in a contradiction, since  $P_K(x)$  minimizes the distance to  $x$  in  $K$  and the points  $y_t$  lie on  $K$ .

□

**A.1.2. Optimization Methods.** In the following paragraphs we will discuss some of the optimization methods that we will use in distance metric learning algorithms. These algorithms will generally try to optimize (we will focus on minimizing without loss of generality) differentiable functions without constraints, or convex functions subject to convex constraints. For the first case, it is well known that the gradient of a differentiable function has the direction of the maximum slope in the function graph, thus by advancing small quantities in the negative gradient direction we manage to reduce the value of our objective function. This iterative method is usually called the gradient descent method. The adaptation rule for this method, for a differentiable function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , is given by  $x_{t+1} = x_t - \eta \nabla f(x_t)$ ,  $t \in \mathbb{N} \cup \{0\}$ , where  $\eta$  is the quantity we advance in the negative gradient direction, and it is called the *learning rate*. This value can be either constant or adapted according to the evaluations of the objective function. For the first option, the choice of a value of  $\eta$  that is too big or too small can lead to poor results. The second option needs to evaluate the objective function at each iteration, which can be computationally expensive.

Foundations of gradient descent are based on the following ideas. Let us consider an objective function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $x \in \mathbb{R}^d$  and  $v \in \mathbb{R}^d \setminus \{0\}$  an arbitrary direction. We also consider the function  $g: \mathbb{R} \rightarrow \mathbb{R}$  given by  $g(\eta) = f(x + \eta v / \|v\|)$ . The rate of change or directional derivative of  $f$  at  $x$  in the direction of  $v$  is given by  $g'(0) = \langle \nabla f(x), v \rangle / \|v\|$ . Applying Cauchy-Schwarz inequality, we have

$$-\|\nabla f(x)\| \leq \frac{1}{\|v\|} \langle \nabla f(x), v \rangle \leq \|\nabla f(x)\|,$$

and equality in the left inequality holds when  $v = -\nabla f(x)$ , thus obtaining the maximum descent rate. In the same way, the maximum ascent rate is achieved when  $v = \nabla f(x)$ .

If gradient at  $x$  is non zero and we consider the first order Taylor approximation with the points  $x$  and  $x - \eta \nabla f(x)$ , we have that

$$f(x - \eta \nabla f(x)) = f(x) - \eta \|\nabla f(x)\|^2 + o(\eta),$$

with  $\lim_{\eta \rightarrow 0} |o(\eta)|/\eta = 0$ , then there is  $\varepsilon > 0$  so that if  $0 < \delta < \varepsilon$ , we have

$$\frac{o(\delta)}{\delta} < \|\nabla f(x)\|,$$

and therefore

$$f(x - \delta \nabla f(x)) - f(x) = \delta \left( -\|\nabla f(x)\|^2 + \frac{o(\delta)}{\delta} \right) < \delta(-\|\nabla f(x)\|^2 + \|\nabla f(x)\|^2) = 0,$$

thus  $f(x - \delta \nabla f(x)) < f(x)$  for  $0 < \delta < \varepsilon$ , so we are guaranteed that for an accurate learning rate the gradient method performs a descent at each iteration. Let us observe that the gradient direction is not the only valid descent direction, but the above calculations are still true for any direction  $v \in \mathbb{R}^d$  with  $\langle \nabla f(x), v \rangle < 0$ . The choice of different descent directions, even if they are not the maximum slope direction, may provide better results in certain situations.

Now we will discuss the constrained convex optimization problems. When we work with constrained problems, gradient descent method cannot be applied directly, as the gradient descent adaptation rule,  $x_{t+1} = x_t - \eta \nabla f(x_t)$ , does not guarantee  $x_{t+1}$  to be a feasible point, that is, a point that fulfills all the constraints. When the optimization problem is convex, the set determined by the constraints is closed and convex, so we can take projections onto this feasible set. The projected gradient method tries to fix the gradient descent problem by adding a projection onto the feasible set in the gradient descent adaptation rule, that is, if  $C$  is the feasible set, and  $P_C$  is the projection onto this set, the projected gradient adaptation rule becomes  $x_{t+1} = P_C(x_t - \eta \nabla f(x_t))$ . To confirm that this method is successful, we have to show that the direction  $v = P_C(x - \eta \nabla f(x)) - x$  is a descent direction, which is attained, thanks to the reasons given above, if  $\langle \nabla f(x), v \rangle < 0$ .

We name  $x_1 = x - \eta \nabla f(x)$ . Then,  $v = P_C(x_1) - x$ . Note that  $\langle \nabla f(x), v \rangle < 0 \iff \langle x_1 - x, P_C(x_1) - x \rangle = -\eta \langle \nabla f(x), v \rangle > 0$ . If gradient is not null and  $x_1 \in C$ , we get  $\langle x_1 - x, P_C(x_1) - x \rangle = \langle x_1 - x, x_1 - x \rangle = \|x_1 - x\|^2 > 0$ . If  $x_1 \notin C$ , then the convex projection theorem (Theorem A.3) ensures that the half-space  $H = \{y \in \mathbb{R}^d : \langle x_1 - P_C(x_1), y - P_C(x_1) \rangle \leq 0\}$  contains  $C$ . In particular,

$$0 \geq \langle x_1 - P_C(x_1), x - P_C(x_1) \rangle = \langle x_1 - x, x - P_C(x_1) \rangle + \|x - P_C(x_1)\|^2.$$

Consequently,  $\langle x_1 - x, P_C(x_1) - x \rangle \geq \|x - P_C(x_1)\|^2 \geq 0$ . In addition, equality holds if and only if  $x = P_C(x_1)$ , in which case the iterative algorithm will have converged (observe that this happens when  $x \in \partial C$  and the gradient descent direction points out of  $C$  and orthogonally to the supporting hyperplane). Therefore, as long as the projected gradient iterations produce changes in the obtained points, an appropriate learning rate will ensure the descent in the objective function. Figure 8 visually compares the gradient descent method and the projected gradient method.

Another problem we can find when trying to optimize constrained convex problems is that we may have multiple constraints, but we only know the projection onto each single restriction, without knowing the projection onto the intersection, which makes up the feasible set. In these cases, a popular method to find a point in the intersection is the so-called *iterated projections method*, which consists of taking successive projections onto each constraint set, and repeating this procedure cyclically. We will analyze the simplest case, that is, let us suppose that we have a feasible set determined by two convex constraints. The following theorem states that, if the intersection of the sets determined by each constraint is not empty, then the sequence of iterated projections converge to a point in the intersection.

**THEOREM A.4 (CONVERGENCE OF THE ITERATED PROJECTIONS METHOD).** *Let  $C, D \subset \mathbb{R}^d$  be closed convex sets, and let  $P_C, P_D : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be the projections onto  $C$  and  $D$ , respectively. Suppose that  $x_0 \in C$  and we build the sequences  $\{x_n\}$  and  $\{y_n\}$  given by  $y_n = P_D(x_n)$  and  $x_{n+1} = P_C(y_n)$ , for each  $n \in \mathbb{N} \cup \{0\}$ .*

*Then, if  $C \cap D \neq \emptyset$ , both sequences converge to a point  $x^* \in C \cap D$ .*

Proof of this result is provided by Boyd and Dattorro [2003]. The extension to the general case can be made following a similar argument, and it is discussed by Bregman [1967]. That is why the general case is also called the *Bregman projections method*. Figure 9 shows a graphical example of the iterated projections method.

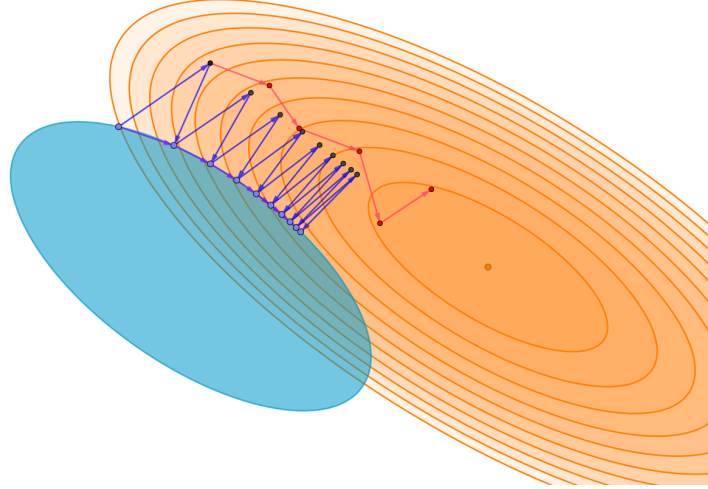


Fig. 8. Orange shaded areas represent contour lines of the function  $f(x, y) = 2(x + y)^2 + 2y^2$  for natural values between 0 and 10. The red path shows the behaviour of the unconstrained gradient descent method applied to  $f$ . The blue path shows the behaviour of the projected gradient descent, with the blue ellipse as the feasible set. In both cases we observe that we are obtaining descent directions.

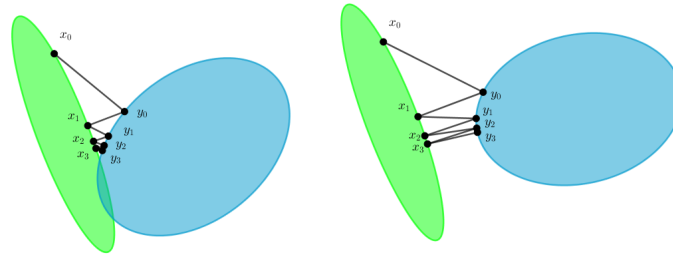


Fig. 9. The iterated projections method. The second image shows how the algorithm works if the sets do not intersect.

To conclude this section, we need to make a last remark. Recall that a convex and differentiable function  $f: \Omega \rightarrow \mathbb{R}$  defined on a convex open set verifies that  $f(x) \geq \langle \nabla f(x_0), x - x_0 \rangle$ , for every  $x, x_0 \in \Omega$ . Let  $x_0 \in \Omega$  be fix. When we work with convex but non differentiable functions, there are still vectors  $v \in \mathbb{R}^d$  for which  $f(x) \geq f(x_0) + \langle v, x - x_0 \rangle$ , for every  $x \in \Omega$ . This is a consequence of the supporting hyperplane theorem applied to the epigraph of  $f$  (recall that  $f$  is convex iff its epigraph is too). In this case we say that  $v$  is a *subgradient* of  $f$  at  $x_0$  and we note it as  $\partial f(x_0)$ , or  $\partial f(x_0)/\partial x$ , if we need to specify the variable.

Subgradients and gradients have similar behaviours, although we cannot always guarantee that subgradients are descent directions. Subgradient methods work in a similar way to gradient methods, replacing the gradient in the adaptation rule by a



subgradient. In subgradient methods it is useful to keep track of the best value obtained, as some subgradients may not be descent directions. In the situations we will handle, subgradient computations are easy: if  $f$  is differentiable at  $x_0$ , then  $\nabla f(x_0)$  is a subgradient (in fact, this is the only subgradient at  $x_0$ ); if  $f$  is a maximum of convex differentiable functions, then a subgradient at  $x_0$  is the gradient of any of the differentiable functions that attains the maximum at  $x_0$ .

## A.2. Matrix Analysis

In distance metric learning, matrices will play a key role, as they will be the structure over which distances will be defined and over which the optimization methods studied in the previous section will be applied. Within the set of all matrices, positive semidefinite matrices will be of even greater importance, so, in order to better understand the learning problems we will be dealing with, it will be necessary to delve into some of their numerous properties.

This section examines in depth the study of matrices, on a basis of the best-known results of diagonalization in linear algebra. From this basis, we will show how to give the set of matrices a Hilbert space structure, in order to be able to apply the convexity results and optimization methods from the previous chapter. In particular, we will be interested in how to obtain projections onto the set of positive semidefinite matrices. Also related to positive semidefinite matrices, we will present several results regarding decomposition that we will need in future sections. Finally, we will study some matrix optimization problems that can be solved via eigenvectors. Table X shows the notations we will use for matrices. We will restrict the study to the real case, since the problem we will deal with is in real variables, although many of the results we will see can be extended to the complex case.

Table X. Matrices notations.

Notation	Concept
$\mathcal{M}_{d' \times d}(\mathbb{R})$	Matrices of order $d' \times d$ .
$\mathcal{M}_d(\mathbb{R})$	Square matrices of order $d$ .
$A_{ij}$	The value of the matrix $A$ at the $i$ -th row and $j$ -th column.
$A_{.j}$ (resp. $A_{i.}$ )	The $j$ -th column (resp. the $i$ -th row) of the matrix $A$ .
$v = (v_1, \dots, v_d)$	A vector in $\mathbb{R}^d$ . Vectors will be treated as column matrices.
$A^T$	The transpose of the matrix $A$ .
$S_d(\mathbb{R})$	Symmetric matrices of order $d$ .
$GL_d(\mathbb{R})$	Invertible matrices of order $d$ .
$r(A), \text{tr}(A), \det(A)$	The rank, trace and determinant of the matrix $A$ .
$O_d(\mathbb{R})$	Orthogonal matrices of order $d$ .
$S_d(\mathbb{R})_0^+$	Positive semidefinite matrices of order $d$ .
$S_d(\mathbb{R})^+$	Positive definite matrices of order $d$ .
$S_d(\mathbb{R})_0^-$	Negative semidefinite matrices of order $d$ .
$S_d(\mathbb{R})^-$	Negative definite matrices of order $d$ .

**A.2.1. Matrices as a Hilbert Space. Projections.** Over the set of matrices we have defined a sum operation, and a matrix product, between matrices of orders  $d \times r$  and  $r \times n$ . When working with square matrices, this sum and product give the matrix set a non-commutative ring structure. These operations only allow us to obtain algebraic properties of matrices, but we also want to obtain geometric and topological properties. That is why we need to introduce a matrix inner product. We will introduce this inner product in the simplest way, that is, we will view matrices as vectors where we add the matrix rows one after the other, and we will consider the usual vector inner product. This matrix product is known as *Frobenius inner product*.

**Definition A.5.** We define the *Frobenius inner product* over the matrices space of order  $d' \times d$  as the mapping  $\langle \cdot, \cdot \rangle_F: \mathcal{M}_{d' \times d}(\mathbb{R}) \times \mathcal{M}_{d' \times d}(\mathbb{R}) \rightarrow \mathbb{R}$  given by

$$\langle A, B \rangle_F = \sum_{i=1}^{d'} \sum_{j=1}^d A_{ij} B_{ij} = \text{tr}(A^T B).$$

We define the *Frobenius norm* over the matrices space of order  $d' \times d$  as the mapping  $\| \cdot \|_F: \mathcal{M}_{d' \times d}(\mathbb{R}) \rightarrow \mathbb{R}_0^+$  given by

$$\|A\|_F = \sqrt{\langle A, A \rangle} = \sqrt{\sum_{i=1}^{d'} \sum_{j=1}^d A_{ij}^2} = \sqrt{\text{tr}(A^T A)}.$$

Frobenius norm is therefore identical to the euclidean norm in  $\mathbb{R}^{d' \times d}$  identifying matrices with vectors as mentioned before. Viewing this norm as a matrix norm, we have to remark that Frobenius norm is sub-multiplicative, but it is not induced by any vector norm. Some interesting properties about Frobenius norm can be deduced from the definition. They are listed below.

**PROPOSITION A.6.**

- (1) For each  $A \in \mathcal{M}_{d' \times d}(\mathbb{R})$ ,  $\|A\|_F = \|A^T\|_F$ .
- (2) For each  $A \in \mathcal{M}_{d' \times d}(\mathbb{R})$ ,  $\|A\|_F = \sqrt{\text{tr}(AA^T)}$ .
- (3) If  $U \in O_d(\mathbb{R})$ ,  $V \in O_{d'}(\mathbb{R})$  and  $A \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , then  $\|AU\|_F = \|VA\|_F = \|VAU\|_F = \|A\|_F$ .
- (4) If  $A \in S_d(\mathbb{R})$ , then  $\|A\|_F^2 = \sum_{i=1}^d \lambda_i^2$ , where  $\lambda_1, \dots, \lambda_d$  are the eigenvalues of  $A$ .

With Frobenius inner product we can apply the convex analysis theory studied in the previous section. The positive semidefinite matrix set has a convex cone structure, that is, it is closed under non-negative linear combinations. That is why  $S_d(\mathbb{R})_0^+$  is usually called the positive semidefinite cone. Under the topology induced by symmetric matrices, we can also see that  $S_d(\mathbb{R})_0^+$  is closed, as it is the intersection of closed sets:

$$S_d(\mathbb{R})_0^+ = \{M \in S_d(\mathbb{R}) : x^T M x \geq 0 \forall x \in \mathbb{R}^d\} = \bigcap_{x \in \mathbb{R}^d} \{M \in S_d(\mathbb{R}) : x^T M x \geq 0\}.$$

So we understand, in particular, that  $S_d(\mathbb{R})_0^+$  is a closed convex set over symmetric matrices, and thus we have a well-defined projection onto the positive semidefinite cone. This property is very important for many of the optimization problems we will study, since they will try to optimize functions defined over the positive semidefinite cone. Here, the projected gradient descent method will be of great use, thus constituting one of the most basic algorithms of the paradigm of semidefinite programming. We can calculate the projection onto the positive semidefinite cone explicitly, as we will see below.

**Definition A.7.** Let  $\Sigma \in \mathcal{M}_d(\mathbb{R})$  be a diagonal matrix,  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_d)$ . We define the *positive part* of  $\Sigma$  as  $\Sigma^+ = \text{diag}(\sigma_1^+, \dots, \sigma_d^+)$ , where  $\sigma_i^+ = \max\{\sigma_i, 0\}$ . In a similar way, we define its *negative part* as  $\Sigma^- = \text{diag}(\sigma_1^-, \dots, \sigma_d^-)$ , where  $\sigma_i^- = \max\{-\sigma_i, 0\}$ .

Let  $A \in S_d(\mathbb{R})$  and let  $A = UDU^T$  be a spectral decomposition of  $A$ . We define the *positive part* of  $A$  as  $A^+ = UD^+U^T$ . In a similar way we define its *negative part* as  $A^- = UD^-U^T$ .

**THEOREM A.8 (SEMIDEFINITE PROJECTION).** Let  $A \in S_d(\mathbb{R})$ . Then,  $A^+$  is the projection of  $A$  onto the positive semidefinite cone.

This result has been proven by Higham [1988], and is extended easily to project any square matrix onto the positive semidefinite cone, as we show below, although the most interesting case is that mentioned previously.

**COROLLARY A.9.** *Let  $A \in \mathcal{M}_d(\mathbb{R})$ . Then, the projection of  $A$  onto the positive semidefinite cone is given by  $((A + A^T)/2)^+$ .*

**A.2.2. Decomposition Theorems.** The positive semidefinite cone allows many of the concepts and properties that we already know about the non negative real numbers to be generalized. For example, we can similarly define concepts as the square roots, and modules or absolute values. These concepts play an important role in elaborating numerous decomposition theorems that involve positive semidefinite matrices. We will use these tools in order to prove a specific decomposition theorem that will motivate the ways of modeling the distance metric learning problem. The statement of this theorem is shown below.

**THEOREM A.10.** *Let  $M \in S_d(\mathbb{R})_0^+$ . Then,*

- (1) *There is a matrix  $L \in \mathcal{M}_d(\mathbb{R})$  so that  $M = L^T L$ .*
- (2) *If  $K \in \mathcal{M}_d(\mathbb{R})$  is any other matrix with  $M = K^T K$ , then  $K = UL$ , where  $U \in O_d(\mathbb{R})$  (that is,  $L$  is unique up to isometries).*

To prove this theorem, we will start with a characterization of positive semidefinite matrices by decomposition, which will also allow us to introduce the concept of square root. We will rely on several previous lemmas.

**LEMMA A.11.** *Let  $A, B \in \mathcal{M}_d(\mathbb{R})$  be two commuting matrices, that is,  $AB = BA$ . Then,  $Ap(B) = p(B)A$ , where  $p$  denotes any polynomial over matrices (that is, a expression of the form  $p(C) = a_0I + a_1C + a_2C^2 + \dots + a_nC^n$ , with  $a_0, \dots, a_n \in \mathbb{R}$ ).*

**PROOF.** Observe that

$$AB^n = (AB)B^{n-1} = B(AB)B^{n-2} = \dots = B^{n-1}(AB) = B^n A,$$

and  $Ap(B) = p(B)A$  is deduced by linearity.  $\square$

**LEMMA A.12.** *Let  $D \in S_d(\mathbb{R})_0^+$  be a diagonal matrix. Then, there is a polynomial over matrices  $p$  so that  $p(D^2) = D$ .*

**PROOF.** Suppose  $D = \text{diag}(\lambda_1, \dots, \lambda_d)$ , with  $0 \leq \lambda_1 \leq \dots \leq \lambda_d$ . Then,  $D^2 = \text{diag}(\lambda_1^2, \dots, \lambda_d^2)$ . We take  $p$  as an interpolation polynomial over the points  $(\lambda_i^2, \lambda_i)$ , for  $i = 1, \dots, d$ . If we evaluate it on  $D^2$  we obtain

$$p(D^2) = p(\text{diag}(\lambda_1^2, \dots, \lambda_d^2)) = \text{diag}(p(\lambda_1^2), \dots, p(\lambda_d^2)) = \text{diag}(\lambda_1, \dots, \lambda_d) = D.$$

$\square$

**THEOREM A.13.** *Let  $M \in \mathcal{M}_d(\mathbb{R})$ . Then,*

- (1)  *$M \in S_d(\mathbb{R})_0^+$  if, and only if, there is  $L \in \mathcal{M}_d(\mathbb{R})$  so that  $M = L^T L$ .*
- (2) *If  $M \in S_d(\mathbb{R})_0^+$ , there is a single matrix  $N \in S_d(\mathbb{R})_0^+$  with  $N^2 = M$ . In addition,  $M \in S_d(\mathbb{R})^+ \iff N \in S_d(\mathbb{R})^+$ .*

**PROOF.** , First we will see that  $L^T L$  is a positive semidefinite, for any  $L \in \mathcal{M}_d(\mathbb{R})$ . Indeed, given  $x \in \mathbb{R}^d$ ,

$$x^T L^T L x = (Lx)^T (Lx) = \|Lx\|_2^2 \geq 0.$$

We will prove the second implication of the first statement finding directly the matrix  $N$  of the second statement. Consider the spectral decomposition  $M = U D U^T$ , with

$U \in O_d(\mathbb{R})$  and  $D = \text{diag}(\lambda_1, \dots, \lambda_d)$ , with  $0 \leq \lambda_1 \leq \dots \leq \lambda_d$  the eigenvalues of  $M$ . We define  $D^{1/2} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_d})$  and construct the matrix  $N = UD^{1/2}U^T$ .  $N$  is positive semidefinite, because its eigenvalues are those of  $D^{1/2}$ , which are all positive, and besides,

$$N^2 = UD^{1/2}U^TUD^{1/2}U^T = UD^{1/2}D^{1/2}U^T = UDU^T = M.$$

Furthermore, the strict positivity of the eigenvalues of  $M$  is equivalent to that of the eigenvalues of  $N$ , then  $M \in \mathcal{M}_d(\mathbb{R})^+ \iff N \in \mathcal{M}_d(\mathbb{R})^+$ . Let us finally see that  $N$  is unique.

Suppose that we have  $N_1, N_2 \in S_d(\mathbb{R})_0^+$  with  $N_1^2 = M = N_2^2$ . Observe that  $N_1$  and  $N_2$  must have the same eigenvalues, since they are necessarily the positive square roots of the eigenvalues of  $M$ . Therefore,  $N_1$  and  $N_2$  are similar to a same diagonal matrix, that is, there are matrices  $U, V \in O_d(\mathbb{R})$  with  $N_1 = UDU^T$  and  $N_2 = VDV^T$ . From  $N_1^2 = N_2^2$  we have

$$UD^2U^T = VD^2V^T \implies V^TUD^2 = D^2V^TU,$$

so for  $W = V^TU \in O_d(\mathbb{R})$  we obtain that  $D^2$  and  $W$  commute. Combining Lemmas A.12 and A.11, we obtain that  $D$  and  $W$  also commute. Therefore,

$$WD = DW \implies V^TUD = DV^TU \implies UDU^T = VDV^T \implies N_1 = N_2,$$

obtaining the uniqueness.  $\square$

As we had anticipated, this theorem motivates the definition of square roots for positive semidefinite matrices.

**Definition A.14.** Let  $M \in S_d(\mathbb{R})_0^+$ . We define the *square root* of  $M$  as the unique matrix  $N \in S_d(\mathbb{R})_0^+$  with  $N^2 = M$ . We denote it as  $N = M^{1/2}$ .

We can also extend other concepts defined over the non-negative real numbers to the positive semidefinite matrices. For example, the square root allows us to define the concept of module for any matrix.

**Definition A.15.** Let  $A \in \mathcal{M}_{d \times d'}(\mathbb{R})$ . We define the *module* of  $A$  as

$$|A| = (A^T A)^{1/2} \in S_{d'}(\mathbb{R})_0^+.$$

With the module we can state a polar decomposition theorem, which shows a decomposition that can be seen as an extension of the polar form for complex numbers.

**THEOREM A.16 (POLAR DECOMPOSITION).** Let  $A \in \mathcal{M}_{d \times d'}(\mathbb{R})$ , with  $d' \leq d$ . Then, there is a matrix  $U \in \mathcal{M}_{d \times d'}(\mathbb{R})$  with  $U^T U = I$ , so that  $A = U|A|$ . This decomposition is called the polar decomposition of  $A$ , and it is not necessarily unique, unless  $A$  is square and invertible.

**PROOF.** First, observe that, given  $x \in \mathbb{R}^{d'}$ , we have

$$\|Ax\|_2^2 = (Ax)^T(Ax) = x^T A^T A x = x^T |A|^2 x = x^T |A| |A| x = (|A|x)^T(|A|x) = \||A|x\|_2^2.$$

This means that  $A$  and  $|A|$  have the same effect on the length of any vector. As an immediate consequence, we can observe that  $\ker A = \ker |A|$ , since

$$x \in \ker A \iff Ax = 0 \iff \|Ax\| = 0 \iff \||A|x\| = 0 \iff |A|x = 0 \iff x \in \ker |A|.$$

As  $d' = \dim \ker A + \dim \text{im } A = \dim \ker |A| + \dim \text{im } |A|$ , we also conclude that  $\dim \text{im } A = \dim \text{im } |A|$ , and then  $r(A) = r(|A|)$ . We will denote this rank as  $r \leq d$ .

$|A|$  is positive semidefinite, so there is an orthonormal basis  $\{w_1, \dots, w_{d'}\} \subset \mathbb{R}^{d'}$  consisting of eigenvectors of  $|A|$ , with corresponding eigenvalues  $\lambda_1, \dots, \lambda_{d'}$ . We can

assume that  $\lambda_1, \dots, \lambda_r > 0$  and  $\lambda_{r+1} = \dots = \lambda_{d'} = 0$ , or equivalently,  $\{w_{r+1}, \dots, w_{d'}\}$  is an orthonormal basis of  $\ker |A| = \ker A$ .

We consider the set of vectors  $\{Aw_1/\lambda_1, \dots, Aw_r/\lambda_r\}$ . Note that

$$\left\langle \frac{1}{\lambda_i} Aw_i, \frac{1}{\lambda_j} Aw_j \right\rangle = \frac{1}{\lambda_i \lambda_j} \langle Aw_i, Aw_j \rangle = \frac{1}{\lambda_i \lambda_j} w_i^T |A|^2 w_j = \frac{1}{\lambda_i \lambda_j} w_i^T \lambda_j^2 w_j = \frac{\lambda_j}{\lambda_i} w_i^T w_j,$$

which equals 1 if  $i = j$ , and 0 otherwise, so this set is also orthonormal. In fact, this set is an orthonormal basis of  $\text{im } A$ .

We extend the previous set to an orthonormal set of size  $d'$  in  $\mathbb{R}^d$ ,

$$\left\{ \frac{1}{\lambda_1} Aw_1, \dots, \frac{1}{\lambda_r} Aw_r, v_{r+1}, \dots, v_{d'} \right\}.$$

Finally, we construct the matrix  $V \in \mathcal{M}_{d \times d'}(\mathbb{R})$  by adding as columns the vectors in the previous set, and the matrix  $W \in \mathcal{M}_{d'}(\mathbb{R})$  by adding as rows the vectors  $w_1, \dots, w_{d'}$ . We define  $U$  as  $U = VW \in \mathcal{M}_{d \times d'}(\mathbb{R})$ . Observe that both  $V$  and  $W$  have orthonormal columns, and then  $V^T V = I = W^T W$ , obtaining that  $U^T U = I$  as well. We can also observe that  $W w_i = e_i$ , where  $\{e_1, \dots, e_{d'}\}$  is the canonical basis of  $\mathbb{R}^{d'}$ . Therefore, we obtain

$$U w_i = \begin{cases} \frac{1}{\lambda_i} w_i, & 1 \leq i \leq r \\ v_i, & r < i \leq d' \end{cases},$$

and finally,

$$U|A|w_i = \begin{cases} \lambda_i U w_i, & 1 \leq i \leq r \\ 0, & r < i \leq d' \end{cases} = \begin{cases} A w_i, & 1 \leq i \leq r \\ 0, & r < i \leq d' \end{cases} = A w_i,$$

where the last equality holds, since  $\{w_{r+1}, \dots, w_{d'}\} \subset \ker A$ . So, we have the equality  $A = U|A|$  on the basis  $\{w_1, \dots, w_{d'}\}$ , concluding the proof. The uniqueness of  $U$  when  $A$  is square and invertible is due to the fact that  $|A|$  is also invertible in that case, and then  $U = A|A|^{-1}$ .

□

**Remark A.17.** When  $A \in \mathcal{M}_d(\mathbb{R})$  is a square matrix, the polar decomposition can be stated as  $A = U|A|$ , where  $U \in O_d(\mathbb{R})$  is an orthogonal matrix.

We are now in a position to prove Theorem A.10. We recall its statement below.

**Theorem** Let  $M \in S_d(\mathbb{R})_0^+$ . Then,

- (1) There is a matrix  $L \in \mathcal{M}_d(\mathbb{R})$  so that  $M = L^T L$ .
- (2) If  $K \in \mathcal{M}_d(\mathbb{R})$  is any other matrix with  $M = K^T K$ , then  $K = UL$ , where  $U \in O_d(\mathbb{R})$  (that is,  $L$  is unique up to isometries).

**PROOF.** The first statement was proved in Theorem A.13. Suppose then that  $L, K \in \mathcal{M}_d(\mathbb{R})$  verify that  $M = L^T L = K^T K$ . Let  $L = V|L|$ ,  $K = W|K|$ , with  $V, W \in O_d(\mathbb{R})$ , be polar decompositions of  $L$  and  $K$ . Then, we have

$$\begin{aligned} L^T L = K^T K &\implies |L|^T V^T V |L| = |K|^T W^T W |K| \\ &\implies |L|^T |L| = |K|^T |K| \implies |L|^2 = |K|^2. \end{aligned}$$

As  $|L|$  and  $|K|$  are positive semidefinite, they must be the only square root of  $|L|^2 = |K|^2$ , that is,  $|L| = |K|$ . We call  $N = |L| = |K|$ . Returning to the polar decompositions

of  $L$  and  $K$ , it follows that

$$N = V^T L = W^T K \implies K = W V^T L.$$

Therefore, taking  $U = W V^T \in O_d(\mathbb{R})$ , we obtain the desired equality.

□

**A.2.3. Matrix Optimization Problems.** To conclude the section about matrix analysis, we consider that the analysis of several specific optimization problems based on eigenvectors is necessary. These problems can be expressed as the maximization of a trace, and they do not need analytical methods, like gradient methods, to find a solution to them. It can be solved only via algebraic methods, specifically by calculating the eigenvectors of the matrices involved in the problem. These problems appear in most of the dimensionality reduction distance metric learning algorithms. We state these problems, together with their solutions, in the lines below.

**THEOREM A.18.** *Let  $d', d \in \mathbb{N}$ , with  $d' \leq d$ . Let  $A \in S_d(\mathbb{R})$ , and we consider the optimization problem*

$$\begin{aligned} \max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \quad & \text{tr}(L A L^T) \\ \text{s.t.} \quad & L L^T = I. \end{aligned}$$

*The problem attains a maximum if  $L = \begin{pmatrix} - & v_1 & - \\ & \dots & \\ - & v_{d'} & - \end{pmatrix}$ , where  $v_1, \dots, v_{d'}$  are orthonormal eigenvectors of  $A$  corresponding to its  $d'$  largest eigenvalues. In addition, the maximum value is the sum of the  $d'$  largest eigenvalues of  $A$ .*

**THEOREM A.19.** *Let  $d', d \in \mathbb{N}$ , with  $d' \leq d$ . Let  $A \in S_d(\mathbb{R})$  and  $B \in S_d(\mathbb{R})^+$ , and we consider the optimization problem*

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \text{tr}((L B L^T)^{-1} (L A L^T))$$

*The problem attains a maximum if  $L = \begin{pmatrix} - & v_1 & - \\ & \dots & \\ - & v_{d'} & - \end{pmatrix}$ , where  $v_1, \dots, v_{d'}$  are eigenvectors of  $B^{-1}A$  corresponding to its  $d'$  largest eigenvalues.*

**THEOREM A.20.** *Let  $d', d \in \mathbb{N}$ , with  $d' \leq d$ . Let  $A, B \in S_d(\mathbb{R})^+$ , and we consider the optimization problem*

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \text{tr}((L B L^T)^{-1} (L A L^T) + (L A L^T)^{-1} (L B L^T))$$

*The problem attains a maximum if  $L = \begin{pmatrix} - & v_1 & - \\ & \dots & \\ - & v_{d'} & - \end{pmatrix}$ , where  $v_1, \dots, v_{d'}$  are the  $d'$  eigenvectors of  $B^{-1}A$  with the highest values for the expression  $\lambda_i + 1/\lambda_i$ , where  $\lambda_i$  is the eigenvalue associated with  $v_i$ .*

These theorems can be proven using tools such as the Rayleigh quotient and the Courant-Fischer theorem and its consequences. First, we will introduce the Rayleigh quotient, and we will see its relationship with the eigenvalues and eigenvectors.

**Definition A.21.** Let  $A \in S_d(\mathbb{R})$ . We define the *Rayleigh quotient* associated with  $A$  as the mapping  $\rho_A: \mathbb{R}^d \setminus \{0\} \rightarrow \mathbb{R}$  given by

$$\rho_A(x) = \frac{x^T A x}{x^T x} = \frac{\langle Ax, x \rangle}{\|x\|_2^2} \quad \forall x \in \mathbb{R}^d \setminus \{0\}.$$

If  $B \in S_d(\mathbb{R})^+$ , we define the *generalized Rayleigh quotient* associated with  $A$  and  $B$  as the mapping  $\mathcal{R}_{A,B}: \mathbb{R}^d \setminus \{0\} \rightarrow \mathbb{R}$  given by

$$\mathcal{R}_{A,B}(x) = \frac{x^T A x}{x^T B x} = \frac{\langle Ax, x \rangle}{\|x\|_B^2} \quad \forall x \in \mathbb{R}^d \setminus \{0\}.$$

Throughout this section we will assume that  $A \in S_d(\mathbb{R})$  and  $B \in S_d(\mathbb{R})^+$  are fixed, and we will refer to Rayleigh quotients as  $\rho = \rho_A$  and  $\mathcal{R} = \mathcal{R}_{A,B}$ . A first observation about  $\rho$  and  $\mathcal{R}$  is that, for  $x \in \mathbb{R}^d \setminus \{0\}$  and  $\lambda \in \mathbb{R}^*$ , it is verified that

$$\mathcal{R}(\lambda x) = \frac{(\lambda x)^T A (\lambda x)}{(\lambda x)^T B (\lambda x)} = \frac{\lambda^2 (x^T A x)}{\lambda^2 (x^T B x)} = \mathcal{R}(x).$$

Therefore,  $\mathcal{R}$  takes all its values over the  $(d-1)$ -dimensional unit sphere, that is,  $\mathcal{R}(\mathbb{R} \setminus \{0\}) = \mathcal{R}(\mathbb{S}^{d-1}) \subset \mathbb{R}$ . Since  $\mathcal{R}$  is continuous and the sphere is compact, it follows that  $\mathcal{R}$  achieves a maximum and a minimum in  $\mathbb{R}^d \setminus \{0\}$ . The same follows with  $\rho$ . These maxima and minima are closely related with the problems we want to analyze. We start studying the extremes of  $\rho$ .

**THEOREM A.22 (RAYLEIGH-RITZ).** Let  $\lambda_{\min}$  and  $\lambda_{\max}$  be the minimum and maximum eigenvalues of  $A$ , respectively. Then,

- (1) For every  $x \in \mathbb{R}^d$ ,  $\lambda_{\min} \|x\|^2 \leq x^T A x \leq \lambda_{\max} \|x\|^2$ .
- (2)  $\lambda_{\max} = \max_{x \in \mathbb{R}^d \setminus \{0\}} \frac{x^T A x}{x^T x} = \max_{\|x\|_2=1} x^T A x$ .
- (3)  $\lambda_{\min} = \min_{x \in \mathbb{R}^d \setminus \{0\}} \frac{x^T A x}{x^T x} = \min_{\|x\|_2=1} x^T A x$ .

Therefore, the maximum and minimum values of  $\rho$  are  $\lambda_{\min}$  and  $\lambda_{\max}$ , respectively. These values are attained in the corresponding eigenvectors.

**PROOF.** Let  $A = U D U^T$ , with  $U \in O_d(\mathbb{R})$  and  $D = \text{diag}(\lambda_1, \dots, \lambda_d)$ , where  $\lambda_1 \leq \dots \leq \lambda_d$ , be a spectral decomposition of  $A$ . Let  $x \in \mathbb{R}^d \setminus \{0\}$  and we take  $y = U^T x$ . Then,

$$\rho(x) = \frac{x^T A x}{x^T x} = \frac{x^T U D U^T x}{x^T x} = \frac{y^T U^T U D U^T U y}{y^T U^T U y} = \frac{y^T D y}{\|y\|_2^2} = \frac{\sum_{i=1}^d \lambda_i y_i^2}{\|y\|_2^2}. \quad (3)$$

In addition, it is clear that

$$\lambda_1 \|y\|_2^2 = \lambda_1 \sum_{i=1}^d y_i^2 \leq \sum_{i=1}^d \lambda_i y_i^2 \leq \lambda_d \sum_{i=1}^d y_i^2 = \lambda_d \|y\|_2^2.$$

Applying this inequality over Eq. 3, it follows that

$$\lambda_1 \leq \rho(x) \leq \lambda_d.$$

Furthermore, if  $u_1$  and  $u_d$  are the corresponding eigenvectors of  $\lambda_1$  and  $\lambda_d$ , we get

$$\rho(u_1) = \frac{u_1^T A u_1}{u_1^T u_1} = \frac{\lambda_1 u_1^T u_1}{u_1^T u_1} = \lambda_1, \quad \rho(u_d) = \frac{u_d^T A u_d}{u_d^T u_d} = \frac{\lambda_d u_d^T u_d}{u_d^T u_d} = \lambda_d.$$

Therefore, the equality is attained, and the three statements of the theorem follow from this equality.  $\square$

Rayleigh-Ritz theorem shows us that  $\rho(\mathbb{R}^d \setminus \{0\}) = [\lambda_{\min}, \lambda_{\max}]$ , obtaining the extreme values in the corresponding eigenvectors. However, these are not the only eigenvalues that can act as an optimal for a Rayleigh quotient. If we restrict ourselves to lower dimensional spaces, we can obtain any eigenvalue of  $A$  as an optimal for the Rayleigh quotient, as we will see below.

**THEOREM A.23 (COURANT-FISCHER).** *Let  $\lambda_1 \leq \dots \leq \lambda_d$  the eigenvectors of  $A$ , and we denote by  $S_k$  a vector subspace of  $\mathbb{R}^d$  of dimension  $k$ . Then, for each  $k \in \{1, \dots, d\}$ , we get*

$$\lambda_k = \min_{S_k \subset \mathbb{R}^d} \max_{\substack{x \in S_k \\ \|x\|_2=1}} x^T A x, \quad (4)$$

$$\lambda_k = \max_{S_{d-k+1} \subset \mathbb{R}^d} \min_{\substack{x \in S_{d-k+1} \\ \|x\|_2=1}} x^T A x. \quad (5)$$

This result extends the Rayleigh-Ritz statement, and this theorem is proven by Horn and Johnson [1990] (chap. 4). There we can also find the proof of an important consequence of Courant-Fischer theorem, usually known as the Cauchy's interlace theorem.

**THEOREM A.24 (CAUCHY'S INTERLACE).** *Suppose that  $\lambda_1 \leq \dots \leq \lambda_d$  are the eigenvalues of  $A$ . Let  $J \subset \{1, \dots, d\}$  be a set of cardinal  $|J| = d'$ , and let  $A_J \in S_{d'}(\mathbb{R})$  be the matrix given by  $A_J = (A_{ij})_{i,j \in J}$ , that is, the submatrix of  $A$  with the entries of  $A$  whose indices are in  $J \times J$ . Then, if  $\tau_1 \leq \dots \leq \tau_{d'}$  are the eigenvalues of  $A_J$ , for each  $k \in \{1, \dots, d'\}$ ,*

$$\lambda_k \leq \tau_k \leq \lambda_{k+d-d'}.$$

The next result follows from Cauchy's interlace theorem, and the inequality it states will help us to solve our optimization problems.

**COROLLARY A.25.** *Let  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$  with  $LL^T = I$ . If  $\mu_1 \geq \dots \geq \mu_d$  are the eigenvalues of  $A$  and  $\sigma_1 \geq \dots \geq \sigma_{d'}$  are the eigenvalues of  $LAL^T$  (now we are considering eigenvalues in decreasing order), then  $\sigma_k \leq \mu_k$ , for  $k = 1, \dots, d'$ .*

**PROOF.** Since  $LL^T = I$ , the rows of  $L$  are orthonormal eigenvectors. We can extend  $L$  to an orthogonal matrix  $\hat{L} \in O_d(\mathbb{R})$  by adding  $d - d'$  orthonormal eigenvectors, and orthonormal to the rows of  $L$ , in its rows. We have then that  $\hat{L}A\hat{L}^T$  and  $A$  have the same eigenvalues, and  $LAL^T$  is a submatrix of  $\hat{L}A\hat{L}^T$  obtained by deleting the last  $d - d'$  rows and columns. The assertion now follows from Cauchy's interlace Theorem A.24, considering eigenvalues in the opposite order.  $\square$

We are now in a position to prove the theorems A.18, A.19 and A.20 proposed at the beginning of this section.

**Theorem** *Let  $d', d \in \mathbb{N}$ , with  $d' \leq d$ . Let  $A \in S_d(\mathbb{R})$ , and we consider the optimization problem*

$$\begin{aligned} \max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \quad & \text{tr}(LAL^T) \\ \text{s.t.} \quad & LL^T = I. \end{aligned} \quad (6)$$

*The problem attains a maximum if  $L = \begin{pmatrix} - & v_1 & - \\ & \dots & \\ - & v_{d'} & - \end{pmatrix}$ , where  $v_1, \dots, v_{d'}$  are orthonormal eigenvectors of  $A$  corresponding to its  $d'$  largest eigenvalues. In addition, the maximum value is the sum of the  $d'$  largest eigenvalues of  $A$ .*



PROOF. Let  $\mu_1 \geq \dots \geq \mu_d$  the eigenvalues of  $A$  in decreasing order, and  $\sigma_1 \geq \dots \geq \sigma_{d'}$  the eigenvalues of  $LAL^T$ . By Corollary A.25, for any  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$  with  $LL^T = I$ ,

$$\text{tr}(LAL^T) = \sum_{i=1}^{d'} \sigma_i \leq \sum_{i=1}^{d'} \mu_i.$$

In addition, when the rows of  $L$  are orthonormal eigenvectors  $v_1, \dots, v_{d'}$  of  $A$  corresponding to  $\mu_1, \dots, \mu_{d'}$ , we get  $LL^T = I$  and  $\text{tr}(LAL^T) = \sum_{i=1}^{d'} \mu_i$ , thus equality holds for these vectors.  $\square$

**LEMMA A.26 (SIMULTANEOUS DIAGONALIZATION).** *Let  $A \in S_d(\mathbb{R})$  and  $B \in S_d(\mathbb{R})^+$ . Then, there is an invertible matrix  $P \in \text{GL}_d(\mathbb{R})$  and a diagonal matrix  $D \in \mathcal{M}_d(\mathbb{R})$  with  $P^T AP = D$  and  $P^T BP = I$ .*

PROOF. We consider the matrix  $C = B^{-1/2}AB^{-1/2}$ .  $C$  is symmetric, since  $A$  and  $B$  are symmetric, thus there is a matrix  $U \in O_d(\mathbb{R})$  so that  $U^T CU$  is diagonal. We call  $D = U^T CU$  and we consider  $P = B^{-1/2}U \in \text{GL}_d(\mathbb{R})$ . We get

$$\begin{aligned} P^T AP &= P^T B^{1/2}CB^{1/2}P = (B^{-1/2}U)^T B^{1/2}CB^{1/2}(B^{-1/2}U) = U^T CU = D, \\ P^T BP &= (B^{-1/2}U)^T B(B^{-1/2}U) = U^T B^{-1/2}BB^{-1/2}U = U^T U = I. \end{aligned}$$

$\square$

**Theorem** *Let  $d', d \in \mathbb{N}$ , with  $d' \leq d$ . Let  $A \in S_d(\mathbb{R})$  and  $B \in S_d(\mathbb{R})^+$ , and we consider the optimization problem*

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \text{tr}((LBL^T)^{-1}(LAL^T)) \quad (7)$$

*The problem attains a maximum if  $L = \begin{pmatrix} - & v_1 - \\ & \dots \\ - & v_{d'} - \end{pmatrix}$ , where  $v_1, \dots, v_{d'}$  are eigenvectors of  $B^{-1}A$  corresponding to its  $d'$  largest eigenvalues.*

PROOF. We denote  $U = L^T \in \mathcal{M}_{d \times d'}(\mathbb{R})$ . If we take the matrix  $P$  from Lemma A.26 and the matrix  $V \in \mathcal{M}_{d \times d'}(\mathbb{R})$  with  $U = PV$  (it exists and it is unique, since  $P$  is regular), we have

$$\begin{aligned} \text{tr}((LBL^T)^{-1}(LAL^T)) &= \text{tr}((U^T BU)^{-1}(U^T AU)) = \text{tr}((V^T P^T BPV)^{-1}(V^T P^T APV)) \\ &= \text{tr}((V^T V)^{-1}(V^T DV)). \end{aligned}$$

Therefore, maximizing Eq. 7 is equivalent to maximize with respect to  $V$  the expression  $\text{tr}((V^T V)^{-1}(V^T DV))$ , because the parameter change is bijective. Now we consider a polar decomposition  $V = Q|V|$ , with  $Q \in \mathcal{M}_{d \times d'}(\mathbb{R})$  verifying  $Q^T Q = I$ . It follows that

$$\begin{aligned} \text{tr}((V^T V)^{-1}(V^T DV)) &= \text{tr}((|V|^T Q^T Q |V|)^{-1}(|V|^T Q^T DQ |V|^T)) \\ &= \text{tr}(|V|^{-1} |V|^{-T} (|V|^T Q^T DQ |V|)) \\ &= \text{tr}(|V|^{-1} Q^T DQ |V|) = \text{tr}(Q^T DQ |V| |V|^{-1}) = \text{tr}(Q^T DQ). \end{aligned}$$

If we call  $W = Q^T$ , what we have obtained is that the maximization of Eq. 7 is equivalent to maximizing in  $W$   $\text{tr}(WDW^T)$ , subject to  $WW^T = I$ , thus obtaining the optimization problem given in Eq. 6. We can suppose the diagonal of  $D$  ordered in descending order, and then a matrix  $W$  that solves the optimization problem can be

obtained adding as rows the vectors  $e_1, \dots, e_{d'}$  of the canonical basis of  $\mathbb{R}^d$ . Then,  $Q$  has contains the same vectors, but added by columns. Observe that the quotient trace  $T(X) = \text{tr}((X^T B X)^{-1}(X^T A X))$ , with  $X \in \mathcal{M}_{d \times d'}(\mathbb{R})$ , is invariant with respect to right multiplications by invertible matrices. Indeed, if  $R \in \text{GL}_{d'}(\mathbb{R})$ ,

$$\begin{aligned} T(XR) &= \text{tr}((R^T X^T B X R)^{-1}(R^T X^T A X R)) = \text{tr}(R^{-1}(X^T B X)^{-1}R^{-T}R^T(X^T A X)R) \\ &= \text{tr}((X^T B X)^{-1}(X^T A X)RR^{-1}) = T(X). \end{aligned}$$

Since  $U$  maximizes  $T$  and  $U = PQ|V|$ , then  $PQ$  also maximizes  $T$ . In addition, as from  $P^T A P = D$  and  $P^T B P = I$  we obtain that

$$D = P^T A P = (P^T B P)^{-1}(P^T A P) = P^{-1}B^{-1}P^{-T}P^T A P = P^{-1}B^{-1}A P,$$

we conclude that  $P$  diagonalizes  $B^{-1}A$ , and then, it contains as columns the eigenvectors of this matrix. Since  $Q$  contains the  $d'$  first eigenvectors of the canonical basis by columns,  $PQ$  contains as columns the  $d'$  first eigenvectors of  $B^{-1}A$ , corresponding to its  $d'$  largest eigenvalues. This ends the proof, because a solution for the problem given by Eq. 7, which is equal to maximizing  $T$  except for a transposition, consists in adding those vectors as rows.  $\square$

**Theorem** *Let  $d', d \in \mathbb{N}$ , with  $d' \leq d$ . Let  $A, B \in S_d(\mathbb{R})^+$ , and we consider the optimization problem*

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \text{tr}((LBL^T)^{-1}(LAL^T) + (LAL^T)^{-1}(LBL^T)) \quad (8)$$

*The problem attains a maximum if  $L = \begin{pmatrix} - & v_1 - \\ & \dots \\ - & v_{d'} - \end{pmatrix}$ , where  $v_1, \dots, v_{d'}$  are the  $d'$  eigen-*

*vectors of  $B^{-1}A$  with the highest values for the expression  $\lambda_i + 1/\lambda_i$ , where  $\lambda_i$  is the eigenvalue associated with  $v_i$ .*

**PROOF.** First of all, given  $C \in S_d(\mathbb{R})^+$  we consider the optimization problem

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \text{tr}((LCL^T + LC^{-1}L^T)) = \max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \text{tr}(L(C + C^{-1})L^T) \quad (9)$$

Using Theorem 6, a solution to this problem can be found by taking as rows of  $L$  the eigenvectors of  $C + C^{-1}$  corresponding to its  $d'$  largest eigenvalues. Observe that the eigenvectors of  $C$  and  $C^{-1}$  are the same, and each one's eigenvalues are the inverse of the other. Therefore,  $C + C^{-1}$  also has the same eigenvectors, and its eigenvalues have the form  $\lambda + 1/\lambda$ , for each  $\lambda$  eigenvalue of  $C$ . Then, the previous solution for Eq. 9 is equivalent to taking the eigenvectors of  $C$  for which  $\lambda + 1/\lambda$  is maximized.

Finally, we only have to realize that we can follow the same proof as in Theorem A.19, considering Eqs. 9 and 8 instead of Eqs. 6 and 7.  $\square$

### A.3. Information Theory

Information theory is a branch of mathematics and computer theory, with the purpose of establishing a rigorous measure to quantify the information and disorder contained in a communication message. It was developed with the aim of finding limits in signal processing operations such as compression, storage and communication. Today, its applications extend to most fields of science and engineering.

Many concepts associated with information theory have been defined, such as entropy, which measures the amount of uncertainty or information expected in an event,

mutual information, which measures the amount of information that one random variable contains about another random variable, or relative entropy, which is a way of measuring the closeness between different random variables. We will focus on the relative entropy, and the concepts derived from it. To do this, we will first define the concept of divergence. Divergence is a magnitude to measure the closeness between certain objects in a set. We should not confuse divergences with distances (we will revisit this concept in Section 2.1), because the magnitudes we will consider may not verify some of the properties required for distances, such as symmetry or triangle inequality.

**Definition A.27.** Let  $X$  be a set. A map  $D(\cdot\|\cdot): X \times X \rightarrow \mathbb{R}$  is said to be a *divergence* if it verifies the following properties:

- (1) Non negativity:  $D(x\|y) \geq 0$ , for every  $x, y \in X$ .
- (2) Coincidence:  $D(x\|y) = 0$  if, and only if,  $x = y$ .

We will use divergences to measure the closeness between probability distributions. The divergences we will use will be presented in the following paragraphs.

**Definition A.28.** Let  $(\Omega, \mathcal{A}, P)$  be a probability space and  $X: \Omega \rightarrow \mathbb{R}$  be a random variable, discrete or continuous, in that space. Suppose that  $p$  is the corresponding probability mass function or density function. Suppose that  $q$  is another probability mass function or density function. Then, we define the *relative entropy* or the *Kullback-Leibler divergence* between  $p$  and  $q$ , as

$$\text{KL}(p\|q) = \mathbb{E}_p \left[ \log \frac{p(X)}{q(X)} \right],$$

as long as such expectation exists. For the discrete case, if  $p$  and  $q$  are valued over the same points, we have

$$\text{KL}(p\|q) = \sum_{x \in X(\Omega)} p(x) \log \frac{p(x)}{q(x)},$$

and for the continuous case, as long as the absolute integral is finite, we have

$$\text{KL}(p\|q) = \int_{-\infty}^{+\infty} p(x) \log \frac{p(x)}{q(x)} dx.$$

For continuity reasons, we assume that  $0 \log(0/0) = 0$ .

The first step is to check that, indeed, Kullback-Leibler divergence is a divergence. This result is known as the information inequality.

**THEOREM A.29 (INFORMATION INEQUALITY).** *Kullback-Leibler divergence is a divergence, that is,  $\text{KL}(p\|q) \geq 0$  and the equality holds if, and only if,  $p(x) = q(x)$  a.e. in  $X(\Omega)$  (the equality is at every point in the discrete case).*

**PROOF.** This result is an immediate consequence of Jensen's inequality [Rudin 1987] applied to the  $-\log$  function, which is strictly convex. We have

$$\begin{aligned} \text{KL}(p\|q) &= \mathbb{E}_p \left[ \log \frac{p(X)}{q(X)} \right] = \mathbb{E}_p \left[ -\log \frac{q(X)}{p(X)} \right] \\ &\geq -\log \mathbb{E}_p \left[ \frac{q(X)}{p(X)} \right] = -\log \int p(x) \frac{q(x)}{p(x)} dx \\ &= -\log \int q(x) dx = -\log 1 = 0. \end{aligned}$$

The proof for the discrete case is similar. In addition, the strict convexity implies that equality holds iff  $p/q$  is constant a.e., iff  $p = q$  a.e., since they are probability density functions or mass functions. And, as in the discrete case  $p$  and  $q$  are valued over sets with no null probabilities, we have equality at every point.  $\square$

As we have already mentioned, Kullback-Leibler divergence is useful to measure closeness between probability distributions and can be used to bring the distributions closer. However, it is not all that useful to put the distributions away, since, as Kullback-Leibler divergence is not symmetric, the values of  $\text{KL}(p\|q)$  and  $\text{KL}(q\|p)$  may differ significantly when  $p$  and  $q$  are not near. That is why it is sometimes helpful to work with a symmetrization of the Kullback-Leibler divergence known as the Jeffrey divergence.

*Definition A.30.* The *Jeffrey divergence* between two probability distributions  $p$  and  $q$  for which  $\text{KL}(p\|q)$  and  $\text{KL}(q\|p)$  exist is defined by

$$\text{JF}(p\|q) = \text{KL}(p\|q) + \text{KL}(q\|p).$$

In the discrete case we have

$$\text{JF}(p\|q) = \sum_{x \in X(\Omega)} (p(x) - q(x))(\log p(x) - \log q(x)).$$

And, for the continuous case,

$$\text{JF}(p\|q) = \int_{-\infty}^{\infty} (p(x) - q(x))(\log p(x) - \log q(x)) dx.$$

It is clear that Jeffrey divergence is a divergence, as a consequence of the information inequality, and it is also symmetric. Observe that both divergences are functions only of the probability distributions, that is, they only depend on the values set on the distributions. This fact allows divergence to be extended to random vectors, as long as we know its probability density functions or mass functions.

A case of special interest in the algorithms we will discuss in subsequent sections is the calculation of divergences between multivariate gaussian distributions. Recall that, if  $\mu \in \mathbb{R}^d$  and  $\Sigma \in S_d(\mathbb{R})^+$ , a random vector  $X = (X_1, \dots, X_d)$  follows a multivariate gaussian distribution with mean  $\mu$  and covariance  $\Sigma$ , if it has the following probability density function:

$$p(x|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} \det(\Sigma)^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right).$$

It is well-known that  $\mathbb{E}[X] = \mu$  and  $\text{Cov}(X) = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T] = \Sigma$ , thus gaussian distributions are completely defined by its mean and covariance. We want to establish an easy way to compute the calculation of divergences between gaussian distributions. To do this, we will find relationships between the studied divergences and matrix divergences. Matrix divergences are an alternative to the Frobenius norm for measuring the closeness between matrices. We are interested in the ones known as Bregman divergences.

*Definition A.31.* Let  $K \subset \mathcal{M}_d(\mathbb{R})$  be an open convex set, and  $\phi: K \rightarrow \mathbb{R}$  a strictly convex and differentiable function. The *Bregman divergence* corresponding to  $\phi$  is the map  $D_\phi(\cdot\|\cdot): K \times K \rightarrow \mathbb{R}$  given by

$$D_\phi(A\|B) = \phi(A) - \phi(B) - \text{tr}(\nabla\phi(B)^T(A - B)).$$

Effectively, Bregman divergences are also divergences, as we can write the expression above as  $D_\phi(A\|B) = \phi(A) - \phi(B) - \langle \nabla\phi(B), A - B \rangle_F$ , which is known to be non

negative when  $\phi$  is strictly convex, and to take the zero value if and only if  $A = B$ . In our situation, we are interested in choosing the *log-det* function to construct a Bregman divergence, that is, the function  $\phi_{ld}: S_d(\mathbb{R})^+ \rightarrow \mathbb{R}$  given by

$$\phi_{ld}(M) = -\log \det(M).$$

This function is known to be strictly convex and its gradient is  $\nabla f(M) = M^{-1}$ , for each  $M$  in  $S_d(\mathbb{R})^+$  [Boyd and Vandenberghe 2004], hence we can construct the known as *log-det divergence* through the expression

$$D_{ld}(A\|B) = \log \det(B) - \log \det(A) - \text{tr}(B^{-1}(A - B)) = \text{tr}(AB^{-1}) - \log \det(AB^{-1}) - d.$$

Once defined the log-det divergence, we are able to express the Kullback-Leibler and Jeffrey divergences between gaussian distributions in terms of this new matrix divergence.

**THEOREM A.32.** *Kullback-Leibler divergence between two multivariate gaussian distributions defined by the probability density functions  $p_1(x|\mu_1, \Sigma_1)$  and  $p_2(x|\mu_2, \Sigma_2)$ , with  $\mu_1, \mu_2 \in \mathbb{R}^d$  and  $\Sigma_1, \Sigma_2 \in S_d(\mathbb{R})^+$ , verifies that*

$$\text{KL}(p_1\|p_2) = \frac{1}{2}D_{ld}(\Sigma_1\|\Sigma_2) + \frac{1}{2}\|\mu_1 - \mu_2\|_{\Sigma_1^{-1}}^2,$$

where  $\|\cdot\|_{\Sigma}$  denotes the norm defined by the positive definite matrix  $\Sigma$ , that is,  $\|v\|_{\Sigma} = \sqrt{v^T \Sigma v}$ , for every  $v \in \mathbb{R}^d$ .

Proof of this result can be found in Davis and Dhillon [2007] (Section 3.1). A simpler version of this theorem can be stated immediately, when we consider equal-mean gaussian distributions.

**COROLLARY A.33.** *Kullback-Leibler divergence between two multivariate gaussian distributions defined by the probability density functions  $p_1$  and  $p_2$  with equal means and covariances  $\Sigma_1$  and  $\Sigma_2$ , verifies that*

$$\text{KL}(p_1\|p_2) = \frac{1}{2}D_{ld}(\Sigma_1\|\Sigma_2).$$

Using these results, we can also express the Jeffrey divergence between gaussian distributions in terms of its mean vectors and covariance matrices. The following expressions can be easily deduced from the theorems above. For more details, see also Nguyen et al. [2017] (Appendix B).

**COROLLARY A.34.** *Jeffrey divergence between two multivariate gaussian distributions defined by the probability density functions  $p_1(x|\mu_1, \Sigma_1)$  and  $p_2(x|\mu_2, \Sigma_2)$  with  $\mu_1, \mu_2 \in \mathbb{R}^d$  and  $\Sigma_1, \Sigma_2 \in S_d(\mathbb{R})^+$ , verifies that*

$$\text{JF}(p_1\|p_2) = \frac{1}{2} \text{tr}(\Sigma_1 \Sigma_2^{-1} + \Sigma_1^{-1} \Sigma_2) - d + \frac{1}{2} \|\mu_1 - \mu_2\|_{\Sigma_1^{-1} + \Sigma_2^{-1}}^2.$$

**COROLLARY A.35.** *Jeffrey divergence between two multivariate gaussian distributions defined by the probability density functions  $p_1$  and  $p_2$  with equal means and covariances  $\Sigma_1$  and  $\Sigma_2$ , verifies that*

$$\text{JF}(p_1\|p_2) = \frac{1}{2} \text{tr}(\Sigma_1 \Sigma_2^{-1} + \Sigma_1^{-1} \Sigma_2) - d.$$

## B. ALGORITHMS FOR DISTANCE METRIC LEARNING: DETAILED EXPLANATION

This appendix describes some of the most popular techniques currently being used in supervised distance metric learning. We also add a review of the principal component analysis, although not supervised, because of its importance for other distance metric

learning algorithms. Some of these techniques, such as PCA or LDA [Cunningham and Ghahramani 2015], are statistical procedures developed over the last century, which are still of great relevance in many problems nowadays. Other more recent proposals are in the state of the art, as is the case of NCMML [Mensink et al. 2012] or DMLMJ [Nguyen et al. 2017], among others. Several of the most popular classic distance metric learning algorithms, such as LMNN [Weinberger and Saul 2009] or NCA [Goldberger et al. 2005], have also been included.

The analyzed techniques are grouped into six subsections. Each of these subsections describes algorithms that share the main purpose, although the purposes described in each section are not exclusive. In the first section (Section B.1) we will study the techniques oriented specifically to dimensionality reduction. Next, the techniques with the purpose of learning distances that improve the nearest neighbors classifiers will be developed (Section B.2), followed by those techniques that aim to improve classifiers based on centroids (Section B.3). The fourth subsection includes methods based on the information theory concepts studied in Section A.3. Subsequently, several distance metric learning mechanisms with less specific goals are described (Section B.5). Finally, kernel-based versions of some of the above algorithms are analyzed, to be able to work in high-dimensionality spaces (Section B.6).

For each of the techniques we will analyze the problem they try to solve or optimize, the mathematical formulations of those problems and the algorithms proposed to solve them.

## B.1. Dimensionality Reduction Techniques

Dimensionality reduction techniques try to learn a distance by searching for a linear transformation from the dataset space to a lower dimensional euclidean space. These kinds of algorithms share many features. For instance, they are usually efficient and their execution involves the calculation of eigenvectors. It is important to point out that there are other non-linear or unsupervised dimensionality reduction techniques [Lee and Verleysen 2007], but they are beyond the scope of this paper (with the exception of kernel versions in Section B.6). The algorithms we will describe are PCA [Jolliffe 2002], LDA [Fisher 1936] and ANMM [Wang and Zhang 2007].

**B.1.1. PCA.** PCA (*principal component analysis*) [Jolliffe 2002] is one of the most popular dimensionality reduction techniques in unsupervised distance metric learning. Although PCA is an unsupervised learning algorithm, it is necessary to talk about it in our work, firstly because of its great relevance, and more particularly, because when a supervised distance metric learning algorithm does not allow a dimensionality reduction, PCA can be first applied to the data in order to be able to use the algorithm later in the lower dimensional space.

Principal component analysis can be understood from two different points of view, which end up leading to the same optimization problem. The first of these approaches consists of finding two linear transformations, one that compresses the data to a smaller space, and another that decompresses them in the original space, so that in the process of compression and decompression the minimum information is lost.

Let us focus on this first approach. Suppose we have the dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , and fix  $0 < d' < d$ . Let us also assume that data are centered, that is, that the mean of the dataset is zero. If it is not the case, it is enough to apply previously to the data the transformation  $x \mapsto x - \mu$ , where  $\mu = \sum x_i / N$  is the dataset mean. We are looking for a compression matrix  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , and a decompression matrix  $U \in \mathcal{M}_{d \times d'}(\mathbb{R})$ , so that, after compressing and decompressing each data the squares of the euclidean distances to the original data are minimal. In other words, the problem we are trying

to solve is

$$\min_{\substack{L \in \mathcal{M}_{d' \times d}(\mathbb{R}) \\ U \in \mathcal{M}_{d \times d'}(\mathbb{R})}} \sum_{i=1}^N \|x_i - ULx_i\|_2^2. \quad (10)$$

To find a solution to this problem, first of all we will see that  $U$  and  $L$  matrices have to be related in a very particular way.

**LEMMA B.1.** *If  $(U, L)$  is a solution of the problem given in Eq. 10, then  $LL^T = I$  (in  $\mathbb{R}^{d'}$ ) and  $U = L^T$ .*

**PROOF.** We fix  $U \in \mathcal{M}_{d \times d'}(\mathbb{R})$  and  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ . We can assume that both  $U$  and  $L$  are full-rank, otherwise the rank of  $UL$  is lower than  $d'$ . Note that in that case, it is always possible to extend  $U$  and  $L$  matrices to full-rank matrices (by replacing linear combinations in the columns by linear independent vectors as long as the dimension allows it) so that the subspace generated extends the one generated by  $UL$ , and in such a case, the error obtained in Eq. 10 for the extension will be, at most, the error obtained for  $U$  and  $L$ .

We consider the linear map  $x \mapsto ULx$ . The image of this map,  $R = \{ULx : x \in \mathbb{R}^d\}$ , is a vector subspace of  $\mathbb{R}^d$  of dimension  $d'$ . Let  $\{u_1, \dots, u_{d'}\}$  be an orthonormal basis of  $R$ , and let  $V \in \mathcal{M}_{d' \times d}(\mathbb{R})$  the matrix that has, by rows, the vectors  $u_1, \dots, u_{d'}$ . It is verified then that the image of  $V$  has dimension  $d'$  and that  $VV^T = I$ . In addition, if we consider  $V^T$  as a linear map, we see that its image is  $R$  (since  $V^T e_i = u_i, i = 1, \dots, d'$ , where  $\{e_1, \dots, e_{d'}\}$  is the canonical basis of  $\mathbb{R}^{d'}$ ).

Therefore, every vector of  $R$  can be written as  $V^T y$ , with  $y \in \mathbb{R}^{d'}$ . Given  $x \in \mathbb{R}^d, y \in \mathbb{R}^{d'}$ , we have

$$\begin{aligned} \|x - V^T y\|_2^2 &= \langle x - V^T y, x - V^T y \rangle \\ &= \|x\|^2 - 2\langle x, V^T y \rangle + \|V^T y\|^2 \\ &= \|x\|^2 - 2\langle y, Vx \rangle + y^T VV^T y \\ &= \|x\|^2 - 2\langle y, Vx \rangle + y^T y \\ &= \|x\|^2 + \|y\|^2 - 2\langle y, Vx \rangle. \end{aligned}$$

If we calculate the gradient with respect to  $y$  from the last previous expression, we obtain  $\nabla_y \|x - V^T y\|_2^2 = 2y - 2Vx$ , which, by equating to zero, allows us to obtain a single critical point,  $y = Vx$ . The convexity of this function (it is the composition of the euclidean norm with an affine map) assures us that this critical point is a global minimum. Therefore, this tells us that, for each  $x \in \mathbb{R}^d$ , the distance to  $x$  in the set  $R$  achieves its minimum at the point  $V^T Vx$ . In particular, for the dataset  $\mathcal{X}$  we conclude that

$$\sum_{i=1}^N \|x_i - ULx_i\|_2^2 \geq \sum_{i=1}^N \|x_i - V^T Vx_i\|_2^2.$$

Since  $U$  and  $L$  were fixed, we can find a matrix  $V$  with these properties for any  $U$  and  $L$  in the conditions of the problem, which concludes the proof.  $\square$

The above lemma allows us to reformulate our problem in terms of only the matrix  $L$ ,

$$\min_{\substack{L \in \mathcal{M}_{d' \times d}(\mathbb{R}) \\ LL^T = I}} \sum_{i=1}^N \|x_i - L^T Lx_i\|_2^2. \quad (11)$$

Let us note now that, for  $x \in \mathbb{R}^d$  and  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , it is verified that

$$\begin{aligned}
 \|x - L^T L x\|_2^2 &= \langle x - L^T L x, x - L^T L x \rangle \\
 &= \|x\|^2 - 2\langle x, L^T L x \rangle + \langle L^T L x, L^T L x \rangle \\
 &= \|x\|^2 - 2x^T L^T L x + x^T L^T L L^T L x \\
 &= \|x\|^2 - x^T L^T L x \\
 &= \|x\|^2 - \text{tr}(x^T L^T L x) \\
 &= \|x\|^2 - \text{tr}(L x x^T L^T).
 \end{aligned}$$

Thus, if we remove terms that do not depend on  $L$ , we can transform the problem in Eq. 11 into the following equivalent problem:

$$\max_{\substack{L \in \mathcal{M}_{d' \times d}(\mathbb{R}) \\ L L^T = I}} \text{tr}(L \Sigma L^T), \quad (12)$$

where  $\Sigma = \sum_{i=1}^N x_i x_i^T$  is, except for a constant, the covariance matrix corresponding to the data in  $\mathcal{X}$ . This matrix is symmetric, and Theorem A.18 guarantees that we can find a maximum of the problem if we build  $L$  adding the  $d'$  orthonormal eigenvectors corresponding to the  $d'$  largest eigenvalues of  $\Sigma$ . The directions that determine these vectors are the *principal directions*, and the components of the data transformed in the orthonormal system determined by the principal directions are the so-called *principal components*.

To conclude, the second approach from which the principal components problem can be dealt with consists of selecting the orthogonal directions for which the variance is maximized. We know that if  $\Sigma$  is the covariance matrix of  $\mathcal{X}$ , when applying a linear transformation  $L$  to the data, the new covariance matrix is given by  $L \Sigma L^T$ . If we want a transformation that reduces the dimensionality and for which the variance is maximized in each variable, what we are looking for is to take the trace of the previous matrix, which leads us back again to Eq. 12. The symmetry of  $\Sigma$  ensures that we can take the main orthonormal directions that maximize the variance for each possible value of  $d'$ .

Finally, it is important to note that the matrix  $L \in \mathcal{M}_d(\mathbb{R})$  (taking all dimensions) that is constructed by adding  $\Sigma$  eigenvectors row by row is the orthogonal matrix that diagonalizes  $\Sigma$ , and therefore, when  $L$  is applied to the data, the transformed data have as the covariance matrix the diagonal matrix  $L \Sigma L^T = \text{diag}(\lambda_1, \dots, \lambda_d)$ , where  $\lambda_1 \geq \dots \geq \lambda_d$  are the eigenvalues of  $\Sigma$ . This tells us that the eigenvalues of the covariance matrix represent the amount of variance explained by each of the principal directions. This provides an additional advantage to PCA, since it allows the percentage of variance that explains each principal component to be analyzed, in order to be able to later choose a dimension that adjusts to the amount of variance that we want to keep in the transformed data.

Figure 10 graphically exemplifies how principal component analysis works.

**B.1.2. LDA.** LDA (*linear discriminant analysis*) [Fisher 1936] is a classical distance metric learning technique with the purpose of learning a projection matrix that maximizes the separation between classes in the projected space, that is, it tries to find the directions that best distinguish the different classes, as shown in Figure 11.

Figure 11 also allows us to compare the results of the projections obtained by PCA and LDA, showing the most remarkable difference between the two techniques: PCA does not take into account the labels information, while LDA does use it. We can observe that the directions obtained by PCA and LDA do not present any type of rela-



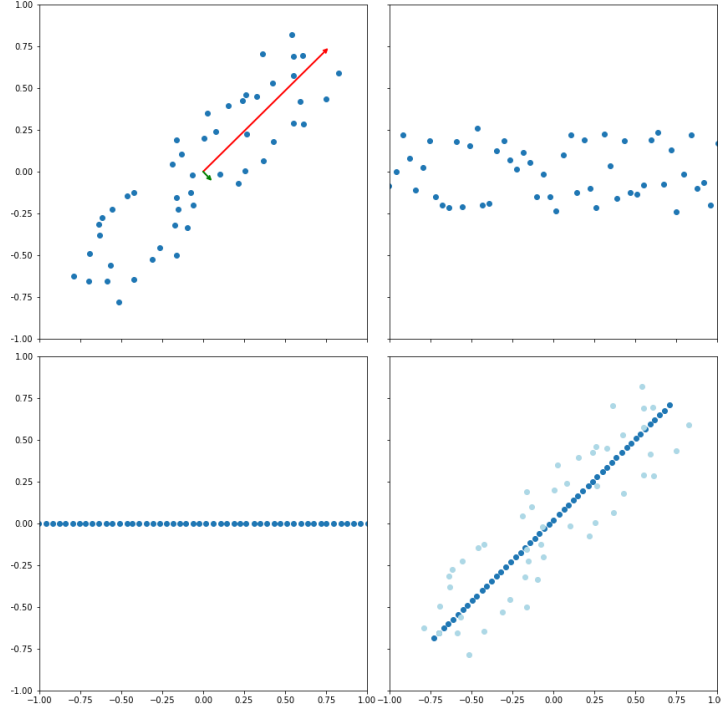


Fig. 10. A graphical example of PCA. The first image shows a dataset, along with the principal directions (proportional according to the explained variance) learned by PCA. To the right, the data is projected at maximum dimension. We observe that this projection consists of rotating the data making the axes coincide with the principal directions. At the bottom left, data is projected onto the first principal component. Finally, to the right, the data recovered through the decompression matrix, along with the original data. We can see that the PCA projection is the one that minimizes the quadratic decomposition error. In this particular case the decompressed data is on the regression line of the original data, due to the dimensions of the problem.

tionship, the latter being the only one of them that provides a data projection oriented to supervised learning.

It is also possible to observe in Figure 11 that it makes no sense to look for a second independent direction that continues to maximize class separation, while in PCA it always makes sense to look inductively for orthogonal directions that maximize variance. If the dataset shown in the figure had a third class, we could find a second direction that maximizes the separation between classes, thus offering the possibility of projecting onto a plane. In general, we will see that if we have  $r$  classes we will be able to find at most (and as long as the dimension of the original space allows it)  $r - 1$  directions that maximize the separation. This indicates that the projections that LDA is going to learn will be, in general, towards a quite low dimension, and always limited by the number of classes in the dataset.

Suppose we have the labeled dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , where  $\mathcal{C}$  is the set of all the classes in the problem and  $y_1, \dots, y_N \in \mathcal{C}$  are the corresponding labels. Suppose that the number of classes in the problem is  $|\mathcal{C}| = r$ . For each  $c \in \mathcal{C}$  we define the set

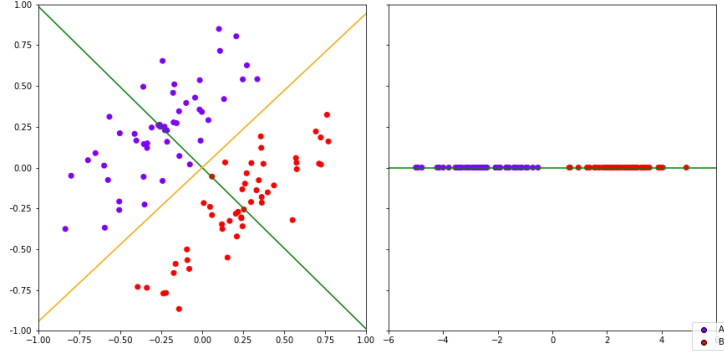


Fig. 11. Graphical example of LDA and comparison with PCA. The first image shows a dataset, with the first principal direction determined by PCA, in orange, and the direction determined by LDA, in green. We observe that if we project the data on the direction obtained by LDA they separate, as it is shown in the right image. In contrast, the direction obtained by PCA only allows us to maximize the variance of the whole dataset, since it does not consider the information of the labels.

$\mathcal{C}_c = \{i \in \{1, \dots, N\} : y_i = c\}$ , and  $N_c = |\mathcal{C}_c|$ . We consider the mean vector of each class,

$$\mu_c = \frac{1}{N_c} \sum_{i \in \mathcal{C}_c} x_i,$$

and the mean vector for the whole dataset,

$$\mu = \frac{1}{N} \sum_{c \in \mathcal{C}} \sum_{i \in \mathcal{C}_c} x_i = \frac{1}{N} \sum_{i=1}^N x_i.$$

We will define two scatter matrices, one *between-class*, denoted as  $S_b$ , and the other *within-class*, denoted as  $S_w$ . The between-class scatter matrix is defined as

$$S_b = \sum_{c \in \mathcal{C}} N_c (\mu_c - \mu)(\mu_c - \mu)^T.$$

And the within-class scatter matrix is defined as

$$S_w = \sum_{c \in \mathcal{C}} \sum_{i \in \mathcal{C}_c} (x_i - \mu_c)(x_i - \mu_c)^T.$$

Note that these matrices represent, except multiplicative constants, the covariances between the data of different classes, taking the class means as representatives for each class in the first case, and the sum, for each class, of the covariances of that class data, in the second case. Since we want to maximize the separation between classes we will formulate the problem of optimization as the search for a projection  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$  that maximizes the quotient of the between-class variances and within-class variances determined by the previous matrices. The problem is established as

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \text{tr} \left( (LS_w L^T)^{-1} (LS_b L^T) \right). \quad (13)$$

Theorem A.19 assures us that, in order to maximize the problem given in Eq. 13,  $L$  has to be composed by the eigenvectors corresponding to the largest eigenvalues of  $S_w^{-1} S_b$ , as long as  $S_w$  is invertible. In practice, this happens in most problems where

$N \gg d$ , because  $S_w$  is the sum of  $N$  outer products, each of which may add a new dimension to the matrix rank. If  $N \gg d$  it is likely that  $S_w$  is full-rank. This, together with the fact that  $S_w$  is positive semidefinite, would guarantee  $S_w$  to be positive definite, thus entering into the theorem hypothesis.

It is interesting to remark the similarity between the optimization problem in Eq. 13 and the expression of the Calinski-Harabasz index [Caliński and Harabasz 1974], an index used in clustering to measure the separation of the established clusters, and that uses the same scatter matrices, and a similar quotient formulation.

Furthermore, let us note, as it was already mentioned at the beginning of this section, that at most we can get  $r - 1$  eigenvectors with a non zero corresponding eigenvalue. This is because the maximum rank of  $S_b$  is  $r - 1$ , because its rank coincides with the rank of the matrix  $A$  that has as columns the vectors  $\mu_c - \mu$  (we get  $S_b = A \text{diag}(N_{c_1}, \dots, N_{c_r}) A^T$ ), which can have as maximum rank  $r$ , and this matrix also includes the linear combination  $\sum N_c(\mu_c - \mu) = 0$ , so at least one column is linearly dependent of the others. Therefore,  $S_w^{-1} S_b$  also has a maximum rank of  $r - 1$ . Consequently, the projection matrix that maximizes Eq. 13 is also going to have, at most, this rank, thus the projection will be contained in a space of this dimension. Therefore, the choice of a dimension  $d' > r - 1$  will not provide any additional information to that provided by the projection onto dimension  $r - 1$ .

To conclude, although we have seen that LDA allows us to reduce dimensionality by adding supervised information as opposed to the non supervised PCA, it can also present some limitations:

- If the size of the dataset is too small, the within-class scatter matrix may be singular, preventing the calculation of  $S_w^{-1} S_b$ . In this situation, several mechanisms are proposed to keep this technique going. One of the most used consists of regularizing the problem, considering, instead of  $S_w$ , the matrix  $S_w + \varepsilon I$ , where  $\varepsilon > 0$ , making  $S_w + \varepsilon I$  be positive definite. The problem of the singularity of  $S_w$  also arises if there are correlated attributes. This case can be avoided by eliminating redundant attributes in a preprocessing prior to learning.
- The definition of the scatter matrices assumes, to some extent, that the data in each class are distributed according to a multivariate gaussian distribution. Therefore, if the data presented other distributions, the projection learned might not be of enough quality.
- As already mentioned, LDA only allows the extraction of  $r - 1$  attributes, which may be suboptimal in some cases, as a lot of information could be lost.

**B.1.3. ANMM.** ANMM (*average neighborhood margin maximization*) [Wang and Zhang 2007] is a distance metric learning technique specifically oriented to dimensionality reduction. It therefore follows the same path as the aforementioned PCA and LDA, trying to solve some of their limitations.

The objective of ANMM is to learn a linear transformation  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , with  $d' \leq d$ , that projects the data onto a lower dimensional space, so that the similarity between the elements of the same class and the separation between classes is maximized, following the criterion of maximization of margins that we will show next.

We consider the training dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , with corresponding labels  $y_1, \dots, y_N$ , and we fix  $\xi, \zeta \in \mathbb{N}$ , and euclidean distance as the initial distance. From these variables we will create two types of neighborhoods.

**Definition B.2.** Let  $x_i \in \mathcal{X}$ .

We define the  $\xi$ -nearest homogeneous neighborhood of  $x_i$  as the set of the  $\xi$  samples in  $\mathcal{X} \setminus \{x_i\}$  nearest to  $x_i$  that belong to its same class. We will denote it by  $\mathcal{N}_i^\xi$ .

We define the  $\zeta$ -nearest heterogeneous neighborhood of  $x_i$  as the set of the  $\zeta$  samples in  $\mathcal{X}$  nearest to  $x_i$  that belong to a different class. We will denote it by  $\mathcal{N}_i^e$ .

ANMM is intended to maximize the concept of *average neighborhood margin*, which we define below.

*Definition B.3.* Given  $x_i \in \mathcal{X}$ , its *average neighborhood margin*  $\gamma_i$  is defined as

$$\gamma_i = \sum_{k: x_k \in \mathcal{N}_i^e} \frac{\|x_i - x_k\|^2}{|\mathcal{N}_i^e|} - \sum_{j: x_j \in \mathcal{N}_i^o} \frac{\|x_i - x_j\|^2}{|\mathcal{N}_i^o|}.$$

The (global) average neighborhood margin  $\gamma$  is defined as

$$\gamma = \sum_{i=1}^N \gamma_i.$$

Note that, for each  $x_i \in \mathcal{X}$ , its average neighborhood margin represents the difference between the average distance from  $x_i$  to its heterogeneous neighbors, and the average distance from  $x_i$  to its homogeneous neighbors. Therefore, maximizing this margin allows, locally, to move data from different classes away, and pulling those of the same class. Figure 12 graphically describes the concept of average neighborhood margin.

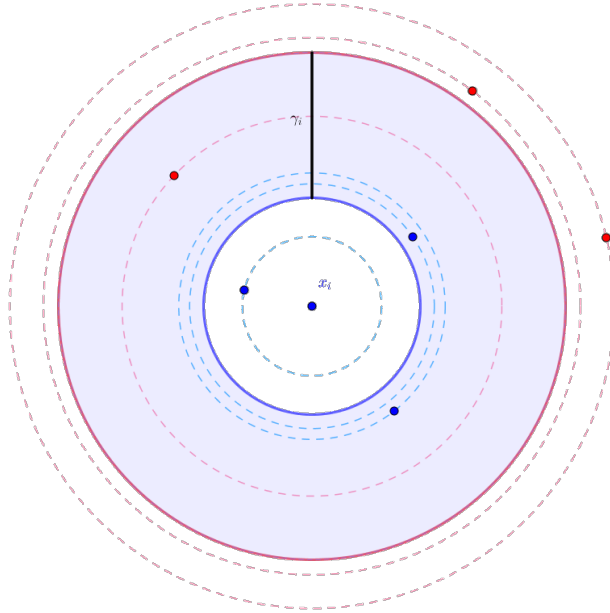


Fig. 12. Graphical description of the average neighborhood margin, for the sample  $x_i$ , for  $\xi = \zeta = 3$ . The blue and red circumferences determine the average distance from  $x_i$  to data of the same and different classes, respectively.

We are now looking for a linear transformation  $L$  that maximizes the margin associated with the projected data,  $\{Lx_i: i = 1, \dots, N\}$ . For such data, we have the average

neighborhood margin corresponding to that transformation,

$$\gamma^L = \sum_{i=1}^N \gamma_i^L = \sum_{i=1}^N \left( \sum_{k: x_k \in \mathcal{N}_i^e} \frac{\|Lx_i - Lx_k\|^2}{|\mathcal{N}_i^e|} - \sum_{j: x_j \in \mathcal{N}_i^o} \frac{\|Lx_i - Lx_j\|^2}{|\mathcal{N}_i^o|} \right).$$

Observe that, thanks to the linearity of the trace operator, we can express

$$\begin{aligned} \sum_{i=1}^n \sum_{k: x_k \in \mathcal{N}_i^e} \frac{\|Lx_i - Lx_k\|^2}{|\mathcal{N}_i^e|} &= \text{tr} \left( \sum_{i=1}^N \sum_{k: x_k \in \mathcal{N}_i^e} \frac{(Lx_i - Lx_k)(Lx_i - Lx_k)^T}{|\mathcal{N}_i^e|} \right) \\ &= \text{tr} \left[ L \left( \sum_{i=1}^N \sum_{k: x_k \in \mathcal{N}_i^e} \frac{(x_i - x_k)(x_i - x_k)^T}{|\mathcal{N}_i^e|} \right) L^T \right] \\ &= \text{tr}(LSL^T), \end{aligned}$$

where  $S = \sum_i \sum_{k: x_k \in \mathcal{N}_i^e} \frac{(x_i - x_k)(x_i - x_k)^T}{|\mathcal{N}_i^e|}$  is called the *scatter matrix*. In a similar way, if we define  $C = \sum_i \sum_{j: x_j \in \mathcal{N}_i^o} \frac{(x_i - x_j)(x_i - x_j)^T}{|\mathcal{N}_i^o|}$ , which we will call the *compactness matrix*, we get

$$\sum_{i=1}^n \sum_{j: x_j \in \mathcal{N}_i^o} \frac{\|Lx_i - Lx_j\|^2}{|\mathcal{N}_i^o|} = \text{tr}(LCL^T).$$

And therefore, combining both expressions,

$$\gamma^L = \text{tr}(L(S - C)L^T). \quad (14)$$

The maximization of  $\gamma^L$  as presented in Eq. 14 is not restrictive enough, because it is enough to multiply  $L$  by positive constants to get a value of  $\gamma^L$  as large as we want. That is why the constraint  $LL^T = I$  is added, so we end up with the next optimization problem:

$$\begin{aligned} \max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} \quad & \text{tr}(L(S - C)L^T) \\ \text{s.t.:} \quad & LL^T = I. \end{aligned}$$

Observe that  $S - C$  is symmetric, as it is the difference between two positive semidefinite matrices (each of them is the sum of outer products). Theorem A.18 tells us that the matrix  $L$  we are looking for can be built by adding, by rows, the  $d'$  eigenvectors of  $S - C$  corresponding to its  $d'$  largest eigenvalues.

To conclude, note that ANMM solves some of the issues of the previously mentioned PCA and LDA. On the one hand, it is a supervised learning algorithm, hence it uses the class information that is ignored by PCA. On the other hand, faced with the shortcomings of LDA, we can see that:

- It does not have computational problems with small samples, for which scatter or compactness matrices may be singular, because it does not have to calculate their inverse matrices.
- It does not make any assumption about the class distributions. The formulation of the problem is purely geometric.
- It admits any size for dimensionality reduction. It does not impose that this size must be lower than the number of classes.

Finally, we can also observe that, if we keep the maximum dimension  $d$ , the condition  $LL^T = I$  implies that  $L$  is orthogonal and  $L^TL = I$ , thus we are just learning an isometry, as already happened with PCA. Therefore, distance-based classifiers will only be able to experience improvements when the chosen dimension is strictly smaller than the original one.

## B.2. Algorithms to Improve Nearest Neighbors Classifiers

In the following paragraphs we will analyze algorithms specifically designed to work with nearest neighbors classifiers. The algorithms we will study are known as LMNN [Weinberger and Saul 2009] and NCA [Goldberger et al. 2005].

**B.2.1. LMNN.** LMNN (*large margin nearest neighbors*) [Weinberger and Saul 2009] is a distance metric learning algorithm aimed specifically at improving the accuracy of the  $k$ -nearest neighbors classifier. It is based on the premise that this classifier will label a sample more reliably if its  $k$  neighbors share the same label, and to do so it tries to learn a distance that maximizes the number of samples that share its label with as many neighbors as possible.

In this way, the LMNN algorithm tries to minimize an error function that penalizes, on the one hand, the large distance between each sample and those considered its ideal neighbors, and on the other hand, the small distances between examples of different classes.

Suppose we have a dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$  with corresponding labels  $y_1, \dots, y_N$ . To work, the algorithm makes use of the concept of *target neighbors*. Given a sample  $x_i \in \mathcal{X}$ , its  $k$  target neighbors are those examples of the same class as  $x_i$  and different from this, for which it is desired to be considered as neighbors in the nearest neighbors classification. If  $x_j$  is a target neighbor of  $x_i$ , then we will write it as  $j \rightsquigarrow i$ . Observe that the relationship given by  $\rightsquigarrow$  may not be symmetric. Target neighbors are fixed during the learning process. If we have some prior information about our dataset we can use it to determine the target neighbors. Otherwise, a good option is to use the nearest neighbors for the euclidean distance as target neighbors.

Once the target neighbors have been established, for each distance and for each sample in  $\mathcal{X}$  we can create a perimeter determined by the the furthest target neighbor. We are looking for distances for which there are no samples of other classes in this perimeter. It is necessary to emphasize that with this perimeter there are not enough separation guarantees, because a feasible distance could have collapsed all the target neighbors in a point, and then the perimeter would have radius zero. For this reason, a margin determined by the radius of the perimeter is considered, to which a positive constant is added. We will see that there is no loss of generality, because of the function that we will define, in supposing that this constant is 1. Any sample of a different class that invades this margin will be called an *impostor*. Our objective, therefore, will be, in addition to bringing each sample as close as possible to its target neighbors, to try to keep impostors as far away as possible.

In mathematical terms, if our distance is determined by the linear transformation  $L \in \mathcal{M}_d(\mathbb{R})$ , and  $x_i, x_j \in \mathcal{X}$  with  $j \rightsquigarrow i$ , we will say that  $x_l$  is an impostor for these samples if  $y_l \neq y_i$  and  $\|L(x_i - x_j)\|^2 \leq \|L(x_i - x_j)\|^2 + 1$ . In Figure 13 the concepts of target neighbor and impostor are graphically described. Finally, note that the margin is defined in terms of the squared distances, instead of considering only the distance. This will make the problem formulation easy to solve.

We now proceed to define accurately the terms of the objective function. As already mentioned, it will be composed of two terms. The first one will penalize distant target neighbors and the second one will penalize nearby impostors. The first term is defined

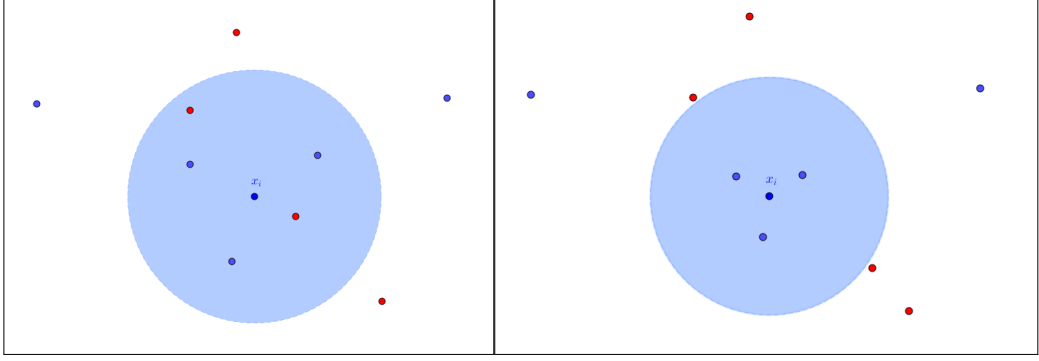


Fig. 13. Graphical description of target neighbors and impostors (with  $k = 3$ ) for the sample  $x_i$ . The blue circle represents the margin determined by the target neighbors. All the points of different classes in this circle are impostors. LMNN's goal will be to bring the target neighbors as close as possible and to remove the impostors from the circle. Therefore, data of the same class that are not target neighbors will not have any influence, and impostors will no longer be penalized as soon as they leave the margin, as shown in right image. This gives a local nature to this learning technique.

as

$$\varepsilon_{pull}(L) = \sum_{i=1}^N \sum_{j \rightsquigarrow i} \|L(x_i - x_j)\|^2.$$

The minimization of this error causes a pulling force between the data samples. The second term is defined as

$$\varepsilon_{push}(L) = \sum_{i=1}^N \sum_{j \rightsquigarrow i} \sum_{l=1}^N (1 - y_{il}) [1 + \|L(x_i - x_j)\|^2 - \|L(x_i - x_l)\|^2]_+,$$

where  $y_{il}$  is a binary variable which takes the value 1 if  $y_i = y_l$ , and 0 if  $y_i \neq y_l$ , and the operator  $[\cdot]_+ : \mathbb{R} \rightarrow \mathbb{R}_0^+$  is defined as  $[z]_+ = \max\{z, 0\}$ . Thus, this error adds up when  $y_{il} = 0$  (that is,  $x_l$  is in different class to  $x_i$ ), and the second factor is strictly positive (that is, the margin defined by the target neighbors is exceeded). The minimization of this second term causes a pushing force between the data samples.

Finally, the objective function results from combining these two terms. After fixing  $\mu \in ]0, 1[$ , we define

$$\varepsilon(L) = (1 - \mu)\varepsilon_{pull}(L) + \mu\varepsilon_{push}(L). \quad (15)$$

The authors state that, experimentally, the choice of  $\mu$  does not cause great differences in results, so it is usually taken  $\mu = 1/2$ . Minimizing this function will lead us to learn the distance we were looking for. Note that this function is sub-differentiable, but not convex, so if we use a subgradient descent method under this approach we may be stuck in a local optimal. However, we can reformulate the objective function in order to make it act over the positive semidefinite cone. If for every  $L \in \mathcal{M}_d(\mathbb{R})$  we take  $M = L^T L \in S_d(\mathbb{R})_0^+$ , we know that  $\|x_i - x_j\|_M^2 = \|L(x_i - x_j)\|_2^2$ , and consequently,

$$\varepsilon(M) = (1 - \mu) \sum_{i=1}^N \sum_{j \rightsquigarrow i} \|x_i - x_j\|_M^2 + \mu \sum_{i=1}^N \sum_{j \rightsquigarrow i} \sum_{l=1}^N [1 + \|x_i - x_j\|_M^2 - \|x_i - x_l\|_M^2]_+ \quad (16)$$

is a convex function in  $M$  that takes the same values as  $\varepsilon(L)$ . The minimization of  $\varepsilon(M)$  in this case is subject to the constraint  $M \in S_d(\mathbb{R})_0^+$ , so the projected subgradient method, with projections onto the positive semidefinite cone, can be used to optimize this function. In addition, we can easily calculate a subgradient  $G \in \partial\varepsilon/\partial M$  given by

$$G = (1 - \mu) \sum_{i,j \rightsquigarrow i} O_{ij} + \mu \sum_{(i,j,l) \in \mathcal{N}} (O_{ij} - O_{il}),$$

where  $\mathcal{N}$  is the set of triplets  $(i, j, l)$  for which  $x_l$  is an impostor over  $x_i$  with the margin determined by  $x_j$ , and  $O_{ij} = (x_i - x_j)(x_i - x_j)^T$  are the outer products obtained from the distances differentiation. The first term of the gradient is constant, while the second term only varies in each iteration with the changes of the impostors that enter or leave the set  $\mathcal{N}$ . These considerations allow a fairly efficient gradient calculation.

As for dimensionality reduction, two different alternatives are presented. If we keep the optimization with respect to  $M$ , it is not feasible to add rank restrictions, as we saw in Example 2.5. Therefore, the use of PCA is suggested prior to the algorithm execution, to project the data onto its first principal components, and then apply LMNN on the projected data. The other alternative is to optimize the objective function with respect to  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , with  $d' < d$  using a gradient descent algorithm. In this case the optimization is not convex, but we learn directly a linear transformation that reduces the dimensionality without making changes in the optimization of Eq. 15. Authors also state, based on empirical results, that this non-convex optimization gives good results.

Other proposals made for the improvement of this algorithm consist of applying LMNN multiple times, learning new metrics each time, and using these metrics to determine increasingly accurate target neighbors, or learning different metrics locally. Finally, although the distance learned by LMNN is designed to be used by the  $k$ -neighbors classifier, it is possible to use the objective function itself as a classification method. These classification models are called *energy-based*. Thus, to classify a test sample  $x_t$ , for each possible label value  $y_t$ , we look for  $k$  target neighbors in the training set for class  $y_t$ , and evaluate the *energy* for the metric learned, finally assigning to  $x_t$  the value of  $y_t$  that provides the lowest energy. According to the objective function, energy will penalize large distances between  $x_t$  and its target neighbors, impostors on the  $x_t$  perimeter, and perimeters of other classes invaded by  $x_t$ . Therefore,

$$\begin{aligned} y_t^{pred} = \arg \min_{y_t} & \left\{ (1 - \mu) \sum_{j \rightsquigarrow t} \|x_t - x_j\|_M^2 \right. \\ & + \mu \sum_{j \rightsquigarrow t, l} (1 - y_{tl}) [1 + \|x_t - x_j\|_M^2 - \|x_t - x_l\|_M^2]_+ \\ & \left. + \mu \sum_{i, j \rightsquigarrow i} (1 - y_{it}) [1 + \|x_i - x_j\|_M^2 - \|x_i - x_t\|_M^2]_+ \right\}. \end{aligned}$$

**B.2.2. NCA.** NCA (*neighborhood components analysis*) [Goldberger et al. 2005] is another distance metric learning algorithm aimed specifically at improving the accuracy of the nearest neighbors classifiers. Its aim is to learn a linear transformation with the goal of minimizing the leave-one-out error expected by the nearest neighbor classification. Additionally, this transformation could be used to reduce the dimensionality of the dataset, and thus make the classifier more efficient.

We consider the training set  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , labeled by  $y_1, \dots, y_N$ . We want to learn a distance, determined by a linear transformation  $L \in \mathcal{M}_d(\mathbb{R})$ , that optimizes the accuracy of the nearest neighbors classifier. Ideally, we would optimize the per-



formance of the classifier over the test dataset, but we only have the training set. Therefore, our goal will be to try to optimize the classification leave-one-out error on the training set. The choice of the leave-one-out error is due to the nature of the nearest neighbors classifier: as we will learn and evaluate over the same set, the nearest neighbor of each sample would be the sample itself, which would not allow the results to be interpreted correctly if the sample is kept while evaluating it.

However, the function that maps each transformation  $L$  to the leave-one-out error for the distance corresponding to  $L$  has no guarantee of differentiability, not even continuity, so it is not easy to deal with it for optimization (observe that the image of this function is a finite set, and its domain is a connected set, so it cannot be continuous unless it is constant, which does not happen in non-trivial examples).

To do this, NCA tries to approach the problem in a stochastic way, that is, instead of operating with the leave-one-out error directly, it operates with its expected value for the probability that we will define below.

Given two samples  $x_i, x_j \in \mathcal{X}$ , we define the probability that  $x_i$  has  $x_j$  as its nearest neighbor, for the distance determined by the mapping  $L$ , as follows:

$$p_{ij}^L = \frac{\exp(-\|Lx_i - Lx_j\|^2)}{\sum_{k \neq i} \exp(-\|Lx_i - Lx_k\|^2)} \quad (j \neq i), \quad p_{ii}^L = 0.$$

Notice that, indeed,  $p_{i*}$  defines a probability measure on the set  $\{1, \dots, N\}$ , for each  $i \in \{1, \dots, N\}$ . Under this probability law, we can define the probability that the sample  $x_i$  is correctly classified as the sum of the probabilities that  $x_i$  has as its nearest neighbor each sample of its same class, that is

$$p_i^L = \sum_{j \in C_i} p_{ij}^L, \text{ where } C_i = \{j \in \{1, \dots, N\} : y_j = y_i\}.$$

Finally, the expected number of correctly classified samples, and the function we will try to maximize, is obtained as

$$f(L) = \sum_{i=1}^N p_i^L = \sum_{i=1}^N \sum_{j \in C_i} p_{ij}^L = \sum_{i=1}^N \sum_{\substack{j \in C_i \\ j \neq i}} \frac{\exp(-\|Lx_i - Lx_j\|^2)}{\sum_{k \neq i} \exp(-\|Lx_i - Lx_k\|^2)}.$$

This function is differentiable, and its derivative can be computed as

$$\nabla f(L) = 2L \sum_{i=1}^N \left( p_i^L \sum_{k=1}^N p_{ik}^L O_{ik} - \sum_{j \in C_i} p_{ij}^L O_{ij} \right),$$

where  $O_{ij} = (x_i - x_j)(x_i - x_j)^T$  represent again the outer products between the differences of the samples in  $\mathcal{X}$ . Once the gradient is known, we can optimize the objective function using a gradient ascent method. Note that the objective function is not concave, and can therefore be trapped in local optima. Another issue for this algorithm is the possibility of overfitting, if the expected leave-one-out error of the learned distance is too low. Authors affirm, based on the experimental results, that normally there is no overfitting, even if we ascend a lot in the objective function.

### B.3. Algorithms to Improve Nearest Centroids Classifiers

In this block we will analyze, following the previous lines, algorithms specifically oriented to improve distance-based classifiers, focusing in this case on the classifiers based on centroids. The algorithms we will study are NCML and NCMC [Mensink et al. 2012].

**B.3.1. NCMML.** NCMML (*nearest class mean metric learning*) [Mensink et al. 2012] is a distance metric learning algorithm specifically designed to improve the nearest class mean (NCM) classifier. To do this, it uses a probabilistic approach similar to that used by NCA to improve the accuracy of the nearest neighbors classifier.

Nearest class mean classifier, during learning process, calculates the mean vectors of each class subset. Then, when predicting a new sample, it assigns the class of the nearest mean vector found. It is a very efficient and simple classifier, although its simplicity makes it a rather weak classifier against datasets that are not grouped around their mean. We will learn in the following lines how to learn a distance for this classifier.

We consider the training set  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , with labels  $y_1, \dots, y_N \in \mathcal{C}$ , where  $\mathcal{C} = \{c_1, \dots, c_r\}$  is the set of available classes. For each  $c \in \mathcal{C}$ , we call  $\mu_c \in \mathbb{R}^d$  the mean vector of the samples belonging to the class  $c$ , that is,  $\mu_c = \frac{1}{N_c} \sum_{i: y_i=c} x_i$ , where  $N_c$  is the number of elements of  $\mathcal{X}$  that belong to class  $c$ . Given a linear transformation  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , we will define, for each  $x \in \mathcal{X}$  and each  $c \in \mathcal{C}$ , the probability that  $x$  will be labeled with the class  $c$  (according to the nearest class mean criterion) as follows:

$$p_L(c|x) = \frac{\exp\left(-\frac{1}{2}\|L(x - \mu_c)\|^2\right)}{\sum_{c' \in \mathcal{C}} \exp\left(-\frac{1}{2}\|L(x - \mu_{c'})\|^2\right)}.$$

Note that  $p_L(\cdot|x)$  effectively defines a probability in the set  $\mathcal{C}$ . Once the above probability is defined, the objective function that NCMML tries to maximize is the log-likelihood for the labeled data in the training set, that is,

$$\mathcal{L}(L) = \frac{1}{N} \sum_{i=1}^N \log p_L(y_i|x_i).$$

This function is differentiable and its gradient is given by

$$\nabla \mathcal{L}(L) = \frac{1}{N} \sum_{i=1}^N \sum_{c \in \mathcal{C}} \alpha_{ic} L(\mu_c - x_i)(\mu_c - x_i)^T,$$

where  $\alpha_{ic} = p_L(c|x_i) - \mathbb{I}[y_i = c]$  and  $\mathbb{I}[R]$  denotes the indicator function for the condition  $R$ . The maximization of this function using gradient methods is the task carried out by NCMML.

**B.3.2. NCMC.** Although nearest class mean classifier is a simple, intuitive and efficient classifier in both learning and prediction processes, it has one major drawback, and that is that it assumes that classes are grouped around their center, which is an overly restrictive hypothesis. In Figure 14 we can see an example where NCM is unable to give good results.

One way to solve this problem is, instead of considering the center of each class to classify new samples, to find subgroups within each class that present a quality grouping, and to consider the center for each of its subgroups. In this way we would have a set of centroids for each class, and at the time of classifying a new sample, it would suffice to select the nearest centroid and assign it the class of which it is centroid.

In this new classifier, which we will call NCMC (*nearest class with multiple centroids*), the clustering algorithms come into play. There are numerous algorithms [Xu and Wunsch 2005] to obtain a set of clusters from a dataset, each with its advantages and disadvantages. Due to the form of our problem, in which we are interested not only in obtaining a set of clusters for each class, but also a center for each cluster, the

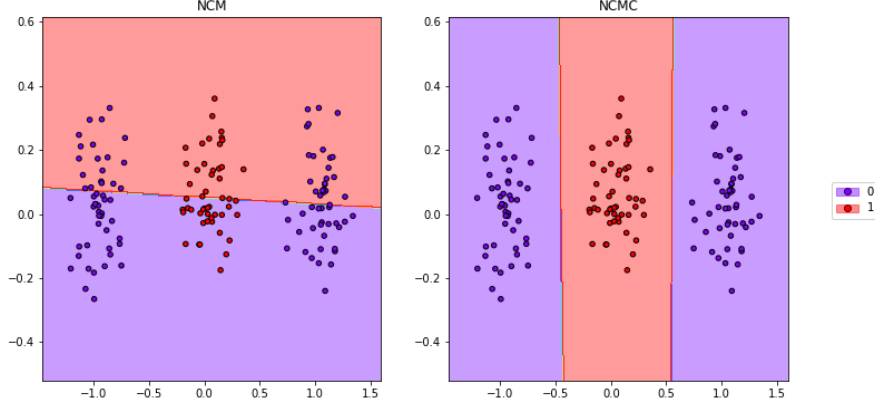


Fig. 14. Dataset where the NCM classifier does not provide good results, because the centroids of both classes are very close and both fall between the points of class 1. We will see that, by choosing more than one centroid in an appropriate way, we can classify this set as shown in right image.

algorithm that meets the most suitable conditions, besides being simple and efficient is  $k$ -Means.

To classify with NCMC, the use  $k$ -Means reduces to applying the segmentation algorithm within each subset of data associated with each of the classes of the problem. In this way, we obtain in a simple way the set of centroids we wanted for each class, and on which we can carry out the classification of new data simply by searching for the nearest centroid. For this algorithm, as it happens with  $k$ -Means, it is necessary to previously establish the number of centroids for each class. These numbers can be estimated by cross validation.

Once the NCMC classifier is defined, the distance learning process [Mensink et al. 2012] is similar to NCM. Following the notation used in NCMML, in this case, instead of a set of class centers  $\{\mu_c\}$ , with  $c \in \mathcal{C}$ , we have a set of centroids,  $\{m_{c_j}\}_{j=1}^{k_c}$ , with  $k_c \in \mathbb{N}$ , for each  $c \in \mathcal{C}$ . In this case, the probabilities associated with each class for the correct prediction of  $x \in \mathcal{X}$  are given by  $p_L(c|x) = \sum_{j=1}^{k_c} p_L(m_{c_j}|x)$ , where the centroids are those whose probability is defined by the softmax function

$$p_L(m_{c_j}|x) = \frac{\exp\left(-\frac{1}{2}\|L(x - m_{c_j})\|^2\right)}{\sum_{c \in \mathcal{C}} \sum_{i=1}^{k_c} \exp\left(-\frac{1}{2}\|L(x - m_{c_i})\|^2\right)}.$$

Again, we maximize the log-likelihood function  $\mathcal{L}(L) = \frac{1}{N} \sum_{i=1}^N p_L(y_i|x_i)$ , whose gradient is given by

$$\nabla \mathcal{L}(L) = \frac{1}{N} \sum_{i=1}^N \sum_{c \in \mathcal{C}} \sum_{j=1}^{k_c} \alpha_{ic_j} L(m_{c_j} - x_i)(m_{c_j} - x_i)^T,$$

where

$$\alpha_{ic_j} = p_L(m_{c_j}|x_i) - \mathbb{I}[y_i = c] \frac{p_L(m_{c_j}|x_i)}{\sum_{j'=1}^{k_c} p_L(m_{c_{j'}}|x_i)}.$$

The log-likelihood maximization by gradient methods is the task carried out by the distance learning technique for NCMC classifier, which we will call with the same name as the classifier.

#### B.4. Information Theory Based Algorithms

In this section we will study several distance metric learning algorithms based on information theory, specifically, in the Kullback-Leibler and Jeffrey divergences. Their working scheme is similar. First of all, they establish different probability distributions on the data, and then they try to bring them closer or further away by using the divergences. The algorithms we will study are ITML [Davis et al. 2007], DMLMJ [Nguyen et al. 2017] and MCML [Globerson and Roweis 2006].

**B.4.1. ITML.** ITML (*information theoretic metric learning*) [Davis et al. 2007] is a distance metric learning technique whose objective is to find a metric as close as possible to an initial distance, understanding this closeness from the point of view of relative entropy, as we will formulate later, making that metric satisfy certain similarity constraints for the trained data.

ITML starts with a dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , not necessarily labeled, but for which it is known that certain pairs of samples considered similar must be at a distance lower than or equal to  $u$ , and other pairs of samples considered not similar must be at a distance greater than or equal to  $l$ , where  $u, l \in \mathbb{R}^+$  are pre-defined constants, with relative small and large values, respectively, with respect to the dataset.

From the data with the indicated restrictions, ITML considers an initial distance corresponding to a positive definite matrix  $M_0$ , and tries to find a positive definite matrix  $M$ , as similar as possible to  $M_0$ , and that respects the imposed similarity constraints. The way to measure the similarity between  $M$  and  $M_0$  is done using information theory tools.

As we saw in Section A.3, there is a correspondence between positive definite matrices and multivariate gaussian distributions, if we fix the same mean vector  $\mu$  for every distribution. Given  $M \in S_d(\mathbb{R})^+$  we can then construct a normal distribution through its density function,

$$p(x|M) = \frac{1}{(2\pi)^{n/2} \det(M)^{1/2}} \exp(-(x - \mu)^T M^{-1}(x - \mu)).$$

Reciprocally, from this distribution, if we calculate the covariance matrix, we recover the matrix  $M$ . Using this correspondence, we will measure the closeness between  $M_0$  and  $M$  through the Kullback-Leibler divergence between their corresponding gaussian distributions, that is,

$$\text{KL}(p(x|M_0)||p(x|M)) = \int p(x|M_0) \log \frac{p(x|M_0)}{p(x|M)} dx.$$

Once we have defined the mechanism to measure the proximity between the metrics, we can formulate the optimization problem of the technique ITML. If we call  $S$  and  $D$  to the sets of pairs of indices on the elements of  $\mathcal{X}$  that represent the samples considered similar and not similar, respectively, and we start from the initial metric  $M_0$ , the problem is

$$\begin{aligned} \min_{M \in S_d(\mathbb{R})^+} \quad & \text{KL}(p(x|M_0)||p(x|M)) \\ \text{s.t.:} \quad & d_M(x_i, x_j) \leq u, \quad (i, j) \in S \\ & d_M(x_i, x_j) \geq l, \quad (i, j) \in D. \end{aligned} \tag{17}$$

We have seen in Theorem A.33 that the Kullback-Leibler divergence between two gaussian distributions with the same mean can be expressed in terms of the *log-det* matrix divergence. This allows us to reformulate Eq. 17 in a way that is easier to deal with computationally:

$$\begin{aligned} \min_{M \in S_d(\mathbb{R})^+} \quad & D_{ld}(M_0 \| M) \\ \text{s.t.:} \quad & \text{tr}(M(x_i - x_j)(x_i - x_j)^T) \leq u, \quad (i, j) \in S \\ & \text{tr}(M(x_i - x_j)(x_i - x_j)^T) \geq l, \quad (i, j) \in D. \end{aligned} \quad (18)$$

We may not be able to find a metric  $M$  that simultaneously satisfies every constraint, so the problem may not have a solution. Therefore, ITML introduces in Eq. 18 slack variables through which we obtain a problem whose optimization establishes a trade-off between the minimization of the divergence and the fulfillment of the constraints, in order to arrive to an approximate solution of the original problem, in case there is no solution for this. Finally, the computational technique used in the resolution of this optimization problem is the *Bregman projections method* discussed in Section A.1.2.

**B.4.2. DMLMJ.** DMLMJ (*distance metric learning through the maximization of the Jeffrey divergence*) [Nguyen et al. 2017] is another distance metric learning technique based on information theory. In this case, the tool that is used by DMLMJ is the Jeffrey divergence, to separate as much as possible the distribution associated with similar points from that associated to dissimilar points, in the sense that we will see below.

We consider the training set  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$  with corresponding labels  $y_1, \dots, y_N$ , and we set  $k \in \mathbb{N}$ . As we have already commented, DMLMJ tries to maximize, with respect to the Jeffrey divergence, the separation between distributions of similar and not similar points. To do this, we will introduce several concepts.

**Definition B.4.** Given  $x_i \in \mathcal{X}$ , the *k-positive neighborhood* of  $x_i$  is defined as the set of the  $k$  nearest neighbors of  $x_i$  in  $\mathcal{X} \setminus \{x_i\}$  whose class is the same as  $x_i$ . It is denoted by  $V_k^+(x_i)$ .

The *k-negative neighborhood* of  $x_i$  is defined as the set of the  $k$  nearest neighbors of  $x_i$  in  $\mathcal{X}$  whose class is different from that of  $x_i$ . It is denoted by  $V_k^-(x_i)$ .

The *k-positive difference space* of the labeled dataset is defined as the set

$$S = \{x_i - x_j : x_i \in \mathcal{X}, x_j \in V_k^+(x_i)\}.$$

Similarly, the *k-negative difference space* of the labeled dataset is defined as the set

$$D = \{x_i - x_j : x_i \in \mathcal{X}, x_j \in V_k^-(x_i)\}.$$

Sets  $S$  and  $D$  represent, therefore, the vectors with the differences between the samples in  $\mathcal{X}$  and its  $k$  nearest neighbors, from the same or a different class, respectively. We refer to  $P$  and  $Q$  as the probability distributions in the spaces  $S$  and  $D$ , respectively, assuming that they are multivariate gaussians. We will also assume that both distributions have zero mean. This assumption is reasonable, since in practice, in most cases, if  $x_i$  is a neighbor of  $x_j$ ,  $x_j$  is also a neighbor of  $x_i$ , then both differences will appear in the difference space, averaging zero. Finally, we will call the corresponding covariance matrices  $\Sigma_S$  and  $\Sigma_D$ , respectively.

If we now apply a linear transformation to the data,  $x \mapsto Lx$ , with  $L \in \mathcal{M}_{d' \times d}(\mathbb{R})$ , the transformed distributions will still have mean zero, and covariances  $L\Sigma_S L^T$  and  $L\Sigma_D L^T$ , respectively. We will call these distributions  $P_L$  and  $Q_L$ . The goal of DMLMJ is to find a transformation that maximizes the Jeffrey divergence between  $P_L$  and  $Q_L$ ,

that is, the problem to optimize is:

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} f(L) = \text{JF}(P_L \| Q_L) = \text{KL}(P_L \| Q_L) + \text{KL}(Q_L \| P_L).$$

As it was shown in Proposition A.35, Jeffrey divergence between the gaussian distributions  $P_L$  and  $Q_L$  can be rewritten as

$$f(L) = \frac{1}{2} \text{tr}((L\Sigma_S L^T)^{-1}(L\Sigma_D L^T) + (L\Sigma_D L^T)^{-1}(L\Sigma_S L^T)) - d'.$$

Since  $d'$  is constant, we obtain the equivalent problem

$$\max_{L \in \mathcal{M}_{d' \times d}(\mathbb{R})} J(L) = \text{tr}((L\Sigma_S L^T)^{-1}(L\Sigma_D L^T) + (L\Sigma_D L^T)^{-1}(L\Sigma_S L^T)).$$

Theorem A.20 tells us that, to maximize  $J(L)$ , we can choose the  $d'$  eigenvectors of  $\Sigma_S^{-1}\Sigma_D$ ,  $v_1, \dots, v_{d'}$  corresponding to the largest values of  $\lambda_i + 1/\lambda_i$ , with  $\lambda_i$  being the eigenvalue of  $\Sigma_S^{-1}\Sigma_D$  associated with  $v_i$ , and add this eigenvectors to the rows of  $L$ . The transformation  $L$  constructed from these eigenvectors determines the distance that is learned by the DMLMJ technique.

Finally, the only additional requirement necessary to complete the construction of  $L$  is the calculation of the covariance matrices  $\Sigma_S$  and  $\Sigma_D$ . Bearing in mind that it has been assumed that the mean of the distributions of  $S$  and  $D$  is 0, we can obtain these matrices quite simply from the difference vectors, as shown below:

$$\Sigma_S = \frac{1}{|S|} \sum_{i=1}^N \left[ \sum_{x_j \in V_k^+(x_i)} (x_i - x_j)(x_i - x_j)^T \right],$$

$$\Sigma_D = \frac{1}{|D|} \sum_{i=1}^N \left[ \sum_{x_j \in V_k^-(x_i)} (x_i - x_j)(x_i - x_j)^T \right].$$

Let us observe that we can also see this algorithm as a dimensionality reduction algorithm and even as an algorithm oriented to improve the nearest neighbors classifier, due to its local character.

**B.4.3. MCML.** MCML (*maximally collapsing metric learning*) [Globerson and Roweis 2006] is a supervised distance metric learning technique, based on the idea that if all the samples of the same class were projected to the same point, and data of different classes were projected to different points and sufficiently far away, we would have, over the projected data, an ideal class separation. Its purpose is to learn a distance metric that allows to collapse as much as possible, within the limitations of the metric, all the samples of the same class in a single point, arbitrarily far from the points where the samples of the remaining classes will collapse.

We consider the dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , with corresponding labels  $y_1, \dots, y_N$ . We want to learn a metric determined by  $M \in S_d(\mathbb{R})^+$  that tries to collapse the classes as much as possible according to the approach of the previous paragraph. The way to deal with this problem will consist once again in using the tools provided by the information theory. To do this, we first introduce a conditional distribution on the points of the dataset, analogous to that established in the case of NCA. If  $i, j \in \{1, \dots, N\}$ , with  $i \neq j$ , we define the probability that  $x_j$  will be classified with the class of  $x_i$  according to the distance between  $x_i$  and  $x_j$  as follows:

$$p^M(j|i) = \frac{\exp(-\|x_i - x_j\|_M^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|_M^2)}.$$

Furthermore, the ideal distribution we are looking for is a binary distribution for which the probability that a sample is correctly classified is 1, and 0 otherwise, that is,

$$p_0(j|i) \propto \begin{cases} 1, & y_i = y_j \\ 0, & y_i \neq y_j \end{cases}.$$

Note that during the training process we know the real classes of the data, therefore we can deal with this last probability. Besides, we can observe that if we get a metric  $M$  whose associated distribution  $p^M$  coincides with  $p_0$ , then, under very mild sufficiency conditions on the data, we will be able to collapse the classes in infinitely distant points.

Indeed, suppose there are at least  $r + 2$  samples in each class, where  $r$  is the rank of  $M$ , and that  $p^M(j|i) = p^0(j|i)$  for any  $i, j \in \{1, \dots, N\}$ . Then, on the one hand, from  $p^M(j|i) = 0$  for  $y_i \neq y_j$ , it follows that  $\exp(-\|x_i - x_j\|_M^2) = 0$ , which undoubtedly leads to  $x_i$  and  $x_j$  being infinitely distant when their classes are different. On the other hand, from  $p^M(j|i) \propto 1$  for any  $x_i, x_j$  with  $y_i = y_j$ , it follows that the value  $\exp(-\|x_i - x_j\|_M^2)$  is constant for all the members of the same class, and consequently, all the points in the same class are equidistant. As  $M$  has rank  $r$ , it is inducing a distance on a subspace of dimension  $r$ , where it is known that at most there can be  $r + 1$  different points and equidistant between them. Since we are assuming that there are at least  $r + 2$  points per class, all the points of the same class must have a distance of 0 between them with respect to  $M$ , thus collapsing into a single point.

Once both distributions are set, the objective of MCML is, as we have already commented, to approximate  $p^M(\cdot|i)$  to  $p_0(\cdot|i)$  as much as possible, for each  $i$ , using the relative entropy between both distributions. The optimization problem is, therefore, to minimize this divergence,

$$\min_{M \in S_d(\mathbb{R})_0^+} f(M) = \sum_{i=1}^N \text{KL} [p_0(\cdot|i) \| p^M(\cdot|i)].$$

We can rewrite the objective function in terms of elementary functions:

$$\begin{aligned} f(M) &= \sum_{i=1}^N \sum_{j=1}^N p_0(j|i) \log \frac{p_0(j|i)}{p^M(j|i)} = \sum_{i=1}^N \sum_{j: y_i=y_j} \log \frac{1}{p^M(j|i)} \\ &= \sum_{i=1}^N \sum_{j: y_i=y_j} -\log p^M(j|i) \\ &= -\sum_{i=1}^N \sum_{j: y_i=y_j} \left( -\|x_i - x_j\|_M^2 - \log \sum_{k \neq i} \exp(-\|x_i - x_k\|^2) \right) \\ &= \sum_{i=1}^N \sum_{j: y_i=y_j} \|x_i - x_j\|_M^2 + \sum_{i=1}^N \log \sum_{k \neq i} \exp(-\|x_i - x_k\|^2). \end{aligned} \tag{19}$$

This function is differentiable, and each summand of the previous expression is convex in  $M$ , the first because it is a distance function in  $M$  (which is affine), and the second because it is a *log-sum-exp* function (see Boyd and Vandenberghe [2004], sec. 3.1.5) composed with a distance function. In addition, the restriction  $M \in S_d(\mathbb{R})_0^+$  is convex, so we can use the projected gradient descent algorithm with projections onto the positive semidefinite cone to optimize the objective function. This requires an expression of the gradient of the objective function, which can be calculated from its expression in

Eq. 19:

$$\nabla f(M) = \sum_{i,j: y_i=y_j} (x_i - x_j)^T (x_i - x_j) - \sum_i \frac{-\sum_{k \neq i} (x_i - x_k)^T (x_i - x_k) \exp(-\|x_i - x_k\|_M^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|_M^2)}.$$

### B.5. Other Distance Metric Learning Techniques

In this section we will study some different proposals for distance metric learning techniques. The algorithms we will analyze are LSI [Xing et al. 2003], DML-eig [Ying and Li 2012] and LDML [Guillaumin et al. 2009].

**B.5.1. LSI.** LSI (*learning with side information*) [Xing et al. 2003], also sometimes referred to as MMC (*Mahalanobis metric for clustering*) is a distance metric learning technique that works with a dataset that is not necessarily labeled, which contains certain pairs of samples that are known to be similar and, optionally, pairs of samples that are known not to be similar. It is possibly one of the first algorithms that has helped make the concept of distance metric learning more well known.

LSI tries to learn a metric  $M$  that respects this additional information. This is why it can be used both in supervised learning, where similar pairs will correspond to data with the same label, and in unsupervised learning with similarity constraints, such as, for example, clustering problems where it is known that certain samples must be grouped in the same cluster.

We now formulate the problem to be optimized by LSI. Suppose we have the dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , and we know additionally the set  $S = \{(x_i, x_j) \in \mathcal{X} \times \mathcal{X} : x_i \text{ and } x_j \text{ are similar}\}$ . In addition, we may know the set  $D = \{(x_i, x_j) \in \mathcal{X} \times \mathcal{X} : x_i \text{ and } x_j \text{ are dissimilar}\}$ . If we do not have the latter, we can take  $D$  as the complement of  $S$  in  $\mathcal{X} \times \mathcal{X}$ .

The first intuition to address this problem, given the information we have, is to minimize the distances between pairs of similar points, that is, to minimize  $\sum_{(x_i, x_j) \in S} \|x_i - x_j\|_M^2$ , where  $M \in S_d(\mathbb{R})_0^+$ . However, this will lead us to the solution  $M = 0$ , which would not give us any productive information. That is why LSI adds the additional constraint  $\sum_{(x_i, x_j) \in D} \|x_i - x_j\|_M \geq 1$ , which leads us to the optimization problem

$$\begin{aligned} \min_M \quad & \sum_{(x_i, x_j) \in S} \|x_i - x_j\|_M^2 \\ \text{s.t.} \quad & \sum_{(x_i, x_j) \in D} \|x_i - x_j\|_M \geq 1 \\ & M \in S_d(\mathbb{R})_0^+. \end{aligned}$$

Note several observations regarding this formula. First, the choice of constant 1 in the constraint is irrelevant; if we choose any constant  $c > 0$  we get a metric proportional to  $M$ . Secondly, the optimization problem is convex, because the sets determined by the restrictions are convex and the function to optimize is also convex. Finally, we may consider a restriction on the set  $D$  of the form  $\sum_{(x_i, x_j) \in D} \|x_i - x_j\|_M^2 \geq 1$ . However, it is possible to rewrite that problem into a formulation similar to that used on the 2-class LDA, where the metric learned would have a rank of 1, which may not be optimal.



To easily optimize this problem, authors propose the equivalent problem

$$\begin{aligned}
& \max_M \quad \sum_{(x_i, x_j) \in D} \|x_i - x_j\|_M \\
& \text{s.a.:} \quad \sum_{(x_i, x_j) \in S} \|x_i - x_j\|_M^2 \leq 1 \\
& \quad M \in S_d(\mathbb{R})_0^+.
\end{aligned} \tag{20}$$

This problem with two convex constraints can be solved by a projected gradient ascent method. In this problem, constraints are easy to satisfy separately. The first constraint consists of a projection onto an affine half-space, while the second constraint consists of a projection onto the positive semidefinite cone. The method of iterated projections makes it possible to fulfill both restrictions by repeatedly projecting onto both sets until convergence is obtained.

**B.5.2. DML-eig.** DML-eig (*distance metric learning with eigenvalue optimization*) [Ying and Li 2012] is a distance metric learning algorithm inspired by the LSI algorithm of the previous section, proposing a very similar optimization problem but offering a completely different resolution method, based on eigenvalue optimization.

We consider, as in the previous case, a training dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , for which we know two sets of pairs,  $S$  and  $D$ , of data considered similar and dissimilar, respectively. In the previous section, in order to optimize Eq. 20 an ascending gradient method with iterated projections was proposed, which may take a long time to converge. DML-eig proposal consists of a slight modification of the objective function, keeping the same constraints, which leads us to the problem

$$\begin{aligned}
& \max_M \quad \min_{(x_i, x_j) \in D} \|x_i - x_j\|_M^2 \\
& \text{s.t.:} \quad \sum_{(x_i, x_j) \in S} \|x_i - x_j\|_M^2 \leq 1 \\
& \quad M \in S_d(\mathbb{R})_0^+.
\end{aligned} \tag{21}$$

To address this problem, it is useful to introduce a notation that simplifies the indexing of the data. First, we will denote  $X_{ij} = (x_i - x_j)(x_i - x_j)^T$  to the outer products between the differences of the elements in  $\mathcal{X}$ . To access pairs of elements  $(i, j)$  we will use a single index  $\tau \equiv (i, j)$ . This index can be assumed ordered when necessary, to access the components of a vector of appropriate size. The previous outer product  $X_{ij}$  can also be written as  $X_\tau$ . Finally, for sets  $S$  and  $D$ , we also assume that they are made by indexes  $\tau$  associated with a pair  $(i, j)$  such that  $x_i$  and  $x_j$  are similar or dissimilar, respectively. Thus, if we denote  $X_S = \sum_{(i, j) \in S} X_{ij}$ , Eq. 21 can be rewritten in terms of Frobenius dot product as

$$\begin{aligned}
& \max_M \quad \min_{\tau \in D} \langle X_\tau, M \rangle \\
& \text{s.t.:} \quad \langle X_S, M \rangle \leq 1 \\
& \quad M \in S_d(\mathbb{R})_0^+.
\end{aligned} \tag{22}$$

Let us see how the formulation of the problem we are looking for is established in terms of eigenvalue optimization. For each symmetric matrix  $X \in S_d(\mathbb{R})$  we denote its highest eigenvalue as  $\lambda_{\max}(X)$ . Associated with the set  $D$  of dissimilar pairs we will

define the simplex

$$\Delta = \left\{ u \in \mathbb{R}^{|D|} : u_\tau \geq 0 \ \forall \tau \in D, \sum_{\tau \in D} u_\tau = 1 \right\}.$$

We also consider the set

$$\mathcal{P} = \{M \in \mathcal{M}_d(\mathbb{R})_0^+ : \text{tr}(M) = 1\}.$$

$\mathcal{P}$  is the intersection of the positive semidefinite cone with an affine subspace of  $\mathcal{M}_d(\mathbb{R})$ . Sets with this structure are known as *spectrahedra*.

So, if  $X_S$  is positive semidefinite, and we define, for each  $\tau \in D$ ,  $\tilde{X}_\tau = X_S^{-1/2} X_\tau X_S^{-1/2}$ , we can prove [Ying and Li 2012] that the problem given by Eq. 22 is equivalent to the following problem:

$$\max_{S \in \mathcal{P}} \min_{u \in \Delta} \sum_{\tau \in D} u_\tau \langle \tilde{X}_\tau, S \rangle,$$

which in turn can be rewritten as an eigenvalue optimization problem:

$$\min_{u \in \Delta} \max_{S \in \mathcal{P}} \left\langle \sum_{\tau \in D} u_\tau \tilde{X}_\tau, S \right\rangle = \min_{u \in \Delta} \lambda_{\max} \left( \sum_{\tau \in D} u_\tau \tilde{X}_\tau \right). \quad (23)$$

The problem of minimizing the largest eigenvalue of a symmetric matrix is well-known and there are some iterative methods that allow this minimum to be reached [Overton 1988]. Furthermore, Ying and Li [2012] also propose an algorithm to solve the problem  $\max_{S \in \mathcal{P}} \min_{u \in \Delta} \sum_{\tau \in D} u_\tau \langle \tilde{X}_\tau, S \rangle + \mu \sum_{\tau \in D} u_\tau \log u_\tau$ , where  $\mu > 0$  is a smoothing parameter, by means of which the problem in Eq. 23 can be approximated.

**B.5.3. LDML.** LDML (*logistic discriminant metric learning*) [Guillaumin et al. 2009] is a distance metric learning algorithm in which the optimization model makes use of the logistic function. Authors affirm that this technique is quite useful to learn distances on sets of labeled images, being able to be used therefore in problems like face identification.

Recall that the *logistic* or *sigmoid* function is the map  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  given by

$$\sigma(x) = \frac{1}{1 + e^{-x}}.$$

This function presents a graph with a sigmoidal shape, is differentiable, strictly increasing and takes values between 0 and 1, reaching these values in their limits at infinity. These properties allow the logistic function to be the cumulative distribution function of a random variable, which gives it an important probabilistic utility. Its graph presents an asymptotic behaviour from small values (in absolute value), with an exponential growth in zones close to zero. This makes logistic function very useful for modeling binary signals. It also presents a derivative that is easy to calculate, and can be expressed in terms of the logistic function itself,  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$ .

Suppose we have the dataset  $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ , with corresponding labels  $y_1, \dots, y_N$ . In LDML, logistic function is used to define a probability, which will assign the greater probability the smaller the distance between points. To measure the distance, LDML will use a positive semidefinite matrix, resulting in the expression of the probability as

$$p_{ij,M} = \sigma(b - d_M(x_i, x_j)^2),$$

where  $b$  is a positive threshold value that will determine the maximum value achievable by the logistic function, and that can be estimated by cross validation. Associated

with this probability, we can define a random variable that follows a Bernoulli distribution, and that takes the values 0 and 1, according to whether the pair  $(x_i, x_j)$  belongs to the same class. This distribution is determined by the probability mass function

$$f_{ij,M}(x) = (p_{ij,M})^x (1 - p_{ij,M})^{1-x}, \quad x \in \{0, 1\}.$$

The function that LDML tries to maximize is the log-likelihood of the previous distribution for the given dataset, that is,

$$\mathcal{L}(M) = \sum_{i,j=1}^N y_{ij} \log p_{ij,M} + (1 - y_{ij}) \log(1 - p_{ij,M}),$$

where  $y_{ij}$  is a binary variable that takes the value 1 if  $y_i = y_j$  and 0 otherwise. This function is differentiable and concave (it is a positive combination of functions that can be expressed as a minus log-sum-exp function, which is concave), so we have a convex maximization problem. Keeping in mind the properties of the logistic function, if  $x_{ij} \equiv (x_i - x_j)(x_i - x_j)^T$  and  $p_{ij} \equiv p_{ij,M}$ , the gradient has the expression

$$\begin{aligned} \nabla \mathcal{L}(M) &= \sum_{i,j=1}^N y_{ij} \frac{-x_{ij} p_{ij} (1 - p_{ij})}{p_{ij}} + (1 - y_{ij}) \frac{x_{ij} p_{ij} (1 - p_{ij})}{1 - p_{ij}} \\ &= \sum_{i,j=1}^N -y_{ij} x_{ij} (1 - p_{ij}) + (1 - y_{ij}) x_{ij} p_{ij} \\ &= \sum_{i,j=1}^N x_{ij} ((1 - y_{ij}) p_{ij} - (1 - p_{ij}) y_{ij}) \\ &= \sum_{i,j=1}^N x_{ij} (p_{ij} - y_{ij}), \end{aligned}$$

The projected gradient method with projections onto the positive semidefinite cone is the semidefinite programming algorithm that is used in LDML to obtain the metric that optimizes its objective function.

## B.6. Kernel Distance Metric Learning

In this part we will analyze some of the kernelized versions of the algorithms presented throughout this section. An introduction to the use of the kernel trick for distance metric learning was already made in Section 3.6. Below we will study the kernel algorithms for LMNN, ANMM, DMLMJ and LDA.

**B.6.1. KLMNN.** KLMNN [Torresani and Lee 2007; Weinberger and Saul 2009] is the kernelized version of LMNN. In it, the data in  $\mathcal{X}$  is sent to the feature space to learn in that space a distance that minimizes the objective function set in the LMNN problem.

Although the problem formulated in the non-kernelized version was made with respect to a positive semidefinite matrix  $M$ , using the error function given in Eq. 16, when working in feature spaces we are more interested in dealing with a linear map, even if the convexity of the problem is lost, in order to be able to use the representer theorem. Therefore, adapting the error function proposed in Eq. 15 to the feature

space, the LMNN problem for the kernelized version consists of

$$\begin{aligned} \min_{L \in \mathcal{L}(\mathcal{F}, \mathbb{R}^d)} \quad \varepsilon(L) = & (1 - \mu) \sum_{i=1}^N \sum_{j \rightsquigarrow i} \|L(\phi(x_i) - \phi(x_j))\|^2 \\ & + \mu \sum_{i=1}^N \sum_{j \rightsquigarrow i} \sum_{l=1}^N (1 - y_{il}) [1 + \|L(\phi(x_i) - \phi(x_j))\|^2 - \|L(\phi(x_i) - \phi(x_l))\|^2]_+. \end{aligned}$$

As a consequence of the representer theorem, it follows that, for each  $x_i \in \mathcal{X}$ ,  $L\phi(x) = AK_{\cdot i}$ , where  $A \in \mathcal{M}_{d' \times N}(\mathbb{R})$  is the matrix given by the representer theorem, and  $K_{\cdot i}$  represents the  $i$ -th column of the kernel matrix for the training set. Using this in the error expression, we obtain

$$\begin{aligned} & (1 - \mu) \sum_{i=1}^N \sum_{j \rightsquigarrow i} \|L(\phi(x_i) - \phi(x_j))\|^2 \\ & + \mu \sum_{i=1}^N \sum_{j \rightsquigarrow i} \sum_{l=1}^N (1 - y_{il}) [1 + \|L(\phi(x_i) - \phi(x_j))\|^2 - \|L(\phi(x_i) - \phi(x_l))\|^2]_+ \\ & = (1 - \mu) \sum_{i=1}^N \sum_{j \rightsquigarrow i} \|A(K_{\cdot i} - K_{\cdot j})\|^2 \\ & + \mu \sum_{i=1}^N \sum_{j \rightsquigarrow i} \sum_{l=1}^N (1 - y_{il}) [1 + \|A(K_{\cdot i} - K_{\cdot j})\|^2 - \|A(K_{\cdot i} - K_{\cdot l})\|^2]_+. \end{aligned}$$

The above expression depends only on  $A$  and kernel functions, and minimizing it as a function of  $A$  (we will denote it  $\varepsilon(A)$ ) we get the same value as minimizing  $\varepsilon(L)$ . Note also that the expression  $\varepsilon(A)$  also requires the calculation of target neighbors and impostors, but these depend only on the distances in the feature space, which, as we have already seen, are computable, as shown in Eq. 1. Therefore, all the components of  $\varepsilon(A)$  are computationally manipulable, so if we apply a gradient descent method on  $\varepsilon(A)$  we can reduce the value of the objective function, always keeping in mind that we can be stuck in a local optimum, because the problem is not convex. Finally, once a matrix  $A$  that minimizes  $\varepsilon(A)$  is found, we will have determined the corresponding map  $L$  thanks to the representer theorem, and we can use  $A$  together with the kernel functions to transform new data.

**B.6.2. KANMM.** KANMM [Wang and Zhang 2007] is the kernelized version of ANMM. In it, the data in  $\mathcal{X}$  is sent to the feature space via the map  $\phi: \mathbb{R}^d \rightarrow \mathcal{F}$ , where ANMM is applied to obtain the linear map we are looking for.

Recall that the first step for the application of ANMM was to obtain the homogeneous and heterogeneous neighborhoods for each sample  $x_i \in \mathcal{X}$ . Note that for this calculation it is only necessary to compare distances in the feature space, which we have seen can be done thanks to the kernel function, through Eq. 1. We will denote the neighborhoods in the feature space as  $N_{\phi(x_i)}^o$  y  $N_{\phi(x_i)}^e$ , respectively, for each  $x_i$ .

The scatter and compactness matrices (or endomorphisms, more in general) in the feature space are given by

$$S^\phi = \sum_{i,k: \phi(x_k) \in N_{\phi(x_i)}^e} \frac{(\phi(x_i) - \phi(x_k))(\phi(x_i) - \phi(x_k))^T}{|N_{\phi(x_i)}^e|}$$

$$C^\phi = \sum_{i,j: \phi(x_j) \in N_{\phi(x_i)}^o} \frac{(\phi(x_i) - \phi(x_j))(\phi(x_i) - \phi(x_j))^T}{|N_{\phi(x_i)}^o|}.$$

The problem to be optimized is therefore expressed as

$$\begin{aligned} \max_{L \in \mathcal{L}(\mathcal{F}, \mathbb{R}^{d'})} \quad & \text{tr}(L(S^\phi - C^\phi)L^T) \\ \text{s.t.:} \quad & LL^T = I. \end{aligned} \quad (24)$$

According to the representer theorem,  $L\varphi(x_i) = AK_{\cdot i}$ , where  $A$  is the matrix of coefficients of the representation theorem and  $K_{\cdot i}$  represents the  $i$ -th column of the kernel matrix for the training set. Then,

$$L(\phi(x_i) - \phi(x_j))(\phi(x_i) - \phi(x_j))^T L^T = A(K_{\cdot i} - K_{\cdot j})(K_{\cdot i} - K_{\cdot j})^T A^T,$$

and if we consider the matrices

$$\tilde{S}^\phi = \sum_{i,k: \phi(x_k) \in N_{\phi(x_i)}^e} \frac{(K_{\cdot i} - K_{\cdot k})(K_{\cdot i} - K_{\cdot k})^T}{|N_{\phi(x_i)}^e|}$$

$$\tilde{C}^\phi = \sum_{i,j: \phi(x_j) \in N_{\phi(x_i)}^o} \frac{(K_{\cdot i} - K_{\cdot j})(K_{\cdot i} - K_{\cdot j})^T}{|N_{\phi(x_i)}^o|},$$

it follows that the average neighborhood margin is given by

$$\gamma^L = \text{tr}(L(S^\phi - C^\phi)L^T) = \text{tr}(LS^\phi L^T - LC^\phi L^T) = \text{tr}(A\tilde{S}^\phi A^T - A\tilde{C}^\phi A^T) = \text{tr}(A(\tilde{S}^\phi - \tilde{C}^\phi)A^T).$$

If we impose the restriction  $AA^T = I$ , Theorem A.18 tells us again that we can take matrix  $A$  that which contains as rows the eigenvectors of  $\tilde{S}^\phi - \tilde{C}^\phi$  corresponding to its  $d'$  largest eigenvalues. Observe that we can calculate both matrices from the kernel function, and the matrix  $A$  we obtain determines the linear map, as a consequence of the representer theorem. Therefore, we have finally obtained a kernel-based method for applying ANMM in feature spaces.

**B.6.3. KDMLMJ.** KDMLMJ [Nguyen et al. 2017] is the kernelized version of DMLMJ. In it, the data in  $\mathcal{X}$  is sent to the feature space, where a distance is learned after applying DMLMJ.

Again, it is possible to calculate the  $k$ -positive and  $k$ -negative neighborhoods,  $V_k^+(\phi(x_i))$  and  $V_k^-(\phi(x_i))$ , for each  $x_i \in \mathcal{X}$ , thanks to Eq. 1. It is not the same with the endomorphisms associated with the difference spaces,

$$\Sigma_S^\phi = \frac{1}{|S|} \sum_{i=1}^N \left[ \sum_{\phi(x_j) \in V_k^+(\phi(x_i))} (\phi(x_i) - \phi(x_j))(\phi(x_i) - \phi(x_j))^T \right]$$

$$\Sigma_D^\phi = \frac{1}{|D|} \sum_{i=1}^N \left[ \sum_{\phi(x_j) \in V_k^-(\phi(x_i))} (\phi(x_i) - \phi(x_j))(\phi(x_i) - \phi(x_j))^T \right].$$

The optimization problem is given by

$$\max_{L \in \mathcal{L}(\mathcal{F}, \mathbb{R}^{d'})} J(L) = \text{tr} \left( (L\Sigma_S^\phi L^T)^{-1} (L\Sigma_D^\phi L^T) + (L\Sigma_D^\phi L^T)^{-1} (L\Sigma_S^\phi L^T) \right).$$

Again we have, as a consequence of the representer theorem, that  $L\phi(x_i) = AK_{\cdot i}$  for each  $x_i \in \mathcal{X}$ , where  $A$  is the matrix provided by the representer theorem, and  $K_{\cdot i}$  is the  $i$ -th column of the kernel matrix for the training set. If, reasoning as in the previous section, we define the matrices

$$U = \frac{1}{|S|} \sum_{i=1}^N \left[ \sum_{\phi(x_j) \in V_k^+(\phi(x_i))} (K_{\cdot i} - K_{\cdot j})(K_{\cdot i} - K_{\cdot j})^T \right]$$

$$V = \frac{1}{|D|} \sum_{i=1}^N \left[ \sum_{\phi(x_j) \in V_k^-(\phi(x_i))} (K_{\cdot i} - K_{\cdot j})(K_{\cdot i} - K_{\cdot j})^T \right],$$

we obtain that

$$\text{tr} \left( (L\Sigma_S^\phi L^T)^{-1} (L\Sigma_D^\phi L^T) + (L\Sigma_D^\phi L^T)^{-1} (L\Sigma_S^\phi L^T) \right) =$$

$$\text{tr} \left( (AUA^T)^{-1} (AVA^T) + (AVA^T)^{-1} (AUA^T) \right).$$

As with DMLMJ, Theorem A.20 tells us that we can find a matrix  $A$  that maximizes this last equality by taking the eigenvectors of  $U^{-1}V$  for which the value  $\lambda + 1/\lambda$  is maximized, where  $\lambda$  is the associated eigenvalue. As matrices  $U$  and  $V$  can be obtained from the kernel function, and  $A$  determines  $L$  by the representer theorem, we have obtained an algorithm for the application of DMLMJ in the feature space.

**B.6.4. KDA.** KDA (*kernel discriminant analysis*) [Mika et al. 1999] is the kernelized version of linear discriminant analysis. The kernelization of this algorithm will make it possible to find non-linear directions that nicely separate the data according to the criteria established in the discriminant analysis. Once again, we send the data in  $\mathcal{X}$  to the feature space using the mapping  $\phi: \mathbb{R}^d \rightarrow \mathcal{F}$ . On that space we will apply linear discriminant analysis.

Suppose, as in LDA, that the set of possible classes is  $\mathcal{C}$ , of cardinal  $r$ , and for each  $c \in \mathcal{C}$  we define  $\mathcal{C}_c = \{i \in \{1, \dots, N\} : y_i = c\}$  and  $N_c = |\mathcal{C}_c|$ , with  $\mu_c^\phi$  the mean vector of the class  $c$ , and  $\mu^\phi$  the mean vector of the whole dataset, considering it within the feature space. The problem we want to solve in this case is

$$\max_{L \in \mathcal{L}(\mathcal{F}, \mathbb{R}^{d'})} \text{tr} \left( (LS_w^\phi L^T)^{-1} (LS_b^\phi L^T) \right), \quad (25)$$

where  $S_b^\phi$  and  $S_w^\phi$  are the operators that measure the between-class and within-class scatter, respectively, and are given by

$$S_b^\phi = \sum_{c \in \mathcal{C}} (\mu_c^\phi - \mu^\phi)(\mu_c^\phi - \mu^\phi)^T$$

$$S_w^\phi = \sum_{c \in \mathcal{C}} \sum_{i \in \mathcal{C}_c} (\phi(x_i) - \mu_c^\phi)(\phi(x_i) - \mu_c^\phi)^T.$$

Again, we use the representer theorem, so that if  $L \in \mathcal{L}(\mathcal{F}, \mathbb{R}^{d'})$ , then, for each  $x \in \mathbb{R}^d$ ,

$$L\phi(x) = A \begin{pmatrix} K(x_1, x) \\ \vdots \\ K(x_N, x) \end{pmatrix},$$

where  $A$  is in the conditions of the representer theorem. Let us look again for an expression of the problem given in Eq. 25 that depends only on the kernel function and the matrix  $A$ . To do this, we have to observe that for the mean vectors of each class we have

$$L\mu_c^\phi = L \left( \frac{1}{N_c} \sum_{i \in \mathcal{C}_c} \phi(x_i) \right) = \frac{1}{N_c} \sum_{i \in \mathcal{C}_c} L\phi(x_i) = \frac{1}{N_c} \sum_{i \in \mathcal{C}_c} AK_{.i},$$

where  $K_{.i}$  is the  $i$ -th column of the kernel matrix for the training set. Similarly, for the global mean vector, we have

$$L\mu^\phi = \frac{1}{N} \sum_{i=1}^N AK_{.i}.$$

Consequently,

$$\begin{aligned} L(\mu_c^\phi - \mu^\phi)(\mu_c^\phi - \mu^\phi)^T L^T &= (L\mu_c^\phi - L\mu^\phi)(L\mu_c^\phi - L\mu^\phi)^T \\ &= \left( \frac{1}{N_c} \sum_{i \in \mathcal{C}_c} AK_{.i} - \frac{1}{N} \sum_{i=1}^N AK_{.i} \right) \left( \frac{1}{N_c} \sum_{i \in \mathcal{C}_c} AK_{.i} - \frac{1}{N} \sum_{i=1}^N AK_{.i} \right)^T. \end{aligned}$$

Note that the last expression depends only on  $A$  and the kernel function. Moreover, for  $x_i \in \mathcal{X}$  with  $y_i = c$ , we have

$$\begin{aligned} L(\phi(x_i) - \mu_c^\phi)(\phi(x_i) - \mu_c^\phi)^T L^T &= (L\phi(x_i) - L\mu_c^\phi)(L\phi(x_i) - L\mu_c^\phi)^T \\ &= \left( AK_{.i} - \frac{1}{N_c} \sum_{j \in \mathcal{C}_c} AK_{.j} \right) \left( AK_{.i} - \frac{1}{N_c} \sum_{j \in \mathcal{C}_c} AK_{.j} \right)^T \\ &= \left( AK_{.i} - \frac{1}{N_c} \sum_{j \in \mathcal{C}_c} AK_{.j} \right) \left( K_{.i}^T A^T - \frac{1}{N_c} \sum_{j \in \mathcal{C}_c} K_{.j}^T A^T \right) \\ &= AK_{.i} K_{.i}^T A^T - \frac{1}{N_c} \sum_{j \in \mathcal{C}_c} AK_{.i} K_{.j}^T A^T - \frac{1}{N_c} \sum_{j \in \mathcal{C}_c} AK_{.j} K_{.i}^T A^T + \frac{1}{N_c^2} \sum_{j \in \mathcal{C}_c} \sum_{l \in \mathcal{C}_c} AK_{.j} K_{.l}^T A^T. \end{aligned}$$

By summing in  $i \in \mathcal{C}_c$ , we obtain

$$\begin{aligned}
& \sum_{i \in \mathcal{C}_c} L(\phi(x_i) - \mu_c^\phi)(\phi(x_i) - \mu_c^\phi)^T L^T \\
&= \sum_{i \in \mathcal{C}_c} \left[ AK_{.i} K_{.i}^T A^T - \frac{1}{N_c} \sum_{j \in \mathcal{C}_c} AK_{.i} K_{.j}^T A^T - \frac{1}{N_c} \sum_{j \in \mathcal{C}_c} AK_{.j} K_{.i}^T A^T + \frac{1}{N_c^2} \sum_{j \in \mathcal{C}_c} \sum_{l \in \mathcal{C}_c} AK_{.j} K_{.l}^T A^T \right] \\
&= \sum_{i \in \mathcal{C}_c} AK_{.i} K_{.i}^T A^T - \frac{2}{N_c} \sum_{i \in \mathcal{C}_c} \sum_{j \in \mathcal{C}_c} AK_{.i} K_{.j}^T A^T + \frac{1}{N_c^2} \sum_{i \in \mathcal{C}_c} \sum_{j \in \mathcal{C}_c} \sum_{l \in \mathcal{C}_c} AK_{.j} K_{.l}^T A^T \\
&= \sum_{i \in \mathcal{C}_c} AK_{.i} K_{.i}^T A^T - \frac{2}{N_c} \sum_{i \in \mathcal{C}_c} \sum_{j \in \mathcal{C}_c} AK_{.i} K_{.j}^T A^T + \frac{N_c}{N_c^2} \sum_{j \in \mathcal{C}_c} \sum_{l \in \mathcal{C}_c} AK_{.j} K_{.l}^T A^T \\
&= \sum_{i \in \mathcal{C}_c} AK_{.i} K_{.i}^T A^T - \frac{1}{N_c} \sum_{i \in \mathcal{C}_c} \sum_{j \in \mathcal{C}_c} AK_{.i} K_{.j}^T A^T \\
&= AK_c K_c^T A^T - AK_c \left( \frac{1}{N_c} \mathbb{1} \right) K_c^T A^T \\
&= AK_c \left( I - \frac{1}{N_c} \mathbb{1} \right) K_c^T A^T,
\end{aligned}$$

where  $\mathbb{1} \in \mathcal{M}_{N_c}(\mathbb{R})$  is a square matrix with the value 1 in all its entries, and  $K_c \in \mathcal{M}_{N \times N_c}$  is a kernel matrix whose entries are the values of the kernel function between all the samples in  $\mathcal{X}$  and the samples with class  $c$ . Again, this last expression depends only on  $A$  and the kernel function.

If we finally define

$$\begin{aligned}
U_c &= \frac{1}{N_c} \sum_{i \in \mathcal{C}_c} K_{.i} \in \mathbb{R}^N, c \in \mathcal{C} \\
U_\mu &= \frac{1}{N} \sum_{j=1}^N K_{.i} \in \mathbb{R}^N \\
U &= \sum_{c \in \mathcal{C}} N_c (U_c - U_\mu)(U_c - U_\mu)^T \in S_N(\mathbb{R}) \\
V &= \sum_{c \in \mathcal{C}} K_c \left( I - \frac{1}{N_c} \mathbb{1} \right) K_c^T \in S_N(\mathbb{R}),
\end{aligned}$$

we can conclude that

$$\text{tr} \left( (LS_w^\phi L^T)^{-1} (LS_b^\phi L^T) \right) = \text{tr} \left( (AV A^T)^{-1} (AU A^T) \right),$$

where  $U$  and  $V$  are computable using the kernel function. Therefore, we obtain a problem equivalent to the original given in Eq. 25, but in terms of  $A$ , for which Theorem A.19 states that, if  $U$  is positive definite, we can maximize the value of the trace by taking as rows of  $A$  the eigenvectors of  $V^{-1}U$  corresponding to its  $d'$  largest eigenvalues. In this way, since  $A$  determines  $L$  thanks to the representer theorem, we obtain a kernel-based method for the application of discriminant analysis in feature spaces.