

Kernel-based Distance Metric Learning in the Output Space

Cong Li, Michael Georgiopoulos and Georgios C. Anagnostopoulos

CL, MG: EE & CS Dept., University of Central Florida; GCA: ECE Dept., Florida Institute of Technology

congli@eecs.ucf.edu, michaelg@ucf.edu and georgio@fit.edu

Abstract

In this paper we present two related, kernel-based Distance Metric Learning (DML) methods. Their respective models non-linearly map data from their original space to an output space, and subsequent distance measurements are performed in the output space via a Mahalanobis metric. The dimensionality of the output space can be directly controlled to facilitate the learning of a low-rank metric. Both methods allow for simultaneous inference of the associated metric and the mapping to the output space, which can be used to visualize the data, when the output space is 2- or 3-dimensional. Experimental results for a collection of classification tasks illustrate the advantages of the proposed methods over other traditional and kernel-based DML approaches.

Keywords: Distance Metric Learning, Kernel Methods, Reproducing Kernel Hilbert Space for Vector-valued Functions

1 Introduction

Distance Metric Learning (DML) has become an active research area due to the fact that many machine learning models and algorithms depend on metric calculations. Considering plain Euclidean distances between samples may not be a suitable approach for some practical problems, *e.g.*, for k -Nearest Neighbors (KNN) classification, where a metric other than the Euclidean may yield higher recognition rates. Hence, it may be important to learn an appropriate metric for the learning problem at hand. DML aims to address this problem, *i.e.*, to infer a parameterized metric from the available training data that maximizes the performance of a model.

Most of past DML research focuses specifically on learning a weighted Euclidean metric, also known as the Mahalanobis distance (*e.g.* see [13]), or generalizations of it, where the weights are inferred from the data. For elements \mathbf{x}, \mathbf{x}' of a finite-dimensional Euclidean space \mathbb{R}^m , the Mahalanobis distance is defined as $d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_{\mathbf{A}} \triangleq \sqrt{(\mathbf{x} - \mathbf{x}')^T \mathbf{A} (\mathbf{x} - \mathbf{x}')}$, where $\mathbf{A} = \mathbf{A}^T \succeq \mathbf{0}$, *i.e.* $\mathbf{A} \in \mathbb{R}^{m \times m}$ is symmetric positive semi-definite matrix of weights to be determined. Note that when \mathbf{A} is not strictly positive definite, it defines a pseudo-metric in \mathbb{R}^m . An obvious DML approach is to learn this metric in the data's native space, which is tantamount to first linearly transforming the data via a matrix \mathbf{L} , such that $\mathbf{A} = \mathbf{L}^T \mathbf{L}$, and then measuring distances using the standard Euclidean metric $\|\cdot\|_2$.

One possible alternative worth exploring is to search for a non-linear transform prior to measuring Mahalanobis distances, so that performance may improve over the case, where a linear transformation is used. Towards this end, efforts have been recently made to develop kernel-based DML approaches. If \mathcal{X} is the original (native) data space, most of these methods choose an appropriate (positive definite) scalar kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, which gives rise to a Reproducing Kernel Hilbert Space (RKHS) \mathcal{H} of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ with inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$. This inner product satisfies the (reproducing) property that, for any $x, x' \in \mathcal{X}$, there are functions $\phi_x, \phi_{x'} \in \mathcal{H}$, such that $\langle \phi_x, \phi_{x'} \rangle_{\mathcal{H}} = k(x, x')$. The mapping $\phi : x \mapsto \phi_x$ is referred to

as the *feature map* and \mathcal{H} is referred to as the (transformed) *feature space* of \mathcal{X} , both of which are implied by the chosen kernel. Notice that the feature map may be highly non-linear. Subsequently, these methods learn a metric in the feature space \mathcal{H} : $d(\phi_x, \phi_{x'}) = \sqrt{\langle (\phi_x - \phi_{x'}), A(\phi_x - \phi_{x'}) \rangle_{\mathcal{H}}}$, where $A : \mathcal{H} \rightarrow \mathcal{H}$ is a self-adjoint, bounded, positive-definite operator, preferably, of low rank. Since any element ϕ_x of \mathcal{H} may be of infinite dimension, operator A may be described by an infinite number of parameters to be inferred from the data. Obviously, learning A is not feasible by following direct approaches and, therefore, needs to be learned in some indirect fashion. For example, the authors in [8] pointed out an equivalence between kernel learning and metric learning in the feature space. In specific, they showed that learning A in \mathcal{H} is implicitly achieved by learning a finite-dimensional matrix.

In this paper, we propose a different DML kernelization strategy, according to which a kernel-based, non-linear transform f maps \mathcal{X} into a Euclidean output space \mathbb{R}^m , in order to learn a Mahalanobis distance in that output space. This strategy gives rise to two new models that simultaneously learn both the mapping and the output space metric. Leveraged by the Representer Theorem proposed in [14], all computations of both methods involve only kernel calculations. Unlike previous kernel-based approaches, whose mapping from input to feature space \mathcal{H} cannot be cast into an explicit form, the relevant mappings from input to output space are explicit for both of our methods. Thus, we can access the transformed data in the output space, and this feature can be even used to visualize the data [20], when the output space is 2- or 3-dimensional. Furthermore, by specifying the dimensionality of the output space, the rank of the learned metric can be easily controlled to facilitate dimensionality reduction of the original data.

Our first approach uses an appropriate, but otherwise arbitrary, matrix-valued kernel function and, hence, provides maximum flexibility in specifying the mapping f . Furthermore, in this approach, Mahalanobis distances are explicitly parameterized by a weight matrix to be learned. Our second method is similar to the first one, but assumes a specific parameterized matrix-valued kernel function that can be inferred from the data. We show that the Mahalanobis distance is implicitly determined by the kernel function, and thus eliminates the need of learning a weight matrix for the Mahalanobis distances. To demonstrate the merit of our methods, we compare them to standard k -NN classification (without DML) and other recent kernelized DML algorithms, including Large Margin Nearest Neighbor (LMNN) [22], Information-Theoretic Metric Learning (ITML) [4] and kernelized LMNN (KLMNN) [3]. The comparisons are drawn using eight UCI benchmark data sets in terms of recognition performance and show that the novel methods can achieve higher classification accuracy.

Related Work Several previous works have been focused on DML. Xing, et. al. [23] proposed an early DML method, which minimizes the distance between similar points, while enlarging the distance between dissimilar points. In [17], relative comparison constraints that involve three points at a time are considered. Neighborhood Components Analysis (NCA) [6] is proposed to learn a Mahalanobis distance for the k -NN classifier by maximizing the leave-one-out k -NN performance. [1] proposed a DML method for clustering. Large Margin Nearest Neighbor (LMNN) DML model [22] aims to produce a mapping, so that the k -nearest neighbors of any given sample belong to the same class, while samples from different classes are separated by large margins. Similarly, a Support Vector-based method is proposed in [15]. Also, LMNN is further extended to a Multi-Task Learning variation [16]. Another multi-task DML model is proposed in [25] that searches for task relationships. In [7], the authors proposed a general framework for sparse DML, such that several previous works are subsumed. Also, some other DML models can be extended to sparse versions by augmenting their formulations. Recently, an eigenvalue optimization framework for DML was developed and presented in [24]. Moreover, the connection between LMNN and Support Vector Machines (SVMs) was discussed in [5].

Besides the problem of learning a metric in the original feature space, there has been increasing interest in kernelized DML methods. In the early work of [19], the Lagrange dual problem of the proposed DML formulation is derived, and the DML method is kernelized in the dual domain. Information-Theoretic Metric Learning (ITML) [4] is another kernelized method, which is based on minimizing the Kullback-Leibler divergence between two distributions. The kernelization of LMNN is discussed in [18] and [10]. Moreover, a Kernel Principal Component Analysis (KPCA)-based kernelized algorithm is developed in [3], such that many DML methods, such as LMNN, can be kernelized. In [12], the Mahalanobis matrix and kernel matrix are learned simultaneously. In [8] and its extended work [9], the authors proposed a framework that builds connections between kernel learning and DML in the kernel-induced feature space. Several kernelized models, such as ITML, are covered by this framework. Finally, Multiple Kernel Learning (MKL)-based metric DML

is discussed in [21].

2 RKHS for Vector-Valued Functions

Before introducing our methods, in this section we will briefly review the concept of Reproducing Kernel Hilbert Space (RKHS) for vector-valued functions as presented in [14]. Let \mathcal{X} be an arbitrary set, which we will refer to as *input space*, although it may not actually be a vector space per se. A matrix function $\mathbf{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^{m \times m}$ is called a *positive-definite matrix-valued kernel*, or simply *matrix kernel*, iff it satisfies the following conditions:

$$\mathbf{K}(x, x') = \mathbf{K}^T(x', x) \quad \forall x, x' \in \mathcal{X} \quad (1)$$

$$\mathbf{K}(x, x) \succeq 0 \quad \forall x \in \mathcal{X} \quad (2)$$

$$\bar{\mathbf{K}}(X) \succeq 0 \quad \forall X \subseteq \mathcal{X} \quad (3)$$

where $X = \{x_i\}_{i=1}^n$ and $\bar{\mathbf{K}}(X) \in \mathbb{R}^{mn \times mn}$ is a $n \times n$ block matrix, whose (i, j) block is given as $\bar{\mathbf{K}}_{i,j} = \mathbf{K}(x_i, x_j) \in \mathbb{R}^{m \times m}$, where $i, j \in \{1, 2, \dots, n\}$. According to [14, Theorem 1], if \mathbf{K} is a matrix kernel, then there exists a unique (up to an isometry) RKHS \mathcal{H} of vector-valued functions $f : \mathcal{X} \rightarrow \mathbb{R}^m$ equipped with an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ that admits \mathbf{K} as its reproducing kernel, *i.e.* $\forall x, x' \in \mathcal{X}$ and $\forall \mathbf{y}, \mathbf{y}' \in \mathbb{R}^m$, there are vector-valued functions $K_x \mathbf{y}, K_{x'} \mathbf{y}' \in \mathcal{H}$ that depend on x, \mathbf{y} and x', \mathbf{y}' respectively, such that it holds

$$\langle K_x \mathbf{y}, K_{x'} \mathbf{y}' \rangle_{\mathcal{H}} = \mathbf{y}^T \mathbf{K}(x, x') \mathbf{y}' \quad (4)$$

Note that $K_x : \mathbb{R}^m \rightarrow \mathcal{H}$ is a bounded linear operator parameterized by $x \in \mathcal{X}$ and that the function $K_x \mathbf{y} \in \mathcal{H}$ is such that, when evaluated on $x' \in \mathcal{X}$, it yields

$$(K_x \mathbf{y})(x') = \mathbf{K}(x', x) \mathbf{y} \quad (5)$$

3 Fixed Matrix Kernel DML Formulation

In this section, we propose our first kernelized DML method based on a RKHS for vector-valued functions. Again, let \mathcal{X} be an arbitrary set. Assume we are provided with a training set $\mathcal{T} = \{(x_i, \mathbf{y}_i)\}_{i=1, \dots, n}$, where $x_i \in \mathcal{X}$ and $\mathbf{y}_i \in \mathbb{R}^m$, and we are considering the supervised learning task that seeks to infer a distance metric in \mathbb{R}^m along with a mapping $f : \mathcal{X} \mapsto \mathbb{R}^m$ from \mathcal{T} . In addition to \mathcal{T} , we also assume that we are provided with a real-valued, symmetric *similarity matrix* $\mathbf{S} \in \mathbb{R}^{n \times n}$ with entries $s_{i,j} = s(\mathbf{y}_i, \mathbf{y}_j)$, where $s : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}_+$ is such that $0 \leq s_{i,j} \leq s_{i,i} \quad \forall i, j$. Other than these constraints, the values $s_{i,j}$ can be arbitrary and assigned appropriately with respect to a specific application context. Moreover, let $\mathbf{K}(x, x')$ be a matrix-valued kernel function (*i.e.*, it satisfies Equation (1) through Equation (3)) on \mathcal{X} of given form and let \mathcal{H} be its associated RKHS of \mathbb{R}^m -valued elements. Consider now the following DML formulation:

$$\min_{f, \mathbf{L}} \frac{\gamma}{2} \sum_{i,j} s_{i,j} \|\mathbf{L}[f(x_i) - f(x_j)]\|_2^2 + \frac{\lambda}{2} \sum_i \|\mathbf{L}[f(x_i) - \mathbf{y}_i]\|_2^2 + \rho \operatorname{tr}(\mathbf{L}) + \frac{1}{2} \|f\|_{\mathcal{H}}^2 \quad (6)$$

Notice that $\|\mathbf{L}\mathbf{y}\|_2 = \|\mathbf{y}\|_{\mathbf{A}}$, $\forall \mathbf{y} \in \mathbb{R}^m$, where $\mathbf{A} \triangleq \mathbf{L}^T \mathbf{L} \succeq \mathbf{0}$. In other words, the Euclidean norms of vector differences appearing in (6) are Mahalanobis distances for the output space. Note that if \mathbf{L} is not full-rank, then \mathbf{A} is not strictly positive definite, thus $\|\cdot\|_{\mathbf{A}}$ will be a pseudo-metric in \mathbb{R}^m . The rationale behind this formulation is as follows. The first term, the *collocation* term, forces similar (w.r.t. the similarity measure s) input samples to be mapped closely in the output space (unsupervised learning task). The second term, the *regression* term, forces samples to be mapped close to their target values (supervised learning task). In the context of classification tasks, the combination of these two terms aims to force data that belong to the same class to be mapped close to the same cluster. Closeness in the output space is measured via a Mahalanobis metric that is parameterized via \mathbf{L} . The third term, as we will show later, controls the magnitude of matrix \mathbf{A} and facilitates the derivation of our proposed algorithm. Finally, the fourth term is a regularization term

and is penalizing the complexity of f . Eventually, one can simultaneously learn the output space distance metric and the mapping f through a joint minimization.

The functional of Problem (6) satisfies the conditions stipulated by the Representer Theorem for Hilbert spaces of vector-valued elements (Theorem 5 in [14]) and, therefore, for a fixed value of \mathbf{L} , the unique minimizer \hat{f} is of the form:

$$\hat{f} = \sum_{i=1}^n K_{x_i} \mathbf{c}_i \quad (7)$$

where the m -dimensional vectors $\{\mathbf{c}_i\}_{i=1}^n$ are to be learned. Notice that, due to Equation (5), the explicit input-to-output mapping is given in Equation (8) and is, in general, non-linear in x , if \mathcal{X} is a vector space over the reals.

$$\hat{f}(x) = \sum_{i=1}^n \mathbf{K}(x, x_i) \mathbf{c}_i \quad (8)$$

Proposition 1. *Problem (6) is equivalent to the following minimization problem:*

$$\min_{\mathbf{c}, \mathbf{L}} \frac{1}{2} \mathbf{c}^T \bar{\mathbf{K}} \mathbf{c} + \frac{\gamma}{2} \sum_{i,j} s_{ij} \|\mathbf{L} \Gamma_{ij} \mathbf{c}\|_2^2 + \frac{\lambda}{2} \sum_i \|\mathbf{L}(\bar{\mathbf{K}}_i \mathbf{c} - \mathbf{y}_i)\|_2^2 + \rho \operatorname{tr}(\mathbf{L}) \quad (9)$$

where $\mathbf{c} \triangleq [\mathbf{c}_1^T, \dots, \mathbf{c}_n^T]^T \in \mathbb{R}^{mn}$, $\bar{\mathbf{K}} \in \mathbb{R}^{mn \times mn}$ is the kernel matrix for the training set (as defined for Equation (3)), $\bar{\mathbf{K}}_i = \bar{\mathbf{K}}(x_i) \triangleq [\mathbf{K}(x_i, x_1), \dots, \mathbf{K}(x_i, x_n)] \in \mathbb{R}^{m \times mn}$, and $\Gamma_{ij} = \Gamma(x_i, x_j) \triangleq \bar{\mathbf{K}}_i - \bar{\mathbf{K}}_j$.

The above proposition can be proved by directly substituting Equation (7) into Problem (6) and then using Equation (4). Given two samples $x, x' \in \mathcal{X}$, the inferred metric will be of the form

$$d(x, x') = \|\mathbf{L} \Gamma(x, x') \mathbf{c}\|_2 = \|\Gamma(x, x') \mathbf{c}\|_{\mathbf{A}} \quad (10)$$

with $\mathbf{A} = \mathbf{L}^T \mathbf{L}$. Next, we state a result that facilitates the solution of Problem (9).

Proposition 2. *Problem (9) is convex with respect to each of the two variables \mathbf{c} and \mathbf{L} individually.*

Proof. The convexity of the objective function, denoted as $Q(\mathbf{c}, \mathbf{L})$, with respect to \mathbf{c} is guaranteed by the positive semi-definiteness of the corresponding Hessian matrix of Q :

$$\frac{\partial^2 Q(\mathbf{c}, \mathbf{L})}{\partial \mathbf{c} \partial \mathbf{c}^T} = \bar{\mathbf{K}} + \gamma \sum_{i,j} s_{ij} \Gamma_{ij}^T \mathbf{L}^T \mathbf{L} \Gamma_{ij} + \lambda \sum_i \bar{\mathbf{K}}_i^T \mathbf{L}^T \mathbf{L} \bar{\mathbf{K}}_i \succeq \mathbf{0} \quad (11)$$

To show the convexity with respect to \mathbf{L} , we consider each term separately. The convexity of $\|\mathbf{L} \Gamma_{ij} \mathbf{c}\|_2^2$ stems from the conclusion in [2, p. 110], which states that $\|\mathbf{X} \mathbf{z}\|_2^2$ is convex with respect to any matrix \mathbf{X} for any \mathbf{z} . For the same reason, $\|\mathbf{L}(\bar{\mathbf{K}}_i \mathbf{c} - \mathbf{y}_i)\|_2^2$ is also convex. Finally, $\operatorname{tr}(\mathbf{L})$ is convex in \mathbf{L} , as shown in [2, p. 109]. Thus, the objective function is also convex with respect to \mathbf{L} . \square

Based on Proposition 2, we can perform the joint minimization Problem (9) by block coordinate descent with respect to \mathbf{c} and \mathbf{L} . We set the partial derivatives of Q with respect to the two variables to zero and obtain

$$\frac{\partial Q(\mathbf{c}, \mathbf{L})}{\partial \mathbf{c}} = \mathbf{0} \Rightarrow \mathbf{c} = \lambda \left(\frac{\partial^2 Q(\mathbf{c}, \mathbf{L})}{\partial \mathbf{c} \partial \mathbf{c}^T} \right)^\dagger \sum_i \bar{\mathbf{K}}_i^T \mathbf{L}^T \mathbf{L} \mathbf{y}_i \quad (12)$$

$$\frac{\partial Q(\mathbf{c}, \mathbf{L})}{\partial \mathbf{L}} = \mathbf{0} \Rightarrow \mathbf{L} = -\rho \left(\gamma \sum_{i,j} s_{ij} \Gamma_{ij} \mathbf{c} \mathbf{c}^T \Gamma_{ij}^T + \lambda \sum_i (\bar{\mathbf{K}}_i \mathbf{c} - \mathbf{y}_i) (\bar{\mathbf{K}}_i \mathbf{c} - \mathbf{y}_i)^T \right)^\dagger \quad (13)$$

where \dagger stands for Moore-Penrose pseudo-inversion. One can update \mathbf{c} via Equation (12) by holding \mathbf{L} fixed to its current estimate and then update \mathbf{L} via Equation (13) by using the most current value of \mathbf{c} . Repeating these steps until convergence would constitute the basis for the block-coordinate descent to train this model.

Due to the calculation of the pseudo-inverse, the time complexity of each iteration, in the worst case scenario, is $O((mn)^3)$.

As we can observe from Equation (13), since $\mathbf{A} = \mathbf{L}^T \mathbf{L}$, the parameter ρ that appears in the term $\rho \text{tr}(\mathbf{L})$ of Problem (6) directly controls the norm of \mathbf{A} . Although other regularization terms on \mathbf{L} may be utilized in place of $\rho \text{tr}(\mathbf{L})$, they may not lead to a simple update equation for \mathbf{L} , such as the one given in Equation (13). The potential appeal of this formulation stems from the simplicity of the training algorithm combined with the flexibility of choosing a matrix kernel function that is suitable to the application at hand.

4 Parameterized Matrix Kernel DML Formulation

Our next formulation shares all assumptions with the previous one with the exception that the matrix kernel function \mathbf{K} is now parameterized. We shall show that, even though the matrix kernel function is somewhat restricted, it has the property that is able to implicitly determine the output space Mahalanobis metric. To start, we assume a matrix kernel of the form:

$$\mathbf{K}(x, x') = k(x, x') \mathbf{B} \quad (14)$$

where k is a scalar kernel function that is predetermined by the user and $\mathbf{B} \in \mathbb{R}^{m \times m}$ is a symmetric, positive semi-definite matrix, which will be learned from \mathcal{T} . Because of this facts, \mathbf{K} satisfies Equation (1) through Equation (3) and, therefore is a legitimate matrix kernel function. The formulation for the alternative DML model reads

$$\min_{f, \mathbf{B}} \frac{\gamma}{2} \sum_{i,j} s_{ij} \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2^2 + \frac{\lambda}{2} \sum_i \|f(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 + \frac{\rho}{2} \|\mathbf{B}\|_F^2 + \frac{1}{2} \|f\|_{\mathcal{H}}^2 \quad (15)$$

where $\|\mathbf{B}\|_F^2 \triangleq \text{tr}\{\mathbf{B}^T \mathbf{B}\} = \text{tr}\{\mathbf{B}^2\}$ is the squared Frobenius norm of \mathbf{B} and $\text{tr}\{\cdot\}$ is the matrix trace operator. Problem (15) differs from Problem (6) in a regularization term and in that the former seems to use Euclidean distances in the output space, while the latter uses Mahalanobis distances in the output space with weight matrix $\mathbf{A} = \mathbf{L}^T \mathbf{L}$. As was the case with the formulation of Section 3, the functional of Problem (15) also satisfies the conditions of the Representer Theorem for Hilbert spaces of vector-valued elements and, for fixed value of \mathbf{B} , the unique minimizer \hat{f} has the same form as the one of Equation (7) and the explicit input-to-output mapping is given as

$$\hat{f}(x) = \sum_{i=1}^n k(x, x_i) \mathbf{B} \mathbf{c}_i \quad (16)$$

which, in all but trivial cases, is again non-linear in x , if \mathcal{X} is a vector space over the reals. In a derivation similar to the one found in Section 3, one can show that Problem (15) is equivalent to the following constrained joint minimization problem:

$$\min_{\mathbf{C}, \mathbf{B} \geq \mathbf{0}} \frac{\gamma}{2} \text{tr}\{\mathbf{C} \widetilde{\mathbf{K}}_{\Delta} \mathbf{C}^T \mathbf{B}^2\} + \frac{\lambda}{2} \|\mathbf{B} \mathbf{C} \widetilde{\mathbf{K}} - \mathbf{Y}\|_F^2 + \frac{\rho}{2} \|\mathbf{B}\|_F^2 + \frac{1}{2} \text{tr}\{\mathbf{C}^T \mathbf{B} \mathbf{C} \widetilde{\mathbf{K}}\} \quad (17)$$

where $\mathbf{C} \triangleq [\mathbf{c}_1, \dots, \mathbf{c}_n] \in \mathbb{R}^{m \times n}$, $\widetilde{\mathbf{K}} \in \mathbb{R}^{n \times n}$ is the kernel matrix with $k(x_i, x_j)$ as its (i, j) element, $\widetilde{\mathbf{K}}_{\Delta} \triangleq \widetilde{\mathbf{K}}[\text{diag}\{\mathbf{S} \mathbf{1}_n\} - \mathbf{S}] \widetilde{\mathbf{K}} \in \mathbb{R}^{n \times n}$, where $\text{diag}\{\cdot\}$ is the operator producing a diagonal matrix with the same diagonal as the operator's argument, $\mathbf{1}_n \in \mathbb{R}^n$ is the all-ones vector and $\mathbf{Y} \triangleq [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathbb{R}^{m \times n}$. The learned metric will be of the form

$$d(x, x') = \|\mathbf{B} \mathbf{C} [\tilde{\mathbf{k}}(x) - \tilde{\mathbf{k}}(x')]\|_2 = \|\mathbf{C} [\tilde{\mathbf{k}}(x) - \tilde{\mathbf{k}}(x')]\|_{\mathbf{A}} \quad (18)$$

where $\tilde{\mathbf{k}}(x) \triangleq [k(x, x_1), \dots, k(x, x_n)]^T$ and, in this case, $\mathbf{A} = \mathbf{B}^2$. It is readily seen that the matrix \mathbf{B} specifying the matrix kernel function also determines the Mahalanobis distance in the output space \mathbf{R}^m . Therefore, this model implicitly learns the Mahalanobis distance by learning the \mathbf{B} matrix in the kernel function.

Proposition 3. *Problem (17) is convex with respect to each of the two variables \mathbf{C} and \mathbf{B} .*

Proof. The proof is based on the following facts outlined in [2, sec. 3.6]: (a) A matrix-valued function g is *matrix convex* if and if for any \mathbf{z} , $\mathbf{z}^T g \mathbf{z}$ is convex. (b) Suppose a matrix-valued function g is matrix convex and a real-valued function h is convex and non-decreasing. Then, $h \circ g$ is convex, where \circ denotes function composition. (c) The function $\text{tr}\{\mathbf{W}\mathbf{X}\}$ is convex and non-decreasing in \mathbf{X} , if $\mathbf{W} \succeq \mathbf{0}$. In what follows, we show convexity for each term in Problem (17). Since $\text{tr}\{\mathbf{C}^T \mathbf{B} \widetilde{\mathbf{C}} \widetilde{\mathbf{K}}\} = \text{tr}\{\mathbf{C} \widetilde{\mathbf{K}} \mathbf{C}^T \mathbf{B}\}$ and $\mathbf{C} \widetilde{\mathbf{K}} \mathbf{C}^T \succeq \mathbf{0}$, therefore $\text{tr}\{\mathbf{C} \widetilde{\mathbf{K}} \mathbf{C}^T \mathbf{B}\}$ is convex with respect to \mathbf{B} based on facts (b) and (c). To show the convexity with respect to \mathbf{C} , note that the matrix-valued function $g(\mathbf{C}) = \mathbf{C}^T \mathbf{B} \mathbf{C}$ is matrix convex with respect to \mathbf{C} based on $\mathbf{B} \succeq \mathbf{0}$ and fact (a). Thus, with $\widetilde{\mathbf{K}} \succeq \mathbf{0}$ and fact (b) and (c), we achieve the convexity. The same method is employed to prove the convexity of the other three terms (note that $\widetilde{\mathbf{K}}_\Delta \succeq \mathbf{0}$). \square

Based on Proposition 3, we can again apply a block coordinate descent algorithm to solve Problem (17). If $\tilde{Q}(\mathbf{C}, \mathbf{B})$ is the relevant objective function, we set the partial derivative of $\tilde{Q}(\mathbf{C}, \mathbf{B})$ with respect to \mathbf{C} zero and obtain:

$$\frac{\partial \tilde{Q}(\mathbf{C}, \mathbf{B})}{\partial \mathbf{C}} = \mathbf{0} \Rightarrow \mathbf{C} + \gamma \mathbf{B} \widetilde{\mathbf{C}} \widetilde{\mathbf{K}}_\Delta \widetilde{\mathbf{K}}^{-1} + \lambda \mathbf{B} \widetilde{\mathbf{C}} \widetilde{\mathbf{K}} = \lambda \mathbf{Y} \quad (19)$$

As noted in [11], this matrix equation can be solved for \mathbf{C} as follows:

$$\text{vec}(\mathbf{C}) = \lambda (\mathbf{I} + \gamma (\widetilde{\mathbf{K}}_\Delta \widetilde{\mathbf{K}}^{-1}) \otimes \mathbf{B} + \lambda \widetilde{\mathbf{K}} \otimes \mathbf{B})^{-1} \text{vec}(\mathbf{Y}) \quad (20)$$

To find the optimum \mathbf{B} for fixed \mathbf{C} , due to the constraint $\mathbf{B} \succeq \mathbf{0}$, we use a projected gradient descent method. In each iteration, we update \mathbf{B} using the traditional gradient descent rule: $\mathbf{B} \leftarrow \mathbf{B} - \alpha \nabla_{\mathbf{B}} \tilde{Q}(\mathbf{C}, \mathbf{B})$, where $\alpha > 0$ is the step length, followed by projecting the updated \mathbf{B} onto the cone of positive semi-definite matrices. Since $\tilde{Q}(\mathbf{C}, \mathbf{B})$ is convex with respect to \mathbf{B} for fixed \mathbf{C} , this procedure is able to find the optimum solution for \mathbf{B} . The gradient with respect to \mathbf{B} is given as

$$\frac{\partial \tilde{Q}(\mathbf{C}, \mathbf{B})}{\partial \mathbf{B}} = \mathbf{G} + \mathbf{G}^T - \mathbf{G} \odot \mathbf{I} \quad (21)$$

where \odot is the Hadamard matrix product and \mathbf{G} is defined as

$$\mathbf{G} \triangleq \mathbf{B} [\mathbf{C} (\gamma \widetilde{\mathbf{K}}_\Delta + \lambda \mathbf{K}^2) \mathbf{C}^T + \rho \mathbf{I}] - (\lambda \mathbf{Y} - \frac{1}{2} \mathbf{C}) \mathbf{K} \mathbf{C}^T \quad (22)$$

Therefore, for each iteration, the time complexity of updating \mathbf{C} is $O((mn)^3)$, due to the calculation of a matrix inverse. When updating \mathbf{B} , the time complexity is determined by the convergence speed of the projected gradient descent method.

5 Experiments

In this section, we evaluate the performance of our two kernelized DML methods on classification problems. Towards this purpose, we opt to set $\mathbf{y}_i = \mathbf{y}^{k(i)}$, $\forall i \in \{1, 2, \dots, n\}$, where $k(i) \in \{1, 2, \dots, c\}$ is the class label of the i^{th} sample and \mathbf{y}^k is an appropriately chosen prototype target vector for the k^{th} class. Additionally, we choose to evaluate the pair-wise sample similarities as $s_{i,j} = [\mathbf{y}_i = \mathbf{y}_j]$, where $[predicate]$ denotes the result of the Iversonian bracket, *i.e.* it equals 1, if *predicate* evaluates to true, and 0, if otherwise. After training each of these models, we employ a KNN classifier to label samples in the range space of f ; the classifier uses the models' learned metrics (given by Equation (10) and Equation (18)) to establish nearest neighbors.

We compare our methods with several other approaches. The first one labels samples of the original feature space via the k -NN classification rule using Euclidean distances and, provides a baseline for the accuracy that can be achieved for each classification problem we considered. The second one relies on a popular DML method, namely the Large Margin Nearest Neighbor (LMNN) DML method [22]. We also selected two kernelized approaches for comparison, namely, Information-Theoretic Metric Learning (ITML) [4] and kernelized LMNN (KLMNN) [3].

We evaluated all approaches on eight datasets from the UCI repository, namely, White Wine Quality (*Wine*), Wall-Following Robot Navigation (*Robot*), Statlog Vehicle Silhouettes (*Vehicle*), Molecular Biology

Splice-junction Gene Sequences (*Molecular*), Waveform Database Generator Version 1 (*Wave*), Ionosphere (*Iono*), Cardiotocography (*Cardio*), Pima Indians Diabetes (*Pima*). For all datasets, each class was equally represented in number of samples. An exception is the original *Wine* dataset that has eleven classes, eight of which are poorly represented; for this dataset we only chose data from the other three classes.

For our model with general matrix kernel function \mathbf{K} , we chose the diagonal matrix $\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) = \text{diag}\{[k_1(\mathbf{x}_i, \mathbf{x}_j), \dots, k_m(\mathbf{x}_i, \mathbf{x}_j)]^T\}$, where k_1 through k_m were Gaussian kernel functions with different spreads. For the second model, where $\mathbf{K} = k \cdot \mathbf{B}$, we also chose k to be a Gaussian kernel. During the test phase for all experiments, the parameters γ , λ , ρ , the output dimension m , the Gaussian kernel’s spread parameter σ and the number of nearest neighbors κ to be used by the KNN classifier are selected through cross-validation. Training of the models was performed using 10% and 50% of each data set. In the sequel, we provide the experimental results in figures, which display the average classification accuracies over 20 runs. Also, the error bars correspond to a 95% confidence interval of the estimated accuracies.

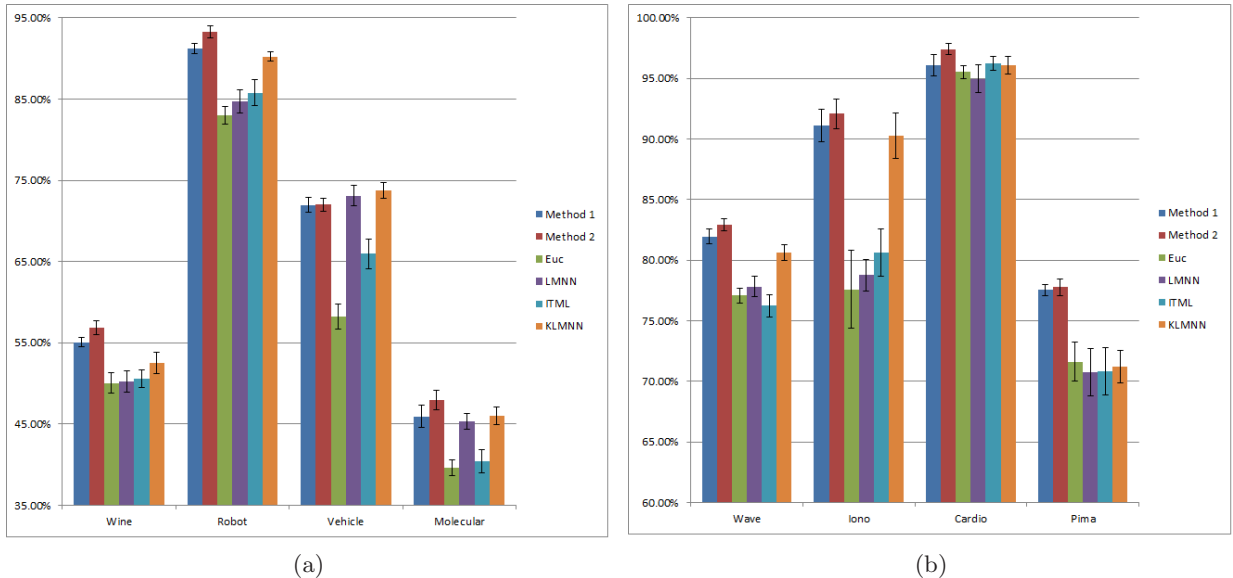


Figure 1: Experimental results for 10% training data. Average classification performance over 20 runs for each data set and each method is shown. Error bars indicate 95% confidence intervals.

We first discuss the results in the case where we used only 10% of the training data; they are depicted in Figure 1. Our first model with general kernel function \mathbf{K} is named as “Method 1”, and the second model with specified kernel function $\mathbf{K} = k \cdot \mathbf{B}$ is called “Method 2”. For almost all datasets, we observe that all five DML methods outperform the scheme involving no transformation of the original feature space (*i.e.*, the output space coincided with the original feature space) and labeling samples via Euclidean-distance KNN classification. This remarkable fact underlines the potential benefits of DML methods. Moreover, we observe that kernelized methods usually outperform LMNN. This observation may partly justify the use of a nonlinear mapping for DML. Furthermore, we observe from the figure that both of our methods typically outperform the other four approaches. More specifically, the proposed two models achieve the highest accuracy across all datasets with the only exception on the *Vehicle* dataset, where ITML and KLMNN outperform slightly. It is worth mentioning that, for the *Pima* data set, none of the other three DML methods can enhance the performance compared to the baseline KNN classification, while our methods achieve significant improvements.

Similar conclusions can be drawn regarding the results generated by using 50% of the training data. These results are depicted in Figure 2. Our methods outperform all the other four methods for most datasets. An exception occurs for the *Molecular* dataset, where KLMNN achieves higher performance than ours. In the case of *Robot* and *Cardio* datasets, all methods perform similarly well. The reason might be that, with enough data, all of the models can be trained well enough to achieve close to optimal performance. For

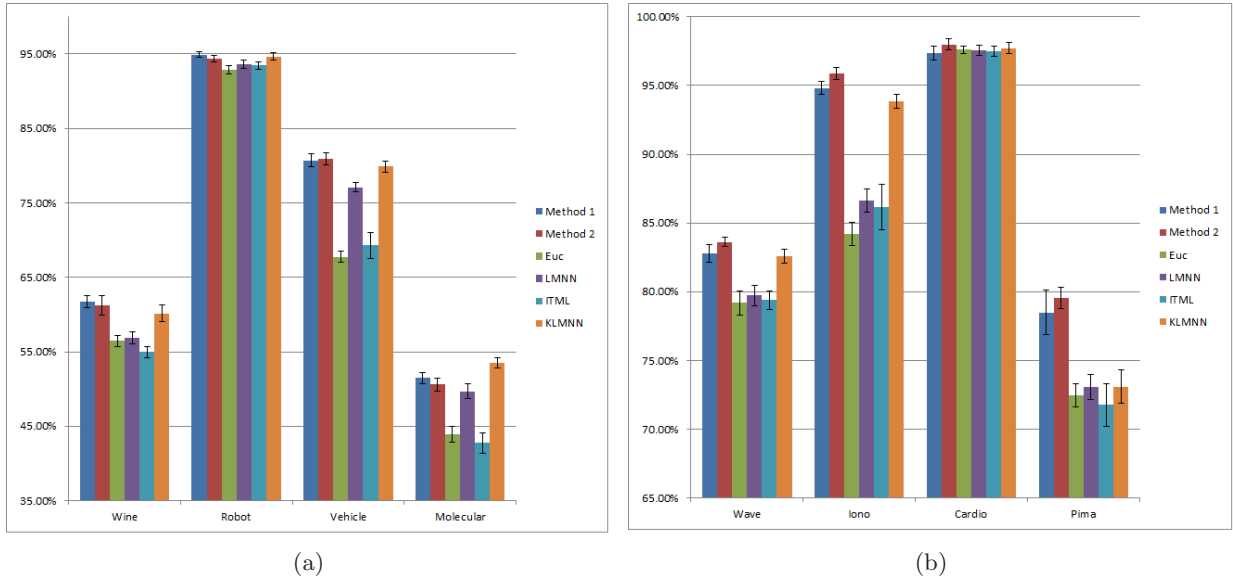


Figure 2: Experimental results for 50% training data. Average classification performance over 20 runs for each data set and each method is shown. Error bars indicate 95% confidence intervals.

the *Pima* data set, again, our methods achieve much better results than all other four methods. It is also important to note that, for our Method 1, despite the relatively simple form of the matrix kernel function we opted for, the resulting model demonstrated very competitive classification accuracy across all datasets. One would likely expect even better performance, if a more sophisticated matrix kernel function is used.

For the sake of visualizing the distribution of the transformed *Robot* data via our models in 2 dimensions, we provide Figure 3 and Figure 4. Similar to [8], we compare the produced mappings of our methods to Kernel Principal Component Analysis (KPCA). KPCA’s 2-dimensional principal subspace was identified based on 10% of the available training data, *i.e.*, 100 training patterns, and the test points were projected onto that subspace. The same training samples were also used for training our two models, which used a Gaussian kernel function and a spread parameter value σ that maximized KNN’s classification accuracy.

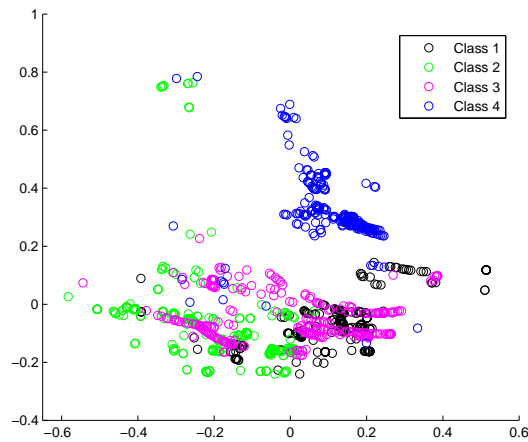


Figure 3: Visualization of the *Robot* data set by applying KPCA.

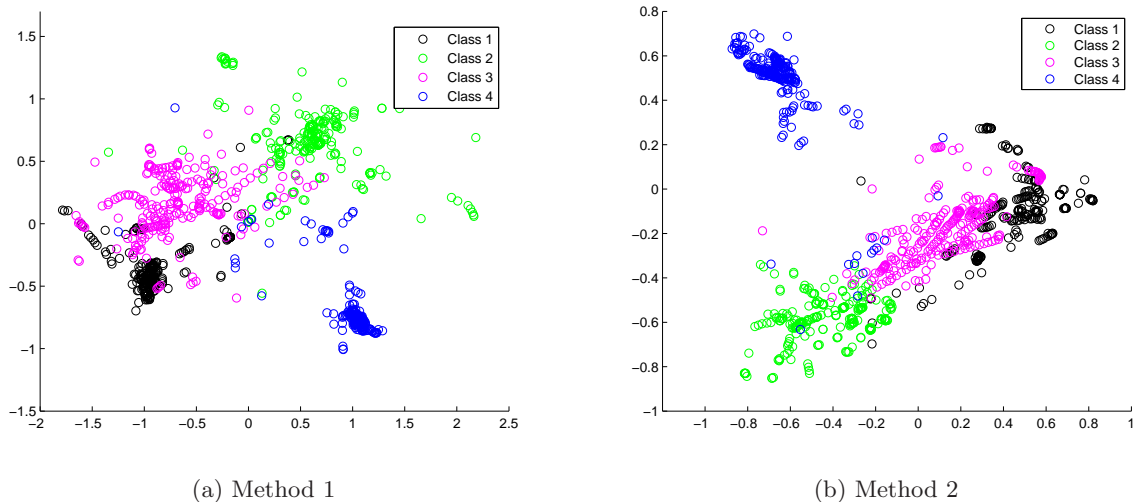


Figure 4: Visualization of the *Robot* data set by applying our methods.

From Figure 3 we observe that KPCA’s projection may only promote good discrimination between samples drawn from class 4 versus the rest. On the other hand, in Figure 4a and Figure 4b, all four classes are reasonably well-clustered in the output space obtained by our two methods. This may explain why our methods are able to achieve high classification accuracy, even when only 10% of the available data are used for training.

6 Conclusions

In this paper, we proposed two new kernel-based Distance Metric Learning (DML) methods, which rely on Reproducing Kernel Hilbert Spaces (RKHSs) of vector-valued functions. Via a mapping f , the two methods map data from their original space to an output space, whose dimension can be directly controlled. Subsequent distance measurements are performed in the output space via a Mahalanobis metric. The first proposed model uses a general matrix kernel function and, thus, provides significant flexibility in modeling the input-to-output space mapping. On the other hand, the second proposed method uses a more restricted matrix kernel function, but has the advantage of implicitly determining the Mahalanobis metric. Furthermore, its matrix kernel function can be learned from data. Unlike previous kernel-based approaches, the relevant f mappings are explicit for both of our two methods. Combined with the fact that the output space dimensionality can be directly specified, the models can also be used for dimensionality reduction purposes, such as for visualizing the data in 2 or 3 dimensions. Experimental results on eight UCI benchmark data sets show that both of the proposed methods can achieve higher performance in comparison to other traditional and kernel-based DML techniques.

Acknowledgements

C. Li acknowledges partial support from National Science Foundation (NSF) grant No. 0806931. Also, M. Georgiopoulos acknowledges partial support from NSF grants No. 0525429, No. 0963146, No. 1200566 and No. 1161228. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF.

References

- [1] Mikhail Bilenko, Sugato Basu, and Raymond J. Mooney. Integrating constraints and metric learning in semi-supervised clustering. In *International Conference on Machine Learning*, 2004. Available from: <http://doi.acm.org/10.1145/1015330.1015360>.
- [2] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [3] Ratthachat Chatpatanasiri, Teesid Korsrilabutr, Pasakorn Tangchanachaianan, and Boonserm Kijirikul. A new kernelization framework for mahalanobis distance learning algorithms. *Neurocomputing*, 73:1570–1579, 2010. Available from: <http://dx.doi.org/10.1016/j.neucom.2009.11.037>.
- [4] Jason V. Davis, Brian Kulis, Prateek Jain, and Inderjit S. Dhillon. Information-theoretic metric learning. In *Proceedings of the 24th international conference on Machine learning*, ICML '07, pages 209–216, 2007. Available from: <http://doi.acm.org/10.1145/1273496.1273523>.
- [5] Huyen Do, Alexandros Kalousis, Jun Wang, and Adam Woznica. A metric learning perspective of svm: On the relation of *LMNN* and SVM. In *JMLR W&CP 5: Proceedings of Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 5, pages 308–317, 2012. Available from: <http://jmlr.csail.mit.edu/proceedings/papers/v22/do12/do12.pdf>.
- [6] Jacob Goldberger, Sam Roweis, Geoff Hinton, and Ruslan Salakhutdinov. Neighbourhood components analysis. In *Neural Information Processing Systems*, 2004. Available from: http://books.nips.cc/papers/files/nips17/NIPS2004_0121.pdf.
- [7] Kaizhu Huang, Yiming Ying, and Colin Campbell. Gsmf: A unified framework for sparse metric learning. In *Data Mining, 2009. ICDM '09. Ninth IEEE International Conference on*, pages 189–198, Dec. 2009. Available from: <http://dx.doi.org/10.1109/ICDM.2009.22>.
- [8] Prateek Jain and Brian Kulis. Inductive regularized learning of kernel functions. In *Neural Information Processing Systems*, 2010. Available from: http://books.nips.cc/papers/files/nips23/NIPS2010_0603.pdf.
- [9] Prateek Jain, Brian Kulis, Jason V. Davis, and Inderjit S. Dhillon. Metric and kernel learning using a linear transformation. *Journal of Machine Learning Research*, 13:519–547, 2012.
- [10] Brian Kulis, Suvrit Sra, Dhillon, and Inderjit. Convex perturbations for scalable semidefinite programming. In *JMLR W&CP 5: Proceedings of Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 5, pages 296–303, 2012. Available from: <http://jmlr.csail.mit.edu/proceedings/papers/v5/kulis09a/kulis09a.pdf>.
- [11] Peter Lancaster. Explicit solutions of linear matrix equations. *SIAM Review*, 12:pp. 544–566, 1970. Available from: <http://www.jstor.org/stable/2028490>.
- [12] Zhengdong Lu, Prateek Jain, and Inderjit S. Dhillon. Geometry-aware metric learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, ICML '09, pages 673–680, 2009. Available from: <http://doi.acm.org/10.1145/1553374.1553461>.
- [13] R. De Maesschalck, D. Jouan-Rimbaud, and D.L. Massart. The mahalanobis distance. *Chemometrics and Intelligent Laboratory Systems*, 50(1):1 – 18, 2000. Available from: <http://www.sciencedirect.com/science/article/pii/S0169743999000477>, doi:10.1016/S0169-7439(99)00047-7.
- [14] Charles A. Micchelli and Massimiliano Pontil. On learning vector-valued functions. *Neural Computation*, 17:177–204, 2005. Available from: <http://dx.doi.org/10.1162/0899766052530802>.
- [15] Nam Nguyen and Yunsong Guo. Metric learning: A support vector approach. In *Proceedings of the European conference on Machine Learning and Knowledge Discovery in Databases - Part II*, ECML PKDD '08, pages 125–136, Berlin, Heidelberg, 2008. Springer-Verlag. Available from: http://dx.doi.org/10.1007/978-3-540-87481-2_9, doi:10.1007/978-3-540-87481-2_9.

- [16] Shibin Parameswaran and Kilian Q. Weinberger. Large margin multi-task metric learning. In *Neural Information Processing Systems*, 2010. Available from: http://books.nips.cc/papers/files/nips23/NIPS2010_0510.pdf.
- [17] Matthew Schultz and Thorsten Joachims. Learning a distance metric from relative comparisons. In *Neural Information Processing Systems*, 2004. Available from: http://books.nips.cc/papers/files/nips16/NIPS2003_AA06.pdf.
- [18] Lorenzo Torresani and Kuang-chih Lee. Large margin component analysis. In *Neural Information Processing Systems*, 2007. Available from: http://books.nips.cc/papers/files/nips19/NIPS2006_0791.pdf.
- [19] Ivor W. Tsang and James T. Kwok. Distance metric learning with kernels. In *Proceedings of International Conference on Artificial Neural Networks (ICANN)*, pages 126–129, 2003.
- [20] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(2579-2605):85, 2008.
- [21] Jun Wang, Huyen Do, Adam Woznica, and Alexandros Kalousis. Metric learning with multiple kernels. In *Neural Information Processing Systems*, 2011. Available from: http://books.nips.cc/papers/files/nips24/NIPS2011_0683.pdf.
- [22] Kilian Q. Weinberger and Laurence K. Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, 10:207–244, 2009. Available from: <http://dl.acm.org/citation.cfm?id=1577069.1577078>.
- [23] Eric P. Xing, Andrew Y. Ng, Michael I. Jordan, and Stuart Russell. Distance metric learning, with application to clustering with side-information. In *Neural Information Processing Systems*, 2002. Available from: <http://books.nips.cc/papers/files/nips15/AA03.pdf>.
- [24] Yiming Ying and Peng Li. Distance metric learning with eigenvalue optimization. *Journal of Machine Learning Research*, 13:1–26, 2012. Available from: <http://dl.acm.org/citation.cfm?id=2188385.2188386>.
- [25] Yu Zhang and Dit-Yan Yeung. Transfer metric learning by learning task relationships. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '10, pages 1199–1208, 2010. Available from: <http://doi.acm.org/10.1145/1835804.1835954>.