

# A Robust and Efficient Doubly Regularized Metric Learning Approach

Meizhu Liu\* and Baba C. Vemuri

Siemens Corporate Research & Technology, Princeton, NJ, 08540  
CISE, University of Florida, Gainesville, FL, 32611

**Abstract.** A proper distance metric is fundamental in many computer vision and pattern recognition applications such as classification, image retrieval, face recognition and so on. However, it is usually not clear what metric is appropriate for specific applications, therefore it becomes more reliable to learn a task oriented metric. Over the years, many metric learning approaches have been reported in literature. A typical one is to learn a Mahalanobis distance which is parameterized by a positive semidefinite (PSD) matrix  $\mathbf{M}$ . An efficient method of estimating  $\mathbf{M}$  is to treat  $\mathbf{M}$  as a linear combination of rank-one matrices that can be learned using a boosting type approach. However, such approaches have two main drawbacks. First, the weight change across the training samples may be non-smooth. Second, the learned rank-one matrices might be redundant. In this paper, we propose a doubly regularized metric learning algorithm, termed by DRMetric, which imposes two regularizations on the conventional metric learning method. First, a regularization is applied on the weight of the training examples, which prevents unstable change of the weights and also prevents outlier examples from being weighed too much. Besides, a regularization is applied on the rank-one matrices to make them independent. This greatly reduces the redundancy of the rank-one matrices. We present experiments depicting the performance of the proposed method on a variety of datasets for various applications.

**Keywords:** Regularized metric learning, boosting, PSD matrix.

## 1 Introduction

The choice of an appropriate distance or similarity measure over the input space is critical to many computer vision and pattern recognition applications, including but not limited to clustering and classification [1], image retrieval [2], shape detection [3], face recognition [4–9], tracking [10]. There are many commonly used distance metrics, e.g. Euclidean distance,  $L_1$ -norm distance,  $\chi^2$  distance, and Mahalanobis distance etc. However, it is usually very hard to predict which

---

\* This work was done when Meizhu Liu was studying for her PhD degree at University of Florida. This research was in part supported by the NIH grant EB007082 to Vemuri.

distance measure is appropriate for a certain application with specific inputs. Therefore, it is more apt to develop a task-dependent metric based on the available knowledge of the inputs. It was shown [1, 9] that a properly designed distance metric, compared with the standard distances, can significantly improve the performance for many applications.

There are a lot of metric learning algorithms in the literature. A good metric learning algorithm should be able to learn a metric that can amplify informative dimensions (feature) and squash non-informative dimensions. This is unlike Euclidean distance, which treats every dimension equally and does not consider the correlation between them.

In most cases, metric learning algorithms are derived from the labeled training datasets, and the goal of the algorithm is to learn a metric which can separate the instances of different classes apart, and bring together the instances belonging to the same class. To be specific, the labeling of the inputs can be provided mainly in three ways. First, the input constraint is  $(\mathbf{x}_i, y_i)$  where  $\mathbf{x}_i \in \mathbb{R}^D$  is an instance and  $y_i$  is its label. Second, the input constraint is  $((\mathbf{x}_i, \mathbf{x}_j), y_{ij})$  where  $y_{ij}$  indicates whether  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are “similar“ or “dissimilar“ [11]. An even weaker representation often used in information retrieval [12] is the proximity relationship over triplets  $(i, j, k)$ , meaning that  $\mathbf{x}_i$  is closer to  $\mathbf{x}_j$  than to  $\mathbf{x}_k$ . Proximity relationships are the most natural constraint for learning a metric, and are of the weakest representation because proximity triplets can be derived from the other kinds of constraints, but not vice versa.

In this paper, we propose a doubly regularized metric learning algorithm, termed by DRMetric. Our goal is to learn a Mahalanobis distance metric which tries to preserve the proximity relationships over the input. Mahalanobis distance metric is parameterized by a positive semidefinite (PSD) matrix [13, 14]<sup>1</sup>. It has been well studied and advantages were shown over some other metrics such as multidimensional scaling (MDS) [15], ISOMAP [16], and locally linear embedding (LLE) [17].

Several aspects of our DRMetric are novel. First, we use the total Kullback-Leibler (tKL) divergence [18] to regularize the evolution of the weights on the training triplets. tKL is a recently proposed divergence which has been proved to be statistically robust [18]. The regularization automatically ensures that the weight of the examples is upper bounded, therefore, the weight can not be extremely large for outliers. Note that without regularization, the weight of an outlier example will keep increasing, which may lead to serious problems such as overfitting and inefficiency. Furthermore, for some noisy examples, their weight may depict severe oscillations. This not only hampers the convergence rate of the metric learning algorithm, but also leads to overfitting, and lowers the accuracy. Second, we regularize the rank-one PSD matrices to minimize the dependence between them. This regularization makes the rank-one matrices least correlated,

---

<sup>1</sup> Strictly speaking, this matrix should be symmetric positive definite (SPD) in order for it to be a metric. However, we relax the requirement and allow two different instances to have zero distance.

and therefore least redundant, which greatly decreases the number of rank-one matrices needed and improves the efficiency.

The rest of the paper is organized as follows. In Section 2, we briefly review the metric learning literature. In Section 3, we present the doubly regularized linear programming metric learning algorithm, termed by DLMetric. In Section 4, we investigate DLMetric empirically by evaluating our algorithm on a number of datasets for various applications. We also compare our method with the state-of-the-art metric learning and other algorithms. Finally, we conclude the paper in Section 5.

## 2 Literature Review

A good task dependent metric has attracted extensive attention recently. The machine learning community has done many researches to automatically learn a distance function from available knowledge of the dataset [13, 19, 20, 11]. Most existing works assume the metrics to be Mahalanobis distance, which are parameterized by PSD matrices.

Various techniques have been proposed to learn a PSD matrix from the dataset. Some techniques force the negative eigenvalues in the learned symmetric matrix to be zero as in [11]. Some others set the matrix to be the inverse of the covariance matrix of the centered data points in small subsets of points with known relevant information [13]. In [20], the matrix exponential gradient update was used which preserves symmetry and positive definiteness due to the fact that the matrix exponential of a symmetric matrix is always an SPD matrix. In [21], Iwasawa factorization was used to ensure the positive definiteness [21]. Most of these techniques are limited from a scalability or a computational complexity view point.

More recently, some researchers [1, 12] adapted the boosting technique to metric learning. This kind of metric learning is based on an important theorem that a PSD matrix with trace one can always be represented as a convex combination of multiple rank-one PSD matrices. This is a generalization of boosting [22] in the sense that the weak learner in these metric learning algorithms is a rank-one matrix instead of a classifier. The main idea behind these boosting-based metric learning algorithms is that at each iteration, they will learn a rank-one matrix from the training examples that follow a distribution. The weighted rank-one matrix is then added to the PSD matrix. This weight is typically related to the rank-one matrix's ability to discriminate the examples from different classes. The higher the discriminatory power, the larger the weight, and vice versa. After learning the rank-one matrix, the distribution of the examples is updated. The examples are reweighted according to the rule that misclassified examples tend to gain weight and correctly classified examples tend to lose weight. Therefore, the rank-one matrices to be learned will be focused more on the examples that were misclassified previously.

However, these methods are not statistically robust i.e., the learning process is sensitive to noisy data and outliers [23, 24]. The reason is that the weight of

noisy examples might switch between severe increase and severe decrease frequently, which seriously slows down the convergence of the learning process. Furthermore, the weight of outliers might keep increasing, which largely affects the metric to be learned. To avoid these issues and inspired by the regularized boosting [23, 24] techniques, we propose a regularized metric learning algorithm, which regularizes the weight updating process involved in the training stage. Furthermore, in order to reduce the redundancy of the learned rank-one matrices, we add another regularization term to make the dependence between the learned rank-one matrices as small as possible. In this way, we can use much fewer number of rank-one matrices (i.e. much fewer number of iterations) to form the PSD matrix which parameterizes a suitable metric. Experimental results illustrate that for a dataset of  $D$  dimensions, DRMetric is able to learn a relatively good metric in  $D$  iterations.

### 3 Proposed Method

Given a dataset  $\mathbf{X} = \{\mathbf{x}_i\}$ , with  $\mathbf{x}_i \in \mathbb{R}^D$ , and its associated triplet set  $\mathcal{T} = \{(i, j, k)\}$ , with  $(i, j, k)$  meaning that  $\mathbf{x}_i$  is more similar to  $\mathbf{x}_j$  than to  $\mathbf{x}_k$ . Let  $N = |\mathcal{T}|$  denote the number of triplets in  $\mathcal{T}$ . The goal is to learn a Mahalanobis distance which preserves the relationship in  $\mathcal{T}$ .

A Mahalanobis distance is parameterized by a PSD matrix  $\mathbf{M} \in \mathbb{R}^{D \times D}$ . The Mahalanobis distance between  $\mathbf{x}_i \in \mathbf{X}$  and  $\mathbf{x}_j \in \mathbf{X}$  based on  $\mathbf{M}$  is

$$d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j) \tag{1}$$

To remove the scalability effect of the distance resulting from  $\mathbf{M}$ , we require  $tr(\mathbf{M}) = 1$ . Since any trace-one PSD matrix can be decomposed as a convex combination of rank-one trace-one PSD matrices, i.e.,

$$\mathbf{M} = \sum_{l=1}^D w_l \mathbf{u}_l \mathbf{u}_l^T, \mathbf{u}_l \in \mathbb{R}^D, \|\mathbf{u}_l\| = 1, \mathbf{w} \in \Delta_D. \tag{2}$$

To avoid notation clutter in later computations, we introduce a vector  $\mathbf{v}_n = [v_{nl}]$ , where  $v_{nl}$  corresponds to  $\mathbf{u}_l$  and the  $n$ th triplet  $(i, j, k)$ , and is defined as

$$v_{nl} = (\mathbf{x}_i - \mathbf{x}_k)^T \mathbf{u}_l \mathbf{u}_l^T (\mathbf{x}_i - \mathbf{x}_k) - (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{u}_l \mathbf{u}_l^T (\mathbf{x}_i - \mathbf{x}_j). \tag{3}$$

A potentially appropriate  $\mathbf{M}$  should be able to maximize the soft margin defined in the following linear programming,

$$\begin{aligned} & \max_{\mathbf{w}, \rho, \zeta} \rho - \alpha \sum_{n=1}^N \zeta_n \\ & \text{s.t.} \quad \sum_{l=1}^t w_l v_{nl} \geq \rho - \zeta_n, n = 1, \dots, N, \\ & \quad \mathbf{w} \in \Delta_t, \zeta \geq \mathbf{0}, \end{aligned} \tag{4}$$

where  $\zeta_n$  is the slack variable, and  $\alpha$  is a constant factor which penalizes the slack variables.

The Lagrangian dual problem of (4) is

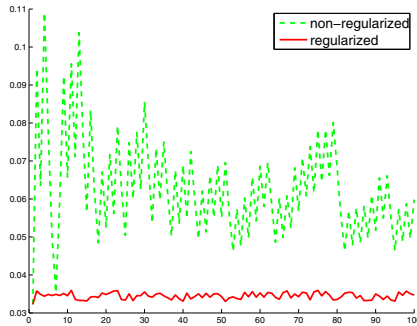
$$\max_{\mathbf{d}, c, \mathbf{q}} \min_{\mathbf{w}, \rho, \zeta} \mathcal{L}(\mathbf{w}, \rho, \zeta, \mathbf{d}, c, \mathbf{q}) = -\rho + \alpha \sum_{n=1}^N \zeta_n - \sum_{n=1}^N d_n \left( \sum_{l=1}^t w_l v_{nl} - \rho + \zeta_n \right) + c(\mathbf{1}^T \mathbf{w} - 1) - \mathbf{q}^T \zeta, \tag{5}$$

where  $\mathbf{d}$ ,  $c$ , and  $\mathbf{q}$  are non-negative regularizers. After some simple algebraic manipulation we arrive at the dual problem of (4) given by,

$$\min_{\mathbf{d} \in \Delta_N, \mathbf{d} \leq \alpha \mathbf{1}} \max_{l=1, \dots, t} \sum_{n=1}^N d_n v_{nl}. \tag{6}$$

### 3.1 Regularization on $\mathbf{d}$

The regularization on  $\mathbf{d}$  is very important because for the non-regularized metric learning algorithm, the weight of the training examples might change very severely, i.e., the weight of a training example might oscillate significantly when it is misclassified or correctly classified by the weak learners (rank-one matrices) as shown in Fig. 1. This instability will seriously affect the learning efficiency, accuracy, and also lead to overfitting. With regularization, severe oscillations and instabilities can be prevented, which makes the algorithm converge faster, i.e. need fewer number of rank-1 PSD matrices. Fig. 1 depicts that, using regularization, the resulting weight change of the training data is stable.



**Fig. 1.** Change in the weight of a training example in the Heart disease dataset from the UCI repository, under metric learning without regularization and with regularization

To overcome the aforementioned instabilities, we add a regularization term to the update of  $\mathbf{d}$  in (6), i.e.,

$$\min_{\mathbf{d} \in \Delta_N, \mathbf{d} \leq \alpha \mathbf{1}} \max_{l=1, \dots, t} \sum_{n=1}^N d_n v_{nl} + \eta \delta(\mathbf{d}, \hat{\mathbf{d}}), \tag{7}$$

where  $\eta$  is the regularization coefficient that balances the margin and the smoothness.  $\eta$  is set to be a fixed number<sup>2</sup> as in [23] to make the number of iterations upper bounded by a constant without hurting the accuracy.  $\delta(\mathbf{d}, \hat{\mathbf{d}})$  is the tKL divergence [23, 25], and

$$\delta(\mathbf{d}, \hat{\mathbf{d}}) = \frac{\sum_{n=1}^N d_n \log \frac{d_n}{\hat{d}_n}}{\sqrt{1 + \sum_{j=1}^N \hat{d}_j (1 + \log \hat{d}_j)^2}}. \tag{8}$$

Note that the regularization term  $\delta(\mathbf{d}, \hat{\mathbf{d}})$  ensures that the evolution of  $\mathbf{d}$  is smooth.

Here,  $\hat{\mathbf{d}}$  can be chosen in different ways. In this paper, we set  $\hat{\mathbf{d}} = \mathbf{d}^0$ , where  $\mathbf{d}^0$  is the initialized distribution, this means  $\mathbf{d}$  should not be far away from the initialized distribution. Since  $\mathbf{d}^0$  is user defined, it is usually set according to the application problem and the data. One tends to initialize larger weight on the examples with more importance, so we use  $\delta(\mathbf{d}, \mathbf{d}^0)$  as the regularizer. Note that the  $d_n$  is upper bounded by  $\alpha$  as in (6), therefore the weight of the noisy examples and outliers is prevented from being too large leading to possible domination in the learning<sup>3</sup>.

To directly compute  $\mathbf{d}^t$  from (7) is complicated, instead, we will first find its Lagrangian and use it to compute  $\mathbf{d}^t$ . To find the Lagrangian, we rewrite (7) into the following form

$$\begin{aligned} & \min_{\beta, \mathbf{d}} \beta + \eta \delta(\mathbf{d}, \mathbf{d}^0) \\ & \text{s.t. } \sum_{n=1}^N d_n v_{nl} \leq \beta, l = 1, \dots, t \\ & \mathbf{d} \in \Delta_N, \mathbf{d} \leq \alpha \mathbf{1}. \end{aligned} \tag{9}$$

The Lagrangian  $\Psi$  of (9) is given by,

$$\Psi(\mathbf{d}, \beta, \mathbf{w}, \xi, \gamma) = \beta + \eta \delta(\mathbf{d}, \mathbf{d}^0) + \sum_{l=1}^t w_l \left( \sum_{n=1}^N d_n v_{nl} - \beta \right) + \sum_{n=1}^N \xi_n (d_n - \alpha) + \gamma (\mathbf{d} \cdot \mathbf{1} - 1) \tag{10}$$

where,  $w_l, l = 1, \dots, t, \xi_n, n = 1, \dots, N$  and  $\gamma$  are non-negative regularizers. Using some simple calculus and the KKT condition [26], we can simplify (10) and get the partial Lagrangian

$$\Psi(\mathbf{d}, \mathbf{w}) = \eta \delta(\mathbf{d}, \mathbf{d}^t) + \sum_{l=1}^t w_l \sum_{n=1}^N d_n v_{nl}. \tag{11}$$

---

<sup>2</sup>  $\eta = \frac{\epsilon \sqrt{1 + (\log N - 1)^2}}{2 \log(ND)}$ , where  $N$  is the number of training samples,  $D$  is the dimension of each training sample, and  $\epsilon$  is the error tolerance of the margin between different classes based on the learned metric [23].

<sup>3</sup> To make such a  $\mathbf{d}$  exist, we should require  $\alpha \geq 1/N$ . It was shown in [24] that  $\alpha = 1/s$ , and  $s \in \{1, \dots, N\}$  is a favorable choice.

Now differentiating  $\Psi$  with respect to  $\mathbf{d}$ , setting it to 0, and normalizing  $\mathbf{d}$ , we get,

$$d_n^t = \frac{d_n^0 \exp\left(-c \sum_{l=1}^t w_l v_{nl}\right)}{Z_t}, \text{ where } c = \frac{1}{\eta} \sqrt{1 + \sum_{n=1}^N d_n^0 (1 + \log d_n^0)^2}, \quad (12)$$

and  $Z_t$  is the normalization parameter to make  $\sum_{n=1}^N d_n^t = 1$ . Here if  $d_n^t > \alpha$ , then we manually set  $d_n^t = \alpha$ .

### 3.2 Regularization on $\mathbf{u}$

We put two constraints on  $\mathbf{u}$ . First, we want it to maximize the margin. Second, we require  $\mathbf{u}$  to be independent of the previously learned  $\mathbf{u}_l, l = 1, 2, \dots, t$ , so that the learned  $\{\mathbf{u}\}$  will not be redundant. Therefore, the number of rank-one matrices needed to form a good metric is reduced. The dependence between  $\mathbf{u}$  and  $\mathbf{u}_l$  is measured by  $\|\mathbf{u}^T \mathbf{u}_l\|^2 \in [0, 1]$ . The larger  $\|\mathbf{u}^T \mathbf{u}_l\|^2$  is, the more dependent they are. When  $\|\mathbf{u}^T \mathbf{u}_l\|^2 = 0$ ,  $\mathbf{u}$  and  $\mathbf{u}_l$  are independent.

The two constraints on  $\mathbf{u}$  are described as

$$\max_{\mathbf{u}} \sum_{n=1}^N d_n^t [(\mathbf{x}_i - \mathbf{x}_k)^T \mathbf{u} \mathbf{u}^T (\mathbf{x}_i - \mathbf{x}_k) - (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{u} \mathbf{u}^T (\mathbf{x}_i - \mathbf{x}_j)] - \lambda \sum_{l=1}^{t-1} \|\mathbf{u}^T \mathbf{u}_l\|^2, \quad (13)$$

where  $\lambda$  is the regularization coefficient to penalize the dependence. (13) can be rewritten as

$$\max_{\mathbf{u}} \mathbf{u}^T \left\{ \sum_{n=1}^N d_n^t [(\mathbf{x}_i - \mathbf{x}_k)(\mathbf{x}_i - \mathbf{x}_k)^T - (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T] - \lambda \sum_{l=1}^t \mathbf{u}_l \mathbf{u}_l^T \right\} \mathbf{u}^T \quad (14)$$

Let matrix  $\mathbf{A}^t = \sum_{n=1}^N d_n^t [(\mathbf{x}_i - \mathbf{x}_k)(\mathbf{x}_i - \mathbf{x}_k)^T - (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T] - \lambda \sum_{l=1}^t \mathbf{u}_l \mathbf{u}_l^T$ , then  $\mathbf{u}_{t+1}$  is the eigenvector corresponding to the largest eigenvalue of  $\mathbf{A}^t$ .

The weight vector  $\mathbf{w}$  for the rank-one matrices should satisfy the linear programming problem (4) which can be solved using column generation [27] or a gradient based method.

### 3.3 Building the Triplets

For each  $\mathbf{x}_i \in \mathbf{X}$ , we first find the  $a$  instances  $\{\mathbf{x}_j\}_{j=1}^a$  which are in the same category as  $\mathbf{x}_i$  but are most different from  $\mathbf{x}_i$ . After that, we find the  $a$  nearest neighbors  $\{\mathbf{x}_k\}_{k=1}^a$  in a different category, then  $(i, j, k)$  will form a triplet. If the size of  $\mathbf{X}$  is small, we will use a larger  $a$ , otherwise, we will use a smaller  $a$ . Furthermore, if the number of triplets is very large, we will randomly select 10 ~ 50% of the triplets for training.

As a summary, the algorithm for the proposed DRMetric is presented in Algorithm 1. The proof of the convergence is very similar to the proof from Schapire and Singer [28].

---

**Algorithm 1.** Doubly Regularized Metric Learning

---

**Input:** Dataset  $\mathbf{X} = \{\mathbf{x}_i\}$ ,  $\mathbf{x}_i \in \mathbb{R}^D$ Triplet set  $\mathcal{T} = \{(i, j, k) \mid \mathbf{x}_i \text{ is closer to } \mathbf{x}_j \text{ than to } \mathbf{x}_k\}$ .  $N = |\mathcal{T}|$ , the number of triplets in  $\mathcal{T}$ .**Output:**  $\mathbf{M} = \sum_{l=1}^t w_l \mathbf{u}_l \mathbf{u}_l^T$ ,  $\mathbf{u}_l \in \mathbb{R}^D$ ,  $\|\mathbf{u}_l\| = 1$ ,  $\mathbf{w} \in \Delta_t$ ,  $t$  is the number of iterations.**Initialization:** Initialize  $d_n^0$ , the weight of the  $n$ th triplet,  $n = 1, \dots, N$ , according to the importance, or set  $d_n^0 = 1/N$  by default.**for**  $l = 1$  to  $t$  **do**    Find the optimal  $\mathbf{u}_l$  according to (13);    Update the distribution  $\mathbf{d}$  according to (12);    Update the weight  $\mathbf{w}$  according to (4)**end for**Return  $\mathbf{M} = \sum_{l=1}^t w_l \mathbf{u}_l \mathbf{u}_l^T$ .

---

## 4 Experimental Results

The proposed algorithm is evaluated on a number of public domain datasets for a variety of applications. We use the UCI machine learning repository [29] for classification, use the COREL image dataset for content based image retrieval, and use the Labeled Faces in the Wild (LFW) [30] dataset for face recognition. We compare our method with many state-of-the-art metric learning and other techniques. The results show that our proposed metric learning method is very promising for many applications.

### 4.1 Classification

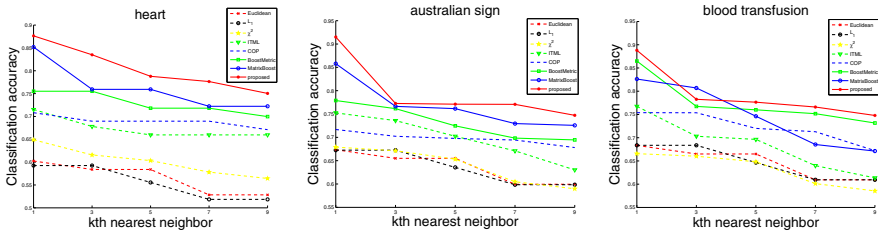
The classification experiments are performed on the UCI machine learning repository [29], which is a collection of datasets that have been extensively used for analyzing machine learning techniques. The repository contains a large variety of datasets, including very noisy datasets (e.g. the Optical Recognition of Handwritten Digits dataset, the wine dataset) as well as relatively clean datasets, which is optimal for testing the robustness and accuracy of classification algorithms. We selected 9 datasets from the UCI repository. The selected datasets include noisy and clean datasets, cover small size to large size datasets in terms of number of instances in the datasets, and range from low dimension to high dimension in terms of number of attributes per instance of the datasets. The description of the selected datasets is shown in Table 1.

We use 5-fold cross validation to evaluate the proposed algorithm. The regularization parameters  $\alpha$ ,  $\eta$  and  $\lambda$  are determined during the training and validation stage, and they are set to be the numbers which maximize the performance on the training dataset. The final result is the average of the results obtained over the 5 runs. The proposed DRMetric is compared with many other non-metric learning and metric learning algorithms, including Euclidean distance,  $L_1$ -norm



**Table 1.** Description of the selected UCI datasets

dataset	# instances	# attributes
Heart disease	303	74
Australian sign	6650	14
Blood transfusion service center	748	5
Artificial characters	6000	7
Glass identification dataset	214	10
Adult dataset	48842	14
Handwritten digits	5620	64
Wine dataset	178	13

**Fig. 2.** The neighbor accuracy curves from different metrics on the Heart disease, Australian sign and Blood transfusion service center datasets in the UCI repository

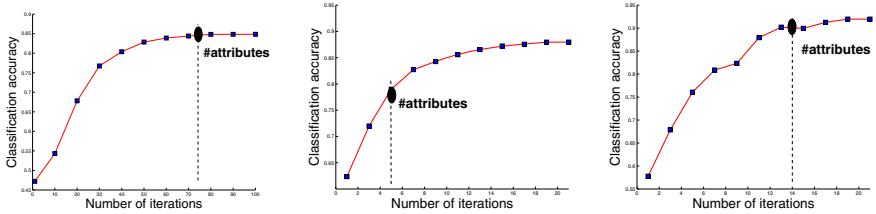
distance,  $\chi^2$  distance, BoostMetric [9], MatrixBoost [1], ITML [14], and COP [11]. The code for metric learning methods is obtained directly from the corresponding authors or downloaded from the authors' webpage. The classification performance is measured based on neighbor accuracy curves. The neighbor accuracy measures the percentage of correctly classified instances based on the  $k$ th ( $k = 1, 3, 5, 7, 9$ ) nearest neighbor. The average neighbor accuracy is shown in Fig. 2. The comparison depicts that in general DRMetric yields higher classification accuracy.

We also evaluate the 3-nearest-neighbor voting classification accuracy on several UCI datasets including the Glass Identification dataset, the Adult dataset, the Optical Recognition of Handwritten Digits dataset, and the Wine dataset. The classification results are shown in Table 2, which reflects that the proposed method outperforms the other methods.

Besides, for DRMetric, we examined the relationship between the classification accuracy change and the number of iterations. The results are shown in Fig. 3, which implies that when the number of iterations is less than  $D$  (the dimension of the dataset), the classification accuracy increases at a higher rate. However, when the number of iterations is larger than  $D$ , the classification accuracy improves very slowly. This means that, using our method,  $D$  rank-one matrices can form a relatively high quality Mahalanobis distance.

**Table 2.** Classification accuracy using different metrics on selected datasets from the UCI repository

dataset	Euclidean	$L_1$	$\chi^2$	ITML	COP	BoostMetric	MatrixBoost	proposed
Characters	0.7235	0.7452	0.7651	0.9114	0.8889	0.9147	0.9049	0.9288
Glass	0.6114	0.6404	0.6479	0.7975	0.7850	0.8135	0.7991	0.8204
Adult	0.6017	0.6249	0.6284	0.7760	0.7752	0.7981	0.7894	0.8075
Digits	0.6865	0.7107	0.7284	0.7352	0.7473	0.8014	0.8148	0.8290
Wine	0.7240	0.7261	0.7602	0.8625	0.8958	0.9074	0.9152	0.9161



**Fig. 3.** The change of classification accuracy with related to the number of iterations using DRMetric on the Heart disease (left), Australian sign (middle), and Blood transfusion service center (right) datasets. The black disk corresponds the classification accuracy when the number of iterations equals to the number of attributes.

### 4.2 Content Based Image Retrieval

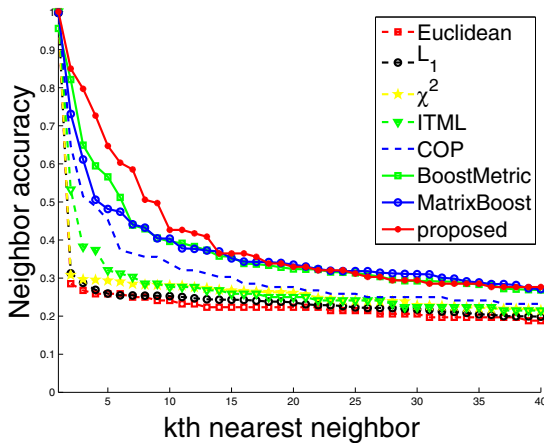
The task for image retrieval is that given one image in a category, find the images in the same category. We use the COREL image database [31] to evaluate our method on content based image retrieval. The database contains 3400 real-world images with 34 different categories, and 100 images per category.

Each image is represented as a 33 dimensional feature vector, which is a combination of low level features including color features, edge features and texture features. For color features, we first represent the images in the HSV color space, and then compute the mean, variance, skewness of the HSV color to get a 9 dimensional feature vector. For edge features, the Canny edge detector [32] is first applied to images to detect the edges, and the histogram for edge direction was quantized into 9 bins of every 40 degrees, which resulted in 9 different edge features. For texture features, we use the multi-resolution simultaneous autoregressive (MASAR) model [33] to get 15 features. In total, there are 33 features for each image.

We use 10-fold cross validation to evaluate the proposed algorithm, i.e., 90% images are used to learn the metric, and 10% images are used for evaluation. We use every image in the test dataset as a query, if the retrieved image belongs to the same category as the query image, the retrieval is correct. We measure the retrieval performance based on the neighbor accuracy curves. Neighbor accuracy measures the percentage of correctly retrieved images in the  $k$ th nearest neighbors of the query images ( $k = 1, \dots, 40$  in our experiments).

We compared our method to many algorithms, including Euclidean distance,  $L_1$ -norm distance,  $\chi^2$  distance, BoostMetric [9], MatrixBoost [1], ITML [14], and COP [11]. The retrieval results are shown in Fig. 4. The results illustrate that our proposed method achieves a higher neighbor accuracy when using 1st  $\sim$  25th and 34th  $\sim$  40th nearest neighbors. However, it's a little worse than MatrixBoost when using the 26th  $\sim$  33rd nearest neighbors.

We compare the computational time of BoostMetric [9], MatrixBoost [1], and DRMetric to learn the distance metric on the COREL database. All algorithms are run on a laptop with Intel(R) Core(TM)2 CPU L7500 @1.6GHz, 4GB memory, GNU Linux and MATLAB (Version R2011a). The average CPU time taken to converge for our algorithm is 167.68s, while BoostMetric takes 244.81s, and MatrixBoost takes 239.28s.



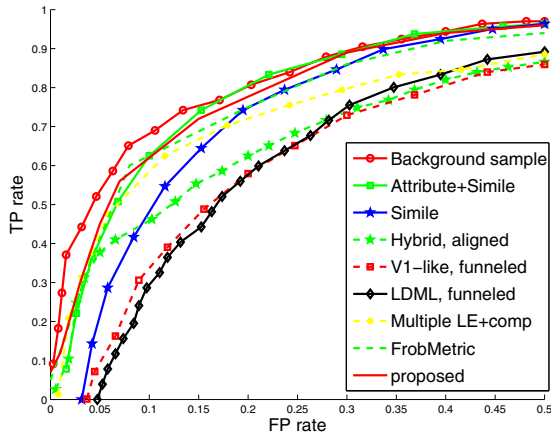
**Fig. 4.** Comparison of using different metric learning methods for content based image retrieval on the COREL dataset

### 4.3 Face Recognition

In this scenario, the goal for face recognition is to do pair matching: given two face images, determine if these two images belong to the same person. We use the Labeled Faces in the Wild (LFW) [30] dataset. This is a fairly difficult dataset for face recognition, because it has a large range of the variation (varying pose: straight, left, right, up; expression: neutral, happy, sad, angry; eyes: wearing glasses or not; clothes: wearing different clothes; size: small, medium or large) seen in real life. It includes 13233 images of 5749 people collected from news articles on the Internet. The number of images per person ranges from 1 to 530, and 1680 people have two or more distinct images in the dataset. This is a popular dataset which has been used by many researchers [4, 6, 34, 8, 35] to evaluate their face recognition frameworks.

In this experiment, we have compared the proposed DRMetric to the state-of-the-art methods for the task of face pair-matching problem. To ensure fairness,

we used the same features as used in the literature [4–9, 36]. Features of face images are extracted by computing the 3-scale, 128-dimensional SIFT descriptors [37], centered on 9 points of facial features extracted by a facial feature descriptor, as described in [5]. In this way, we get  $3 \times 128 \times 9 = 3456$  features in total for each image. PCA is then performed on the feature vectors to reduce the dimension to 400 (because the result in [9] showed that dimension 400 is a good compromise between performance and efficiency) for training. The triplets are built according to Section 3.3. The number of generated triplets is 44794, out of which, we use 20% (i.e. 8960) for training. We compared our method with LDML funneled [38], Hybrid aligned [36], V1-like funneled [8], Simile [6], Attribute + Simile [6], Background sample [39] Multiple LE + comp [4], and FrobMetric [9]. All these methods except FrobMetric are more complicated than our method, because they either use additional information, hybrid descriptors or combination of classifiers. The performance is described using an ROC curve<sup>4</sup> on which each point represents the average over the 10 runs of (false positive rate, true positive rate) for a fixed distance threshold. The results from all other techniques were taken from their latest published results. The comparison is shown in Fig. 5, which depicts that our method is only slightly worse than the two leading techniques (Attribute + Simile [6] and Background sample [39]) that are much more complicated. Furthermore, our method is comparable to or better than other state-of-the-art techniques on face recognition.



**Fig. 5.** False positive (FP) rate versus true positive (TP) rate on face recognition using different metric learning methods on the LFW dataset

<sup>4</sup> If the distance, based on the learned metric, is above some threshold, the two images will be declared as not belonging to the same person, and vice versa. For each threshold, we get the corresponding FP/TP rate. By changing the thresholds, we get a set of {FP/TP rate}, which forms the ROC curve. Using ROC curve to evaluate face recognition methods is widely used in literature [38, 36, 6, 39, 8].

## 5 Conclusions

We proposed an efficient and robust doubly regularized metric learning algorithm DRMetric. It has two regularization parts. First, we use tKL to regularize the update of the weight of the training examples. This avoids instabilities in the weight change, and consequently avoids overfitting and make it more robust to noisy data as well as outliers. Second, we add a regularization to the rank-one matrices enforcing them to be as independent as possible. In this way, the redundancy of the learned rank-one matrices as well as the number of necessary rank-one matrices are significantly reduced, which leads to higher efficiency. Furthermore, DRMetric is robust and capable of handling a variety of datasets for different applications. Though the idea behind DRMetric seems simple, its robustness and applicability can not be undervalued.

## References

1. Bi, J., et al.: AdaBoost on low-rank PSD matrices for metric learning with applications in Computer Aided Diagnosis. In: IEEE CVPR, pp. 1049–1056 (2011)
2. Yang, L., et al.: A Boosting framework for visually-preserving distance metric learning and its application to medical image retrieval. TPAMI, 30–44 (2010)
3. Ong, E., Bowden, R.: A boosted classifier tree for hand shape detection. In: IEEE Int. Conf. Automatic Face & Gesture Recogn., pp. 889–894 (2004)
4. Cao, Z., Yin, Q., Tang, X., Sun, J.: Face recognition with learning-based descriptor. In: IEEE CVPR, pp. 2707–2714 (2010)
5. Guillaumin, M., Verbeek, J., Schmid, C.: Multimodal semi-supervised learning for image classification. In: IEEE CVPR, pp. 902–909 (2010)
6. Kumar, N., Berg, A.C., Belhumeur, P.N., Nayar, S.K.: Attribute and simile classifiers for face verification. In: IEEE ICCV, pp. 365–372 (2009)
7. Nowak, E., Jurie, F.: Learning visual similarity measures for comparing never seen objects. In: IEEE CVPR (2007)
8. Pinto, N., DiCarlo, J.J., Cox, D.D.: How far can you get with a modern face recognition test set using only simple features? In: IEEE CVPR (2009)
9. Shen, C., Kim, J., Wang, L.: A scalable dual approach to semidefinite metric learning. In: IEEE CVPR, pp. 2601–2608 (2011)
10. Jiang, N., Liu, W., Wu, Y.: Adaptive and Discriminative Metric Differential Tracking. In: IEEE CVPR (2011)
11. Xing, E., Ng, A., Jordan, M., Russell, S.: Distance metric learning, with application to clustering with side-information. NIPS 15, 505–512 (2002)
12. Shen, C., Kim, J., Wang, L., van den Hengel, A.: Positive semidefinite metric learning with boosting. In: NIPS (2009)
13. Bar-Hillel, A., Hertz, T., Shental, N., Weinshall, D.: Learning a mahalanobis metric from equivalence constraints. JMLR 6, 937–965 (2005)
14. Davis, J., Kulis, B., Jain, P., Sra, S., Dhillon, I.S.: Information-theoretic metric learning. In: ICML, pp. 209–216 (2007)
15. Cox, M., Cox, T.: Multidimensional Scaling. In: Handbooks Comp. Statistics, pp. 315–347. Springer (2008)
16. Tenenbaum, J., Silva, V., Langford, J.: A Global Geometric Framework for Non-linear Dimensionality Reduction. Science 290, 2319–2323 (2000)

17. Roweis, S., Saul, L.: Nonlinear dimensionality reduction by locally linear embedding. *Science* 290, 2323–2326 (2000)
18. Vemuri, B.C., Liu, M., Amari, S.I., Nielsen, F.: Total Bregman divergence and its applications to DTI analysis. *IEEE TMI* 30, 475–483 (2011)
19. Shalev-Shwartz, S., Singer, Y., Ng, A.Y.: Online and batch learning of pseudo-metrics. In: *ICML* (2004)
20. Tsuda, K., Rsch, G., Warmuth, M.K.: Matrix exponentiated gradient updates for on-line learning and bregman projection. *JMLR* 6, 995–1018 (2005)
21. Jian, B., Vemuri, B.C.: Metric learning using iwasawa decomposition. In: *IEEE ICCV*, pp. 1–6 (2007)
22. Saberian, M.J., Vasconcelos, N.: Multiclass Boosting: Theory and Algorithms. In: *NIPS* (2011)
23. Liu, M., Vemuri, B.C.: Robust and efficient regularized boosting using total bregman divergence. In: *IEEE CVPR*, pp. 2897–2902 (2011)
24. Warmuth, M.K., Glocer, K.A., Vishwanathan, S.V.: Entropy regularized LPBoost. In: *Int. Conf. Alg. Learn. Theory*, pp. 256–271 (2008)
25. Liu, M., Vemuri, B.C., Amari, S., Nielsen, F.: Shape retrieval using hierarchical total bregman soft clustering. *IEEE TPAMI* (2012)
26. Boyd, S., Vandenberghe, L.: *Convex optimization*. Cambridge University Press, Cambridge (2004)
27. Demiriz, A., Bennett, K.P., Shawe-Taylor, J.: Linear programming boosting via column generation. *Mach. Learn.* 46, 225–254 (2002)
28. Schapire, R., Singer, Y.: Improved boosting algorithms using confidence-rated predictions. *Mach. Learn.* 37, 297–336 (1999)
29. Frank, A., Asuncion, A.: *UCI machine learning repository* (2010)
30. Huang, G.B., Ramesh, M., Berg, T., Learned-Miller, E.: Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In: *ECCV* (2008)
31. French, J., Watson, J., Jin, X., Martin, W.: An exogenous approach for adding multiple image representations to content-based image retrieval systems. In: *Int. Sym. Signal Processing App.*, vol. 1, pp. 201–204 (2003)
32. Canny, J.: A computational approach to edge detection. *IEEE TPAMI* 8, 679–698 (1986)
33. Manjunath, B.S., Ma, W.: Texture features for browsing and retrieval of image data. *IEEE TPAMI* 18, 837–842 (1996)
34. Nguyen, H.V., Bai, L.: Cosine Similarity Metric Learning for Face Verification. In: Kimmel, R., Klette, R., Sugimoto, A. (eds.) *ACCV 2010, Part II. LNCS*, vol. 6493, pp. 709–720. Springer, Heidelberg (2011)
35. Yin, Q., Tang, X., Sun, J.: An Associate-Predict Model for Face Recognition. In: *IEEE CVPR*, pp. 497–504 (2011)
36. Taigman, Y., Wolf, L., Hassner, T., Tel-Aviv, I.: Multiple One-Shots for utilizing class label information. In: *BMVC* (2009)
37. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. *IJCV* 60, 91–110 (2004)
38. Guillaumin, M., Verbeek, J., Schmid, C.: Is that you? metric learning approaches for face identification. In: *IEEE ICCV*, pp. 498–505 (2009)
39. Wolf, L., Hassner, T., Taigman, Y.: Similarity Scores Based on Background Samples. In: Zha, H., Taniguchi, R.-I., Maybank, S. (eds.) *ACCV 2009, Part II. LNCS*, vol. 5995, pp. 88–97. Springer, Heidelberg (2010)