

# Sparse Compositional Metric Learning

Yuan Shi, Aurélien Bellet, Fei Sha

► **To cite this version:**

Yuan Shi, Aurélien Bellet, Fei Sha. Sparse Compositional Metric Learning. AAI Conference on Artificial Intelligence (AAAI 2014), Jul 2014, Quebec City, Canada. hal-01430847

**HAL Id: hal-01430847**

**<https://hal.inria.fr/hal-01430847>**

Submitted on 12 Jan 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Sparse Compositional Metric Learning\*

Yuan Shi<sup>† ‡</sup>, Aurélien Bellet<sup>† ‡</sup>, Fei Sha<sup>‡</sup>

## Abstract

We propose a new approach for metric learning by framing it as learning a sparse combination of locally discriminative metrics that are inexpensive to generate from the training data. This flexible framework allows us to naturally derive formulations for global, multi-task and local metric learning. The resulting algorithms have several advantages over existing methods in the literature: a much smaller number of parameters to be estimated and a principled way to generalize learned metrics to new testing data points. To analyze the approach theoretically, we derive a generalization bound that justifies the sparse combination. Empirically, we evaluate our algorithms on several datasets against state-of-the-art metric learning methods. The results are consistent with our theoretical findings and demonstrate the superiority of our approach in terms of classification performance and scalability.

## 1 Introduction

The need for measuring distance or similarity between data instances is ubiquitous in machine learning and many application domains. However, each problem has its own underlying semantic space for defining distances that standard metrics (e.g., the Euclidean distance) often fail to capture. This has led to a growing interest in *metric learning* for the past few years, as summarized in two recent surveys (Bellet et al., 2013; Kulis, 2012). Among these methods, learning a globally linear Mahalanobis distance is by far the most studied setting. Representative methods include (Xing et al., 2002; Goldberger et al., 2004; Davis et al., 2007; Jain et al., 2008; Weinberger and Saul, 2009; Shen et al., 2012; Ying and Li, 2012). This is equivalent to learning a linear projection of the data to a feature space where constraints on the training set (such as “ $\mathbf{x}_i$  should be closer to  $\mathbf{x}_j$  than to  $\mathbf{x}_k$ ”) are better satisfied.

Although the performance of these learned metrics is typically superior to that of standard metrics in practice, a single linear metric is often unable to accurately capture the complexity of the task, for instance when the data are multimodal or the decision boundary is complex. To overcome this limitation, recent work has focused on learning *multiple locally linear metrics* at several locations of the feature space (Frome et al., 2007; Weinberger and Saul, 2009; Zhan et al., 2009; Hong et al., 2011; Wang et al., 2012), to the extreme of learning one metric per training instance (Noh et al., 2010). This line of research is motivated by the fact that locally, simple linear metrics perform well (Ramanan and Baker, 2011; Hauberg et al., 2012). The main challenge is to integrate these metrics into a meaningful global one while keeping the number of learning parameters to a reasonable level in order to avoid heavy computational burden and severe overfitting. So far, existing methods are not able to compute valid (smooth) global metrics from the local metrics they learn and do not provide a principled way of generalizing to new regions of the space at test time. Furthermore,

---

\*This document is an extended version of a conference paper (Shi et al., 2014) that provides additional details and results.

<sup>†</sup>Equal contribution.

<sup>‡</sup>Department of Computer Science, University of Southern California, {yuanshi, bellet, feisha}@usc.edu.

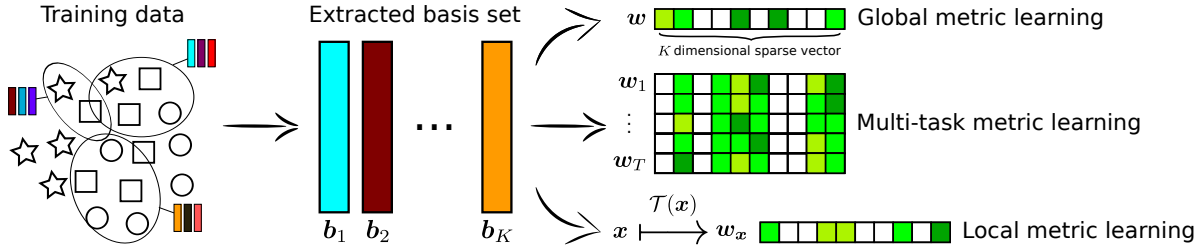


Figure 1: Illustration of the general framework and its applications. We extract locally discriminative basis elements from the training data and cast metric learning as learning sparse combinations of these elements. We formulate global metric learning as learning a single sparse weight vector  $w$ . For multi-task metric learning, we learn a vector  $w_t$  for each task where all tasks share the same basis subset. For local metric learning we learn a function  $\mathcal{T}(x)$  that maps any instance  $x$  to its associated sparse weight vector  $w_x$ . Shades of grey encode weight magnitudes.

they scale poorly with the dimensionality  $D$  of the data: typically, learning a Mahalanobis distance requires  $O(D^2)$  parameters and the optimization involves projections onto the positive semidefinite cone that scale in  $O(D^3)$ . This is expensive even for a single metric when  $D$  is moderately large.

In this paper, we study metric learning from a new perspective to efficiently address these key challenges. We propose to learn metrics as *sparse compositions of locally discriminative metrics*. These “basis metrics” are low-rank and extracted efficiently from the training data at different local regions, for instance using Fisher discriminant analysis. Learning higher-rank linear metrics is then formulated as learning the combining weights, using sparsity-inducing regularizers to select only the most useful basis elements. This provides a unified framework for metric learning, as illustrated in Figure 1, that we call SCML (for Sparse Compositional Metric Learning). In SCML, the number of parameters to learn is much smaller than existing approaches and projections onto the positive semidefinite cone are not needed. This gives an efficient and flexible way to learn a single global metric when  $D$  is large.

The proposed framework also applies to multi-task metric learning, where one wants to learn a global metric for several related tasks while exploiting commonalities between them (Caruana, 1997; Parameswaran and Weinberger, 2010). This is done in a natural way by means of a group sparsity regularizer that makes the task-specific metrics share the same basis subset. Our last and arguably most interesting contribution is a new formulation for local metric learning, where we learn a transformation  $\mathcal{T}(x)$  that takes as input any instance  $x$  and outputs a sparse weight vector defining its metric. This can be seen as learning a smoothly varying metric tensor over the feature space (Ramanan and Baker, 2011; Hauberg et al., 2012). To the best of our knowledge, it is the first discriminative metric learning approach capable of computing, in a principled way, an instance-specific metric for *any* point in the feature space. All formulations can be solved using scalable optimization procedures based on stochastic subgradient descent with proximal operators (Duchi and Singer, 2009; Xiao, 2010).

We present both theoretical and experimental evidence supporting the proposed approach. We derive a generalization bound which provides a theoretical justification to seeking sparse combinations and suggests that the basis set  $B$  can be large without incurring overfitting. Empirically, we evaluate our algorithms against state-of-the-art global, local and multi-task metric learning methods on several datasets. The results strongly support the proposed framework.

The rest of this paper is organized as follows. Section 2 describes our general framework and illustrates how it can be used to derive efficient formulations for global, local and multi-task metric learning. Section 3

provides a theoretical analysis supporting our approach. Section 4 reviews related work. Section 5 presents an experimental evaluation of the proposed methods. We conclude in Section 6.

## 2 Proposed Approach

In this section, we present the main idea of sparse compositional metric learning (SCML) and show how it can be used to unify several existing metric learning paradigms and lead to efficient new formulations.

### 2.1 Main Idea

We assume the data lie in  $\mathbb{R}^D$  and focus on learning (squared) Mahalanobis distances  $d_M(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \mathbf{x}')^T \mathbf{M} (\mathbf{x} - \mathbf{x}')$  parameterized by a positive semidefinite (PSD)  $D \times D$  matrix  $\mathbf{M}$ . Note that  $\mathbf{M}$  can be represented as a nonnegative weighted sum of  $K$  rank-1 PSD matrices:<sup>1</sup>

$$\mathbf{M} = \sum_{i=1}^K w_i \mathbf{b}_i \mathbf{b}_i^T, \quad \text{with } \mathbf{w} \geq 0, \quad (1)$$

where the  $\mathbf{b}_i$ 's are  $D$ -dimensional column vectors.

In this paper, we use the form (1) to cast metric learning as learning a *sparse combination of basis elements* taken from a basis set  $B = \{\mathbf{b}_i\}_{i=1}^K$ . The key to our framework is the fact that such a  $B$  is made readily available to the algorithm and consists of rank-one metrics that are *locally discriminative*. Such basis elements can be easily generated from the training data at several local regions — in the experiments, we simply use Fisher discriminant analysis (see the corresponding section for details). They can then be combined to form a single global metric, multiple global metrics (in the multi-task setting) or a metric tensor (implicitly defining an infinite number of local metrics) that varies smoothly across the feature space, as we will show in later sections.

We use the notation  $d_{\mathbf{w}}(\mathbf{x}, \mathbf{x}')$  to highlight our parameterization of the Mahalanobis distance by  $\mathbf{w}$ . Learning  $\mathbf{M}$  in this form makes it PSD by design (as a nonnegative sum of PSD matrices) and involves  $K$  parameters (instead of  $D^2$  in most metric learning methods), enabling it to more easily deal with high-dimensional problems. We also want the combination to be *sparse*, i.e., some  $w_i$ 's are zero and thus  $\mathbf{M}$  only depends on a small subset of  $B$ . This provides some form of regularization (as shown later in Theorem 1) as well as a way to tie metrics together when learning multiple metrics. In the rest of this section, we apply the proposed framework to several metric learning paradigms (see Figure 1). We start with the simple case of global metric learning (Section 2.1.1) before considering more challenging settings: multi-task (Section 2.1.2) and local metric learning (Section 2.1.3). Finally, Section 2.2 discusses how these formulations can be solved in a scalable way using stochastic subgradient descent with proximal operators.

#### 2.1.1 Global Metric Learning

In global metric learning, one seeks to learn a single metric  $d_{\mathbf{w}}(\mathbf{x}, \mathbf{x}')$  from a set of distance constraints on the training data. Here, we use a set of triplet constraints  $C$  where each  $(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \in C$  indicates that the distance between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  should be smaller than the distance between  $\mathbf{x}_i$  and  $\mathbf{x}_k$ .  $C$  may be constructed from label information, as in LMNN (Weinberger and Saul, 2009), or in an unsupervised manner based for instance on implicit users' feedback (such as clicks on search engine results). Our formulation for global

---

<sup>1</sup>Such an expression exists for any PSD matrix  $\mathbf{M}$  since the eigenvalue decomposition of  $\mathbf{M}$  is of the form (1).

metric learning, SCML-Global, is simply to combine the local basis elements into a higher-rank global metric that satisfies well the constraints in  $C$ :

$$\min_{\mathbf{w}} \frac{1}{|C|} \sum_{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \in C} L_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) + \beta \|\mathbf{w}\|_1, \quad (2)$$

where  $L_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) = [1 + d_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_j) - d_{\mathbf{w}}(\mathbf{x}_i, \mathbf{x}_k)]_+$  with  $[\cdot]_+ = \max(0, \cdot)$ , and  $\beta \geq 0$  is a regularization parameter. The first term in (2) is the classic margin-based hinge loss function. The second term  $\|\mathbf{w}\|_1 = \sum_{i=1}^K w_i$  is the  $\ell_1$  norm regularization which encourages sparse solutions, allowing the selection of relevant basis elements. SCML-Global is convex by the linearity of both terms and is bounded below, thus it has a global minimum.

### 2.1.2 Multi-Task Metric Learning

Multi-task learning (Caruana, 1997) is a paradigm for learning several tasks simultaneously, exploiting their commonalities. When tasks are related, this can perform better than separately learning each task. Recently, multi-task learning methods have successfully built on the assumption that the tasks should share a common low-dimensional representation (Argyriou et al., 2008; Yang et al., 2009; Gong et al., 2012). In general, it is unclear how to achieve this in metric learning. In contrast, learning metrics as sparse combinations allows a direct translation of this idea to multi-task metric learning.

Formally, we are given  $T$  different but somehow related tasks with associated constraint sets  $C_1, \dots, C_T$  and we aim at learning a metric  $d_{\mathbf{w}_t}(\mathbf{x}, \mathbf{x}')$  for each task  $t$  while sharing information across tasks. In the following, the basis set  $B$  is the union of the basis sets  $B_1, \dots, B_T$  extracted from each task  $t$ . Our formulation for multi-task metric learning, mt-SCML, is as follows:

$$\min_{\mathbf{W}} \sum_{t=1}^T \frac{1}{|C_t|} \sum_{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \in C_t} L_{\mathbf{w}_t}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) + \beta \|\mathbf{W}\|_{2,1},$$

where  $\mathbf{W}$  is a  $T \times K$  nonnegative matrix whose  $t$ -th row is the weight vector  $\mathbf{w}_t$  defining the metric for task  $t$ ,  $L_{\mathbf{w}_t}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) = [1 + d_{\mathbf{w}_t}(\mathbf{x}_i, \mathbf{x}_j) - d_{\mathbf{w}_t}(\mathbf{x}_i, \mathbf{x}_k)]_+$  and  $\|\mathbf{W}\|_{2,1}$  is the  $\ell_2/\ell_1$  mixed norm used in the group lasso problem (Yuan and Lin, 2006). It corresponds to the  $\ell_1$  norm applied to the  $\ell_2$  norm of the columns of  $\mathbf{W}$  and is known to induce group sparsity at the column level. In other words, this regularization makes most basis elements either have zero weight or nonzero weight *for all tasks*.

Overall, while each metric remains task-specific ( $d_{\mathbf{w}_t}$  is only required to satisfy well the constraints in  $C_t$ ), it is composed of *shared features* (i.e., it potentially benefits from basis elements generated from other tasks) that are regularized to be relevant *across tasks* (as favored by the group sparsity). As a result, all learned metrics can be expressed as combinations of the same basis subset of  $B$ , though with different weights for each task. Since the  $\ell_2/\ell_1$  norm is convex, mt-SCML is again convex.

### 2.1.3 Local Metric Learning

Local metric learning addresses the limitations of global methods in capturing complex data patterns (Frome et al., 2007; Weinberger and Saul, 2009; Zhan et al., 2009; Noh et al., 2010; Hong et al., 2011; Wang et al., 2012). For heterogeneous data, allowing the metric to vary across the feature space can capture the semantic distance much better. On the other hand, local metric learning is costly and often suffers from severe overfitting since the number of parameters to learn can be very large. In the following, we show how our framework can be used to derive an efficient local metric learning method.

We aim at learning a *metric tensor*  $\mathcal{T}(\mathbf{x})$ , which is a smooth function that (informally) maps any instance  $\mathbf{x}$  to its metric matrix (Ramanan and Baker, 2011; Hauberg et al., 2012). The distance between two points should then be defined as the geodesic distance on a Riemannian manifold. However, this requires solving an intractable problem, so we use the widely-adopted simplification that distances from point  $\mathbf{x}$  are computed based on its own metric alone (Zhan et al., 2009; Noh et al., 2010; Wang et al., 2012):

$$\begin{aligned} d_{\mathcal{T}}(\mathbf{x}, \mathbf{x}') &= (\mathbf{x} - \mathbf{x}')^{\top} \mathcal{T}(\mathbf{x})(\mathbf{x} - \mathbf{x}') \\ &= (\mathbf{x} - \mathbf{x}')^{\top} \sum_{i=1}^K w_{\mathbf{x},i} \mathbf{b}_i \mathbf{b}_i^{\top} (\mathbf{x} - \mathbf{x}'), \end{aligned}$$

where  $w_{\mathbf{x}}$  is the weight vector for instance  $\mathbf{x}$ .

We could learn a weight vector for each training point. This would result in a formulation similar to mt-SCML, where each training instance is considered as a task. However, in the context of local metric learning, this is not an appealing solution. Indeed, for a training sample of size  $S$  we would need to learn  $SK$  parameters, which is computationally difficult and leads to heavy overfitting for large-scale problems. Furthermore, this gives no principled way of computing the weight vector of a test instance.

We instead propose a more effective solution by constraining the weight vector for an instance  $\mathbf{x}$  to parametrically depend on some embedding of  $\mathbf{x}$ :

$$\mathcal{T}_{\mathbf{A},\mathbf{c}}(\mathbf{x}) = \sum_{i=1}^K (\mathbf{a}_i^{\top} \mathbf{z}_{\mathbf{x}} + c_i)^2 \mathbf{b}_i \mathbf{b}_i^{\top}, \quad (3)$$

where  $\mathbf{z}_{\mathbf{x}} \in D'$  is an embedding of  $\mathbf{x}$ ,<sup>2</sup>  $\mathbf{A} = [\mathbf{a}_1 \dots \mathbf{a}_K]^{\top}$  is a  $D' \times K$  real-valued matrix and  $\mathbf{c} \in \mathbb{R}^K$ . The square makes the weights nonnegative  $\forall \mathbf{x} \in \mathbb{R}^D$ , ensuring that they define a valid (pseudo) metric. Intuitively, (3) combines the locally discriminative metrics with weights that depend on the position of the instance in the feature space.

There are several advantages to this formulation. First, by learning  $\mathbf{A}$  and  $\mathbf{c}$  we implicitly learn a different metric not only for the training data but for any point in the feature space. Second, if the embedding is smooth,  $\mathcal{T}_{\mathbf{A},\mathbf{c}}(\mathbf{x})$  is a smooth function of  $\mathbf{x}$ , therefore similar instances are assigned similar weights. This can be seen as some kind of manifold regularization. Third, the number of parameters to learn is now  $K(D' + 1)$ , thus independent of both the size of the training sample and the dimensionality of  $\mathbf{x}$ . Our formulation for local metric learning, SCML-Local, is as follows:

$$\min_{\tilde{\mathbf{A}}} \frac{1}{|C|} \sum_{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) \in C} L_{\mathcal{T}_{\mathbf{A},\mathbf{c}}}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) + \beta \|\tilde{\mathbf{A}}\|_{2,1},$$

where  $\tilde{\mathbf{A}}$  is a  $(D' + 1) \times K$  matrix denoting the concatenation of  $\mathbf{A}$  and  $\mathbf{c}$ , and  $L_{\mathcal{T}_{\mathbf{A},\mathbf{c}}}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k) = [1 + d_{\mathcal{T}_{\mathbf{A},\mathbf{c}}}(\mathbf{x}_i, \mathbf{x}_j) - d_{\mathcal{T}_{\mathbf{A},\mathbf{c}}}(\mathbf{x}_i, \mathbf{x}_k)]_+$ . The  $\ell_2/\ell_1$  norm on  $\tilde{\mathbf{A}}$  introduces sparsity at the column level, regularizing the local metrics to use the same basis subset. Interestingly, if  $\mathbf{A}$  is the zero matrix, we recover SCML-Global. SCML-Local is nonconvex and is thus subject to local minima.

## 2.2 Optimization

Our formulations use (nonsmooth) sparsity-inducing regularizers and typically involve a large number of triplet constraints. We can solve them efficiently using stochastic composite optimization (Duchi and Singer,

<sup>2</sup>In our experiments, we use kernel PCA (Schölkopf et al., 1998) as it provides a simple way to limit the dimension and thus the number of parameters to learn. We use RBF kernel with bandwidth set to the median Euclidean distance in the data.

2009; Xiao, 2010), which alternates between a stochastic subgradient step on the hinge loss term and a proximal operator (for  $\ell_1$  or  $\ell_{2,1}$  norm) that explicitly induces sparsity. We solve SCML-Global and mt-SCML using Regularized Dual Averaging (Xiao, 2010), which offers fast convergence and levels of sparsity in the solution comparable to batch algorithms. For SCML-Local, due to local minima, we ensure improvement over the optimal solution  $\mathbf{w}^*$  of SCML-Global by using a forward-backward algorithm (Duchi and Singer, 2009) which is initialized with  $\mathbf{A} = \mathbf{0}$  and  $c_i = \sqrt{w_i^*}$ .

Recall that unlike most existing metric learning algorithms, we do not need to perform projections onto the PSD cone, which scale in  $O(D^3)$  for a  $D \times D$  matrix. Our algorithms thereby have a significant computational advantage for high-dimensional problems.

### 3 Theoretical Analysis

In this section, we provide a theoretical analysis of our approach in the form of a generalization bound based on algorithmic robustness analysis (Xu and Mannor, 2012) and its adaptation to metric learning (Bellet and Habrard, 2012). For simplicity, we focus on SCML-Global, our global metric learning formulation in (2).

Consider the supervised learning setting, where we are given a labeled training sample  $S = \{z_i = (\mathbf{x}_i, y_i)\}_{i=1}^n$  drawn i.i.d. from some unknown distribution  $P$  over  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ . We call a triplet  $(z, z', z'')$  *admissible* if  $y = y' \neq y''$ . Let  $C$  be the set of admissible triplets built from  $S$  and  $L(\mathbf{w}, z, z', z'') = [1 + d_{\mathbf{w}}(\mathbf{x}, \mathbf{x}') - d_{\mathbf{w}}(\mathbf{x}, \mathbf{x}'')]_+$  denote the loss function used in (2), with the convention that  $L$  returns 0 for non-admissible triplets.

Let us define the *empirical loss* of  $\mathbf{w}$  on  $S$  as

$$\mathcal{R}_{emp}^S(\mathbf{w}) = \frac{1}{|C|} \sum_{(z, z', z'') \in C} L(\mathbf{w}, z, z', z''),$$

and its *expected loss* over distribution  $P$  as

$$\mathcal{R}(\mathbf{w}) = \mathbb{E}_{z, z', z'' \sim P} L(\mathbf{w}, z, z', z'').$$

The following theorem bounds the deviation between the empirical loss of the learned metric and its expected loss.

**Theorem 1.** *Let  $\mathbf{w}^*$  be the optimal solution to SCML-Global with  $K$  basis elements,  $\beta > 0$  and  $C$  constructed from  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  as above. Let  $K^* \leq K$  be the number of nonzero entries in  $\mathbf{w}^*$ . Let us assume the norm of any instance bounded by some constant  $R$  and  $L$  uniformly upper-bounded by some constant  $U$ . Then for any  $\delta > 0$ , with probability at least  $1 - \delta$  we have:*

$$|\mathcal{R}(\mathbf{w}^*) - \mathcal{R}_{emp}^S(\mathbf{w}^*)| \leq \frac{16\gamma RK^*}{\beta} + 3U \sqrt{\frac{N \ln 2 + \ln \frac{1}{\delta}}{0.5n}},$$

where  $N$  is the size of an  $\gamma$ -cover of  $\mathcal{Z}$ .

This bound has a standard  $O(1/\sqrt{n})$  asymptotic convergence rate.<sup>3</sup> Its main originality is that it provides a theoretical justification to enforcing sparsity in our formulation. Indeed, notice that  $K^*$  (and not  $K$ )

<sup>3</sup>In robustness bounds, the cover radius  $\gamma$  can be made arbitrarily close to zero at the expense of increasing  $N$ . Since  $N$  appears in the second term, the right hand side of the bound indeed goes to zero when  $n \rightarrow \infty$ . This is in accordance with other similar learning bounds, for example, the original robustness-based bounds in (Xu and Mannor, 2012).

appears in the bound as a penalization term, which suggests that one may use a large basis set  $K$  without overfitting as long as  $K^*$  remains small. This will be confirmed by our experiments (Section 5.3). A similar bound can be derived for mt-SCML, but not for SCML-Local because of its nonconvexity. The details and proofs can be found in Appendix A.

## 4 Related Work

In this section, we review relevant work in global, multi-task and local metric learning. The interested reader should refer to the recent surveys of Kulis (2012) and Bellet et al. (2013) for more details.

**Global methods** Most global metric learning methods learn the matrix  $M$  directly: see (Xing et al., 2002; Goldberger et al., 2004; Davis et al., 2007; Jain et al., 2008; Weinberger and Saul, 2009) for representative papers. This is computationally expensive and subject to overfitting for moderate to high-dimensional problems. An exception is BoostML (Shen et al., 2012) which uses rank-one matrices as weak learners to learn a global Mahalanobis distance via a boosting procedure. However, it is not clear how BoostML can be generalized to multi-task or local metric learning.

**Multi-task methods** Multi-task metric learning was proposed in (Parameswaran and Weinberger, 2010) as an extension to the popular LMNN (Weinberger and Saul, 2009). The authors define the metric for task  $t$  as  $d_t(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \mathbf{x}')^T(\mathbf{M}_0 + \mathbf{M}_t)(\mathbf{x} - \mathbf{x}')$ , where  $\mathbf{M}_t$  is task-specific and  $\mathbf{M}_0$  is shared by all tasks. Note that it is straightforward to incorporate their approach in our framework by defining a shared weight vector  $\mathbf{w}_0$  and task-specific weights  $\mathbf{w}_t$ . However, this assumption of a metric that is common to all tasks can be too restrictive in cases where task relatedness is complex, as illustrated by our experiments.

**Local methods** MM-LMNN (Weinberger and Saul, 2009) is an extension of LMNN which learns only a small number of metrics (typically one per class) in an effort to alleviate overfitting. However, no additional regularization is used and a full-rank metric is learned for each class, which becomes intractable when the number of classes is large. msNCA (Hong et al., 2011) learns a function that splits the space into a small number of regions and then learns a metric per region using NCA (Goldberger et al., 2004). Again, the metrics are full-rank so msNCA does not scale well with the number of metrics. Like SCML-Local, PLML (Wang et al., 2012) is based on a combination of metrics but there are major differences with our work: (i) weights only depend on a manifold assumption: they are not sparse and use no discriminative information, (ii) the basis metrics are full-rank, thus expensive to learn, and (iii) a weight vector is learned explicitly for each training instance, which can result in a large number of parameters and prevents generalization to new instances (in practice, for a test point, they use the weight vector of its nearest neighbor in the training set). As observed by Ramanan and Baker (2011), the above methods make the implicit assumption that the metric tensor is locally constant (at the class, region or neighborhood level), while SCML-Local learns a smooth function that maps any instance to its specific metric. ISD (Zhan et al., 2009) is an attempt to learn the metrics for unlabeled points by propagation, but is limited to the transductive setting. Unlike the above discriminative approaches, GLML (Noh et al., 2010) learns a metric for each point independently in a generative way by minimizing the 1-NN expected error under some assumption for the class distributions.



	Vehicle	Vowel	Segment	Letters	USPS	BBC
# samples	846	990	2,310	20,000	9,298	2,225
# classes	4	11	7	26	10	5
# features	18	10	19	16	256	9,636

Table 1: Datasets for global and local metric learning.

Dataset	Euc	Global-Frob	SCML-Global
Vehicle	29.7±0.6	<b>21.5±0.8</b>	<b>21.3±0.6</b>
Vowel	<b>11.1±0.4</b>	<b>10.3±0.4</b>	<b>10.9±0.5</b>
Segment	5.2±0.2	<b>4.1±0.2</b>	<b>4.1±0.2</b>
Letters	14.0±0.2	<b>9.0±0.2</b>	<b>9.0±0.2</b>
USPS	10.3±0.2	5.1±0.2	<b>4.1±0.1</b>
BBC	8.8±0.3	5.5±0.3	<b>3.9±0.2</b>

Table 2: Global metric learning results (best in bold).

## 5 Experiments

In this section, we compare our methods to state-of-the-art algorithms on global, multi-task and local metric learning.<sup>4</sup> We use a 3-nearest neighbor classifier in all experiments. To generate a set of locally discriminative rank-one metrics, we first divide data into regions via clustering. For each region center, we select  $J$  nearest neighbors from each class (for  $J = \{10, 20, 50\}$  to account for different scales), and apply Fisher discriminant analysis followed by eigenvalue decomposition to obtain the basis elements.<sup>5</sup> Section 5.1 presents results for global metric learning, Section 5.2 for multi-task and Section 5.3 for local metric learning.

### 5.1 Global Metric Learning

We use 6 datasets from UCI<sup>6</sup> and BBC<sup>7</sup> (see Table 1). The dimensionality of USPS and BBC is reduced to 100 and 200 using PCA to speed up computation. We normalize the data as in (Wang et al., 2012) and split into train/validation/test (60%/20%/20%), except for Letters and USPS where we use 3,000/1,000/1,000. Results are averaged over 20 random splits.

#### 5.1.1 Proof of Concept

**Setup** Global metric learning is a convenient setting to study the effect of combining basis elements. To this end, we consider a formulation with the same loss function as SCML-Global but that directly learns the metric matrix, using Frobenius norm regularization to reduce overfitting. We refer to it as Global-Frob. Both algorithms use the same training triplets, generated by identifying 3 target neighbors (nearest neighbors with same label) and 10 imposters (nearest neighbors with different label) for each instance. We tune the regularization parameter on the validation data. For SCML-Global, we use a basis set of 400 elements for Vehicle, Vowel, Segment and BBC, and 1,000 elements for Letters and USPS.

<sup>4</sup>For all compared methods we use MATLAB code from the authors’ website. The MATLAB code for our methods is available at <http://www-bcf.usc.edu/~bellet/>.

<sup>5</sup>We also experimented with a basis set based on local GLML metrics. Preliminary results were comparable to those obtained with the procedure above.

<sup>6</sup><http://archive.ics.uci.edu/ml/>

<sup>7</sup><http://mlg.ucd.ie/datasets/bbc.html>

Dataset	Euc	LMNN	BoostML	SCML-Global
Vehicle	29.7±0.6	23.5±0.7	<b>19.9±0.6</b>	21.3±0.6
Vowel	<b>11.1±0.4</b>	<b>10.8±0.4</b>	11.4±0.4	<b>10.9±0.5</b>
Segment	5.2±0.2	4.6±0.2	<b>3.8±0.2</b>	<b>4.1±0.2</b>
Letters	14.0±0.2	11.6±0.3	10.8±0.2	<b>9.0±0.2</b>
USPS	10.3±0.2	<b>4.1±0.1</b>	7.1±0.2	<b>4.1±0.1</b>
BBC	8.8±0.3	<b>4.0±0.2</b>	9.3±0.3	<b>3.9±0.2</b>
Avg. rank	3.3	2.0	2.3	<b>1.2</b>

Table 3: Comparison of SCML-Global against LMNN and BoostML (best in bold).

Dataset	BoostML	SCML-Global
Vehicle	334	164
Vowel	19	47
Segment	442	49
Letters	20	133
USPS	2,375	300
BBC	3,000	59

Table 4: Average number of basis elements in the solution.

**Results** Table 2 shows misclassification rates with standard errors, where Euc is the Euclidean distance. The results show that SCML-Global performs similarly as Global-Frob on low-dimensional datasets but has a clear advantage when dimensionality is high (USPS and BBC). This demonstrates that learning a sparse combination of basis elements is an effective way to reduce overfitting and improve generalization. SCML-Global is also faster to train than Global-Frob on these datasets (about 2x faster on USPS and 3x on BBC) because it does not require any PSD projection.

### 5.1.2 Comparison to Other Global Algorithms

**Setup** We now compare SCML-Global to two state-of-the-art global metric learning algorithms: Large Margin Nearest Neighbor (LMNN, Weinberger and Saul, 2009) and BoostML (Shen et al., 2012). The datasets, preprocessing and setting for SCML-Global are the same as in Section 5.1.1. LMNN uses 3 target neighbors and all imposters, while these are set to 3 and 10 respectively for BoostML (as in SCML-Global).

**Results** Table 3 shows the average misclassification rates, along with standard error and the average rank of each method across all datasets. SCML-Global clearly outperforms LMNN and BoostML, ranking first on 5 out of 6 datasets and achieving the overall highest rank. Furthermore, its training time is smaller than competing methods, especially for high-dimensional data. For instance, on the BBC dataset, SCML-Global trained in about 90 seconds, which is about 20x faster than LMNN and 35x faster than BoostML. Note also that SCML-Global is consistently more accurate than linear SVM, as shown in Appendix B.

**Number of selected basis elements** Like SCML-Global, recall that BoostML is based on combining rank-one elements (see Section 4). The main difference with SCML-Global is that our method is given a set of locally discriminative metrics and picks the relevant ones by learning sparse weights, while BoostML generates a new basis element at each iteration and adds it to the current combination. Table 4 reports the number of basis elements used in SCML-Global and BoostML solutions. Overall, SCML-Global uses fewer

Task	st-Euc	st-LMNN	st-SCML	u-Euc	u-LMNN	u-SCML	mt-LMNN	mt-SCML
Books	33.5±0.5	29.7±0.4	27.0±0.5	33.7±0.5	29.6±0.4	28.0±0.4	29.1±0.4	25.8±0.4
DVD	33.9±0.5	29.4±0.5	26.8±0.4	33.9±0.5	29.4±0.5	27.9±0.5	29.5±0.5	26.5±0.5
Electronics	26.2±0.4	23.3±0.4	21.1±0.5	29.1±0.5	25.1±0.4	22.9±0.4	22.5±0.4	20.2±0.5
Kitchen	26.2±0.6	21.2±0.5	19.0±0.4	27.7±0.5	23.5±0.3	21.9±0.5	22.1±0.5	19.0±0.4
Avg. accuracy	30.0±0.2	25.9±0.2	23.5±0.2	31.1±0.3	26.9±0.2	25.2±0.2	25.8±0.2	<b>22.9±0.2</b>
Avg. runtime	N/A	57 min	3 min	N/A	44 min	2 min	41 min	5 min

Table 5: Multi-task metric learning results.

elements than BoostML (on two datasets, it uses more but this yields to significantly better performance). The results on USPS and BBC also suggest that the number of basis elements selected by SCML-Global seems to scale well with dimensionality. These nice properties come from its knowledge of the entire basis set and the sparsity-inducing regularizer. On the contrary, the number of elements (and therefore iterations) needed by BoostML to converge seems to scale poorly with dimensionality.

## 5.2 Multi-task Metric Learning

**Dataset** Sentiment Analysis (Blitzer et al., 2007) is a popular dataset for multi-task learning that consists of Amazon reviews on four product types: kitchen appliances, DVDs, books and electronics. Each product type is treated as a task and has 1,000 positive and 1,000 negative reviews. To reduce computational cost, we represent each review by a 200-dimensional feature vector by selecting top 200 words of the largest mutual information with the labels. We randomly split the dataset into training (800 samples), validation (400 samples) and testing (400 samples) sets.

**Setup** We compare the following metrics: st-Euc (Euclidean distance), st-LMNN and st-SCML (single-task LMNN and single-task SCML-Global, trained independently on each task), u-Euc (Euclidean trained on the union of the training data from all tasks), u-LMNN (LMNN on union), u-SCML (SCML-Global on union), multi-task LMNN (Parameswaran and Weinberger, 2010) and finally our own multi-task method mt-SCML. We tune the regularization parameters in mt-LMNN, st-SCML, u-SCML and mt-SCML on validation sets. As in the previous experiment, the number of target neighbors and imposters for our methods are set to 3 and 10 respectively. We use a basis set of 400 elements for each task for st-SCML, the union of these (1,600) for mt-SCML, and 400 for u-SCML.

**Results** Table 5 shows the results averaged over 20 random splits. First, notice that u-LMNN and u-SCML obtain significantly higher error rates than st-LMNN and st-SCML respectively, which suggests that the dataset may violate mt-LMNN’s assumption that all tasks share a similar metric. Indeed, mt-LMNN does not outperform st-LMNN significantly. On the other hand, mt-SCML performs better than its single-task counterpart and than all other compared methods by a significant margin, demonstrating its ability to leverage some commonalities between tasks that mt-LMNN is unable to capture. It is worth noting that the solution found by mt-SCML is based on only 273 basis elements on average (out of a total of 1,600), while st-SCML makes use of significantly more elements (347 elements *per task* on average). Basis elements selected by mt-SCML are evenly distributed across all tasks, which indicates that it is able to exploit meaningful information across tasks to get both more accurate and more compact metrics. Finally, note that our algorithms are about an order of magnitude faster.

Dataset	MM-LMNN	GLML	PLML	SCML-Local
Vehicle	23.1±0.6	23.4±0.6	22.8±0.7	<b>18.0±0.6</b>
Vowel	6.8±0.3	<b>4.1±0.4</b>	8.3±0.4	6.1±0.4
Segment	<b>3.6±0.2</b>	<b>3.9±0.2</b>	<b>3.9±0.2</b>	<b>3.6±0.2</b>
Letters	9.4±0.3	10.3±0.3	<b>8.3±0.2</b>	<b>8.3±0.2</b>
USPS	<b>4.2±0.7</b>	7.8±0.2	4.1±0.1	<b>3.6±0.1</b>
BBC	4.9±0.4	5.7±0.3	<b>4.3±0.2</b>	<b>4.1±0.2</b>
Avg. rank	2.0	2.7	2.0	<b>1.2</b>

Table 6: Local metric learning results (best in bold).

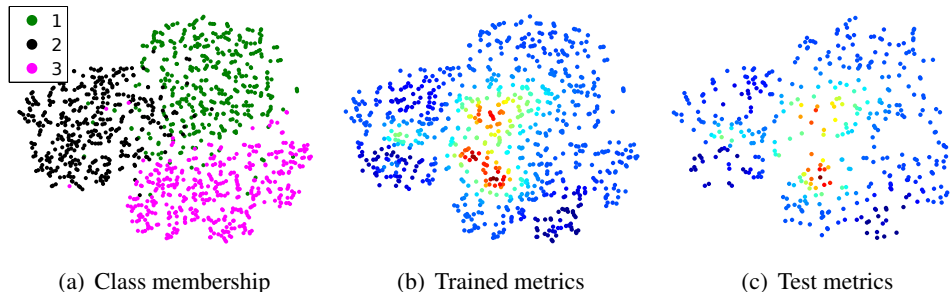


Figure 2: Illustrative experiment on digits 1, 2 and 3 of USPS in 2D. Refer to the main text for details.

### 5.3 Local Metric Learning

**Setup** We use the same datasets and preprocessing as for global metric learning. We compare SCML-Local to MM-LMNN (Weinberger and Saul, 2009), GLML (Noh et al., 2010) and PLML (Wang et al., 2012). The parameters of all methods are tuned on validation sets or set by authors’ recommendation. MM-LMNN use 3 target neighbors and all imposters, while these are set to 3 and 10 in PLML and SCML-Local. The number of anchor points in PLML is set to 20 as done by the authors. For SCML-Local, we use the same basis set as SCML-Global, and embedding dimension  $D'$  is set to 40 for Vehicle, Vowel, Segment and BBC, and 100 for Letters and USPS.

**Results** Table 6 gives the error rates along with the average rank of each method across all datasets. Note that SCML-Local significantly improves upon SCML-Global on all but one dataset and achieves the best average rank. PLML does not perform well on small datasets (Vehicle and Vowel), presumably because there are not enough points to get a good estimation of the data manifold. GLML is fast but has rather poor performance on most datasets because its Gaussian assumption is restrictive and it learns the local metrics independently. Among discriminative methods, SCML-Local offers the best training time, especially for high-dimensional data (e.g. on BBC, it trained in about 8 minutes, which is about 5x faster than MM-LMNN and 15x faster than PLML). Note that on this dataset, both MM-LMNN and PLML perform worse than SCML-Global due to severe overfitting, while SCML-Local avoids it by learning significantly fewer parameters. Finally, SCML-Local achieves accuracy results that are very competitive with those of a kernel SVM, as shown in Appendix B.

**Visualization of the learned metrics** To provide a better understanding of why SCML-Local works well, we apply it to digits 1, 2, and 3 of USPS projected in 2D using t-SNE (van der Maaten and Hinton, 2008),

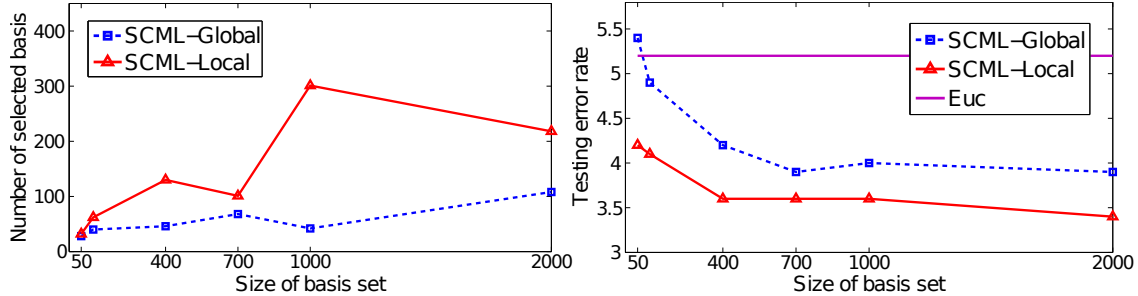


Figure 3: Effect of the number of bases on Segment dataset.

shown in Figure 2(a). We use 10 basis elements and  $D' = 5$ . Figure 2(b) shows the training points colored by their learned metric (based on the projection of the weight vectors in 1D using PCA). We see that the local metrics vary smoothly and are thereby robust to outliers. Unlike MM-LMNN, points within a class are allowed to have different metrics: in particular, this is useful for points that are near the decision boundary. While smooth, the variation in the weights is thus driven by discriminative information, unlike PLML where they are only based on the smoothness assumption. Finally, Figure 2(c) shows that the metrics consistently generalize to test data.

**Effect of the basis set size** Figure 3 shows the number of selected basis elements and test error rate for SCML-Global and SCML-Local as a function of the size of basis set on Segment (results were consistent on other datasets). The left pane shows that the number of selected elements increases sublinearly and eventually converges, while the right pane shows that test error may be further reduced by using a larger basis set without significant overfitting, as suggested by our generalization bound (Theorem 1). Figure 3 also shows that SCML-Local generally selects more basis elements than SCML-Global, but notice that it can outperform SCML-Global even when the basis set is very small.

## 6 Conclusion

We proposed to learn metrics as sparse combinations of rank-one basis elements. This framework unifies several paradigms in metric learning, including global, local and multi-task learning. Of particular interest is our local metric learning algorithm which can compute instance-specific metrics for both training and test points in a principled way. The soundness of our approach is supported theoretically by a generalization bound, and we showed in experimental studies that the proposed methods improve upon state-of-the-art algorithms in terms of accuracy and scalability.

**Acknowledgements** This research is partially supported by the IARPA via DoD/ARL contract # W911NF-12-C-0012 and DARPA via contract # D11AP00278. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, DoD/ARL, DARPA, or the U.S. Government.

## References

- A. Argyriou, T. Evgeniou, and M. Pontil. Convex multi-task feature learning. *Mach. Learn.*, 73(3):243–272, 2008.
- A. Bellet and A. Habrard. Robustness and Generalization for Metric Learning. Technical report, arXiv:1209.1086, 2012.
- Aurélien Bellet, Amaury Habrard, and Marc Sebban. A Survey on Metric Learning for Feature Vectors and Structured Data. Technical report, arXiv:1306.6709, June 2013.
- J. Blitzer, M. Dredze, and F. Pereira. Biographies, Bollywood, Boom-boxes and Blenders: Domain Adaptation for Sentiment Classification. In *ACL*, 2007.
- R. Caruana. Multitask Learning. *Mach. Learn.*, 28(1):41–75, 1997.
- C.-C. Chang and C.-J. Lin. LIBSVM : a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2(3):27–27, 2011.
- J. Davis, B. Kulis, P. Jain, S. Sra, and I. Dhillon. Information-theoretic metric learning. In *ICML*, 2007.
- J. Duchi and Y. Singer. Efficient Online and Batch Learning Using Forward Backward Splitting. *JMLR*, 10:2899–2934, 2009.
- A. Frome, Y. Singer, F. Sha, and J. Malik. Learning Globally-Consistent Local Distance Functions for Shape-Based Image Retrieval and Classification. In *ICCV*, 2007.
- J. Goldberger, S. Roweis, G. Hinton, and R. Salakhutdinov. Neighbourhood Components Analysis. In *NIPS*, 2004.
- P. Gong, J. Ye, and C. Zhang. Robust multi-task feature learning. In *KDD*, 2012.
- S. Hauberg, O. Freifeld, and M. Black. A Geometric take on Metric Learning. In *NIPS*, 2012.
- Y. Hong, Q. Li, J. Jiang, and Z. Tu. Learning a mixture of sparse distance metrics for classification and dimensionality reduction. In *CVPR*, 2011.
- P. Jain, B. Kulis, I. Dhillon, and K. Grauman. Online Metric Learning and Fast Similarity Search. In *NIPS*, 2008.
- A. Kolmogorov and V. Tikhomirov.  $\epsilon$ -entropy and  $\epsilon$ -capacity of sets in functional spaces. *American Mathematical Society Translations*, 2(17):277–364, 1961.
- Brian Kulis. Metric Learning: A Survey. *Foundations and Trends in Machine Learning*, 5(4):287–364, 2012.
- Y.-K. Noh, B.-T. Zhang, and D. Lee. Generative Local Metric Learning for Nearest Neighbor Classification. In *NIPS*, 2010.
- S. Parameswaran and K. Weinberger. Large Margin Multi-Task Metric Learning. In *NIPS*, 2010.
- D. Ramanan and S. Baker. Local Distance Functions: A Taxonomy, New Algorithms, and an Evaluation. *TPAMI*, 33(4):794–806, 2011.

- B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Comput.*, 10(1):1299–1319, 1998.
- C. Shen, J. Kim, L. Wang, and A. van den Hengel. Positive Semidefinite Metric Learning Using Boosting-like Algorithms. *JMLR*, 13:1007–1036, 2012.
- Y. Shi, A. Bellet, and F. Sha. Sparse Compositional Metric Learning. In *AAAI*, 2014.
- L. van der Maaten and G. Hinton. Visualizing Data using t-SNE. *JMLR*, 9:2579–2605, 2008.
- J. Wang, A. Woznica, and A. Kalousis. Parametric Local Metric Learning for Nearest Neighbor Classification. In *NIPS*, 2012.
- K. Weinberger and L. Saul. Distance Metric Learning for Large Margin Nearest Neighbor Classification. *JMLR*, 10:207–244, 2009.
- L. Xiao. Dual Averaging Methods for Regularized Stochastic Learning and Online Optimization. *JMLR*, 11:2543–2596, 2010.
- E. Xing, A. Ng, M. Jordan, and S. Russell. Distance Metric Learning with Application to Clustering with Side-Information. In *NIPS*, 2002.
- H. Xu and S. Mannor. Robustness and Generalization. *Mach. Learn.*, 86(3):391–423, 2012.
- X. Yang, S. Kim, and E. Xing. Heterogeneous multitask learning with joint sparsity constraints. In *NIPS*, 2009.
- Y. Ying and P. Li. Distance Metric Learning with Eigenvalue Optimization. *JMLR*, 13:1–26, 2012.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *J. Roy. Statist. Soc. Ser. B*, 68:49–67, 2006.
- D.-C. Zhan, M. Li, Y.-F. Li, and Z.-H. Zhou. Learning instance specific distances using metric propagation. In *ICML*, 2009.

## Appendix A Detailed Analysis

In this section, we give the details of the derivation of the generalization bounds for the global and multi-task learning formulations given in Section 3.

### A.1 Preliminaries

We start by introducing some notation. We are given a training sample  $S = \{z_i = (\mathbf{x}_i, y_i)\}_{i=1}^n$  drawn i.i.d. from a distribution  $P$  over the labeled space  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ . We assume that  $\|\mathbf{x}\| \leq R$  (for some convenient norm),  $\forall \mathbf{x} \in \mathcal{X}$ . We call a triplet  $(z, z', z'')$  *admissible* if  $y = y' \neq y''$ . Let  $C_S$  be the set of all admissible triplets built from instances in  $S$ .<sup>8</sup>

---

<sup>8</sup>When the training triplets consist of only a subset of all admissible triplets (which is often the case in practice), a relaxed version of the robustness property can be used to derive similar results (Bellet and Habrard, 2012). For simplicity, we focus here on the case when all admissible triplets are used.

Let  $L(h, z, z', z'')$  be the loss suffered by some hypothesis  $h$  on triplet  $(z, z', z'')$ , with the convention that  $L$  returns 0 for non-admissible triplets. We assume  $L$  to be uniformly upper-bounded by a constant  $U$ . The empirical loss  $\mathcal{R}_{emp}^{C_S}(h)$  of  $h$  on  $C_S$  is defined as

$$\mathcal{R}_{emp}^{C_S}(h) = \frac{1}{|C_S|} \sum_{(z, z', z'') \in C_S} L(h, z, z', z''),$$

and its expected loss  $\mathcal{R}(h)$  over distribution  $P$  as

$$\mathcal{R}(h) = \mathbb{E}_{z, z', z'' \sim P} L(h, z, z', z'').$$

Our goal is to bound the deviation between  $\mathcal{R}(\mathcal{A}_{C_S})$  and  $\mathcal{R}_{emp}^{C_S}(\mathcal{A}_{C_S})$ , where  $\mathcal{A}_{C_S}$  is the hypothesis learned by algorithm  $\mathcal{A}$  on  $C_S$ .

## A.2 Algorithmic Robustness

To derive our generalization bounds, we use the recent framework of algorithmic robustness (Xu and Mannor, 2012), in particular its adaptation to pairwise and tripletwise loss functions used in metric learning (Bellet and Habrard, 2012). For the reader’s convenience, we briefly review these main results below.

Algorithmic robustness is the ability of an algorithm to perform “similarly” on a training example and on a test example that are “close”. The proximity of points is based on a partitioning of the space  $\mathcal{Z}$ : two examples are close to each other if they lie in the same region. The partition is based on the notion of covering number (Kolmogorov and Tikhomirov, 1961).

**Definition 1** (Covering number). *For a metric space  $(\mathcal{S}, \rho)$  and  $\mathcal{V} \subset \mathcal{S}$ , we say that  $\hat{\mathcal{V}} \subset \mathcal{V}$  is a  $\gamma$ -cover of  $\mathcal{V}$  if  $\forall \mathbf{t} \in \mathcal{V}, \exists \hat{\mathbf{t}} \in \hat{\mathcal{V}}$  such that  $\rho(\mathbf{t}, \hat{\mathbf{t}}) \leq \gamma$ . The  $\gamma$ -covering number of  $\mathcal{V}$  is*

$$\mathcal{N}(\gamma, \mathcal{V}, \rho) = \min \left\{ |\hat{\mathcal{V}}| : \hat{\mathcal{V}} \text{ is a } \gamma\text{-cover of } \mathcal{V} \right\}.$$

In particular, when  $\mathcal{X}$  is compact,  $\mathcal{N}(\gamma, \mathcal{X}, \rho)$  is finite, leading to a finite cover. Then,  $\mathcal{Z}$  can be partitioned into  $|\mathcal{Y}| \mathcal{N}(\gamma, \mathcal{X}, \rho)$  subsets such that if two examples  $\mathbf{z} = (\mathbf{x}, y)$  and  $\mathbf{z}' = (\mathbf{x}', y')$  belong to the same subset, then  $y = y'$  and  $\rho(\mathbf{x}, \mathbf{x}') \leq \gamma$ . The definition of robustness for tripletwise loss functions (adapted from Xu and Mannor, 2012) is as follows.

**Definition 2** (Robustness for metric learning (Bellet and Habrard, 2012)). *An algorithm  $\mathcal{A}$  is  $(N, \epsilon(\cdot))$  robust for  $N \in \mathbb{N}$  and  $\epsilon(\cdot) : (\mathcal{Z} \times \mathcal{Z})^n \rightarrow \mathbb{R}$  if  $\mathcal{Z}$  can be partitioned into  $N$  disjoint sets, denoted by  $\{Q_i\}_{i=1}^N$ , such that the following holds for all  $S \in \mathcal{Z}^n$ :*

*$\forall (\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3) \in C_S, \forall \mathbf{z}, \mathbf{z}', \mathbf{z}'' \in \mathcal{Z}, \forall i, j \in [N] : \text{if } \mathbf{z}_1, \mathbf{z} \in Q_i, \mathbf{z}_2, \mathbf{z}' \in Q_j, \mathbf{z}_3, \mathbf{z}'' \in Q_k \text{ then}$*

$$|L(\mathcal{A}_{C_S}, \mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3) - L(\mathcal{A}_{C_S}, \mathbf{z}, \mathbf{z}', \mathbf{z}'')| \leq \epsilon(C_S),$$

*where  $\mathcal{A}_{C_S}$  is the hypothesis learned by  $\mathcal{A}$  on  $C_S$ .*

$N$  and  $\epsilon(\cdot)$  quantify the robustness of the algorithm and depend on the training sample. Again adapting the result from (Xu and Mannor, 2012), (Bellet and Habrard, 2012) showed that a metric learning algorithm that satisfies Definition 2 has the following generalization guarantees.



**Theorem 2.** *If a learning algorithm  $\mathcal{A}$  is  $(N, \epsilon(\cdot))$ -robust and the training sample consists of the triplets  $C_S$  obtained from a sample  $S$  generated by  $n$  i.i.d. draws from  $P$ , then for any  $\delta > 0$ , with probability at least  $1 - \delta$  we have:*

$$|\mathcal{R}(\mathcal{A}_{C_S}) - \mathcal{R}_{emp}^{C_S}(\mathcal{A}_{C_S})| \leq \epsilon(C_S) + 3U \sqrt{\frac{N \ln 2 + \ln \frac{1}{\delta}}{0.5n}}.$$

As shown in (Bellet and Habrard, 2012), establishing the robustness of an algorithm is easier using the following theorem, which basically says that if a metric learning algorithm has approximately the same loss for triplets that are close to each other, then it is robust.

**Theorem 3.** *Fix  $\gamma > 0$  and a metric  $\rho$  of  $\mathcal{Z}$ . Suppose that  $\forall z_1, z_2, z_3, z, z', z'' : (z_1, z_2, z_3) \in C_S, \rho(z_1, z) \leq \gamma, \rho(z_2, z') \leq \gamma, \rho(z_3, z'') \leq \gamma$ ,  $\mathcal{A}$  satisfies*

$$|L(\mathcal{A}_{C_S}, z_1, z_2, z_3) - L(\mathcal{A}_{C_S}, z, z', z'')| \leq \epsilon(C_S),$$

and  $\mathcal{N}(\gamma/2, \mathcal{Z}, \rho) < \infty$ . Then the algorithm  $\mathcal{A}$  is  $(\mathcal{N}(\gamma/2, \mathcal{Z}, \rho), \epsilon(C_S))$ -robust.

We now have all the tools we need to prove the results of interest.

### A.3 Generalization Bounds for SCML

#### A.3.1 Bound for SCML-Global

We first focus on SCML-Global where the loss function is defined as follows:

$$L(\mathbf{w}, z, z', z'') = [1 + d_{\mathbf{w}}(\mathbf{x}, \mathbf{x}') - d_{\mathbf{w}}(\mathbf{x}, \mathbf{x}'')]_{+}.$$

We obtain a generalization bound by showing that SCML-Global satisfies Definition 2 using Theorem 3. To establish the result, we need a bound on the  $\ell_2$  norm of the basis elements. Since they are obtained by eigenvalue decomposition, their norm is equal to (and thus bounded by) 1.

Let  $\mathbf{w}^*$  be the optimal solution to SCML-Global. By optimality of  $\mathbf{w}^*$  we have:

$$L(\mathbf{w}^*, z, z', z'') + \beta \|\mathbf{w}^*\|_1 \leq L(\mathbf{0}, z, z', z'') + \beta \|\mathbf{0}\|_1 = 1,$$

thus we get  $\|\mathbf{w}^*\|_1 \leq 1/\beta$ . Let  $\mathbf{M}^* = \sum_{i=1}^K w_i^* \mathbf{b}_i \mathbf{b}_i^T$  be the learned metric. Then using Holder's inequality and the bound on  $\mathbf{w}^*$  and the  $\mathbf{b}$ 's:

$$\|\mathbf{M}^*\|_1 = \left\| \sum_{i=1}^K w_i^* \mathbf{b}_i \mathbf{b}_i^T \right\|_1 = \left\| \sum_{i:w_i^* \neq 0} w_i^* \mathbf{b}_i \mathbf{b}_i^T \right\|_1 \leq \|\mathbf{w}^*\|_1 \sum_{i:w_i^* \neq 0} \|\mathbf{b}_i\|_{\infty} \|\mathbf{b}_i\|_{\infty} \leq K^*/\beta,$$

where  $K^* \leq K$  is the number of nonzero entries in  $\mathbf{w}^*$ .

Using Definition 1, we can partition  $\mathcal{Z}$  into  $|\mathcal{Y}| \mathcal{N}(\gamma, \mathcal{X}, \rho)$  subsets such that if two examples  $z = (\mathbf{x}, y)$  and  $z' = (\mathbf{x}', y')$  belong to the same subset, then  $y = y'$  and  $\rho(\mathbf{x}, \mathbf{x}') \leq \gamma$ . Now, for  $z_1, z_2, z_3, z'_1, z'_2, z'_3 \in \mathcal{Z}$ , if  $y_1 = y'_1, \|\mathbf{x}_1 - \mathbf{x}'_1\|_1 \leq \gamma, y_2 = y'_2, \|\mathbf{x}_2 - \mathbf{x}'_2\|_1 \leq \gamma, y_3 = y'_3, \|\mathbf{x}_3 - \mathbf{x}'_3\|_1 \leq \gamma$ , then  $(z_1, z_2, z_3)$  and  $(z'_1, z'_2, z'_3)$  are either both admissible or both non-admissible triplets.

In the non-admissible case, by definition their respective loss is 0 and so is the deviation between the losses. In the admissible case we have:

$$\begin{aligned}
& \left| [1 + d_{w^*}(\mathbf{x}_1, \mathbf{x}_2) - d_{w^*}(\mathbf{x}_1, \mathbf{x}_3)]_+ - [1 + d_{w^*}(\mathbf{x}'_1, \mathbf{x}'_2) - d_{w^*}(\mathbf{x}'_1, \mathbf{x}'_3)]_+ \right| \\
\leq & |(\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_2) - (\mathbf{x}_1 - \mathbf{x}_3)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_3) + (\mathbf{x}'_1 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_3) - (\mathbf{x}'_1 - \mathbf{x}'_2)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_2)| \\
= & |(\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_2) - (\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_2) + (\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_2) - (\mathbf{x}'_1 - \mathbf{x}'_2)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_2) \\
& + (\mathbf{x}'_1 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_3) - (\mathbf{x}'_1 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_3) + (\mathbf{x}'_1 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_3) - (\mathbf{x}_1 - \mathbf{x}_3)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_3)| \\
= & |(\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_2 - (\mathbf{x}'_1 - \mathbf{x}'_2)) + (\mathbf{x}_1 - \mathbf{x}_2 - (\mathbf{x}'_1 - \mathbf{x}'_2))^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_2) \\
& + (\mathbf{x}'_1 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_3 - (\mathbf{x}_1 - \mathbf{x}_3)) + (\mathbf{x}'_1 - \mathbf{x}'_3 - (\mathbf{x}_1 - \mathbf{x}_3))^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_3)| \\
\leq & |(\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}'_1)| + |(\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}'_2 - \mathbf{x}_2)| + |(\mathbf{x}_1 - \mathbf{x}'_1)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_2)| \\
& + |(\mathbf{x}'_2 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_2)| + |(\mathbf{x}'_1 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}_1)| + |(\mathbf{x}'_1 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}_3 - \mathbf{x}'_3)| \\
& + |(\mathbf{x}'_1 - \mathbf{x}_1)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_3)| + |(\mathbf{x}_3 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_3)| \\
\leq & \|\mathbf{x}_1 - \mathbf{x}_2\|_\infty \|\mathbf{M}^*\|_1 \|\mathbf{x}_1 - \mathbf{x}'_1\|_1 + \|\mathbf{x}_1 - \mathbf{x}_2\|_\infty \|\mathbf{M}^*\|_1 \|\mathbf{x}'_2 - \mathbf{x}_2\|_1 + \|\mathbf{x}_1 - \mathbf{x}'_1\|_1 \|\mathbf{M}^*\|_1 \|\mathbf{x}'_1 - \mathbf{x}'_2\|_\infty \\
& + \|\mathbf{x}'_2 - \mathbf{x}_2\|_1 \|\mathbf{M}^*\|_1 \|\mathbf{x}'_1 - \mathbf{x}'_2\|_\infty + \|\mathbf{x}'_1 - \mathbf{x}'_3\|_\infty \|\mathbf{M}^*\|_1 \|\mathbf{x}'_1 - \mathbf{x}_1\|_1 + \|\mathbf{x}'_1 - \mathbf{x}'_3\|_\infty \|\mathbf{M}^*\|_1 \|\mathbf{x}_3 - \mathbf{x}'_3\|_1 \\
& + \|\mathbf{x}'_1 - \mathbf{x}_1\|_1 \|\mathbf{M}^*\|_1 \|\mathbf{x}_1 - \mathbf{x}_3\|_\infty + \|\mathbf{x}_3 - \mathbf{x}'_3\|_1 \|\mathbf{M}^*\|_1 \|\mathbf{x}_1 - \mathbf{x}_3\|_\infty \\
\leq & \frac{16\gamma RK^*}{\beta},
\end{aligned}$$

by using the property that the hinge loss is 1-Lipschitz, Holder's inequality and bounds on the involved quantities. Thus SCML-Global is  $(|Y|\mathcal{N}(\gamma, X, \|\cdot\|_1), \frac{16\gamma RK^*}{\beta})$ -robust and the generalization bound follows.

### A.3.2 Bound for mt-SCML

In the multi-task setting, we are given a training sample  $S_t = \{\mathbf{z}_i = (\mathbf{x}_i^t, y_i^t)\}_{i=1}^{n_t}$ . Let  $C_{S_t}$  be the set of all admissible triplets built from instances in  $S_t$ .

Let  $\mathbf{W}^*$  be the optimal solution to mt-SCML. Using the same arguments as for SCML-Global, by optimality of  $\mathbf{W}^*$  we have  $\|\mathbf{W}^*\|_{2,1} \leq 1/\beta$ . Let  $\mathbf{M}_t^* = \sum_{i=1}^K W_{ti}^* \mathbf{b}_i \mathbf{b}_i^\top$  be the learned metric for task  $t$  and  $\mathbf{W}_t^*$  be the weight vector for task  $t$ , corresponding to the  $t$ -th row of  $\mathbf{W}^*$ . Then using the fact that  $\|\mathbf{W}_t^*\|_{2,1} \leq \|\mathbf{W}^*\|_{2,1}$  and  $\|\mathbf{b}\|_{2,1} = \|\mathbf{b}\|_2$ , we have:

$$\|\mathbf{M}_t^*\|_{2,1} = \left\| \sum_{i=1}^K W_{ti}^* \mathbf{b}_i \mathbf{b}_i^\top \right\|_{2,1} = \left\| \sum_{i:W_{ti}^* \neq 0} W_{ti}^* \mathbf{b}_i \mathbf{b}_i^\top \right\|_{2,1} \leq \|\mathbf{W}_t^*\|_{2,1} \sum_{i:W_{ti}^* \neq 0} \|\mathbf{b}_i\|_{2,1} \|\mathbf{b}_i\|_{2,1} \leq K_t^* / \beta,$$

where  $K_t^* \leq K$  is the number of nonzero entries in  $\mathbf{W}_t^*$ .

From this we can derive a generalization bound for each task using arguments similar to the global case, using a partition specific to each task defined with respect to  $\|\cdot\|_2$ . Without loss of generality, we focus on task  $t$  and only explicitly write the last derivations as the beginning is the same as above:

$$\begin{aligned}
& \left| [1 + d_{w^*}(\mathbf{x}_1, \mathbf{x}_2) - d_{w^*}(\mathbf{x}_1, \mathbf{x}_3)]_+ - [1 + d_{w^*}(\mathbf{x}'_1, \mathbf{x}'_2) - d_{w^*}(\mathbf{x}'_1, \mathbf{x}'_3)]_+ \right| \\
\leq & |(\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}'_1)| + |(\mathbf{x}_1 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}'_2 - \mathbf{x}_2)| + |(\mathbf{x}_1 - \mathbf{x}'_1)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_2)| \\
& + |(\mathbf{x}'_2 - \mathbf{x}_2)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}'_2)| + |(\mathbf{x}'_1 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}'_1 - \mathbf{x}_1)| + |(\mathbf{x}'_1 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}_3 - \mathbf{x}'_3)| \\
& + |(\mathbf{x}'_1 - \mathbf{x}_1)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_3)| + |(\mathbf{x}_3 - \mathbf{x}'_3)^\top \mathbf{M}^*(\mathbf{x}_1 - \mathbf{x}_3)| \\
\leq & \|\mathbf{x}_1 - \mathbf{x}_2\|_2 \|\mathbf{M}^*\|_{\mathcal{F}} \|\mathbf{x}_1 - \mathbf{x}'_1\|_2 + \|\mathbf{x}_1 - \mathbf{x}_2\|_2 \|\mathbf{M}^*\|_{\mathcal{F}} \|\mathbf{x}'_2 - \mathbf{x}_2\|_2 + \|\mathbf{x}_1 - \mathbf{x}'_1\|_2 \|\mathbf{M}^*\|_{\mathcal{F}} \|\mathbf{x}'_1 - \mathbf{x}'_2\|_2 \\
& + \|\mathbf{x}'_2 - \mathbf{x}_2\|_2 \|\mathbf{M}^*\|_{\mathcal{F}} \|\mathbf{x}'_1 - \mathbf{x}'_2\|_2 + \|\mathbf{x}'_1 - \mathbf{x}'_3\|_2 \|\mathbf{M}^*\|_{\mathcal{F}} \|\mathbf{x}'_1 - \mathbf{x}_1\|_2 + \|\mathbf{x}'_1 - \mathbf{x}'_3\|_2 \|\mathbf{M}^*\|_{\mathcal{F}} \|\mathbf{x}_3 - \mathbf{x}'_3\|_2 \\
& + \|\mathbf{x}'_1 - \mathbf{x}_1\|_2 \|\mathbf{M}^*\|_{\mathcal{F}} \|\mathbf{x}_1 - \mathbf{x}_3\|_2 + \|\mathbf{x}_3 - \mathbf{x}'_3\|_2 \|\mathbf{M}^*\|_{\mathcal{F}} \|\mathbf{x}_1 - \mathbf{x}_3\|_2 \\
\leq & \frac{16\gamma RK_t^*}{\beta},
\end{aligned}$$

Dataset	Linear SVM	Kernel SVM	SCML-Global	SCML-Local
Vehicle	21.4±0.6	<b>16.6±0.8</b>	21.3±0.6	<b>18.0±0.6</b>
Vowel	24.3±0.7	<b>4.4±0.4</b>	10.9±0.5	6.1±0.4
Segment	5.1±0.2	<b>3.6±0.2</b>	4.1±0.2	<b>3.6±0.2</b>
Letters	19.5±0.4	8.8±0.2	9.0±0.2	<b>8.3±0.2</b>
USPS	6.5±0.2	4.2±0.1	4.1±0.1	<b>3.6±0.1</b>
BBC	4.3±0.2	<b>3.8±0.2</b>	<b>3.9±0.2</b>	<b>4.1±0.2</b>
Avg. rank	3.8	1.3	2.3	<b>1.2</b>

Table 7: Comparison of SCML against linear and kernel SVM (best in bold).

where we used the same arguments as above and the inequality  $\|M^*\|_F \leq \|M^*\|_{2,1}$ . Thus mt-SCML is  $(|Y|\mathcal{N}(\gamma, X, \|\cdot\|_2), \frac{16\gamma RK_t^*}{\beta})$ -robust and the bound for task  $t$  follows. Note that the number of training examples in the bound is only that from task  $t$ , i.e.,  $n = n_t$ .

### A.3.3 Comments on SCML-Local

It would be interesting to be able to derive a similar bound for SCML-Local. Unfortunately, as it is nonconvex, we cannot assume optimality of the solution. If a similar formulation can be made convex, the same proof technique should apply: even though each instance has its own metric, it essentially depends on the instance itself (whose norm is bounded) and on the learned parameters shared across metrics (which could be bounded using optimality of the solution). Deriving such a convex formulation and the corresponding generalization bound is left as future work.

## Appendix B Experimental Comparison with Support Vector Machines

In this section, we compare SCML-Global and SCML-Local to Support Vector Machines (SVM) using a linear and a RBF kernel. We used the software LIBSVM (Chang and Lin, 2011) and tuned the parameter  $C$  as well as the bandwidth for the RBF kernel on the validation set. Table 7 shows misclassification rates averaged over 20 random splits, along with standard error and the average rank of each method across all datasets. First, we can see that SCML-Global consistently performs better than linear SVM. Second, SCML-Local is competitive with kernel SVM. These results show that a simple  $k$ -nearest neighbor strategy with a good metric can be competitive (and even outperform) SVM classifiers.