

---

# Active Metric Learning for Supervised Classification

---

Krishnan Kumaran<sup>1</sup> Dimitri Papageorgiou<sup>1</sup> Yutong Chang<sup>2</sup> Minhan Li<sup>2</sup> Martin Takac<sup>2</sup>

## Abstract

Clustering and classification critically rely on distance metrics that provide meaningful comparisons between data points. We present mixed-integer optimization approaches to find optimal distance metrics that generalize the Mahalanobis metric extensively studied in the literature. Additionally, we generalize and improve upon leading methods by removing reliance on pre-designated “target neighbors,” “triplets,” and “similarity pairs.” Another salient feature of our method is its ability to enable active learning by recommending precise regions to sample after an optimal metric is computed to improve classification performance. This targeted acquisition can significantly reduce computational burden by ensuring training data completeness, representativeness, and economy. We demonstrate classification and computational performance of the algorithms through several simple and intuitive examples, followed by results on real image and medical datasets.

## 1. Introduction and Motivation

Selecting an appropriate distance metric is fundamental to many learning algorithms such as k-means, nearest neighbor searches, and others, as observed by (Davis et al., 2007) and other researchers in this field. Further, they observe that choosing such a measure is highly problem-specific and ultimately dictates the success - or failure - of the learning algorithm. Nevertheless, the algorithm used to select the metric based on data can be designed to be more general, and designing such algorithm(s) is indeed the objective of this work, as well as past research on this problem.

In this work, we formulate a general framework for choosing such metrics that improves and extends previous formulations in some important ways. In particular, we attempt to couple the metric learning problem with that of recommending targeted data acquisition, which has not been sufficiently

addressed in past work which has mostly assumed that the data is a given, static collection of N-dimensional vectors. However, in many real-world settings, one does not have the luxury of learning a once-and-for-all distance metric. Rather, an iterative approach is required whereby an initial distance metric is learned, new data is acquired, the distance metric is refined, and so on. A main goal of this work is to present a systematic framework for optimizing this iterative procedure so that an optimal and interpretable metric is learned and is, in turn, used to recommend precise regions to sample in order to acquire new data to be used to further refine the metric and improve classification performance.

### 1.1. Problem Setting

As described in (Xiang et al., 2008), there are two prominent batch distance metric learning settings, both of which assume that we are given a set of  $N$  points  $\mathbf{x}_i \in \mathcal{R}^D$ . In the first setting, a class (or label)  $C_i$  is explicitly given for each point  $i \in \mathcal{N} = \{1, \dots, N\}$ . In the second setting, classes are implicitly furnished through pairwise constraints in the form of must-links and cannot-links. Must-links are given as  $\{(i, j) : i \text{ and } j \text{ are in the same class}\}$ , whereas cannot-links are specified as  $\{(i, k) : i \text{ and } k \text{ are not in the same class}\}$ . For both settings, we let  $\mathcal{C}_i$  and  $\bar{\mathcal{C}}_i$  denote the co-class and non-class neighbors of  $i$ , respectively.

We ask the question: Is there a metric  $\mathbb{D}(\mathbf{x}, \mathbf{y})$  that enforces the condition that the nearest neighbor of each point is a co-class point, i.e.,  $\forall i \in \mathcal{N}$

$$\min_{j \in \mathcal{C}_i} \mathbb{D}(\mathbf{x}_i, \mathbf{x}_j) < \min_{k \in \bar{\mathcal{C}}_i} \mathbb{D}(\mathbf{x}_i, \mathbf{x}_k)? \quad (1)$$

More generally, let  $\mathcal{N}_i^K(\mathbb{D}) = \{j_1, \dots, j_K\}$  be the  $K$  nearest neighbors to  $i$  with respect to a distance metric  $\mathbb{D}$  (assume no ties). Then, we are interested in finding a distance metric  $\mathbb{D}$  satisfying the condition: Given  $K$ , the majority of the  $K$  nearest neighbors are co-class points, i.e.,  $\forall i \in \mathcal{N}$

$$\exists \bar{K} = \lfloor \frac{K}{2} \rfloor + 1 \text{ points } j_1, \dots, j_{\bar{K}} \in \mathcal{N}_i^K(\mathbb{D}) \cap \mathcal{C}_i. \quad (2)$$

Note that the above condition does not *a priori* define target neighbors, as required by most previous work. This is an important distinction because the closest neighbors of a point are not determined unless the metric is specified. Our formulation incorporates variables that compare distances between true neighbors contingent on the distance metric. This property avoids the pitfalls of pre-specified target neighbors as shown below, while preserving the desirable characteris-

<sup>1</sup>ExxonMobil Corporate Strategic Research, Annandale, NJ  
<sup>2</sup>Lehigh University, Bethlehem, PA. Correspondence to: Krishnan Kumaran <krishnan.kumaran@exxonmobil.com>, Dimitri Papageorgiou <dimitri.j.papageorgiou@exxonmobil.com>.

tics of agglomerative and  $K$ -nearest neighbor clustering methods such as permitting multiple disjoint islands of the same class and non-convex class regions while maintaining simplicity and interpretability of the metric.

**Form of Distance Metric.** In general, we allow the metric to be a power series of the form

$$\mathbb{D}(\mathbf{x}, \mathbf{y}) \equiv \mathbf{a} \cdot |\mathbf{x} - \mathbf{y}| + (\mathbf{x} - \mathbf{y}) \cdot \mathbf{B} \cdot (\mathbf{x} - \mathbf{y}) + \sum_{p,q,r=1}^D C_{pqr} |x_p - y_p| \cdot |x_q - y_q| \cdot |x_r - y_r| + \text{higher order terms.} \quad (3)$$

Strictly speaking, the mapping  $\mathbb{D}$  proposed in (3) may not satisfy the four properties - non-negativity, symmetry, triangle inequality, distinguishability - that are required to be a “metric.” Nevertheless, we use this terminology throughout.

Restricting to the first term loosely corresponds to SVM/discriminant analysis, while the second term is commonly known in the literature as the *Mahalanobis metric* if the matrix  $\mathbf{B}$  is symmetric and positive definite. Higher order terms introduce additional parameters at a power law rate, e.g., the fully symmetric tensor  $C_{pqr}$  has  $O(D^3)$  parameters. In principle, almost any dataset can eventually be fitted with a metric with a potentially infinite number of these parameters. However, in practice, a mis-classification trade-off curve and knee-point can be used to prevent over-fitting. Further, the mis-classified points could point to possible outliers, errors in input data, or class boundaries. In all of these cases, the algorithm points to regions in the space where further data acquisition/quality testing would be of most value. This aspect of our method is unique, and provides significant value in selecting the most effective training for supervised classification in general, even applied to SVM/Deep Neural Networks or other algorithms. Further, we will show in our experimental results that the ratio of closest co-class to closest non-class point, defined as

$$R_i = \frac{\min_{j \in \mathcal{C}_i} \mathbb{D}(\mathbf{x}_i, \mathbf{x}_j)}{\min_{k \in \mathcal{C}_i} \mathbb{D}(\mathbf{x}_i, \mathbf{x}_k)} \quad (4)$$

is a useful metric to separate “interior” points from “boundary points” of classes.

Even in that case where we consider only the second-order term, the Mahalanobis distance, our approach differs from past approaches due to condition (1) which we show empirically results in better solutions.

## 1.2. Comparison with Prior Work

(Wang & Sun, 2015) survey distance metric learning in unsupervised and supervised settings. Several recent publications have proposed metric learning methods similar to ours, see (Weinberger & Saul, 2009; Davis et al., 2007; Ying & Li, 2012; Rosales & Fung, 2006; Xing et al., 2003) and references therein.

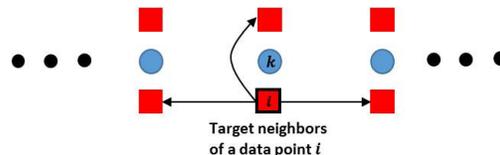


Figure 1. Poor choice of target neighbors (Weinberger & Saul, 2009) can result in infeasible/distorted metrics. No Mahalanobis metric can bring the target neighbors closer while simultaneously pushing the intermediate non-class neighbor  $k$  further. Observe that a simple metric with high vertical weighting and low horizontal weighting will classify correctly in our approach, which does not use target neighbors.

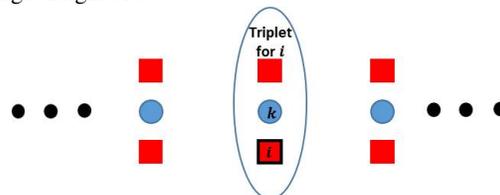


Figure 2. Poorly chosen triplets (Rosales & Fung, 2006) can lead to similar distortions as target neighbors. The non-class point  $k$  can distort the metric, or result in infeasibility of desired metric condition 1. As with target neighbors, we can produce a simple metric here without using triplets.

A fundamental distinction that makes our approach more general than the distance metric learning approaches of these authors is that ours does not rely on auxiliary input information in the form of so-called target co-class neighbors (Weinberger & Saul, 2009; Ying & Li, 2012; Xing et al., 2003), similar and dissimilar point pairs (Davis et al., 2007), nor triplets with one co-class and non-class neighbor (Rosales & Fung, 2006). The target neighbors of a point are co-class points that the user desires to be closest to it. Target neighbor-based methods fix a priori a set of points and attempt to learn a linear transformation of the input space such that the resulting nearest neighbors of the point are indeed its target neighbors. Unfortunately, in many applications target neighbors or triples are not available. In the absence of prior target neighbor knowledge, (Weinberger & Saul, 2009) suggest using the  $K$  nearest neighbors with the same class label, as determined by Euclidean distance. While these requirements appear reasonable, they can be misleading and highly data-dependent as shown in the example below. Further, target neighbors and/or triplets, even if available with an initial data set, may become burdensome and/or error-prone to update when additional data becomes available.

Figures 1 and 2 illustrate the potential for distorted distance metrics when target neighbors or triplets are pre-defined. It is worth emphasizing that (Weinberger & Saul, 2009) and (Ying & Li, 2012) rely on Euclidean distance in their

computational experiments.

### 1.3. Contributions

The contributions of this work are:

1. We present a distance metric learning algorithm that is competitive with other state-of-the-art metric learning algorithms including (Weinberger & Saul, 2009) and (Davis et al., 2007). Moreover, our approach is more general than the aforementioned algorithms since we do not require pre-specification of target neighbors or triplets, which involves a high degree of user and data choice dependency, and hence possible errors.
2. Our method **provides recommendations** for new data acquisition and data quality control to improve classification performance. This is a key “value-of-information” criterion that can significantly improve both classification performance and computational burden by ensuring training data completeness, representativeness and economy, which are not adequately addressed in current applications of DNN and other methods which often depend on very large quantities of training data (e.g. internet cat images).
3. We show that our underlying metric learning problem can be formulated and solved as a **mixed-integer linear optimization (MIO)** problem. To the best of our knowledge, this is the first time such a claim has been made in the metric learning arena. Indeed, this work is also timely as it builds on what (Hastie et al., 2017) call “exciting new work” applying MIO to prominent machine learning problems with great success (Bertsimas et al., 2016; 2017; Bertsimas & Van Parys, 2017; Bertsimas & King, 2015; Bertsimas et al., 2014; Friesen & Domingos, 2017; Wilson & Sahinidis, 2017)

## 2. Mixed-Integer Linear Optimization Formulations for Metric Learning

In this section, we present mixed-integer linear formulations to determine an “optimal” distance metric that satisfies condition (1) or (2), where optimality is governed by an appropriately chosen loss function. For ease of exposition, we describe the formulations for a distance metric (3) with only first- and second-order terms. Specifically, the feasible region  $\mathcal{F}$  for the distance metric is

$$\mathcal{F} = \{(\mathbf{a}, \mathbf{B}, \mathbf{d}) \in \mathbb{R}^D \times \mathbb{R}^{D \times D} \times [0, 1]^{N \times N} : \mathbf{B} = \mathbf{B}^\top, \\ d_{ij} = \mathbf{a}^\top |\boldsymbol{\delta}_{ij}| + \boldsymbol{\delta}_{ij}^\top \mathbf{B} \boldsymbol{\delta}_{ij} \quad \forall i, j \in \mathcal{N} \quad (5) \\ d_{ik} \geq d^{\min} \quad \forall i, k \in \mathcal{N} : i \text{ and } k \text{ are not in same class}\}$$

Here  $\boldsymbol{\delta}_{ij} = \mathbf{x}_i - \mathbf{x}_j$  for all  $i, j \in \mathcal{N}$ .  $d^{\min} > 0$  is a given positive parameter that prevents the degenerate metric  $(\mathbf{a}, \mathbf{B}, \mathbf{d}) = \mathbf{0}$ . Note that  $\mathcal{F}$  is a polyhedron, since we do not enforce  $\mathbf{B}$  to be positive semidefinite; this extension is addressed in Section 2.3. Distances are bounded above by 1

(although any positive upper bound suffices) since our MIO methods require an upper bound.

### 2.1. Metric Learning for Single Nearest Neighbor

Consider first the task of determining an “optimal” distance metric satisfying condition (1) assuming that such a metric exists. We relax this assumption below. Let  $\lambda_i = \min\{\mathbb{D}(\mathbf{x}_i, \mathbf{x}_k) : k \in \bar{\mathcal{C}}_i\} - \min\{\mathbb{D}(\mathbf{x}_i, \mathbf{x}_j) : j \in \mathcal{C}_i\}$  be the separation between the distance to point  $i$ ’s nearest non-class neighbor and the distance to its nearest co-class neighbor. Since condition (1) seeks a distance metric such that  $\lambda_i > 0$  for all  $i \in \mathcal{N}$ , we first consider the loss function

$$L^0(\lambda_1, \dots, \lambda_N) = -\min\{\lambda_i : i \in \mathcal{N}\}, \quad (6)$$

which rewards the minimum separation over all points. The following mixed-integer *nonlinear* formulation attempts to minimize the loss function  $L^0$ , or equivalently, to maximize  $\lambda$ , the minimum separation over all points:

$$\max_{\lambda, \mathbf{a}, \mathbf{B}, \mathbf{d}, \mathbf{y}} \lambda \quad (7a)$$

$$\text{s.t. } \sum_{j \in \mathcal{C}_i} d_{ij} y_{ij} + \lambda \leq d_{ik}, \quad \forall i \in \mathcal{N}, k \in \bar{\mathcal{C}}_i, \quad (7b)$$

$$\sum_{j \in \mathcal{C}_i} y_{ij} = 1, \quad \forall i \in \mathcal{N}, \quad (7c)$$

$$y_{ij} \in \{0, 1\}, \quad \forall i \in \mathcal{N}, j \in \mathcal{C}_i, \quad (7d)$$

$$\lambda \in \mathbb{R}, (\mathbf{a}, \mathbf{B}, \mathbf{d}) \in \mathcal{F}. \quad (7e)$$

Binary decision variables  $y_{ij}$ , which take value 1 if point  $i$  is assigned to co-class point  $j$  (0 otherwise), are required to select a single co-class neighbor as nearest neighbor. Together, constraints (7b) and (7c) attempt to ensure that at least one co-class neighbor is closer than all other non-class neighbors to point  $i$ . Although it is helpful to think that  $y_{ij}$  will take value 1 if point  $j$  is chosen as the nearest co-class point to  $i$  (under the resulting optimal distance metric), it is possible in an optimal solution  $(\lambda^*, \mathbf{a}^*, \mathbf{B}^*, \mathbf{d}^*, \mathbf{y}^*)$  that  $y_{i\hat{j}}^* = 0$  when  $d_{i\hat{j}}^* < d_{ij}^*$  for all  $j \in \mathcal{C}_i, j \neq \hat{j}$  for some point  $i$ . This simply means that, for that point  $i$ , constraint (7b) is not tight in that optimal solution. Nevertheless, an optimal solution to formulation (7) is guaranteed to find the largest minimum separation  $\lambda^*$  over all points. If  $\lambda^* > 0$ , then there exists a distance metric satisfying condition (1).

There are at least two deficiencies with formulation (7). First and most important, although it is guaranteed to be feasible, it is not guaranteed to return a distance metric satisfying condition (1). Second, it contains bilinear terms  $d_{ij} y_{ij}$  (the multiplication of two decision variables), which are non-convex and are undesirable when solving (7) with an off-the-shelf optimization engine. We now discuss how to overcome these two deficiencies.

Enforcing condition (1) to hold for all points  $i \in \mathcal{N}$  could be too stringent. For example, in a sparse data set, it is quite possible that certain classes are both “islanded” and under-represented leading to nearest co-class neighbors that are

far away under most optimized metrics. Furthermore, it is possible that a small number of points severely limit classes from being separated by a large margin. In such cases, we may wish to identify this subset of bottleneck points as “outliers” and only enforce condition (1) for non-outliers.

To this end, let  $\mathcal{O} \subseteq \mathcal{N}$  be a set of outliers and re-define  $\lambda_i$  as  $\lambda_i = \min\{d_{ik} : k \in \bar{\mathcal{C}}_i \setminus \mathcal{O}\} - \min\{d_{ij} : j \in \mathcal{C}_i \setminus \mathcal{O}\}$  for all non-outliers  $i \in \mathcal{N} \setminus \mathcal{O}$ . To obtain a distance metric such that  $\lambda_i > 0$  for all  $i \in \mathcal{N} \setminus \mathcal{O}$ , we adopt the loss function

$$L(\lambda_1, \dots, \lambda_N) = \rho|\mathcal{O}| - \min_{i \in \mathcal{N} \setminus \mathcal{O}} \lambda_i, \quad (8)$$

which penalizes all outliers (with respect to the given distance metric) and rewards the minimum separation over all non-outliers. It is important to emphasize that  $\rho$  is the only user-defined parameter in our approach as a large value of  $\rho$  signals that outliers are highly undesirable, whereas a small value indicates greater tolerance of outliers and more preference for larger margin.

In order to handle outliers within an optimization framework, we introduce binary decision variables  $z_i$  that take value 1 if point  $i$  is deemed an outlier; 0 otherwise. The main interactions between the binary variables  $y_{ij}$  and  $z_i$  are captured through the following set:

$$\mathcal{YZ} = \{(\mathbf{y}, \mathbf{z}) : \sum_{j \in \mathcal{C}_i} y_{ij} = 1 - z_i, \quad \forall i \in \mathcal{N}, \quad (9a)$$

$$y_{ij} \leq 1 - z_j, \quad \forall i \in \mathcal{N}, j \in \mathcal{C}_i, \quad (9b)$$

$$y_{ij} \in \{0, 1\}, \quad \forall i \in \mathcal{N}, j \in \mathcal{C}_i, \quad (9c)$$

$$z_i \in \{0, 1\}, \quad \forall i \in \mathcal{N}. \quad (9d)$$

Constraints (9a) ensure that each point is assigned to exactly one co-class neighbor or is deemed an outlier. Constraints (9b) only allow point  $i$  to be assigned to point  $j$  if  $j$  is not an outlier. With these additions for outliers, we can re-formulate (7) as follows:

$$\max_{\lambda, \mathbf{a}, \mathbf{B}, \mathbf{d}, \mathbf{y}, \mathbf{z}} \lambda - \rho \sum_{i \in \mathcal{N}} z_i \quad (10a)$$

$$\text{s.t. } \sum_{j \in \mathcal{C}_i} d_{ij} y_{ij} + \lambda \leq d_{ik} + M_{ik} z_k, \quad \forall i \in \mathcal{N}, k \in \bar{\mathcal{C}}_i \quad (10b)$$

$$\lambda \geq 0, (\mathbf{a}, \mathbf{B}, \mathbf{d}) \in \mathcal{F}, (\mathbf{y}, \mathbf{z}) \in \mathcal{YZ}. \quad (10c)$$

Here,  $M_{ik}$  is known as a “Big M” parameter and can be set to 1 for all  $(i, k)$  pairs since all distance variables  $d_{ik}$  are bounded above by 1. Besides the additions noted above, constraints (10b) are now active only if point  $i$  and point  $k$  are non-outliers.

Finally, to overcome the computational issues with the non-linear terms  $d_{ij} y_{ij}$ , we apply a standard “trick” in mixed-integer optimization to arrive at the following MILP formu-

lation:

$$\max_{\lambda, \mathbf{w}, \mathbf{y}, \mathbf{z}} \lambda - \rho \sum_{i \in \mathcal{N}} z_i \quad (11a)$$

$$\text{s.t. } \sum_{j \in \mathcal{C}_i} w_{ij} + \lambda \leq d_{ik} + M_{ik} z_k, \quad \forall i \in \mathcal{N}, k \in \bar{\mathcal{C}}_i, \quad (11b)$$

$$w_{ij} \leq d_{ij}, \quad \forall i \in \mathcal{N}, j \in \mathcal{C}_i, \quad (11c)$$

$$w_{ij} \leq y_{ij}, \quad \forall i \in \mathcal{N}, j \in \mathcal{C}_i, \quad (11d)$$

$$w_{ij} \geq d_{ij} + y_{ij} - 1, \quad \forall i \in \mathcal{N}, j \in \mathcal{C}_i, \quad (11e)$$

$$w_{ij} \geq 0, \quad \forall i \in \mathcal{N}, j \in \mathcal{C}_i, \quad (11f)$$

$$\lambda \geq 0, (\mathbf{a}, \mathbf{B}, \mathbf{d}) \in \mathcal{F}, (\mathbf{y}, \mathbf{z}) \in \mathcal{YZ}. \quad (11g)$$

Constraints (11c)-(11e) are known as McCormick envelopes of the bilinear expression  $w_{ij} = d_{ij} y_{ij}$ . Note that, in actual implementation, the decision variables  $d_{ij}$  are replaced by their definition given in (5).

## 2.2. Metric Learning for $K$ Nearest Neighbors

We next turn to the task of determining an “optimal” distance metric satisfying condition (2), i.e., the majority of the  $K$  nearest neighbors are co-class points. For the remainder of this subsection, we assume  $K$  is given and that each class has at least  $K$  points. As above, we first present a formulation that assumes such a metric exists and then relax this assumption to allow for outliers.

Let  $\mathcal{P} = \{(i, j) : i \in \mathcal{N}, j \in \mathcal{N}, i \neq j\}$ . Let  $\kappa_i = |\mathcal{N}_i^K(D) \cap \mathcal{C}_i|$  be the number of co-class points that are among point  $i$ 's  $K$  nearest neighbors (with respect to the distance metric  $\mathbb{D}$ ). Since condition (2) seeks a distance metric such that  $\kappa_i \geq \frac{K}{2} + 1$  for all  $i \in \mathcal{N}$ , analogous to loss function (6), we first consider the loss function

$$L_K^0(\kappa_1, \dots, \kappa_N) = -\min\{\kappa_i : i \in \mathcal{N}\}. \quad (12)$$

The following MIO formulation attempts to minimize the loss function  $L^0$ . Binary decision variables  $u_{ij}$ , which take value 1 if point  $i \in \mathcal{N}$  is assigned to point  $j \in \mathcal{N}$  (0 otherwise), are required to keep track of which points are selected as the  $K$  nearest neighbors of each point  $i \in \mathcal{N}$ .

$$\max_{\kappa, \Delta, \mathbf{u}} \kappa \quad (13a)$$

$$\text{s.t. } \sum_{j \in \mathcal{C}_i} u_{ij} \geq \kappa, \quad \forall i \in \mathcal{N}, \quad (13b)$$

$$\sum_{j \in \mathcal{N}} u_{ij} \leq K, \quad \forall i \in \mathcal{N}, \quad (13c)$$

$$d_{ij} \leq \Delta_i + M_{ij}(1 - u_{ij}), \quad \forall (i, j) \in \mathcal{P}, \quad (13d)$$

$$d_{ik} \geq \Delta_i + \epsilon - M_{ik} u_{ik}, \quad \forall (i, k) \in \mathcal{P} : k \in \bar{\mathcal{C}}_i, \quad (13e)$$

$$u_{ij} \in \{0, 1\}, \quad \forall (i, j) \in \mathcal{P}, \quad (13f)$$

$$\kappa \geq 0, \Delta \geq 0, (\mathbf{a}, \mathbf{B}, \mathbf{d}) \in \mathcal{F}. \quad (13g)$$

Constraints (13b) allow us to maximize the minimum  $\kappa_i$ .

Constraints (13c) ensure that no more than  $K$  points are chosen as point  $i$ 's nearest neighbors. Note that, since we are maximizing  $\kappa$ , constraints (13c) can be written with an inequality " $\leq K$ " rather than an equality " $= K$ ". Constraints (13d) require point  $i$ 's (at most)  $K$  nearest neighbors to be within a distance of  $\Delta_i$ , while constraints (13e) enforce the opposite condition that all of point  $i$ 's non-nearest non-class neighbors be at least  $\Delta_i + \epsilon$  units from  $i$ , where  $\epsilon > 0$  is a user-defined parameter. Assuming there are no co-located points (i.e.,  $\mathbf{x}_i \neq \mathbf{x}_j \forall (i, j) \in \mathcal{P}$ ), setting  $\epsilon = \alpha \min\{\|\mathbf{x}_i - \mathbf{x}_j\| : (i, j) \in \mathcal{P}\}$  for  $\alpha \in (0, 1]$  will guarantee that a feasible solution always exists.

An optimal solution  $(\kappa^*, \mathbf{\Delta}^*, \mathbf{u}^*, \mathbf{a}^*, \mathbf{B}^*, \mathbf{d}^*)$  satisfies condition (2) if  $\kappa^* \geq \frac{K}{2} + 1$  since, together, constraints (13d) and (13e) ensure that there are at least  $\kappa^*$  co-class points among the  $K$  nearest neighbors of every point. Note also that formulation (13) allows for "ties" amongst co-class points, thus it could return a distance metric for which a given point's  $K$  nearest neighbors are not unique.

There are at least two deficiencies with formulation (13). First, there may not exist a distance metric satisfying condition (2), i.e., an optimal solution to (13) could result in  $\kappa^* < \frac{K}{2} + 1$ . Second, when using the loss function (12), there may be a large number of optimal solutions to (13) with objective function value  $\kappa^*$ , even though we would prefer one in which the remaining  $\kappa_i$ 's are maximized. To this end, we modify the loss function to account for outliers and give weight for having more co-class points among the  $K$  nearest neighbors:

$$L_K(\kappa_1, \dots, \kappa_N) = \rho|\mathcal{O}| - \min_{i \in \mathcal{N}} \kappa_i - W \sum_{i \in \mathcal{N} \setminus \mathcal{O}} \kappa_i. \quad (14)$$

Here,  $W$  is a non-negative user-defined scalar that governs the preference between maximizing the minimum  $\kappa_i$  and encouraging more co-class points among the  $K$  nearest neighbors. Setting  $W = (NK/2)^{-1}$  suffices to ensure that maximizing the minimum  $\kappa_i$  is the dominant objective.

This leads to the following MIO formulation:

$$\max_{\substack{\kappa, \mathbf{\Delta}, \mathbf{u}, \mathbf{z} \\ \mathbf{a}, \mathbf{B}, \mathbf{d}}} \kappa + W \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{C}_i} u_{ij} - \rho \sum_{i \in \mathcal{N} \setminus \mathcal{Z}} z_i \quad (15a)$$

$$\text{s.t.} \quad \sum_{j \in \mathcal{C}_i} u_{ij} \geq \kappa, \quad \forall i \in \mathcal{N}, \quad (15b)$$

$$\sum_{j \in \mathcal{N}} u_{ij} \leq K, \quad \forall i \in \mathcal{N}, \quad (15c)$$

$$u_{ij} \leq 1 - z_i, \quad \forall (i, j) \in \mathcal{P}, \quad (15d)$$

$$u_{ij} \leq 1 - z_j, \quad \forall (i, j) \in \mathcal{P}, \quad (15e)$$

$$d_{ij} \leq \Delta_i + M_{ij}(1 - u_{ij}), \quad \forall (i, j) \in \mathcal{P}, \quad (15f)$$

$$d_{ik} \geq \Delta_i + \epsilon - M_{ik}(u_{ik} + z_i + z_k), \quad \forall (i, k) \in \mathcal{P} : k \in \bar{\mathcal{C}}_i, \quad (15g)$$

$$u_{ij} \in \{0, 1\}, \quad \forall (i, j) \in \mathcal{P}, \quad (15h)$$

$$z_i \in \{0, 1\}, \quad \forall i \in \mathcal{N}, \quad (15i)$$

$$\kappa \geq 0, \mathbf{\Delta} \geq \mathbf{0}, (\mathbf{a}, \mathbf{B}, \mathbf{d}) \in \mathcal{F}. \quad (15j)$$

The objective function is the negative loss function  $-L_K(\kappa_1, \dots, \kappa_N)$ , where  $\sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{C}_i} u_{ij}$  plays the role of  $\kappa_i$ . Constraints (15d) and (15e) ensure that point  $j$  is not chosen as one of point  $i$ 's  $K$  nearest neighbors if either  $i$  or  $j$  is deemed an outlier. All other constraints resemble those of formulation (13) except perhaps with modifications to account for outliers.

Note that the  $K$  nearest neighbor approach requires more binary variables than the single nearest neighbor formulation. In particular, the former requires  $(N - 1)^2$  binary variables  $u_{ij}$  for all  $(i, j) \in \mathcal{P}$ , whereas the latter requires far fewer since  $y_{ij}$  is only defined for all  $(i, j) : i \in \mathcal{N}, j \in \mathcal{C}_i$ .

### 2.3. Extensions

Thus far, positive semidefiniteness of the  $\mathbf{B}$  matrix is not enforced. Although there is evidence that positive semidefiniteness is a desirable attribute and may improve interpretability of the resulting metric, it may not be as essential as others have described. Indeed, several highly touted state-of-the-art approaches do not enforce positive semidefiniteness, e.g. NCA (Goldberger et al., 2004) and deep neural nets, yet are still garnering considerable attention. Nevertheless, if positive semidefiniteness is strongly desired, we could trivially extend formulation (11) to include diagonal dominance constraints as done in (Rosales & Fung, 2006). Although such constraints would not allow for fully general psd matrices to be generated, they would keep the formulation mixed-integer linear. On the other hand, if all psd matrices are desired, we would need to adopt a more sophisticated approach as in (Weinberger & Saul, 2009).

While we have thus far extolled the fact that our approach does not rely on a priori information, our MILP formulations can easily accommodate user-provided target neighbors and similarity/dissimilarity pairs. Indeed, linear constraints, like those used in (Davis et al., 2007) and Shavel-Shwartz (Shalev-Shwartz et al., 2004), that require user-specified similar points to be closer than dissimilar points are simple to incorporate.

## 3. Active Learning for Targeted Data Acquisition Using Boundary and Outlier Identification

Our algorithm is particularly suited for continuous, online data acquisition aimed at converging to an optimal metric with the smallest amount of data. This type of *active learning* approach is critical to maintaining economy, completeness and representativeness in data selection. To the best of our knowledge, this work is the first to address the connection between metric learning and active data selection. A similar approach is applicable to alternative metric and other learning paradigms including LMNN, ITML, DNN,

and SVM. For example, for DNN, the current approach often involves using very large quantities of training data, causing high computational burdens and convergence issues. Prioritization of data based on empirical boundary point and outlier determination could significantly improve performance.

Figure 3 shows a summary of our approach for the  $K = 1$  nearest neighbor case<sup>1</sup>. It involves the following steps:

- Compute the  $R$ -ratio (4) for all current points.
- Compute the cumulative histogram of each class separately (bottom of Figure 3).
- Points with  $R \geq 1$  are potential outliers, while points between the "knee-point"  $R^*$  and 1 are designated as *boundary points*. Points with  $R \leq R^*$  are the *interior points*.<sup>2</sup>
- Our recommendation for new data acquisition is to selectively acquire additional data at the outliers and boundary points, prioritized in decreased order of  $R$ .

The above described procedure is motivated by the self-evident observation that interior points typically have many co-class neighbors closer than their nearest non-class neighbor, and are hence less likely to be misclassified, while the reverse is true of the boundary points. Empirical experimental validation of this observation is shown in Figure 4, where the metric learned from all the data is very close to the metric learned from the boundary points only.

#### 4. Numerical Experiments

In this section, we demonstrate the performance of our algorithm on real and synthetic datasets. In all experiments presented here, we restrict our attention to the search of optimal Mahalanobis distances for which the leading method in the literature is LMNN (Weinberger & Saul, 2009)<sup>3</sup> against which we compare our approach.

##### 4.1. Synthetic Data

The purpose of this section is to demonstrate the characteristics of the proposed algorithm on synthetic data sets designed to illustrate complex structure, including classes with multiple "domains/islands" and non-convex class shapes.

<sup>1</sup>Both the formula for  $R$  and the data selection procedure can be extended for  $K > 1$  with small modifications.

<sup>2</sup>Note that the parallel line construction shown in the figure typically results in  $R^* < 1$  in practice as the cumulative histogram is typically convex and monotonically decreasing with high enough sampling density, which results in many more interior points than boundary points. However, when this is not the case, we suggest retaining all data points for the class, i.e.  $R^* = 0$  as this condition is an indication of insufficient sampling and more samples will be needed before discriminating data selection is possible.

<sup>3</sup>We have used LMNN 3.0.0 available at <http://www.cs.cornell.edu/~kilian/code/code.html>.

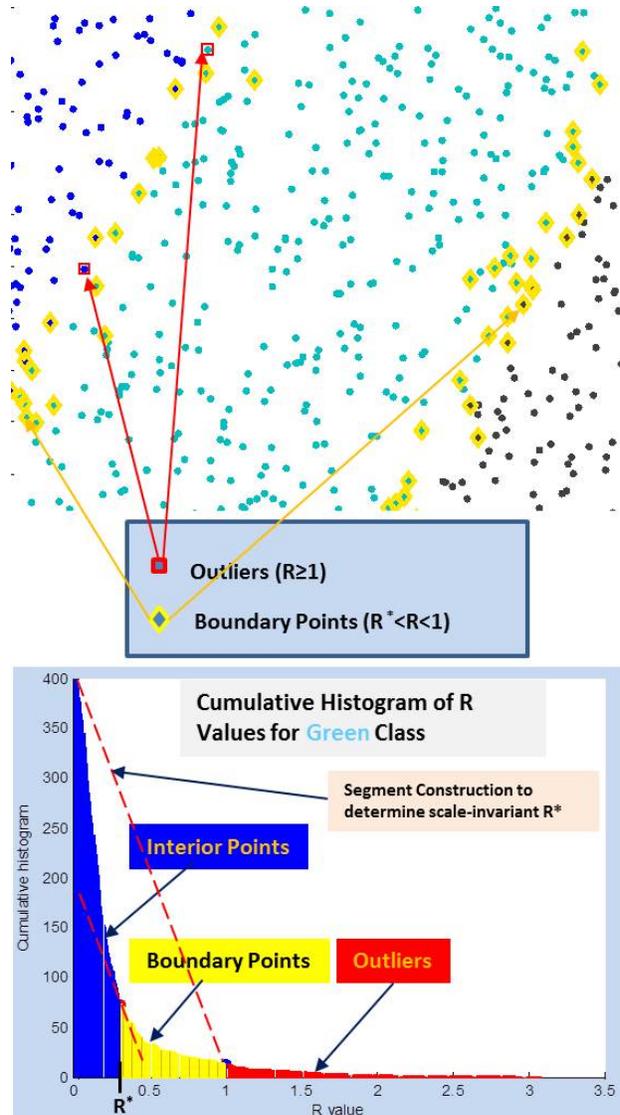


Figure 3. Histogramming the data by  $R$ -ratio (4) allows empirical identification of outlier and boundary points for each class separately using the histogram as shown on the bottom figure. The scatter plot shows the resulting boundary and outlier points. Active learning involves selectively acquiring more data at these outliers and boundaries (as opposed to the interior), in sorted order of  $R$ -value for better data economy and faster convergence to true metric.

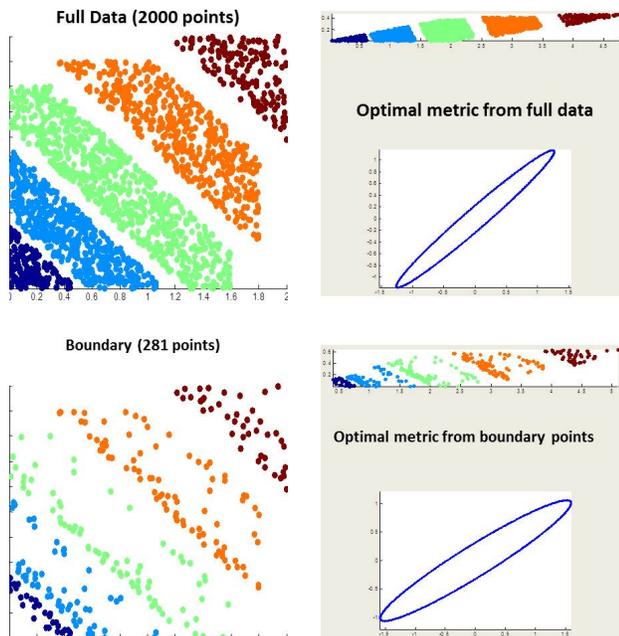


Figure 4. The closest co-class to non-class neighbor ratios  $R_i$  can reveal the class boundaries in combination with a trade-off curve as shown in the top right Figure. Note that the metric inferred from all class points (top ellipse) is very close to the metric from the much smaller set of boundary points only (bottom ellipse). The vertically squished images (top right and center right) are scatter plots of the points in the transformed coordinates from the optimal metric in each case. This example demonstrates our algorithm’s ability to implement data economy. As mentioned in the text, the boundary points also suggest the most desirable regions for further data acquisition to improve classification results.

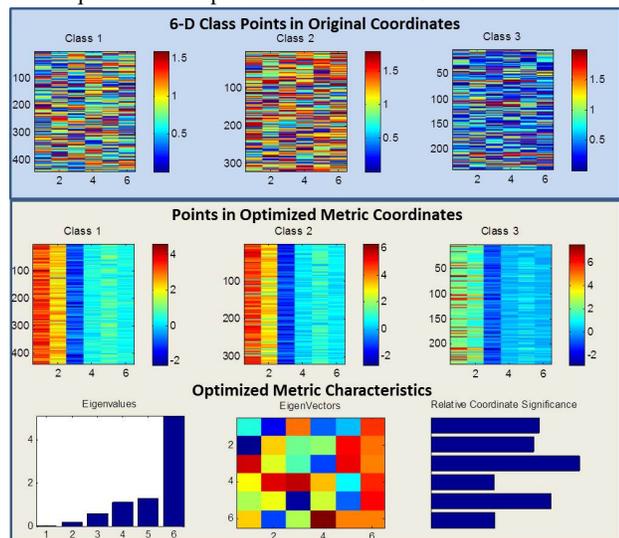


Figure 5. Synthetic example of a 6-D dataset with 3 classes. The pattern of the data within classes, as well as the differences between classes, can be characterized and visualized much more easily on the scaled coordinates from our optimized metric, as shown in the figures titled ‘Points in Optimized Metric Coordinates’. The occurrence of vertical stripes in each class visualization implies that each class has relatively uniform coordinate values in the transformed metric space.

Further, we show how the boundary points of classes can be inferred from the metric for data economy and recommendations of further data acquisition. Also, in the Appendix, we demonstrate that our method produces smoother class boundaries that standard metrics, thus reducing errors due to over-fitting. Finally, we provide a higher dimensional example to illustrate how classes are better represented and visualized using the optimal metric.

**Synthetic Data: No Outliers** We start with 2D examples that permit easy visualization and intuitive geometric explanation. In these examples, we assume that the data has no outliers, and the sampling is sufficiently large for robust classification. The Mahalanobis metric allows for simple and direct interpretation of the metric using eigenvalues and eigenvectors. We use the co-class to non-class distance ratio of (4) to define *relative coordinate significance* for each coordinate as the eigenvalue-weighted absolute sum of the corresponding components of the eigenvectors, i.e., if the Mahalanobis matrix has (eigenvalue, eigenvector) pairs  $(\lambda_i, \mathbf{v}_i), i = 1, \dots, D$ , with each eigenvector  $\mathbf{v}_i \equiv \{v_{ij} : j = 1, \dots, D\}$ , then the relative significance of each coordinate dimension  $j$  is defined as

$$W_j \equiv \sum_{i=1}^D \lambda_i |v_{ij}|. \quad (16)$$

Figures 4, 6, and 5 show how these concepts can be used, in 2D as well as higher dimensions, to produce more intuitive and interpretable metrics even when the sample points are restricted primarily to class boundaries. We deal with identification and exclusion of outliers in the Appendix. Now we show some applications of these methods to real data along with comparisons to competing alternatives.

**4.2. Real Datasets.** In this section, we present results on real datasets drawn from handwritten number recognition MNIST (LeCun et al., 1998) and medical/biological datasets like diabetes and iris data.

**Experiments on MNIST Dataset.** The state-of-the-art classification method for MNIST data uses convolutional neural networks which extracts useful features. In this experiment we trained convolution auto-encoder (with ReLU activation) and represented each image by the compressed feature representation resulting in 16 features per image. The results of our classification – performed on this reduced 16D space – are shown in Figure 8. These results demonstrate how our algorithm can be used to classify the training data into interpretable interiors, boundaries and outliers, with the latter two being visibly similar to other neighboring classes leading to possible mis-classification.

**Experiments on Medical Datasets.** In this experiment we use *cod-rna* and *diabetes* datasets<sup>4</sup> each having 8 features. We trained the optimal metric on 300 random points and tested on larger sets of points (as indicated in table). Note

<sup>4</sup>All datasets used in this Section are available at <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>.

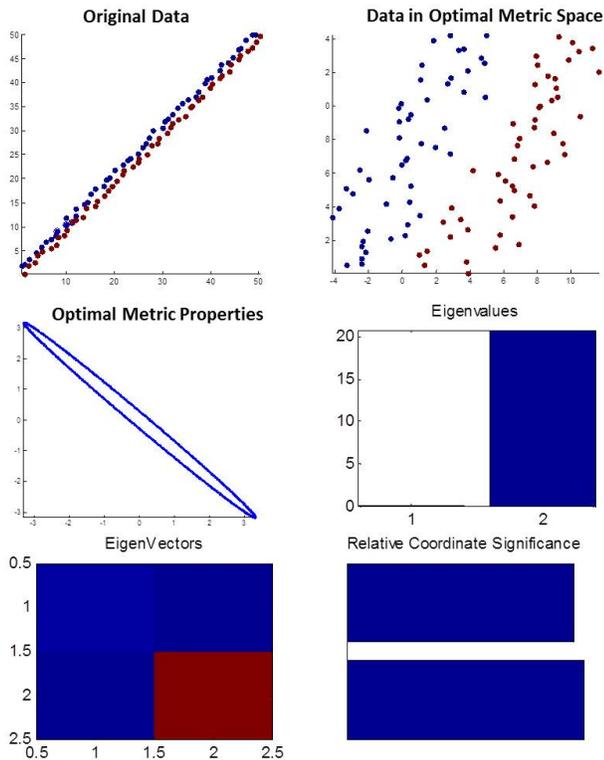


Figure 6. Results on a closely spaced pair of noisy linear class examples show a metric highly weighted orthogonal to the line as expected. This example could be challenging for existing methods like LMNN if target neighbors are chosen using Euclidean metric at the start as normally suggested.

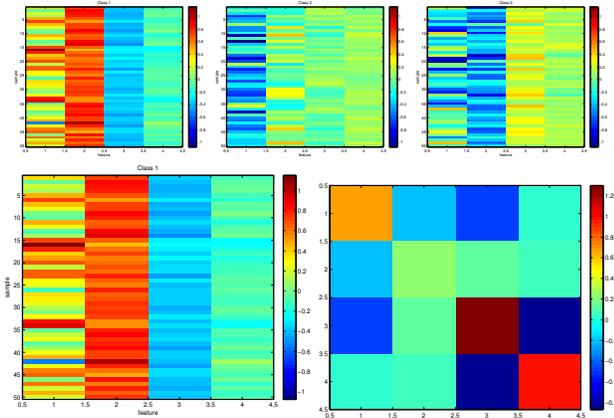


Figure 7. Results on 4D IRIS dataset. This example illustrates visualization of high-dimensional data using our optimized metric. As shown in the synthetic example earlier, note that representation of class features (top row) on the scaled coordinates reveals unique signatures of each class as vertical streaks indicating relatively uniform feature values within each class.

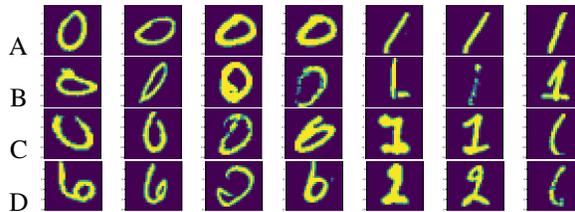


Figure 8. Row A shows interior points of the classes (class 0 left and class 1 right); Row B shows typical boundary points of the classes, while row C shows outliers and row D shows the closest point to the outlier. Note that the results are generally consistent with human judgement and provide guidelines for training data enhancement based on the outliers.

approach	<i>cod-rna</i>	<i>cod-rna</i>	<i>diabetes</i>
# of test	4,000	8,000	769
Euclidean distance	14.47%	13.43%	29.42%
MeLL0	<b>6.27%</b>	<b>6.17%</b>	<b>23.30%</b>
LMNN $K = 1$	8.82%	9.18%	29.81%
LMNN $K = 3$	8.29%	8.61%	30.59%
LMNN $K = 5$	8.19%	8.47%	29.94%
LMNN $K = 7$	8.32%	8.78%	30.07%
LMNN $K = 11$	8.64%	8.88%	30.46%

Table 1. Comparison of classification error for various distance metrics. Note that MeLL0, which uses outlier removal in both training and testing phases, provides some improvement over state-of-the-art alternatives. Further, the outliers, as before, suggest new desirable data to improve performance.

that our algorithm (labeled MeLL0), provides improvement in performance over Euclidean and various LMNN alternatives without any significant tuning of parameters on our part.

**Experiments on Biological Dataset.** In Figure 7 we show results on Iris Flower Dataset (Fisher, 1936) which contains 3 classes of 150 training samples, each represented by 4 features. Note that our method separates the classes reasonably well, while also depicting relatively internally uniform characteristics for each class in terms of the Mahalanobis vectors, as seen by the vertical streaks in the images on the top row in Figure 7. In all cases demonstrated above we observed that our method produces competitive and interpretable classification results as well as outlier identification. This combination is useful for choosing comprehensive and economical training datasets.

### 5. Conclusion

We have presented novel methods for classification based on optimized metric learning. Our methods are designed to overcome limitations of state-of-the-art alternatives such as SVM, DNN by being more interpretable, robust and more general. We have shown promising results on synthetic and real datasets.

## Acknowledgements

The work of MT was partially supported by the U.S. National Science Foundation, under award number NSF:CCF:1618717, NSF:CMMI:1663256 and NSF:CCF:1740796.

## References

- Bertsimas, Dimitris and King, Angela. Or foruman algorithmic approach to linear regression. *Operations Research*, 64(1):2–16, 2015.
- Bertsimas, Dimitris and Van Parys, Bart. Sparse high-dimensional regression: Exact scalable algorithms and phase transitions. *arXiv preprint arXiv:1709.10029*, 2017.
- Bertsimas, Dimitris, Mazumder, Rahul, et al. Least quantile regression via modern optimization. *The Annals of Statistics*, 42(6):2494–2525, 2014.
- Bertsimas, Dimitris, King, Angela, Mazumder, Rahul, et al. Best subset selection via a modern optimization lens. *The Annals of Statistics*, 44(2):813–852, 2016.
- Bertsimas, Dimitris, King, Angela, et al. Logistic regression: From art to science. *Statistical Science*, 32(3):367–384, 2017.
- Bixby, Robert E. A brief history of linear and mixed-integer programming computation. *Documenta Mathematica*, pp. 107–121, 2012.
- Davis, Jason V, Kulis, Brian, Jain, Prateek, Sra, Suvrit, and Dhillon, Inderjit S. Information-theoretic metric learning. In *Proceedings of the 24th international conference on Machine learning*, pp. 209–216. ACM, 2007.
- Fisher, Ronald A. The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2):179–188, 1936.
- Friesen, Abram L and Domingos, Pedro. Deep learning as a mixed convex-combinatorial optimization problem. *arXiv preprint arXiv:1710.11573*, 2017.
- Goldberger, Jacob, Roweis, Sam, Hinton, Geoff, and Salakhutdinov, Ruslan. Neighbourhood components analysis. *NIPS04*, 2004.
- Hastie, Trevor, Tibshirani, Robert, and Tibshirani, Ryan J. Extended comparisons of best subset selection, forward stepwise selection, and the lasso. *arXiv preprint arXiv:1707.08692*, 2017.
- LeCun, Yann, Bottou, Léon, Bengio, Yoshua, and Haffner, Patrick. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Rosales, Rómer and Fung, Glenn. Learning sparse metrics via linear programming. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 367–373. ACM, 2006.
- Shalev-Shwartz, Shai, Singer, Yoram, and Ng, Andrew Y. Online and batch learning of pseudo-metrics. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 94. ACM, 2004.
- Wang, Fei and Sun, Jimeng. Survey on distance metric learning and dimensionality reduction in data mining. *Data Mining and Knowledge Discovery*, 29(2): 534–564, Mar 2015. ISSN 1573-756X. doi: 10.1007/s10618-014-0356-z. URL <https://doi.org/10.1007/s10618-014-0356-z>.
- Weinberger, Kilian Q and Saul, Lawrence K. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, 10(Feb): 207–244, 2009.
- Wilson, Zachary T and Sahinidis, Nikolaos V. The alamo approach to machine learning. *Computers & Chemical Engineering*, 2017.
- Xiang, Shiming, Nie, Feiping, and Zhang, Changshui. Learning a mahalanobis distance metric for data clustering and classification. *Pattern Recognition*, 41(12):3600 – 3612, 2008. ISSN 0031-3203. doi: <https://doi.org/10.1016/j.patcog.2008.05.018>. URL <http://www.sciencedirect.com/science/article/pii/S0031320308002057>.
- Xing, Eric P, Ng, Andrew Y, Jordan, Michael I, and Russell, Stuart. Distance metric learning with application to clustering with side-information. *Advances in neural information processing systems*, 15:505–512, 2003.
- Ying, Yiming and Li, Peng. Distance metric learning with eigenvalue optimization. *Journal of Machine Learning Research*, 13(Jan):1–26, 2012.

## Part I

## Appendix

## A. Class boundaries

Figure 9 demonstrates the characteristics of class boundaries inferred from our optimal metric in comparison to commonly used metrics like Manhattan ( $\ell_1$ ), Euclidean ( $\ell_2$ ) and Maximum ( $\ell_\infty$ ). Smooth and robust boundaries are often practically important for robust online classification. Note that our boundaries are significantly smoother and less prone to overfitting, although the Euclidean metric performs similarly. However, in the presence of outliers and noise, we will later show that the Euclidean metric performs much worse.

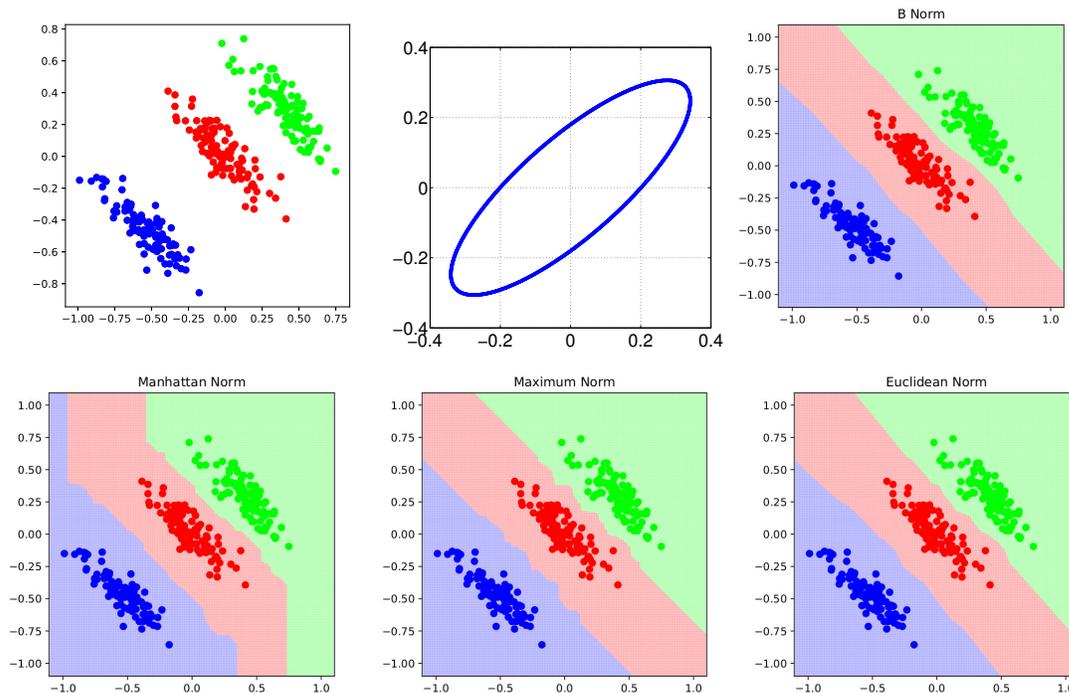


Figure 9. The boundaries learnt by our algorithm are smoother than alternatives. Top left is the original data; top right is the result produced by our algorithm; bottom 3 are the alternative approaches.

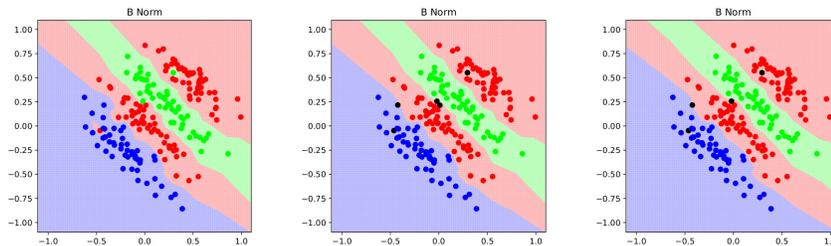
## A.1. Handling Outliers

Handling outliers is a very important task for robust classification. In this section, we show how we can improve the distance metric by iteratively removing outliers using synthetic example. Figure 10 demonstrates our outlier removal procedure, which consists of 3 major iterations.

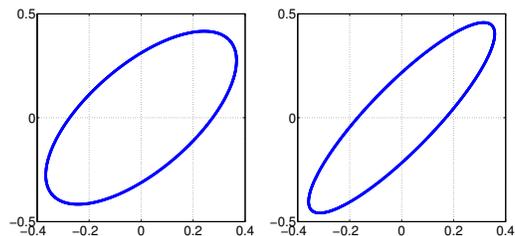
1. determination of optimal metric for the whole dataset;
2. identification of outliers using the distance ratio;
3. reoptimization of the metric after removal of outliers.

Steps 1–3 can be optionally repeated iteratively until satisfactory metric is obtained.<sup>5</sup>

<sup>5</sup>This procedure is a heuristic version of our thorough scheme which incorporates trade-off curves.



(a) shows the original data (b) identification of outliers (c) smooth class regions from consisting of sparse outliers, shown as black dots refined metric e.g. red point in the blue region.



(d) is the optimal metric for (e) refined metric with outliers removed a full data

Figure 10. Outlier removal.

Finally we present a comparison of our optimal metric with outlier removal to the commonly used Euclidean, Manhattan and Maximum metric. As mentioned earlier, Figure 11 shows that our optimal metric produces smoother class boundaries with outliers. In contrast, Euclidean distance gives significantly inferior results.

## B. Training DNN using only boundary points

In this experiment, we explore how our approach can be used in training a deep neural network using far fewer points than typically required. Specifically, after finding a near-optimal metric using our proposed strategy, we identify the “boundary” points and “outliers” of the data set and then feed these points to a DNN for training.

In Figure 12 left, we show a simple datasets with 10 classes containing 10000 points in 2D. In the middle, we show the boundary points which were then used to train a fully-connected DNN with 4 hidden layers and ReLU activation function. We used soft-max cross-entropy loss function and trained it using only boundary points. In the right figure, we show how the DNN predicts classes for different points in the unit square. The testing error for interior points is 0.7836%, i.e. less than 1%. This experiment hence further confirms the value of careful data selection using boundary points based on a suitably chosen metric, by demonstrating substantial computational and data savings even for other more intensive techniques like DNN.

## C. Comparison with LMNN

### C.1. Qualitative comparison with LMNN

Given that the LMNN approach of (Weinberger & Saul, 2009) is arguably the nearest neighbor to our approach, we now highlight the salient differences between the two approaches. Let  $\mathcal{T}_i$  be the pre-defined target co-class neighbors of point  $i$  needed as input for LMNN;  $\mathcal{T} = \{(i, j) : i \in \mathcal{N}, j \in \mathcal{T}_i\}$  be the set of target pairs;  $\mathcal{U} = \{(i, j, k) : i \in \mathcal{N}, j \in \mathcal{T}_i, k \in \bar{\mathcal{C}}_i\}$

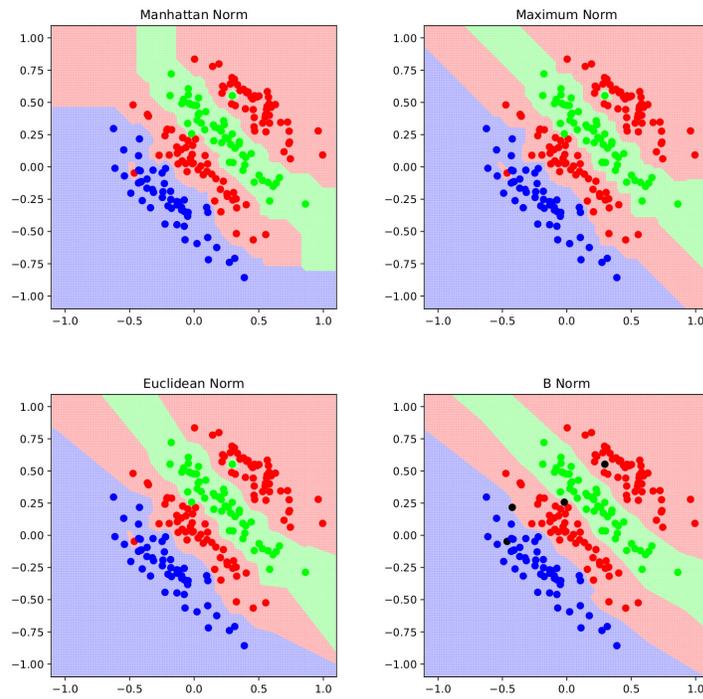


Figure 11. Our result (bottom right) produces smoother class boundaries with outlier removal compared to alternatives.

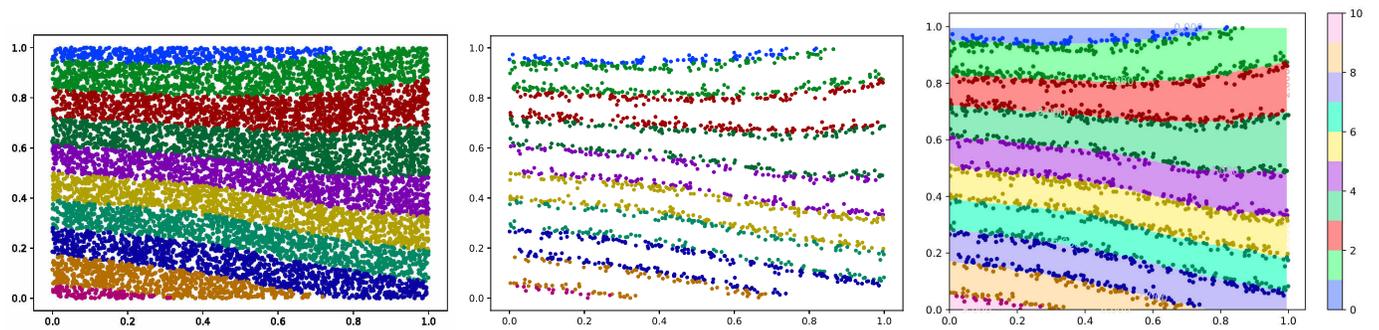


Figure 12. Left figure shows a full, densely sampled 10-class, 10000-point dataset in 2D. Middle figure shows the boundary points computed using the  $R$ -ratio (4). Right figure shows DNN classification of the entire domain using only the boundary points as training set.

be the triplets for which a large margin of separation is desired. In our notation, the LMNN optimization problem becomes

$$\min_{\mathbf{B}, \mathbf{s}} (1 - \mu) \sum_{(i,j) \in \mathcal{T}} \delta_{ij}^{\top} \mathbf{B} \delta_{ij} + \mu \sum_{(i,j,k) \in \mathcal{U}} s_{ijk} \quad (17a)$$

$$\text{s.t. } \delta_{ij}^{\top} \mathbf{B} \delta_{ij} + 1 \leq \delta_{ik}^{\top} \mathbf{B} \delta_{ik} + s_{ijk}, \quad \forall (i, j, k) \in \mathcal{U}, \quad (17b)$$

$$s_{ijk} \geq 0, \quad \forall (i, j, k) \in \mathcal{U}, \quad (17c)$$

$$\mathbf{B} \succeq \mathbf{0}. \quad (17d)$$

Here  $\mu \in [0, 1]$  is a scalar weight to balance the trade-off between imposter violation penalties and the choice of distance metric. Constraint (17d) enforces  $\mathbf{B}$  to be positive semidefinite.

There are several notable differences between our approach and the LMNN. First, (Weinberger & Saul, 2009) rely on a set  $\mathcal{T}$  of pre-defined target co-class neighbors, which can lead to distorted distance metrics if not chosen judiciously (see Figure 1).

Second, we use a 0-1 loss function, as opposed to a hinge loss function, for outliers/imposters. This means that, in contrast to LMNN, which incurs a small penalty for a small violation of condition (1), our approach treats even the slightest violation as a grave infringement. A common argument for avoiding the 0-1 loss function is that the resulting optimization problem is NP-hard. While this is true in theory, it by no means implies that such problems are unsolvable or prohibitively expensive for current methods. On the contrary, mixed-integer optimization solvers have witness tremendous improvements over the last two decades (Bixby, 2012) and challenging MILPs with millions of decision variables and constraints are routinely solved today. Moreover, there are powerful heuristics for solving MILPs that do not guarantee provable optimality, yet are capable of quickly generating high quality solutions.

Third, LMMN and other approaches require a minimum margin between co-class and non-class point while allowing distances to go to infinity. In contrast, we impose an upper bound of 1 on all distance pairs. It is computationally advantageous when using a 0-1 loss function for outliers to have a known finite bound on all distances. Specifically, it allows for smaller values of  $M_{ik}$ , which, in turn, leads to tighter linear relaxations and faster solve times for standard mixed-integer optimization methods.

Fourth, our distance metric is more general than a Mahalanobis metric. However, when we restrict our method to search only for a Mahalanobis metric as done in the LMNN approach, then we see that LMNN strictly enforces positive semidefiniteness, whereas our current formulation does not.

## C.2. Numerical comparison with LMNN

As we argued in Section 1.2, state-of-the-art methods that require target neighbors are susceptible to the choice of target neighbors, which can lead to distorted distance metrics. Figure 13 demonstrates precisely this undesirable behavior for LMNN on a synthetic data set resembling that in Figure 1 with small vertical distance between classes. The optimal metric for various  $K$  from LMNN is inconsistent, thus demonstrating the difficulty of choosing target neighbours effectively – even in this easy example. In contrast, our approach identifies the ideal Mahalanobis metric  $\mathbf{B}^* = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$ <sup>6</sup>. Further, note that the prescribed option in LMNN is the use of Euclidean distance to choose target neighbours, which can be highly misleading if the optimal metric is skewed.

## D. Metric learning modeling extensions

### D.1. Enforcing positive semi-definiteness

We may require the matrix  $\mathbf{B}$  to be positive semi-definite. Gershgorin’s diagonal dominance theorem provides a partial solution to this approach as it restricts  $\mathbf{B}$  to live in a restricted space of positive semi-definite matrices. To accomplish this with our matrix  $\mathbf{B}$ , let  $R_k(\mathbf{B}) = \sum_{\ell \neq k} |b_{k\ell}|$  and note that the condition  $\lambda_G^{\min}(\mathbf{B}) \geq 0$  implies that  $\min_{k=1, \dots, p} \{b_{kk} - R_k(\mathbf{B})\} \geq 0$ , which is equivalent to  $b_{kk} - \sum_{\ell \neq k} |b_{k\ell}| \geq 0$  for all  $k = 1, \dots, p$ . After introducing auxiliary non-negative decision

<sup>6</sup>This means that only if we move in  $y$ -axis, we are measuring some distance, all points in  $x$ -axis have distance 0.

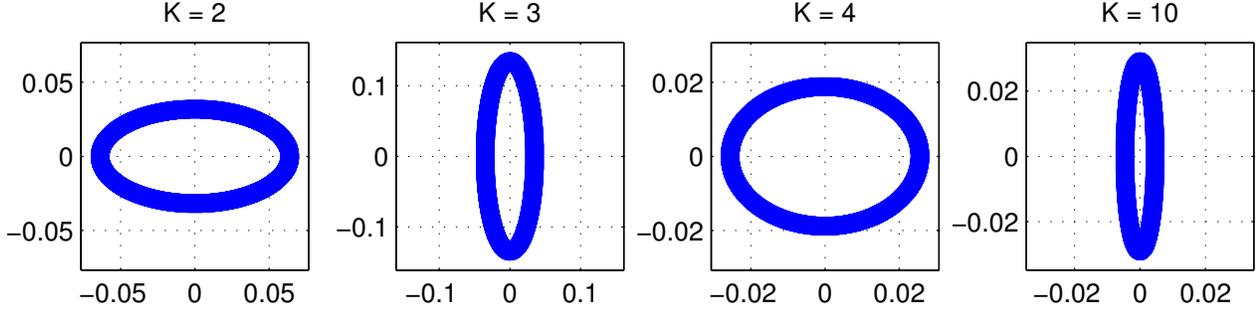


Figure 13. Suboptimal performance of LMNN on dataset described in Figure 1. The optimal metric for various  $K$  from LMNN is inconsistent, thus demonstrating the difficulty of choosing target neighbours effectively.

variables  $b_{k\ell}^+$  to model  $|b_{k\ell}|$  for all  $k \neq \ell$ , the latter can be converted into a set of linear constraint as follows:

$$b_{kk} - \sum_{\ell \neq k} b_{k\ell}^+ \geq 0 \quad \forall k = 1, \dots, p \quad (18a)$$

$$b_{k\ell}^+ \geq b_{k\ell} \quad \forall k = 1, \dots, p, \ell = 1, \dots, p (k \neq \ell) \quad (18b)$$

$$b_{k\ell}^+ \geq -b_{k\ell} \quad \forall k = 1, \dots, p, \ell = 1, \dots, p (k \neq \ell) \quad (18c)$$

$$b_{k\ell}^+ \geq 0 \quad \forall k = 1, \dots, p, \ell = 1, \dots, p (k \neq \ell). \quad (18d)$$

Appending these constraints to the feasible region (5) allows the user to enforce a restricted version of positive semi-definiteness.

## D.2. Sparsification

Regularization terms can easily be incorporated into the objective function to promote a sparse distance metric. Indeed, adding the term  $\mu \sum_{k,\ell:k \leq \ell} |b_{k,\ell}|$  to the objective/loss function as is done in lasso can accomplish this task. Going a step further, one can easily include a cardinality constraint to ensure that no more than  $U$  coefficients are included. Let  $B^{\max}$  be an upper bound on the absolute value of any coefficient  $b_{k,\ell}$  and let  $q_{k,\ell}$  denote a binary variable that takes value 1 if  $b_{k,\ell}$  is non-zero (0 otherwise). Appending the following constraints to the feasible region (5) allows the user to enforce sparsification:

$$-B^{\max} q_{k,\ell} \leq b_{k,\ell} \leq B^{\max} q_{k,\ell} \quad \forall k, \ell \quad (19)$$

$$\sum_{k,\ell:k \leq \ell} q_{k,\ell} \leq U \quad (20)$$