



Unsupervised representation learning with Minimax distance measures

Morteza Haghir Chehreghani¹

Received: 27 June 2019 / Revised: 2 March 2020 / Accepted: 4 June 2020 / Published online: 28 July 2020
© The Author(s) 2020

Abstract

We investigate the use of Minimax distances to extract in a nonparametric way the features that capture the unknown underlying patterns and structures in the data. We develop a general-purpose and computationally efficient framework to employ Minimax distances with many machine learning methods that perform on numerical data. We study both computing the pairwise Minimax distances for all pairs of objects and as well as computing the Minimax distances of all the objects to/from a fixed (test) object. We first efficiently compute the pairwise Minimax distances between the objects, using the equivalence of Minimax distances over a graph and over a minimum spanning tree constructed on that. Then, we perform an embedding of the pairwise Minimax distances into a new vector space, such that their squared Euclidean distances in the new space equal to the pairwise Minimax distances in the original space. We also study the case of having multiple pairwise Minimax matrices, instead of a single one. Thereby, we propose an embedding via first summing up the centered matrices and then performing an eigenvalue decomposition to obtain the relevant features. In the following, we study computing Minimax distances from a fixed (test) object which can be used for instance in K -nearest neighbor search. Similar to the case of all-pair pairwise Minimax distances, we develop an efficient and general-purpose algorithm that is applicable with any arbitrary base distance measure. Moreover, we investigate in detail the edges selected by the Minimax distances and thereby explore the ability of Minimax distances in detecting outlier objects. Finally, for each setting, we perform several experiments to demonstrate the effectiveness of our framework.

Keywords Representation learning · Distance measure · Computational efficiency · Minimax distances

Editors: Larisa Soldatova, Joaquin Vanschoren.

✉ Morteza Haghir Chehreghani
morteza.chehreghani@gmail.com

¹ Department of Computer Science and Engineering (CSE), Chalmers University of Technology, Gothenburg, Sweden

1 Introduction

Data is usually described by a set of objects and a corresponding representation. The representation can be for example the vectors in a vector space or the pairwise dissimilarities between the objects. In real-world applications, the data is often very complicated and a priori unknown. Thus, the basic representation, e.g., *Euclidean* distance, *Mahalanobis* distance, *cosine* similarity and *Pearson correlation*, might fail to correctly capture the underlying patterns or classes. Thereby, the raw data needs to be processed further in order to obtain a more sophisticated representation. Kernel methods are a common approach to enrich the basic representation of the data and model the underlying patterns (Shawe-Taylor and Cristianini 2004; Hofmann et al. 2008). However, the applicability of kernels is confined by several limitations, such as, (1) finding the optimal parameter(s) of a kernel function is often very critical and nontrivial (Nadler and Galun 2007; Luxburg 2007), and (2) as we will see later, kernels assume a global structure which does not distinguish between the different type of classes in the data.

A category of distance measures, called *link-based* measure (Fouss et al. 2007; Chebotarev 2011), takes into account all the *paths* between the objects represented in a graph (where the edge weights indicate the respective pairwise dissimilarities). The *path-specific* distance between nodes i and j is computed by summing the edge weights on this path (Yen et al. 2008). Their link-based distance is then obtained by summing up the *path-specific* measures of all paths between them. Such a distance measure is known to better capture the arbitrarily shaped patterns compared to the basic representations such as Euclidean or Mahalanobis distances. Link-based measures are often obtained by inverting the Laplacian of the distance matrix, in the context of regularized Laplacian kernel and Markov diffusion kernel (Yen et al. 2008; Fouss et al. 2012). However, computing all-pairs link-based distances requires $\mathcal{O}(N^3)$ runtime, where N is the number of objects; thus it is not applicable to large-scale datasets.

A more effective distance measure, called *Minimax* measure, selects the minimum largest gap among all possible paths between the objects. This measure, known also as *Path-based* distance measure, has been first investigated on clustering applications (Fischer and Buhmann 2003; Chehreghani 2016; Pavan and Pelillo 2007). It was also proposed as an axiom for evaluating clustering methods (Zadeh and Ben-David 2009). A straightforward approach to compute all-pairs Minimax distances is to use an adapted variant of the Floyd–Warshall algorithm. The runtime of this algorithm is $\mathcal{O}(N^3)$ (Aho and Hopcroft 1974; Cormen et al. 2001). This distance measure has been also integrated into a variant of K -means providing an agglomerative algorithm whose runtime is $\mathcal{O}(N^2|E| + N^3 \log N)$ (Fischer and Buhmann 2003) ($|E|$ indicates the number of edges in the corresponding graph).

In addition, Minimax distances have been so far applied to a limited type of classification, i.e. to K -nearest neighbor search. The method in Kim and Choi (2007) presents a message passing algorithm with forward and backward steps, similar to the sum-product algorithm (Kschischang et al. 2006). The method takes $\mathcal{O}(N)$ time, which is in theory equal to the standard K -nearest neighbor search, but the algorithm needs several visits of the training dataset. Moreover, this method requires computing a minimum spanning tree (MST) in advance which might require $\mathcal{O}(N^2)$ runtime. Later on, a greedy algorithm (Kim and Choi 2013), proposes to compute the Minimax K nearest neighbors by space partitioning and using Fibonacci heaps whose runtime is $\mathcal{O}(\log N + K \log K)$. However, this method is applicable only to Euclidean spaces and assumes the graph is sparse.

Motivation Minimax distances enable to cope with arbitrarily shaped classes and structures in the data. For example, it has been shown that Minimax K -nearest neighbor classification is effective on non-spherical data, whereas the standard variant, the metric learning approach (Weinberger and Saul 2009), or the shortest path distance (Tenenbaum et al. 2000) might give poor results, since they ignore the underlying geometry. See for example Figure 1 in Kim and Choi (2013). In particular, four properties of Minimax distances are attractive for us:

- They enable to compute the patterns and structures in a non-parametric way, i.e. unlike many kernel methods, they do not require fixing any critical parameter in advance.
- They extract the structures adaptively, i.e., they adapt appropriately whenever the classes differ in shape or type.
- They take into account the transitive relations: if object a is similar to b , b is similar to c , ..., to z , then the Minimax distance between a and z will be small, although their direct distance might be large. This property is particularly useful when dealing with elongated or arbitrarily shaped patterns. Moreover, if a basic pairwise similarity is broken due to noise, then their Minimax distance is able to correct it via taking into account the other paths and relations. For example if the similarity of i and j is broken, then according to Minimax distances, they might still be similar via $i \rightarrow k \rightarrow l \rightarrow \dots \rightarrow j$.
- Many learning methods perform on a vector representation of the objects. However, such a representation might not be available. We might be given only the pairwise distances which do not necessarily induce a *metric*. Minimax distances satisfy the *metric* conditions and enable to compute an embedding, as we will study in this paper.

Contributions Our goal is to develop a generic and computationally efficient framework wherein many different machine learning algorithms can be applied to Minimax distances, beyond e.g., K -nearest neighbor classification or the few clustering methods mentioned before. Within this unified framework, we consider computing both the all-pair pairwise Minimax distances and the Minimax distances of all the objects to/from a fixed (test) object.

1. We first efficiently compute the pairwise Minimax distances between the objects, using the equivalence of Minimax distances over a graph and over a minimum spanning tree constructed on that. This approach reduces the runtime of computing the pairwise Minimax distances to $\mathcal{O}(N^2)$ from $\mathcal{O}(N^3)$.
2. Then, we investigate the possibility of embedding the pairwise Minimax distances into a vector space. This feasibility, allows us to perform an Euclidean embedding of the pairwise Minimax distances, such that the pairwise squared Euclidean distances in the new space equal to the pairwise Minimax distances in the original space. Such an embedding enables us to apply any numerical learning algorithm on the resultant Minimax vectors.
3. We also consider the cases where there are multiple pairwise Minimax matrices instead of a single matrix, which might happen when dealing with multiple pairwise relations, robustness or high-dimensional data. Hence, to obtain a collective embedding, we propose to first center the individual Minimax matrices and then sum them up. This makes an embedding feasible, because the resultant matrix is positive semidefinite.
4. In the following, we study computing Minimax distances from a fixed (test) object which can be used for instance in K -nearest neighbor search. For this case, we propose an efficient and computational optimal Minimax K -nearest neighbor algorithm whose

runtime is $\mathcal{O}(N)$, similar to the standard K -nearest neighbor search, and it can be applied with any arbitrary base dissimilarity measure.

5. Moreover, we investigate in detail the edges selected by the Minimax distances and thereby explore the ability of Minimax distances in detecting outlier objects.
6. Finally, we experimentally study our framework in different machine learning problems (classification, clustering and K -nearest neighbor search) on several synthetic and real-world datasets and illustrate its effectiveness and superior performance in different settings.

The rest of the paper is organized as following. In Sect. 2, we introduce the notations and definitions. In Sect. 3, we develop our framework for computing pairwise Minimax distances and extracting the relevant features applicable to general machine learning methods. Here, we also extend our approach for computing and embedding Minimax vectors for multiple pairwise data relations, i.e., to collective Minimax distances. In Sect. 4, we extend further our framework for computing the Minimax K -nearest neighbor search and outlier detection. In Sect. 5, we describe the experimental studies, and finally, we conclude the paper in Sect. 6.

2 Notations and definitions

A dataset can be modeled by a graph $\mathcal{G}(\mathbf{O}, \mathbf{D})$, where \mathbf{O} and \mathbf{D} respectively indicate the set of N objects (nodes) and the corresponding edge weights such that \mathbf{D}_{ij} shows the pairwise dissimilarity between objects i and j . The base pairwise dissimilarities \mathbf{D}_{ij} can be computed for example according to squared Euclidean distances or cosine (dis)similarities between the vectors that represent i and j . In several applications, \mathbf{D}_{ij} might be given directly.¹ In general, \mathbf{D} might not yield a *metric*, i.e. the triangle inequality may not hold. We have recently studied the application of Minimax distances to correlation clustering, where the triangle inequality does not necessarily hold for the base pairwise dissimilarities (Chehreghani 2020). In our study, \mathbf{D} needs to satisfy three basic conditions:

1. zero self distances, i.e. $\forall i, \mathbf{D}_{ii} = 0$,
2. non-negativity, i.e. $\forall i, j, \mathbf{D}_{ij} \geq 0$, and
3. symmetry, i.e. $\forall i, j, \mathbf{D}_{ij} = \mathbf{D}_{ji}$.

We assume there are no duplicates, i.e. the pairwise dissimilarity between every two distinct objects is positive. For this purpose, we may either remove the duplicate objects or perturb them slightly to make the zero non-diagonal elements of \mathbf{D} positive. The goal is then to use the Minimax distances and the respective features to perform the machine learning task of interest.

The Minimax (MM) distance between objects i and j is defined as

$$\mathbf{D}_{i,j}^{MM} = \min_{r \in \mathcal{R}_{ij}(\mathbf{O})} \{ \max_{1 \leq l \leq |r|-1} \mathbf{D}_{r(l)r(l+1)} \}, \quad (1)$$

¹ For simplicity of explanation, we assume that the graph is full, i.e. the missing edges are filled by a large value. However, our analysis can be easily extended to arbitrary graphs (Chehreghani 2017).

where $\mathcal{R}_{ij}(\mathbf{O})$ is the set of all paths between i and j . Each path r is identified by a sequence of object indices, i.e. $r(l)$ shows the l th object on the path.

3 Features extraction from pairwise Minimax distances

We aim at developing a unified framework for performing arbitrary numerical learning methods with Minimax distances. To cover different algorithms, we pursue the following strategy:

1. We compute the pairwise Minimax distances for all pairs of objects i and j in the dataset.
2. We, then, compute an embedding of the objects into a new vector space such that their pairwise (squared Euclidean) distances in this space equal their Minimax distances in the original space.

Notice that vectors are the most basic way for data representation, since they render a bijective mapping between the objects and the measurements. Hence, any machine learning method which performs on numerical data can benefit from our approach. Such methods perform either on the objects in a vector space, e.g. logistic regression, or on a kernel matrix computed from the pairwise relations between the objects. In the later case, the pairwise Minimax distances can be used for this purpose, for example through an exponential transformation, or the kernel can be computed from the final Minimax vectors.

3.1 Computing pairwise Minimax distances

Previous works for computing Minimax distances (e.g. applied to clustering) use a variant of Floyd–Warshall algorithm whose runtime is $\mathcal{O}(N^3)$ (Aho and Hopcroft 1974; Cormen et al. 2001), or combine it with K -means in an agglomerative algorithm whose runtime is $\mathcal{O}(N^2|E| + N^3 \log N)$ (Fischer and Buhmann 2003). To reduce such a computational demand, we follow a more efficient procedure established in two steps: (1) build a minimum spanning tree (MST) over the graph, and then, (2) compute the Minimax distances over the MST.

I. *Equivalence of Minimax distances over a graph and over a minimum spanning tree on the graph*

We exploit the equivalence of Minimax distances over an arbitrary graph and those obtained from a minimum spanning tree on the graph, as expressed in Theorem 1.²

² This result may lead to several simplifications for computing the base pairwise distance measure \mathbf{D} and the respective graph $\mathcal{G}(\mathbf{O}, \mathbf{D})$, where the graph does not need to be necessarily full. For example, instead of a full graph, we can compute a minimal (connected) graph that sufficiently includes a minimum spanning tree. For instance, one might use the K nearest neighbor graph, which is not necessarily full. On the other hand, we might have no idea about the pairwise dissimilarities. Then we may query them from a user and infer the rest by Minimax distances. In fact, this result proposes an efficient method for this purpose: due to equivalence of Minimax distances on a graph and on any minimum spanning tree on that, we may be able to query significantly smaller number of pairwise dissimilarities that sufficiently represent a MST. This can be significantly faster than querying or inferring all the pairwise distances to obtain \mathbf{D} .

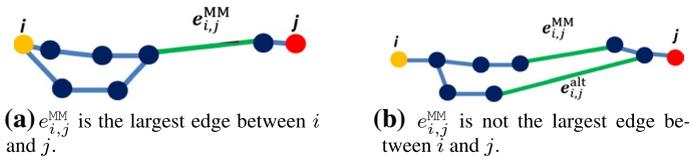


Fig. 1 Equivalence of minimax distances over a graph and over a MST constructed on the graph

Theorem 1 Given graph $\mathcal{G}(\mathbf{O}, \mathbf{D})$, a minimum spanning tree constructed on that provides the sufficient and necessary edges to obtain the pairwise Minimax distances.

Proof The equivalence is concluded in a similar way to the *maximum capacity* problem (Hu 1961), where one can show that picking an edge which does not belong to a minimum spanning tree, yields a larger Minimax distance (i.e., a contradiction occurs).

We here describe a high level proof. Let e_{ij}^{MM} denote the edge representing the Minimax distance between i and j . We investigate two cases:

1. There is only one path between i and j . Then, this path will be part of any minimum spanning tree constructed over the graph. Since, otherwise, if some edge for example e_{ij}^{MM} is not selected, then the tree will lose the connectivity, i.e. there will be less edges in the tree than $N - 1$. The same argument holds when there are several paths between i and j , but all share e_{ij}^{MM} such that e_{ij}^{MM} is the largest edge on all of them (Fig. 1a). Then e_{ij}^{MM} will be selected by any minimum spanning tree, as otherwise the tree will be disconnected.
2. There are several paths between i and j , whose largest edges are different. We need to show that only the path including e_{ij}^{MM} is selected by a minimum spanning tree. To prove, consider two paths, one including e_{ij}^{MM} and the other containing e_{ij}^{alt} which is the largest edge on the alternative path (Fig. 1b). It is easy to see that the minimum spanning tree would choose e_{ij}^{MM} instead of e_{ij}^{alt} : existence of two different paths indicates presence of at least one cycle. According to the definition of a tree, one edge of a cycle should be eliminated to yield a tree. In the cycle including e_{ij}^{MM} , according to the definition of Minimax distances, there is at least one edge not smaller than e_{ij}^{MM} (which is e_{ij}^{alt}), as otherwise e_{ij}^{MM} would not represent the Minimax distance (the minimum of largest gaps) between i and j . Thereby the minimum spanning tree algorithm keeps e_{ij}^{MM} instead of e_{ij}^{alt} , since this choice leads to a tree with a smaller total weight. Notice that the algorithm cannot discard both e_{ij}^{MM} and e_{ij}^{alt} , because then the connectivity will be lost.

On the other hand, it is *necessary* to build a minimum spanning tree, e.g. to have at least $N - 1$ edges selected given that the graph \mathcal{G} is connected. Otherwise, the computed MST will lose its connectivity and thereby for some pairs of objects there will be no path to be computed. □

Therefore, to compute the Minimax distances \mathbf{D}^{MM} , we need to care only about the edges which are present in an MST over the graph. Then, the Minimax distances are written by

$$\mathbf{D}_{ij}^{MM} = \max_{1 \leq l \leq |r_{ij}|-1} \mathbf{D}_{r_{ij}(l)r_{ij}(l+1)}, \tag{2}$$

where r_{ij} indicates the (only) path between i and j . This result does not depend on the particular choice of the algorithm used for constructing the minimum spanning tree. The

graph that we assume is a full graph, i.e. we compute the pairwise dissimilarities between all pairs of objects. For full graphs, the straightforward implementation of the Prim's algorithm (1957) using an auxiliary vector requires $\mathcal{O}(N^2)$ runtime which is an optimal choice.

II. Pairwise Minimax distances over a MST

At the next step, after constructing a minimum spanning tree, we compute the pairwise Minimax distances over that. A naive and straightforward algorithm would perform a Depth First Search (DFS) from each node to compute the Minimax distances by keeping the track of the largest distance between the initial node and each of the traversed nodes. A single run of DFS requires $\mathcal{O}(N)$ runtime and thus the total time will be $\mathcal{O}(N^2)$. However, such a method might lead to visiting some edges multiple times which renders an unnecessary extra computation. For instance, a maximal edge weight might appear in many DFSs such that it is processed several times. Hence, a more elegant approach would first determine all the objects whose pairwise distances are represented by an edge weight and then assigns the weight as their Minimax distances. According to Theorem 1, every edge in the MST represents some Minimax distances. Thus, we first find the edge(s) which represent only few Minimax distances, namely only one pairwise distance, such that the respective objects can be immediately identified. Lemma 1 suggest existence of this kind of edges.

Lemma 1 *In a minimum spanning tree, for the minimal edge weights we have $\mathbf{D}_{ij}^{MM} = \mathbf{D}_{ij}$, where i and j are the two objects that occur exactly at the two sides of the edge.*

Proof Without loss of generality, we assume that the edge weights are distinct. Let e_{min} denote the edge with minimal weight in the MST. We consider the pair of objects p and q such that at least one of them is not directly connected to e_{min} and show that e_{min} does not represent their Minimax distance. In a MST, there is exactly one path between each pair of objects. Thus, on the path between p and q there is at least another edge whose weight is not smaller than the weight of e_{min} , which hence represents the Minimax distance between p and q , instead of e_{min} . \square

Lemma 1 yields a dynamic programming approach to compute the pairwise Minimax distances from a tree. We first sort the edge weights of the MST via for example merge sort or heap sort (Cormen et al. 2001) which in the worst case require $\mathcal{O}(N \log N)$ time. We then consider each object as a separate component. We process the edge weights one by one from the smallest to the largest, and at each step, we set the Minimax distance of the two respective components by this weight. Then, we remove the two components and replace them by a new component constructed from the combination (union) of the two. We repeat these two steps for all edges of the minimum spanning tree. A main advantage of this algorithm is that whenever an edge is processed, all the nodes that it represents their Minimax distances are ready in the components connected to the two sides of the edge. Algorithm 1 describes the procedure in detail.

Algorithm 1 Computation of pairwise Minimax distances on a tree.

Require: A tree specified by, i) *indices*: a $|E| \times 2$ matrix of (endpoints of) edges, and ii) *weights*: vector of edge weights.

Ensure: Matrix of pairwise Minimax distances \mathbf{D}^{MM} .

```

1: sorting_ind = arg sort(weights)
2: weights = weights[sorting_ind]
3: indices = indices[sorting_ind, :]

4: component_list = []
5: for i = 0 to N - 1 do
6:   component_id[i] = i
7:   component_list.append([i])
8: end for
9: cur_id = N - 1
10: for i = 0 to  $|E| - 1$  do
11:   ind1 = indices[i, 0]
12:   ind2 = indices[i, 1]
13:   if component_id[ind1]  $\neq$  component_id[ind2] then
14:     first_side = component_list[component_id[ind1]]
15:     scnd_side = component_list[component_id[ind2]]

16:     new_component = first_side  $\cup$  scnd_side
17:     cur_id = cur_id + 1

18:     component_id[new_component] = cur_id
19:     component_list.append(new_component)

20:      $\mathbf{D}^{MM}$ [first_side, scnd_side] = weights[i]
21:      $\mathbf{D}^{MM}$ [scnd_side, first_side] = weights[i]
22:   end if
23: end for
24: return  $\mathbf{D}^{MM}$ 

```

Algorithm 1 uses the following data structures.

1. *component_list*: is a list of lists, wherein each list (*component*) contains the set of nodes (objects) that are treated similarly, i.e. they have the same Minimax distance to/from an external node or component.
2. *component_id*: a N -dimensional vector containing the ID of the latest component that each object belongs to.

Each object initially constitutes a separate component. The algorithm, at each step i , pops out the unselected edge e_i which has the smallest weight. For this purpose, we assume that the vector of the edge weights (i.e. *weights*) has been sorted in advance. The nodes associated to the edges are rearranged according to the ordering of *weights* and are stored in *indices*. If the edge e_i does not entirely occur inside a single component, then the nodes reachable from each side of e_i (i.e. from *ind1* and *ind2*) are selected and stored respectively in *first_side* and *scnd_side*. For this purpose, we use a vector called *component_id* which keeps the ID (index) of the latest component that each node belongs to. Therefore, the component *first_side* is obtained by *component_list*[*component_id*[*ind1*]] and similarly *scnd_side* by

$component_list[component_id[ind2]]$. Then \mathbf{D}^{MM} is updated by

$$\begin{aligned} \mathbf{D}^{MM}[first_side, scnd_side] &= weights[i], \\ \mathbf{D}^{MM}[scnd_side, first_side] &= weights[i]. \end{aligned} \quad (3)$$

Finally, a new component is constructed and added to $component_list$ by combining the two base components $first_side$ and $scnd_side$. The ID of this new component is used as the ID of its members in $component_id$.

3.2 Embedding of pairwise Minimax distances

In the next step, given the matrix of pairwise Minimax distances \mathbf{D}^{MM} , we obtain an embedding of the objects into a vector space such that their pairwise squared Euclidean distances in this new space are equal to their Minimax distances in the original space. For this purpose, in Theorem 2, we investigate a useful property of Minimax distance measures, called the *ultrametric* property, and use this property to prove the existence of such an embedding.

Theorem 2 *Given the pairwise distances \mathbf{D} , the matrix of Minimax distances \mathbf{D}^{MM} induces an \mathcal{L}_2 embedding, i.e. there exist a new vector space for the set of objects \mathbf{O} wherein the pairwise squared Euclidean distances are equal to the pairwise Minimax distances in the original space.*

Proof First, we investigate that the pairwise Minimax distances \mathbf{D}^{MM} constitute an *ultrametric*. The conditions to be satisfied are:

1. $\forall i, j : \mathbf{D}_{ij}^{MM} = 0$ if and only if $i = j$. We verify each of the conditions separately. (1) If $i = j$, then $\mathbf{D}_{ij}^{MM} = 0$: We have $\mathbf{D}_{ii}^{MM} = \mathbf{D}_{ii} = 0$ because the smallest maximal gap between every object and itself is zero. (2) If $\mathbf{D}_{ij}^{MM} = 0$, then $\mathbf{D}_{ij} = 0$ and $i = j$, because we have assumed that all the distinct pairwise distances are positive, i.e. zero base or Minimax pairwise distances can occur only if $i = j$.
2. $\forall i, j : \mathbf{D}_{ij}^{MM} \geq 0$. All the edge weights, i.e. the elements of \mathbf{D} , are non-negative. Thus the minimum of them, i.e. $\min(\mathbf{D})$, is also non-negative. Moreover, by definition we have $\mathbf{D}_{ij}^{MM} \geq \min(\mathbf{D})$. Hence, we conclude $\mathbf{D}_{ij}^{MM} \geq \min(\mathbf{D}) \geq 0$.
3. $\forall i, j : \mathbf{D}_{ij}^{MM} = \mathbf{D}_{ji}^{MM}$. By assumption \mathbf{D} is symmetric, therefore, any path from i to j will also be a path from j to i , and vice versa. Thereby, their maximal weights and the minimum among different paths are identical.
4. $\forall i, j, k : \mathbf{D}_{ij}^{MM} \leq \max\{\mathbf{D}_{ik}^{MM}, \mathbf{D}_{kj}^{MM}\}$. We show that otherwise a contradiction occurs. Suppose there is a triplet i, j, k such that $\mathbf{D}_{ij}^{MM} > \max\{\mathbf{D}_{ik}^{MM}, \mathbf{D}_{kj}^{MM}\}$. Then, according to the definition of Minimax distance, the path from i to k and then to j must be used for computing the Minimax distance \mathbf{D}_{ij}^{MM} which leads to $\mathbf{D}_{ij}^{MM} \leq \max\{\mathbf{D}_{ik}^{MM}, \mathbf{D}_{kj}^{MM}\}$, i.e. a contradiction occurs.

On the other hand, *ultrametric* matrices are positive definite (Fiedler 1998; Varga and Nabben 1993) and positive definite matrices induce an Euclidean embedding (Schoenberg 1937). \square

Notice that we do not require \mathbf{D} to induce a *metric*, i.e. the triangle inequality is not necessarily fulfilled. After satisfying the feasibility condition, there exist several ways to

compute a squared Euclidean embedding. We exploit a method motivated in Young and Householder (1938) and further analyzed in Torgerson (1958) that is known as *multidimensional scaling* [25]. This method works based on centering \mathbf{D}^{MM} to compute a Mercer kernel which is positive semidefinite, and then performing an eigenvalue decomposition, as following:³

1. Center \mathbf{D}^{MM} by

$$\mathbf{W}^{MM} \leftarrow -\frac{1}{2}\mathbf{A}\mathbf{D}^{MM}\mathbf{A}. \tag{4}$$

\mathbf{A} is defined as $\mathbf{A} = \mathbf{I}_N - \frac{1}{N}\mathbf{e}_N\mathbf{e}_N^T$, where \mathbf{e}_N is a vector of length N with 1's and \mathbf{I}_N is an identity matrix of size N .

2. Under this transformation, \mathbf{W}^{MM} is positive semidefinite. Thus, we decompose \mathbf{W}^{MM} into its eigenbasis, i.e.,

$$\mathbf{W}^{MM} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T, \tag{5}$$

where $\mathbf{V} = (v_1, \dots, v_N)$ contains the eigenvectors v_i and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_N)$ is a diagonal matrix of eigenvalues $\lambda_1 \geq \dots \geq \lambda_d \geq \lambda_{d+1} = 0 = \dots = \lambda_N$. Note that the eigenvalues are nonnegative, since \mathbf{W}^{MM} is positive semidefinite.

3. Calculate the $N \times d$ matrix

$$\mathbf{Y}_d^{MM} = \mathbf{V}_d(\mathbf{\Lambda}_d)^{1/2}, \tag{6}$$

where $\mathbf{V}_d = (v_1, \dots, v_d)$ and $\mathbf{\Lambda}_d = \text{diag}(\lambda_1, \dots, \lambda_d)$.

Here, d shows the dimensionality of the Minimax vectors, i.e. the number of Minimax features. The Minimax dimensions are ordered according to the respective eigenvalues and thereby we might choose only the first most representative ones, instead of taking them all.

3.3 Embedding of collective Minimax pairwise matrices

We extend our generic framework to the cases where multiple pairwise Minimax matrices are available, instead of a single matrix. Then, the goal would be to find an embedding of the objects into a new space wherein their pairwise squared Euclidean distances are related to the collective Minimax distances over different Minimax matrices. Such a scenario might be interesting in several situations:

- There might exist different type of relations between the objects, where each renders a separate graph. Then, we compute several pairwise Minimax distances, each for a specific relation.
- Minimax distances might fail when for example few noise objects connect two compact classes. Then, the inter-class Minimax distances become very small, even if the objects from the two classes are connected via only few outliers. To solve this issue, similar to model averaging, one could use the higher order, i.e. the second, third, ...Minimax distances, e.g., the second smallest maximal gap. Then, there will be multiple pairwise Minimax matrices each representing the k th Minimax distance.

³ One could instead use more efficient methods such as Globerson et al. (2007) and Khoshneshin and Street (2010), or use faster approximate of eigenvectors/eigenvalues (Quarteroni et al. 2007).

- In many real-world applications, we encounter high-dimensional data, the patterns might be hidden in some unknown subspace instead of the whole space, such that they are disturbed in the high-dimensional space due to the curse of dimensionality. Minimax distances rely on the existence of well-connected paths, whereas such paths might be very sparse or fluctuated in high dimensions. Thereby, similar to ensemble methods, it is natural to seek for connectivity paths and thus for Minimax distances in some subspaces of the original space.⁴ However, investigating all the possible subspaces is computationally intractable as the respective cardinality scales exponentially with the number of dimensions. Hence, we propose an approximate approach based on computing Minimax distances for each dimension, which leads to having multiple Minimax matrices.

In all the aforementioned cases, we need to deal with (let say M) different matrices of Minimax distances computed for the same set of objects. Then, the next step is to investigate the existence of an embedding that represents the pairwise collective Minimax distances. To proof the existence of such an embedding for $M = 1$, we used the *ultrametric* property of the respective Minimax distances and then concluded the positive definiteness. However, as shown in Theorem 3, the sum of $M > 1$ Minimax matrices does not necessarily satisfy the *ultrametric* property.

Theorem 3 Given $M > 1$ pairwise Minimax matrices $\mathbf{D}^{MM}(m), m \in \{1, \dots, M\}$ for the same set of objects \mathbf{O} , then, the accumulative Minimax matrix $\mathbf{D}^{aMM} = \mathbf{D}^{MM}(1) + \dots + \mathbf{D}^{MM}(m)$ does not necessarily constitute an ultrametric.

Proof We investigate the *ultrametric* conditions for \mathbf{D}^{aMM} (using the results of the analysis in Theorem 2):

- $\forall i, j : \mathbf{D}_{ij}^{aMM} = 0$, if and only if $i = j$. We investigate each of the conditions.
 - If $i = j$, then $\mathbf{D}_{ij}^{aMM} = 0$: Since we have $\mathbf{D}_{ii}^{MM}(1) = \dots = \mathbf{D}_{ii}^{MM}(M) = 0$, then $\mathbf{D}_{ii}^{aMM} = \sum_{m=1}^M \mathbf{D}_{ii}^{MM}(m) = 0$.
 - If $\mathbf{D}_{ij}^{aMM} = 0$, then $i = j$: According to our assumption, for $i \neq j$ we have $\mathbf{D}_{ij} > 0$ and thus $\mathbf{D}_{ij}^{MM} > 0$. Therefore, $\mathbf{D}_{ij}^{aMM} > 0$ if $i \neq j$, i.e. $\mathbf{D}_{ij}^{aMM} = 0$ implies that $i = j$.
- $\forall i, j : \mathbf{D}_{ij}^{aMM} \geq 0$. We have: $\forall 1 \leq m \leq M, \mathbf{D}_{ij}^{MM}(m) \geq 0$. Thus,

$$\mathbf{D}_{ij}^{aMM} = \sum_{m=1}^M \mathbf{D}_{ij}^{MM}(m) \geq 0.$$
- $\forall i, j : \mathbf{D}_{ij}^{aMM} = \mathbf{D}_{ji}^{aMM}$. We have: $\forall 1 \leq m \leq M, \mathbf{D}_{ij}^{MM}(m) = \mathbf{D}_{ji}^{MM}(m)$. Thus,

$$\sum_{m=1}^M \mathbf{D}_{ij}^{MM}(m) = \sum_{m=1}^M \mathbf{D}_{ji}^{MM}(m), \text{ i.e. } \mathbf{D}_{ij}^{aMM} = \mathbf{D}_{ji}^{aMM}.$$
- $\forall i, j, k : \mathbf{D}_{ij}^{aMM} \leq \max\{\mathbf{D}_{ik}^{aMM}, \mathbf{D}_{kj}^{aMM}\}$. $\forall 1 \leq m \leq M$, we have

⁴ An alternative approach would be to perform dimension reduction via e.g. principle component analysis (PCA) and so on. Our approach can be easily extended to such cases too.

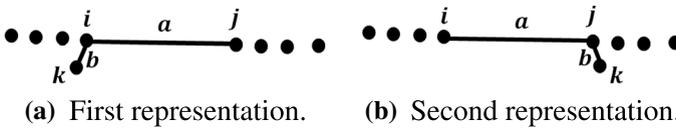


Fig. 2 Two different representation of the same objects. The difference is only in the position of the k th object. Summing up the pairwise Minimax distances of these two datasets does not lead to an *ultrametric*

$$D_{ij}^{MM}(m) \leq \max\{D_{ik}^{MM}(m), D_{kj}^{MM}(m)\}.$$

Then,

$$\sum_{m=1}^M D_{ij}^{MM}(m) \leq \sum_{m=1}^M \max\{D_{ik}^{MM}(m), D_{kj}^{MM}(m)\}.$$

However, we need to show that

$$\sum_{m=1}^M D_{ij}^{MM}(m) \leq \max\left\{\sum_{m=1}^M D_{ik}^{MM}(m), \sum_{m=1}^M D_{kj}^{MM}(m)\right\}.$$

Thus, if we can approximate $\max\{\sum_{m=1}^M D_{ik}^{MM}(m), \sum_{m=1}^M D_{kj}^{MM}(m)\}$ by $\sum_{m=1}^M \max\{D_{ik}^{MM}(m), D_{kj}^{MM}(m)\}$, then, this condition is satisfied too and thereby D^{aMM} induces an *ultrametric*.

However, in general this approximation is not valid and one can find cases where $D_{ij}^{aMM} > \max\{D_{ik}^{aMM}, D_{kj}^{aMM}\}$. Here, we show an example. Let fix $M = 2$ and consider the datasets in Fig. 2. The two datasets represent the same set of objects. They differ only in the position of the k th object and the other objects are fixed among the two datasets. In both datasets, the Minimax distance between i and j is equal to a , i.e. $\sum_{m=1}^M D_{ij}^{MM}(m) = 2a$.

On the other hand, we have $\sum_{m=1}^M D_{ik}^{MM}(m) = b + a$, and similarly $\sum_{m=1}^M D_{kj}^{MM}(m) = a + b$. Thus, $\sum_{m=1}^M D_{ij}^{MM}(m) = 2a > \max\{\sum_{m=1}^M D_{ik}^{MM}(m), \sum_{m=1}^M D_{kj}^{MM}(m)\} = a + b$, i.e., D^{aMM} does not induce an *ultrametric*. \square

Theorem 3 indicates that the accumulative Minimax matrix D^{aMM} does not always yield an *ultrametric*. However, we do not necessarily need the *ultrametric* property to compute an embedding, rather it is only a sufficient condition. Theorem 4 suggests that after computing $D^{MM}(m), m \in \{1, \dots, M\}$, first we center each of them via Eq. 4 and then sum them up. The resultant matrix W^{cMM} is embeddable as it is positive definite.

Theorem 4 For a set of objects O , we are given M pairwise Minimax matrices $D^{MM}(m), m \in \{1, \dots, M\}$ and the respective centered matrices $W^{MM}(m)$ defined as

$$W^{MM}(m) = -\frac{1}{2}AD^{MM}(m)A. \tag{7}$$

Then, the collective matrix $W^{cMM} = \sum_{m=1}^M W^{MM}(m)$ induces an \mathcal{L}_2^2 embedding.

Proof According to Theorem 2, each $D^{MM}(m)$ yields an \mathcal{L}_2^2 embedding, which implies that $W^{MM}(m)$ is positive definite (Torgerson 1958; Young and Householder 1938). On the other

hand, sum of multiple positive definite matrices (i.e. \mathbf{W}^{cMM}) is a positive definite matrix too (Horn and Johnson 1990). Thus, the eigenvalues of \mathbf{W}^{cMM} are non-negative and it induces an \mathcal{L}_2^2 embedding. \square

Thereby, instead of summing the $\mathbf{D}^{MM}(m)$ matrices, we sum up the $\mathbf{W}^{MM}(m)$'s and then perform eigenvalue decomposition, to compute an embedding wherein the pairwise squared Euclidean distances correspond to the collective pairwise Minimax distances. As will be mentioned in the experiments, we perform the above computations on the squared pairwise base dissimilarities.

Efficient calculation of dimension-specific Minimax distances In this paper, we particularly study the use of collective Minimax embedding for high-dimensional data (called the *dimension-specific* variant). The dimension-specific variant may require computing pairwise Minimax distances for each dimension, which can be computationally expensive. However, in this setting, the objects stay in an one-dimensional space. A main property of one-dimensional data is that sorting them immediately gives a minimum spanning tree. We, then, compute the pairwise distances for each pair of consecutive objects in the sorted list to obtain the edge weights of the minimum spanning tree. Finally, we compute the pairwise Minimax distances from the minimum spanning tree.

4 One-to-all Minimax distance measures

In this section, we study computing Minimax distances from a fixed (test) object which can be in particular used in K -nearest neighbor search. In this setup, we are given a graph $\mathcal{G}(\mathbf{O}, \mathbf{D})$ (as the training dataset), and a new object (node) v . \mathbf{O} indicates a set of N (training) objects with the corresponding pairwise dissimilarities \mathbf{D} , i.e., \mathbf{D} denotes the weights of the edges in \mathcal{G} . The goal is to compute the K objects in \mathbf{O} whose Minimax distances from v are not larger than any other object in \mathbf{O} . We assume there is a function $d(v, \mathbf{O})$ which computes the pairwise dissimilarities between v and all the objects in \mathbf{O} . Thereby, by adding v to \mathcal{G} we obtain the graph $\mathcal{G}^+(\mathbf{O} \cup v, \mathbf{D} \cup d(v, \mathbf{O}))$.

4.1 Minimax K -nearest neighbor search

We introduce an incremental approach for computing the Minimax K -nearest neighbors of the new object v , i.e. we obtain the $(T + 1)$ th Minimax neighbor of v given the first T Minimax neighbors. Thereby, first, we define the *partial neighbors* of v as the first Minimax neighbors of v over graph $\mathcal{G}^+(\mathbf{O} \cup v, \mathbf{D} \cup d(v, \mathbf{O}))$, i.e.,

$$\mathcal{N}\mathcal{P}_T(v) = \left\{ \{i\} \subseteq \mathbf{O} : |\{i\}| = T, \nexists j \in \{\mathbf{O} \setminus \mathcal{N}\mathcal{P}_T(v)\} : \mathbf{D}_{vj}^{MM} < \mathbf{D}_{vi}^{MM} \right\}. \quad (8)$$

Theorem 5 provides a way to extend $\mathcal{N}\mathcal{P}_T(v)$ step by step until it contains the K nearest neighbors of v according to Minimax distances.

Theorem 5 *Given the (training) graph $\mathcal{G}(\mathbf{O}, \mathbf{D})$ and a test node v , assume we have already computed the first T Minimax neighbors of v (i.e., the set $\mathcal{N}\mathcal{P}_T(v)$). Then, the node with the*

minimal distance to the set $\{v \cup \mathcal{N}\mathcal{P}_T(v)\}$ gives the $(T + 1)$ th Minimax nearest neighbor of v .⁵

Proof Let us call this new (potential) neighbor u . Let e_{u^*} indicate the edge (with the smallest weight) connecting u to $\{v \cup \mathcal{N}\mathcal{P}_T(v)\}$ and $p \in \{v \cup \mathcal{N}\mathcal{P}_T(v)\}$ denote the node at the other side of e_{u^*} . We consider two cases:

1. $\mathbf{D}_{p,u} \leq \mathbf{D}_{v,p}^{MM}$, i.e. the weight of e_{u^*} is not larger than the Minimax distance between v and p . Then, $\mathbf{D}_{v,u}^{MM} = \mathbf{D}_{v,p}^{MM}$ and therefore u is naturally a *valid* choice.
2. $\mathbf{D}_{p,u} > \mathbf{D}_{v,p}^{MM}$, then we show that there exist no other unselected node u' (i.e. from $\{\mathbf{O} \setminus \mathcal{N}\mathcal{P}_T(v)\}$) which has a smaller Minimax distance to v than u . Thereby:
 - (a) $\mathbf{D}_{v,u}^{MM} > \mathbf{D}_{v,p}^{MM}$ (according to the assumption $\mathbf{D}_{p,u} > \mathbf{D}_{v,p}^{MM}$). Moreover, it can be shown that there is no other path from v to u whose maximum weight is smaller than $\mathbf{D}_{p,u}$. Because, if such a path exists, then there must be a node $p' \in \mathcal{N}\mathcal{P}_T(v)$ which has a smaller distance to an external node like u' . This leads to a contradiction since u is the closest neighbor of $\mathcal{N}\mathcal{P}_T(v)$.
 - (b) For any other unselected node u' we have $\mathbf{D}_{v,u'}^{MM} \geq \mathbf{D}_{p,u}$. Because computing the Minimax distance between v and u' requires visiting an edge whose one side is inside $\mathcal{N}\mathcal{P}_T(v)$, but the other side is outside $\mathcal{N}\mathcal{P}_T(v)$. Among such edges, e_{u^*} has the minimal weight, therefore, $\mathbf{D}_{v,u'}^{MM}$ cannot be smaller than $\mathbf{D}_{p,u}$.

Finally, from (a) and (b) we conclude that $\mathbf{D}_{v,u}^{MM} \leq \mathbf{D}_{v,u'}^{MM}$ for any unselected node $u' \neq u$. Hence, u has the smallest Minimax distance to v among all unselected nodes. □

Theorem 5 proposes a dynamic programming approach to compute the Minimax K -nearest neighbors of v . Iteratively, at each step, we obtain the next Minimax nearest neighbor by selecting the external (unselected) node u which has the minimum distance to the nodes in $\{v \cup \mathcal{N}\mathcal{P}_T(v)\}$. This procedure is repeated for K times. Algorithm 2 describes the steps in detail. Since the algorithm performs based on finding the nearest unselected node, therefore vector *dist* is used to keep the minimum distance of unselected nodes to (one of the members of) $\{v \cup \mathcal{N}\mathcal{P}_T(v)\}$. Thereby, the algorithm finds a new neighbor in two steps: (1) *extension*: it picks the minimum of *dist* and adds the respective node to $\mathcal{N}\mathcal{P}_T(v)$, and (2) *update*: it updates *dist* by checking if an unselected node has a smaller distance to the new $\mathcal{N}\mathcal{P}_T(v)$. Thus *dist* is updated by $\mathit{dist} = \min(\mathit{dist}, \mathbf{D}_{\min_ind, \cdot})$, except for the members of $\mathcal{N}\mathcal{P}_T(v)$.

⁵ Given a set $S \subset \mathbf{O}$, the distance of object $i \in \{\mathbf{O} \setminus S\}$ to S is obtained by the minimum of the distances between i and the objects in S .

Algorithm 2 Calculation of Minimax K -nearest neighbor

Require: Graph $\mathcal{G}^+(\mathbf{O} \cup v, \mathbf{D} \cup d(v, \mathbf{O}))$ including a test object v .
Ensure: A list of K Minimax nearest neighbors of v stored in $\mathcal{NP}(v)$.

- 1: Initialize vector $dist$ by the distance of v to each object in \mathbf{O} .
- 2: $\mathcal{NP}(v) = []$
- 3: **for** $T = 1$ **to** K **do**
- 4: {extension}
- 5: $min_ind = \arg \min(dist)$
- 6: $\mathcal{NP}(v).append(min_ind)$
- 7: {update}
- 8: $dist_{min_ind} = inf$
- 9: **for** $i \notin \mathcal{NP}(v)$ **do**
- 10: **if** $\mathbf{D}_{new_ind,i} < dist_i$ **then**
- 11: $dist_i = \mathbf{D}_{min_ind,i}$
- 12: **end if**
- 13: **end for**
- 14: **end for**
- 15: **return** $\mathcal{NP}(v)$

Computational complexity Each step of Algorithm 2, either the extension or the update, requires an $\mathcal{O}(N)$ running time. Therefore the total complexity is $\mathcal{O}(N)$ which is the same as for the standard K -nearest neighbor method. The standard K -nearest neighbor algorithm computes only the first step. Thus our algorithm only adds a second update step, thereby it is more efficient than the message passing method (Kim and Choi 2007) that requires more visits of the objects and also builds a complete minimum spanning tree in advance on $\mathcal{G}(\mathbf{O}, \mathbf{D})$.

4.2 Minimax K -NN search, the Prim's algorithm and computational optimality

According to Theorem 1, computing pairwise Minimax distances requires first computing a minimum spanning tree over the underlying graph. Thereby, here, we study the connection between Algorithm 2 (i.e., Minimax K -NN search) and minimum spanning trees. For this purpose, we first consider a general framework for constructing minimum spanning trees called *generalized greedy algorithm* (Gabow et al. 1986). Consider a forest (collection) of trees $\{T_p\}$. The distance between the two trees T_p and T_q is obtained by

$$\Delta T_{pq} = \min_{i \in T_p} \min_{j \in T_q} \mathbf{D}_{ij}. \quad (9)$$

The *nearest neighbor tree* of T_p , i.e. T_p^* , is obtained by

$$T_p^* = \arg \min_{T_q} \Delta T_{pq}, \quad q \neq p. \quad (10)$$

Then, the edge e_p^* represents the nearest tree from T_p . It can be shown that e_p^* will belong to a minimum spanning tree on the graph as otherwise it yields a contradiction (Gabow et al.

1986). This result provides a *generic* way to compute a minimum spanning tree. A greedy MST algorithm at each step, (1) picks two candidate (base) trees where one is the nearest neighbor of the other, and (2) combines them via their shortest distance (edge) to build a larger tree.

This analysis guarantees that Algorithm 2 yields a minimum spanning subtree on $\{v \cup \mathcal{N}_{\mathcal{P}_T}(v)\}$ which is additionally a subset of a larger minimum spanning tree on the whole graph $\mathcal{G}^+(\mathbf{O} \cup v, \mathbf{D} \cup d(v, \mathbf{O}))$. In the context of *generalized greedy algorithm*, Algorithm 2 generates the MST via growing only the tree started from v and the other candidate trees are singleton nodes. A complete minimum spanning tree would be constructed, if the attachment continues for N steps, instead of K . This procedure, then, would yield a *Prim* minimum spanning tree (Prim 1957). This analysis reveals an interesting property of the *Prim's* algorithm:

The Prim's algorithm 'sorts' the nodes based on their Minimax distances from/to the initial (test) node v .

Computational optimality This analysis also reveals the '*computational optimality*' of our approach (Algorithm 2) compared to the alternative Minimax search methods, e.g. the method introduced in Kim and Choi (2007). Our algorithm *always* expands the tree which contains the initial node v , whereas the alternative methods might sometimes combine the trees that do not have any impact on the Minimax K -nearest neighbors of v . In particular, the method in Kim and Choi (2007) constructs a *complete* MST, whereas we build a *partial* MST with only and exactly K edges. Any new node that we add to the partial MST, belongs to the K Minimax nearest neighbors of v , i.e., we do not investigate/add any unnecessary node to that.

4.3 One-to-all Minimax distances and minimum spanning trees

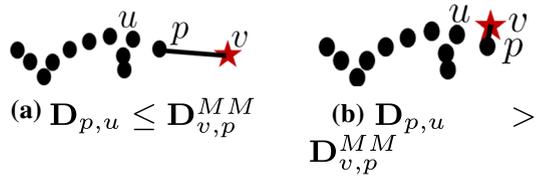
A straightforward generalization of Algorithm 2 gives sorted *one-to-all* Minimax distances from the target (test) v to the all other nodes. We only need to run the steps (extension and update) for N times instead of K . In Theorem 1, we observed that computing *all-pair* Minimax distances is consistent/equivalent with computing a minimum spanning tree, such that a MST provide all the necessary and sufficient edge weights for this purpose. Then later we have proposed an efficient algorithm for the *one-to-all* problem that also yields computing a (partial) minimum spanning tree (according to Prim's algorithm). Therefore, here we study the following question:

Is it necessary to build a minimum spanning tree to compute the 'one-to-all' Minimax distances, similar to 'all-pair' Minimax distances?

We reconsider the proof of Theorem 5 which has led us to Algorithm 2. The proof investigates two cases: (1) $\mathbf{D}_{p,u} \leq \mathbf{D}_{v,p}^{MM}$, (2) $\mathbf{D}_{p,u} > \mathbf{D}_{v,p}^{MM}$. We investigate the first case in more detail.

Let us look at the largest Minimax distance selected up to step T , i.e., the edge with maximal weight whose both sides occur inside $\{v \cup \mathcal{N}_{\mathcal{P}_T}(v)\}$ and its weight represents a Minimax distance. We call this edge e_T^{max} . Among different external neighbors of $\{v \cup \mathcal{N}_{\mathcal{P}_T}(v)\}$, any node $u' \in \{\mathbf{O} \setminus \mathcal{N}_{\mathcal{P}_T}(v)\}$ whose minimal distance to $\{v \cup \mathcal{N}_{\mathcal{P}_T}(v)\}$ does not exceed the weight of e_T^{max} can be selected, even if it does not belong to a minimum spanning tree over \mathcal{G}^+ . Because, anyway the Minimax distance will be still the weight of e_T^{max} .

Fig. 3 Computing Minimax nearest neighbors (of the red object) does not necessarily need to agree with a complete MST on \mathcal{G}^+ , as any remaining node can be replaced with u



In general, any node $u' \in \{\mathbf{O} \setminus \mathcal{N}^{\mathcal{P}_T}(v)\}$, whose Minimax distance to a selected member p , i.e. $\mathbf{D}_{p,u}^{MM}$, is not larger than the weight of e_T^{max} , can be selected as the next nearest neighbor.⁶ This is concluded from the following property of Minimax distances:

$$\mathbf{D}_{v,u}^{MM} \leq \max(\mathbf{D}_{v,p}^{MM}, \mathbf{D}_{p,u}^{MM}), \tag{11}$$

where p is an arbitrary object (node). In this setting we then have,

1. $\mathbf{D}_{v,u}^{MM} \leq \max(\mathbf{D}_{v,p}^{MM}, \mathbf{D}_{p,u}^{MM})$, and
2. $\mathbf{D}_{v,p}^{MM} > \mathbf{D}_{p,u}^{MM}$.

Thus, we conclude $\mathbf{D}_{v,u}^{MM} \leq \max(\mathbf{D}_{v,p}^{MM})$, i.e. $\mathbf{D}_{v,u}^{MM} = \mathbf{D}_{v,p}^{MM}$.

An example is shown in Fig. 3, where K is fixed at 2. After computing the first nearest neighbor (i.e., p), the next one can be any of the remaining objects, as their Minimax distance to v is the weight of the edge connecting p to v (Fig. 3a) or p to u (Fig. 3b). Thereby, one could choose a node such that the *partial* MST on $\{v \cup \mathcal{N}^{\mathcal{P}}(v)\}$ is not a subset of any *complete* MST on the whole graph. Therefore, computing *one-to-all* Minimax distances might not necessarily require taking into account construction of a minimum spanning tree.

Thereby, such an analysis suggests an even more generic algorithm for computing *one-to-all* Minimax distances (including Minimax K -nearest neighbor search) which does not necessarily yield a MST on the entire graph. To expand $\mathcal{N}^{\mathcal{P}_T}(v)$, we add a new u' whose *Minimax distance* (not the direct dissimilarity) to $\{v \cup \mathcal{N}^{\mathcal{P}_T}(v)\}$ is minimal. Algorithm 2 is only one way, but an efficient way, that even sorts one-to-all Minimax distances.

4.4 Outlier detection with Minimax K -nearest neighbor search

While performing K -nearest neighbor search, the test objects might not necessarily comply with the structure in the train dataset, i.e. some of the test objects might be outliers or belong to other classes than those existing in the train dataset. More concretely, when computing the K nearest neighbors of v , we might realize that v is an outlier or irrelevant with respect to the underlying structure in \mathbf{O} . We study this case in more detail and propose an efficient algorithm which while computing Minimax K nearest neighbors, it also detects whether v could be an outlier. Notice that the ability to detect the outliers will be an additional feature of our algorithm, i.e., we do not aim at proposing the best outlier detection method, rather, the goal is to maximize the benefit from performing Minimax K -NN while still having an $\mathcal{O}(N)$ runtime.

⁶ We notice that in the second case, i.e. when $\mathbf{D}_{p,u} > \mathbf{D}_{v,p}^{MM}$, we need to add an unselected node whose Minimax to v (to one of the nodes in $\mathcal{N}^{\mathcal{P}_T}(v)$) is equal to $\mathbf{D}_{p,u}$. This situation might also yield adding an edge which does not necessarily belong to a MST over \mathcal{G}^+ . The argument is similar to the case when $\mathbf{D}_{p,u} \leq \mathbf{D}_{v,p}^{MM}$.

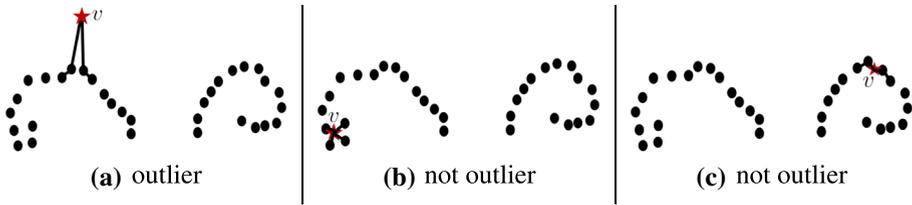


Fig. 4 The edges representing the Minimax distances can provide useful information for detecting outliers

Thereby, we follow our analysis by investigating another special aspect of Algorithm 2: as we extend $\mathcal{N}\mathcal{P}_T(v)$, the edge representing the Minimax distance between v and the new member u (i.e., the edge with the largest weight on the path indicating the Minimax distance between v and u) is always connected to v . For this purpose, we reconsider the proof of Theorem 5. When u is being added to $\mathcal{N}\mathcal{P}_T(v)$, two special cases might occur:

1. u is directly connected to v , i.e. $p = v$.
2. u is not directly connected to v (i.e. $p \neq v$), but its Minimax distance to v is *not* represented by $\mathbf{D}_{p,u}$, i.e., $\mathbf{D}_{p,u} < \mathbf{D}_{v,p}^{MM}$.

Figure 4a illustrates these situations, where $K = 4$. Among the nearest neighbors, two are directly connected to v (i.e., $p = v$), and the two others are connected to v via the early members of $\mathcal{N}\mathcal{P}_T(v)$, but the Minimax distance is still represented by the edges connected to v . In other words, there is no edge in graph $\mathcal{G}(\mathbf{O}, \mathbf{D})$ which represents a Minimax distance, although some edges of $\mathcal{G}(\mathbf{O}, \mathbf{D})$ are involved in Minimax paths. Thus, the type of connection between v and its neighbors is *different* from the type of the distances inside graph $\mathcal{G}(\mathbf{O}, \mathbf{D})$ (without test object v). Thereby, in this case we report v as an *outlier*. Notice that both of the above mentioned conditions should occur, i.e. if we have always $p = v$ (as shown in Fig. 4b where the nearest neighbors are always connected to v), then v might not be labeled as an *outlier*.

On the other hand, as shown in Fig. 4c, if some edges of $\mathcal{G}(\mathbf{O}, \mathbf{D})$ contribute in computing the Minimax K -nearest neighbors of v (in addition to the edges which meet v), then v has the same type of neighborhood as some other nodes in \mathcal{G} . Thereby, it is *not* labeled as an outlier.

This analysis suggests an algorithm for simultaneously computing Minimax K -nearest neighbors and detecting if v is an outlier (Algorithm 3). For this purpose, we use a vector called *updated*, which determines the type of nearest neighborhood of each external object to $\{v \cup \mathcal{N}\mathcal{P}_T(v)\}$, i.e.,

1. $updated_i = 0$, if v is the node representing the nearest neighbor of i in the set $\{v \cup \mathcal{N}\mathcal{P}_T(v)\}$.
2. $updated_i = 1$, otherwise.

At the beginning, $\mathcal{N}\mathcal{P}_T(v)$ is empty. Thereby, *updated* is initialized by a vector of zeros. At each update step, whenever we modify an element of the vector *dist*, we then set the corresponding index in *updated* to 1. At each step of Algorithm 2, a new edge is added to the set of selected edges whose weight is $\mathbf{D}_{p,u}$. As mentioned before, v is labeled as an outlier if, (1) some of the edges are directly connected to v , and some others indirectly, and

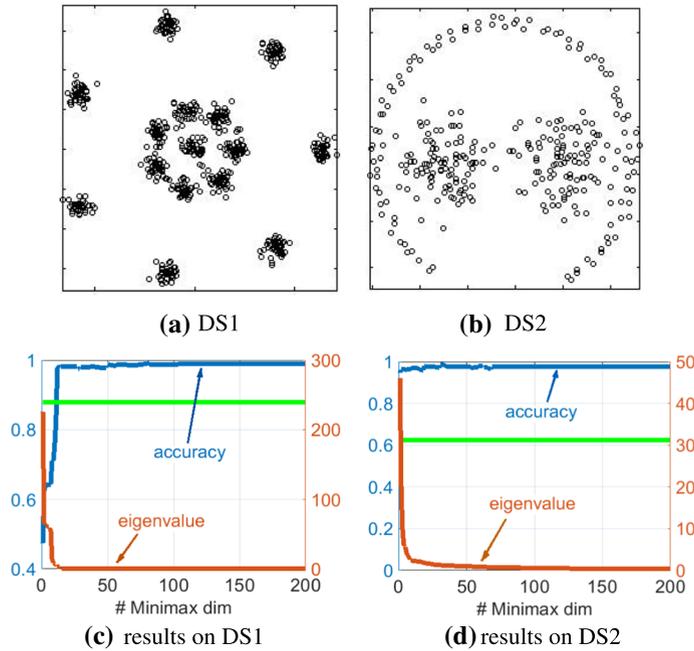


Fig. 5 Illustration of DS1, DS2 and the accuracy. Different classes of DS1 have a similar structure as shown in (a), but the classes are different in DS2 as illustrated in (b). The accuracy scores (shown for LogReg-MM) are stable w.r.t. the dimensionality of the Minimax vectors (c, d). The straight green line shows the accuracy for the base LogReg

(2) no indirect edge weight represents a Minimax distance, i.e. the minimum of the edge weights directly connected to v (stored in min_direct) is larger than the maximum of the edge weights not connected to v (stored in max_indrct). The later corresponds to the condition that $\mathbf{D}_{v,p}^{MM} > \mathbf{D}_{p,u}$. Thereby, whenever we pick the nearest neighbor of $\{v \cup \mathcal{N}\mathcal{P}_T(v)\}$ in the extension step, we then check the *updated* status of the new member:

1. If $updated_{min_ind} = 1$, we then update max_indrct by $\max(max_indrct, dist_{min_ind})$.
2. If $updated_{min_ind} = 0$, we then update min_direct by $\min(min_direct, dist_{min_ind})$.

Finally, if $min_direct > max_indrct$ and $max_indrct \neq -1$ (max_indrct is initialized by -1 ; this condition ensures that at least one indirect edge has been selected), then v is labeled as an outlier. Algorithm 3 describes the procedure in detail. Compared to Algorithm 2, we have added: (1) steps 11–15 for keeping the statistics about the direct and indirect edge weights (i.e., via updating max_indrct and min_direct), (2) an additional step in line 20 for updating the type of the connection of the external nodes to the set $\{v \cup \mathcal{N}\mathcal{P}_T(v)\}$, and (3) finally checking whether v is an outlier.⁷

Computational complexity The extension step is done in $\mathcal{O}(N)$ via a sequential search. Lines 11–15 are constant in time and the loop 17–22 requires an $\mathcal{O}(N)$ time. Thus, the total

⁷ Our goal is not to propose the best outlier detection algorithm, rather, we aim at adding a new feature (i.e., detecting potential outliers) to Minimax K -nearest neighbor search while we require the runtime to be still $\mathcal{O}(N)$. Thus, we improve the gain of performing a linear-time Minimax K -NN search query.

runtime is $\mathcal{O}(N)$ which is identical to standard K -nearest neighbor search over arbitrary graphs.

Algorithm 3 Outlier detection

Require: Graph $\mathcal{G}^+(\mathbf{O} \cup v, \mathbf{D} \cup d(v, \mathbf{O}))$ including a test object v .

Ensure: Minimax K nearest neighbors of v stored in $\mathcal{NP}(v)$, detecting whether v is an outlier.

```

1:  $\mathcal{NP}(v) = \square$ 
2: Initialize vector  $dist$  by the distance of  $v$  to each object in  $\mathbf{O}$ .
3: Initialize vector  $updated$  by 0.
4:  $max\_indrct = -1$ 
5:  $min\_direct = inf$ 
6: for  $T = 1$  to  $K$  do
7:   {extension}
8:    $min\_ind = \arg \min(dist)$ 
9:    $\mathcal{NP}(v).append(min\_ind)$ 
10:  {update}
11:  if  $updated_{min\_ind} = 1$  then
12:     $max\_indrct = \max(max\_indrct, dist_{min\_ind})$ 
13:  else
14:     $min\_direct = \min(min\_direct, dist_{min\_ind})$ 
15:  end if
16:   $dist_{min\_ind} = inf$ 
17:  for  $i \notin \mathcal{NP}(v)$  do
18:    if  $\mathbf{D}_{min\_ind,i} < dist_i$  then
19:       $dist_i = \mathbf{D}_{min\_ind,i}$ 
20:       $updated_i = 1$ 
21:    end if
22:  end for
23: end for
24:  $is\_outlier = 0$ 
25: if  $max\_indrct \neq -1$  and  $min\_direct > max\_indrct$  then
26:    $is\_outlier = 1$ 
27: end if
28: return  $\mathcal{NP}(v), is\_outlier$ 

```

5 Experiments

We experimentally study the performance of Minimax distance measures on a variety of synthetic and real-world datasets and illustrate how the use of Minimax distances as an intermediate step improves the results. In each dataset, each object is represented by a vector.⁸

5.1 Classification with Minimax distances

First, we study classification with Minimax distances. We use logistic regression (LogReg) and support vector machines (SVM) as the baseline methods and investigate how

⁸ When not mentioned, we compute the pairwise squared Euclidean distances between the vectors to construct the base distance matrix \mathbf{D} .

Table 1 Accuracy of different methods on synthetic datasets

Dataset	Standard				Minimax	
	SVM-lin	SVM-rbf	SVM-sig	LogReg	SVM-lin	LogReg
DS1	0.5749	0.9924	0.4260	0.8066	0.9917	0.9918
DS2	0.5654	0.9295	0.3294	0.6252	0.9950	0.9983

Minimax measure particularly improves the results when the classes in the dataset have different shapes and types (e.g. DS2)

Table 2 Accuracy scores of different methods on UCI datasets, when 60% of the data is used for training

Data	Standard				Minimax		dim.spec. Minimax	
	SVM-lin	SVM-rbf	SVM-sig	LogReg	SVM-lin	LogReg	SVM-lin	LogReg
UCI1	0.8709	0.8974	0.4468	0.8687	0.6187	0.6086	0.9211	0.9739
UCI2	0.9876	1.0000	0.5577	0.9872	0.9989	1.0000	0.8847	0.9827
UCI3	0.9988	0.5788	0.5349	0.9988	1.0000	1.0000	1.0000	1.0000
UCI4	0.5190	0.5533	0.4250	0.5107	0.5647	0.5747	0.5392	0.5389
UCI5	0.5924	0.6012	0.3371	0.6053	0.5971	0.6671	0.4918	0.6347
UCI6	0.6893	0.7230	0.7344	0.7426	0.7434	0.7377	0.7418	0.7352
UCI7	0.6058	0.7750	0.3596	0.5365	0.7038	0.7115	0.8635	0.8558
UCI8	0.8779	0.9300	0.6536	0.8621	0.9457	0.9450	0.8843	0.9336
UCI9	0.6917	0.7500	0.7500	0.6917	0.7750	0.7500	0.8333	0.8333
UCI10	0.7783	0.9318	0.1498	0.6933	0.9865	0.9870	0.9798	0.9830
UCI11	0.9287	0.9693	0.7635	0.8980	0.9994	0.9994	0.9906	0.9919
UCI12	0.7835	0.7233	0.3019	0.8612	0.5893	0.6757	0.6010	0.8592

Using Minimax vectors improves the results. The best results are bolded

performing these methods on the vectors induced from Minimax distances improves the accuracy of prediction. With SVM, we examine three different kernels: (1) linear (lin), (2) radial basis function (rbf), and (3) sigmoid (sig), and choose the best result. With Minimax distances, we only use the linear kernel, since we assume that Minimax distances must be able to capture the correct classes, such that they can be then discriminated via a linear separator.

Experiments with synthetic data

We first perform our experiments on two synthetic datasets, called: (1) DS1 (Chang and Yeung 2008), and (2) DS2 (Veenman et al. 2002), which are shown in Fig. 5. The goal is to demonstrate the superior ability of Minimax distances to capture the correct class-specific structures, particularly when the classes have different types (DS2), compared to kernel methods. Table 1 shows the accuracy scores for different methods. The standard SVM is performed with three different kernels (lin, rbf and sig), and the best choice which is the rbf kernel is shown. As mentioned, with Minimax distances, we only use the linear kernel. We observe that performing classification on Minimax vectors yields the best results, since it enables the method to better identify the correct classes. The datasets differ in the type and consistency of the classes. DS1 contains very similar classes which are Gaussian. But DS2

Table 3 Accuracy of different methods when only 10% of the data is used for training, where using Minimax distances often improves the results

Data	Standard				Minimax		dim.spec. Minimax	
	SVM-lin	SVM-rbf	SVM-sig	LogReg	SVM-lin	LogReg	SVM-lin	LogReg
UCI1	0.8459	0.8477	0.4566	0.8694	0.5114	0.6021	0.8270	0.7879
UCI2	0.9859	0.9843	0.5305	0.9870	0.9895	0.9916	0.6632	0.9060
UCI3	0.9607	0.5207	0.5069	0.9849	1.0000	1.0000	1.0000	1.0000
UCI4	0.4746	0.5078	0.4082	0.4871	0.4761	0.5206	0.4714	0.4655
UCI5	0.3417	0.4677	0.3385	0.3958	0.4365	0.4844	0.4100	0.5000
UCI6	0.7011	0.7275	0.7347	0.7311	0.7369	0.7362	0.7336	0.7176
UCI7	0.4258	0.4593	0.3352	0.4775	0.5186	0.5758	0.5958	0.6636
UCI8	0.8033	0.7692	0.6256	0.7943	0.9043	0.9097	0.8000	0.8786
UCI9	0.7045	0.7227	0.7136	0.7205	0.7159	0.6750	0.7386	0.7500
UCI10	0.6937	0.5732	0.1053	0.5129	0.8787	0.8775	0.8236	0.8481
UCI11	0.9161	0.9294	0.6203	0.9166	0.9983	0.9983	0.9855	0.9861
UCI12	0.5179	0.3772	0.3129	0.7200	0.3759	0.5315	0.3681	0.7638

The best results are bolded

consists of classes which differ in shape and type. Therefore, for DS1 we are able to find an optimal kernel (rbf, since the classes are Gaussian) with a *global* form and parametrization, which fits with the data and thus yields very good results. However, in the case of DS2, since classes have different shapes, then a single kernel *is not* able to capture correctly all of them. For this dataset, LogReg and SVM with Minimax vectors perform better, since they enable to adapt to the class-specific structures. Note that in the case of DS1, using Minimax vectors is equally good to using the optimal rbf kernel. Remember that, to Minimax vectors, we apply only SVM with a linear kernel. Figure 5c, d show the accuracy and eigenvalues w.r.t. different dimensionality of Minimax vectors. As mentioned earlier, the dimensions of Minimax vectors are sorted according to the respective eigenvalues, since a larger eigenvalue indicates a higher importance. By choosing only few dimensions, the accuracy attains its maximal value. We will elaborate in more detail on this further on.

We note that after computing the matrix of pairwise Minimax distances, one can in principle apply any of the embedding methods (e.g., PCA, SVD, etc) either to the original pairwise distance matrix or to the pairwise Minimax distance matrix. Therefore, our contribution is orthogonal to such embedding methods. For example on these synthetic datasets, when we apply PCA to the original distance matrix, we do not observe a significant improvement in the results (e.g., 0.7838 on DS1 and 0.6671 on DS2 with LogReg).

Classification of UCI data

We then perform real-world experiments on twelve datasets from different domains, selected from the UCI repository (Dua and Graff 2019):⁹

- (1) *Balance Scale*: contains 625 observations modeling 3 types of psychological experiments.

⁹ The specifications of the datasets can be found at <https://archive.ics.uci.edu/ml/index.php>.

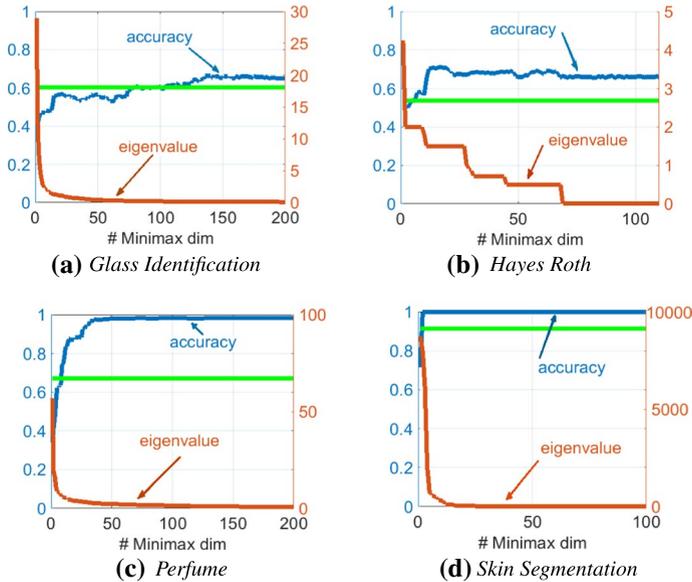


Fig. 6 Accuracy score of LogReg-MM applied to different datasets when we choose different number of dimensions of Minimax vectors. The straight green lines show the base LogReg result

- (2) *Banknote Authentication*: includes 1372 images taken from genuine and forged banknote-like specimens (number of classes is 2).
- (3) *Cloud*: consists of 1024 10-dimensional vectors, each dimension representing a specific parameter.
- (4) *Contraceptive Method*: contains information of 1473 women, where the three classes are about the pregnancy status.
- (5) *Glass Identification*: contains 6 types (classes) of glass w.r.t the oxide content. The number of instances is 214.
- (6) *Haberman Survival*: contains the survival of 306 patients who had surgery for breast cancer. The number of classes is 2.
- (7) *Hayes Roth*: is about a study on human subjects which contains 160 instances and 3 classes.
- (8) *Ionosphere*: includes 351 34-dimensional instances collected from radars and organized into 2 classes.
- (9) *Lung Cancer*: describes 3 types of pathological lung cancer, including 32 instances each with 56 dimensions.
- (10) *Perfume*: consists of odors of 20 different perfumes (classes), where the data is collected via OMX-GR sensor. There are in total 560 measurements.
- (11) *Skin Segmentation*: the original dataset contains 245,057 instances generated using skin textures from face images of different people. However, to make the classification task more difficult, we pick only the first 1000 instances of each class (to decrease the number of objects per class). The target variable is skin or non-skin sample, i.e. the number of classes is 2.
- (12) *User Knowledge*: describes 403 students' knowledge level (4 classes) about the subject of Electrical DC Machines.

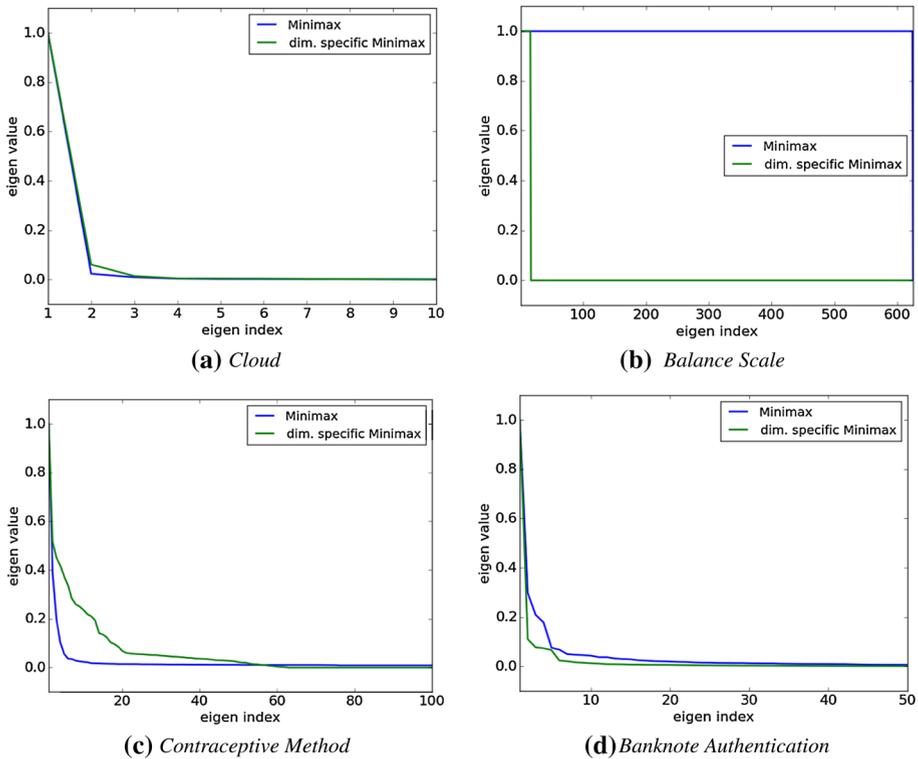


Fig. 7 Eigenspectra for different example datasets. A better Minimax variant often yields a sharper eigenspectrum

Table 4 Runtimes (in s) of FW-MM and MST-MM for computing pairwise Minimax distances on different real-world datasets

Dataset	<i>Glass Iden.</i>	<i>Haberman Surv.</i>	<i>Hayes Roth</i>	<i>Perfume</i>	<i>Cloud</i>	<i>Skin Seg.</i>
Size	160	214	306	560	1024	2000
FW-MM	0.030	0.067	0.230	1.399	18.54	127.12
MST-MM	0.019	0.024	0.080	0.150	0.239	0.540

MST-MM performs significantly faster

We call these datasets respectively UC1, UC12, ..., UC112.

Accuracy scores Table 2 shows the results for different methods applied to the datasets, when 60% of the objects are used for training. We have repeated the random split of the data for 20 times and report the average results. The scores and the ranking of different methods are very consistent among different splits, such that the standard deviations are low (i.e., maximally about 0.016). We observe that often performing the classification methods on Minimax vectors improves the classification accuracy. In only very few cases the standard setup outperforms (slightly). In the rest, either the Minimax vectors or the dimension-specific variant of Minimax vectors yield a better performance. In particular,

the Minimax variant is more appropriate for low dimensional data, whereas the dimension-specific Minimax variant outperforms on high-dimensional data. We elaborate more on the choice between them in the ‘model order selection’ section.

Table 3 shows the results when only 10% of the objects are used for training. For this setting, we observe a consistent performance with the setting where 60% of data is used for training, i.e. the use of Minimax vectors (either the standard variant or the dimension-specific one) improves the accuracy scores. In this setting, only with UCI12 the non-Minimax method performs better. However, even on this dataset, the dimension-specific Minimax performs very closely to the best result. Such a consistency is observed for other ratios of train and test sets too.

Model order selection Choosing the appropriate number of dimensions for Minimax vectors (i.e. their dimensionality) constitutes a model order selection problem. We study in detail how the dimensionality of the Minimax vectors affects the results. Figure 6 shows the accuracy scores for Minimax-LogReg applied to four of the datasets w.r.t. different number of dimensions (the other datasets behave similarly). The dimensions are ordered according to their importance (the value of respective eigenvalue). Choosing a very small number of dimensions might be insufficient since we might lose some informative dimensions, which yields underfitting. By increasing the dimensionality, the method extracts more sufficient information from the data, thus the accuracy improves. We note that this phase occurs for a very wide range of choices of dimensions. However, by increasing the number of dimensions even more, we might include non-informative or noisy features (with very small eigenvalues), where then the accuracy stays constant or decreases slightly, due to overfitting. However, an interesting advantage of this approach is that the overfitting dimensions (if there exists any) have a very small eigenvalue, thus their impact is negligible. This analysis leads to a simple and effective model order selection criteria: Fix a small threshold and pick the dimensions whose respective eigenvalues are larger than the threshold. The exact value of this threshold may not play a critical role or it can be estimated via a validation set in a supervised learning setting.

Model selection According to our experimental results, very often, either the standard Minimax classification or the dimension-specific variant outperform the baselines. In a supervised learning setting, the correct choice between these two Minimax variants (i.e., the model selection task) can be made via measuring the prediction ability, e.g. the accuracy score. However, this approach might not be applicable for unsupervised learning, e.g. clustering, where the ground truth is not given. We investigate a simple heuristics which can be employed in any arbitrary setting. When we perform eigen decomposition of the centered matrix (\mathbf{W}^{MM} or \mathbf{W}^{cMM}), we normalize the eigenvalues by the largest of them such that the largest eigenvalue equals 1. Then, the variant that yields a shaper decay in eigenvalues is supposed to be potentially a better model. A sharper decay of the eigenvalues indicates a tighter confinement to a low dimensional subspace, i.e. lower complexity in data representation. This possibly yields a better learning, and thereby a higher accuracy score.

To analyze this heuristics, we observe that in our experiments on *Cloud*, *Glass Identification*, *Haberman Survival*, *Perfume* and *Skin Segmentation* either the accuracy scores are not significantly distinguishable or there is no consistent ranking of the two different variants (Minimax or dimension-specific Minimax) such that it is difficult to decide which one is better. For these datasets the eigenspectra are very similar as shown for example for *Cloud* dataset in Fig. 7a. Among the remaining datasets, the heuristics works properly on *Balance Scale*, *Contraceptive Method*, *Hayes Roth*, *Lung Cancer* and *User Knowledge*

Table 5 Performance of GMM on original vectors (*base*), on the standard Minimax vectors (*Minimax*), and on the dimension-specific Minimax vectors at subspaces (*DimMM*) in the clustering task

Method	Dataset 1		Dataset 2		Dataset 3	
	Rand score	Mutual info.	Rand score	Mutual info.	Rand score	Mutual info.
Base	0.2539	0.5914	0.3191	0.6478	0.4268	0.7195
Minimax	0.3033	0.6263	0.3995	0.6943	0.4611	0.7407
DimMM050	0.3605(0.016)	0.6489(0.015)	0.4306(0.016)	0.7052(0.018)	0.4556(0.014)	0.7294(0.018)
DimMM075	0.3767(0.013)	0.6431(0.019)	0.4238(0.014)	0.7003(0.019)	0.5254(0.015)	0.7587(0.018)
DimMM100	0.3728(0.014)	0.6451(0.018)	0.4126(0.011)	0.6996(0.012)	0.5002(0.013)	0.7577(0.016)
DimMM125	0.3538(0.015)	0.6470(0.015)	0.4167(0.017)	0.6985(0.019)	0.4731(0.015)	0.7448(0.019)
DimMM150	0.3431(0.012)	0.6418(0.016)	0.3932(0.013)	0.6960(0.014)	0.4917(0.012)	0.7530(0.014)
DimMM175	0.3484(0.017)	0.6459(0.016)	0.3766(0.018)	0.6896(0.019)	0.4512(0.018)	0.7375(0.017)
DimMM200	0.3351(0.015)	0.6374(0.018)	0.4023(0.011)	0.6948(0.011)	0.4607(0.014)	0.7429(0.016)
DimMM250	0.3322(0.012)	0.6413(0.016)	0.3809(0.013)	0.6866(0.0114)	0.5034(0.011)	0.7530(0.012)
DimMM300	0.3243(0.012)	0.6376(0.015)	0.3942(0.014)	0.6954(0.017)	0.4800(0.016)	0.7469(0.015)
DimMM350	0.3332(0.016)	0.6419(0.018)	0.3831(0.012)	0.6798(0.013)	0.4660(0.011)	0.7452(0.011)
DimMM400	0.3385(0.013)	0.6406(0.013)	0.4089(0.014)	0.6944(0.015)	0.4680(0.013)	0.7429(0.012)
DimMM500	0.3272(0.015)	0.6378(0.0179)	0.4008(0.016)	0.6929(0.019)	0.4653(0.016)	0.7470(0.017)

The number at the front of *DimMM* shows the dimensionality of the subspace used to compute the dimension-specific Minimax distances. Different methods are compared w.r.t. adjusted Rand score and adjusted Mutual Information criteria. We repeat the *DimMM* variant 10 times and report the mean(std) results. On all datasets, *DimMM* yields higher scores of the evaluation criteria

datasets, where two sample eigenspectra are shown in Fig. 7b, c. Finally, the heuristics fails only on the two *Banknote Authentication* and *Ionosphere* datasets. However, for these datasets the accuracy scores as well as the decays are rather similar, as shown for example for *Banknote Authentication* in Fig. 7d.

Efficiency As a side study, we also investigate the runtimes of computing pairwise Minimax distance via either based on computing an MST (MST-MM) or the Floyd–Warshall algorithm (FW-MM). The results (in s) have been shown in Table 4 for some of the datasets. We observe that MST-MM yields a significantly faster approach to computer the pairwise Minimax distances. The difference is even more obvious when the dataset is larger. We note that the both methods yield the same Minimax distances. We observe consistent results on the other datasets.

5.2 Experiments on clustering of document scanings

We study the impact of performing Minimax distances on clustering of scanings of the documents collocated by a large document processing corporation. The dataset contains the vectors of 675 documents each represented in a 4096 dimensional space. The vectors contain metadata, the document layout, tags used in the document, the images in the document, the text in the document, the author information, the mathematical expressions and other structural information. This dataset contains 56 clusters some of which have only one single document. We call this dataset *dataset 1*. Then, by removing the clusters with only one or two documents, we obtain *dataset 2* which includes 634 documents and 34 clusters.

Table 6 Runtime (in s) of computing L^+ for different 20 newsgroup datasets

Dataset	COMP	REC	SCI	TALK
Runtime	20.67	12.93	11.84	6.55

L^+ is considerably slower

Finally, we obtain *dataset 3* which contains the clusters that have at least 5 documents. This datasets consists of 592 documents and 21 clusters. We compute the pairwise distances of pairs of documents according to squared Euclidean distance. The goal is to demonstrate the applicability of Gaussian Mixture Models (GMM) to Minimax distances, as well as to show if the use of Minimax distances improves the results.

We apply dimension-specific Minimax distances to different subspaces of the original data. For this purpose, we define the parameter b to be the number of dimensions of the subspace. Then, we randomly partition the original features (dimensions) such that the dimensionality of each subset equals b (the last subset might have less dimensions than b). Note that different subsets have the same number of documents which equals to the number of documents in the original dataset. For each subset, we compute the pairwise Minimax distances and obtain multiple Minimax matrices. Thus, we apply the collective Minimax embedding to compute an embedding in a new space. Finally, we apply GMM to the Minimax vectors and compare with the results of GMM on original (base) vectors and GMM on standard Minimax vectors.

We use adjusted Rand score (Hubert and Arabie 1985) and adjusted Mutual Information (Vinh et al. 2010) to evaluate the performance of different methods. Rand score computes the similarity between the estimated and the true clusterings. Mutual Information measures the mutual information between the two solutions. Note that we compute the adjusted version of these criteria, such that they yield zero for random clustering. Table 5 shows the results on the three datasets. The number at the front of different dimension-specific (DimMM) variants indicates b . We repeat the DimMM variant 10 times and report the mean(std) results. We observe: (1) computing Minimax distances enables GMM to better capture the underlying structures, thus yields improving the results. (2) The dimension-specific variant improves the clusters even further via extracting appropriate structures in different subspaces. (3) However, the first improvement, i.e., the use or not use of Minimax distances, has a more significant impact than the choice between the standard or the dimension-specific Minimax variants. Finally, we note that DimMM is equivalent to the standard Minimax if $b = 4096$.

5.3 Minimax K -NN classification

We then examine the performance of our algorithm for K nearest neighbor search with Minimax distances. We perform our experiments on seven datasets, four selected from 20 newsgroup collection and the others come from image and plant specification.

- (1) COMP: a subset of 20 newsgroup contains 2936 documents around *computers*: ‘comp.graphics’, ‘comp.os.ms-windows.misc’, ‘comp.sys.ibm.pc.hardware’, ‘comp.sys.mac.hardware’, ‘comp.windows.x’.
- (2) REC: a subset of 20 newsgroup with 2389 documents on *sports*: ‘rec.autos’, ‘rec.motorcycles’, ‘rec.sport.baseball’, ‘rec.sport.hockey’.
- (3) SCI: a subset of 20 newsgroup having 2373 documents about *science*: ‘sci.crypt’, ‘sci.electronics’, ‘sci.med’, ‘sci.space’.

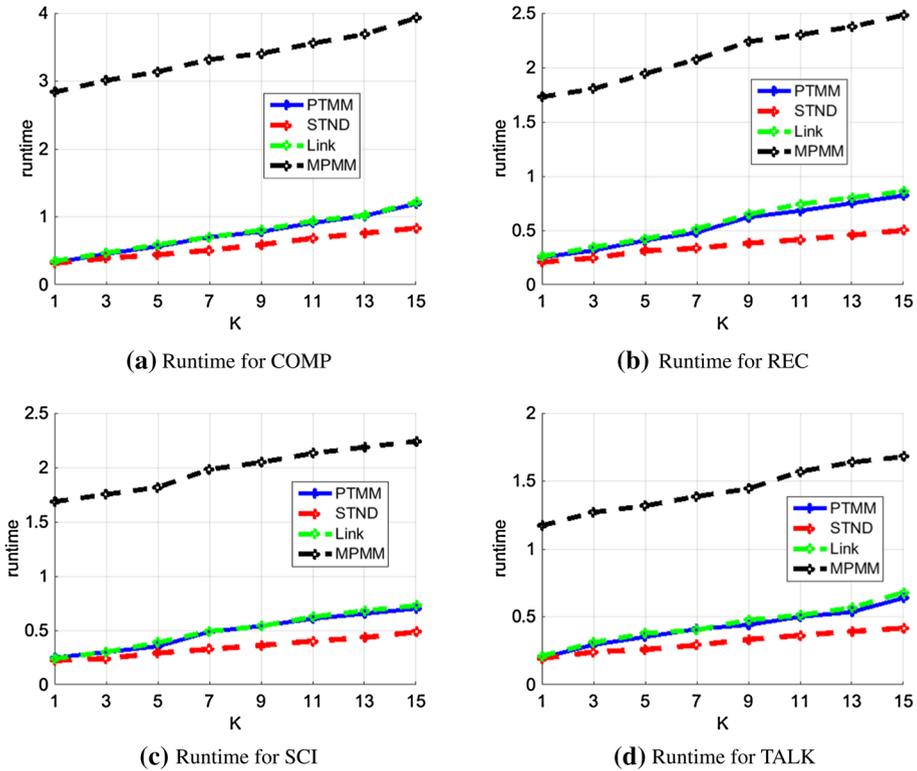


Fig. 8 Comparison of runtime (in s) of different methods on 20 news group datasets. PTMM is significantly faster than MPMM and is only slightly slower than STND. Figure 10a, b show the accuracy scores, where Minimax K -NN gives the best results

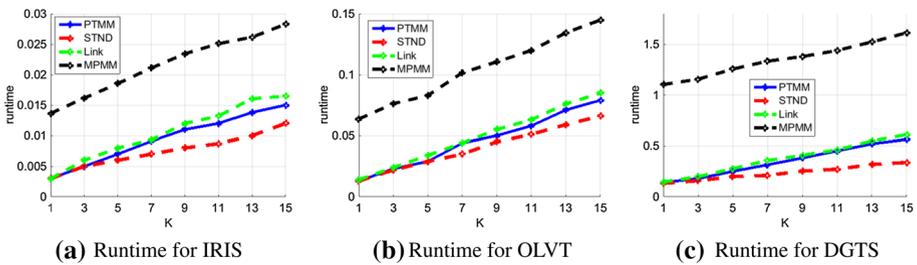


Fig. 9 Comparison of runtime (in s) of different methods on non-textual datasets. PTMM runs much faster than MPMM

- (4) TALK: a subset of 20 newsgroup with 1952 documents related to *talk*: ‘talk.politics.guns’, ‘talk.politics.mideast’, ‘talk.politics.misc’, ‘talk.religion.misc’.
- (5) IRIS: a common dataset with 150 samples from three species ‘setosa’, ‘virginica’ and ‘versicolor’.

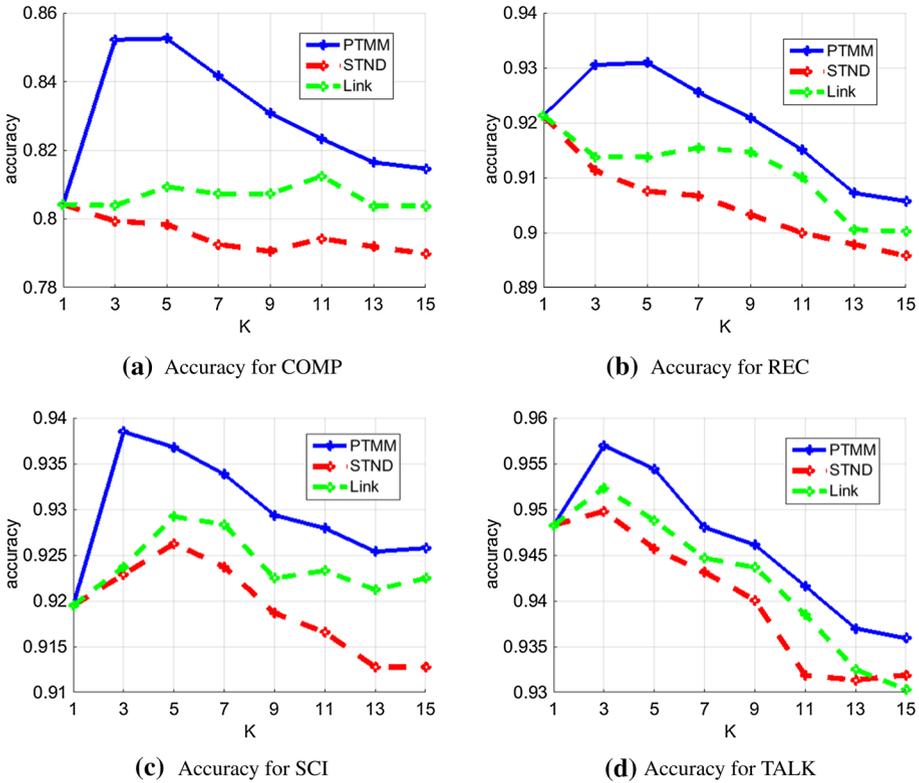


Fig. 10 Accuracy scores of different K -NN methods on 20 new group datasets, where Minimax K -NN gives the best results

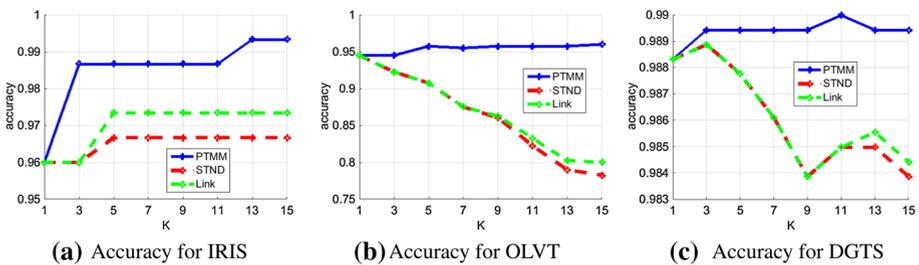


Fig. 11 Comparison of accuracy of different methods on IRIS, OLVT and DGTS datasets. Using Minimax distance measure improves the results

- (6) OLVT: the Olivetti faces dataset from AT&T which contains pictures from 40 individuals from each 10 pictures. The dimensionality is 4096.
- (7) DGTS: images of 10 digits each with 64 dimensions (1797 digits in total).

Table 7 Runtime (in s) of different methods on two-moon datasets

Size	5000	10,000	15,000	20,000
PTMM	1.03	4.42	10.96	19.187
MPMM	2.43	13.15	34.72	58.48
L^+	107.61	1105.13	2855.1	6603.16

PTMM performs significantly faster. PTMM and MPMM produce the same results

These datasets are well-known and publicly accessible, e.g. via *sklearn.datasets*.¹⁰ For the 20 newsgroup datasets, we obtain the vectors after eliminating the stop words. *Cosine* similarity is known to be a better choice for textual data. Thus, for each pair of document vectors \mathbf{x}_i and \mathbf{x}_j , we compute their *cosine* similarity and then their dissimilarity by $1 - \text{cosine}(\mathbf{x}_i, \mathbf{x}_j)$ to construct the pairwise dissimilarity matrix \mathbf{D} . For non-textual datasets, we directly compute the pairwise squared Euclidean distances. Thereby, we construct graph $\mathcal{G}(\mathbf{O}, \mathbf{D})$. Note that the method in Kim and Choi (2013) is not applicable to cosine-based dissimilarities. We compare our algorithm (called PTMM) against the following methods.

1. STND: Standard K -NN with the basic dissimilarities \mathbf{D} .
2. MPMM: Minimax K -NN proposed in Kim and Choi (2007).
3. Link: Link-based K -NN based on shortest distance algorithm (Dijkstra 1959).

L^+ (Fouss et al. 2012; Yen et al. 2008) is a link-based method computed based on the pseudo inverse of Laplacian matrix. However, this method has two main deficiencies: (1) it is very slow; as mentioned earlier, its runtime is $\mathcal{O}(N^3)$ (see Table 6 for runtime results on 20 newsgroup datasets), and (2) in the way that it is used (e.g. Yen et al. 2008; Kim and Choi 2013) L^+ is computed for all data, including training and test objects, i.e. the test objects are not completely separated from training. The correct approach would compute the pseudo inverse of the Laplacian only for the training dataset and then would extend it to contain the out-of-sample object.

We perform the experiments and compute accuracy in *leave-one-out* manner, i.e. the i -th object is left out and its label is predicted via the other objects. In the final voting, the nearest neighbors are weighted according to inverse of their dissimilarities with the test object. We calculate the accuracy by averaging over all N objects. M -fold cross validation could be used instead of the *leave-one-out* approach. However, the respective ranking of the runtimes of different methods is independent of the way the accuracy is computed.

Figure 8 illustrates the runtime of different measures on 20 news group datasets. Our algorithm for computing the Minimax distances (i.e. PTMM) is significantly faster than MPMM. PTMM is only slightly slower than the standard method and is even slightly faster than Link. The main source of computational inefficiency of MPMM comes from requiring a MST which needs an $\mathcal{O}(N^2)$ runtime. PTMM theoretically has the same complexity as STND (i.e. $\mathcal{O}(N)$) and in practice performs only slightly slower which is due to one

¹⁰ These datasets can be read and extracted via Python interfaces as described in <http://scikit-learn.org/stable/modules/classes.html#module-sklearn.datasets>.

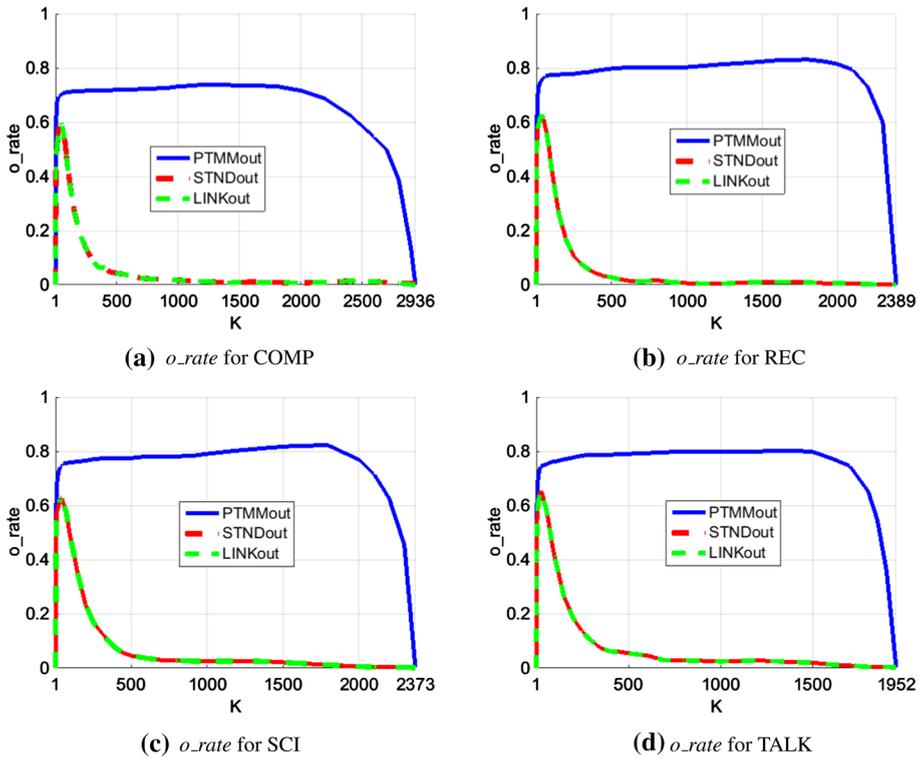


Fig. 12 *o_rate* score for different text datasets when PTMMout, STNDout and LINKout are used to predict outliers. PTMMout yields significantly a higher and more stable *o_rate*. Figure 13a, b show *rate_1* and *rate_2* for COMP

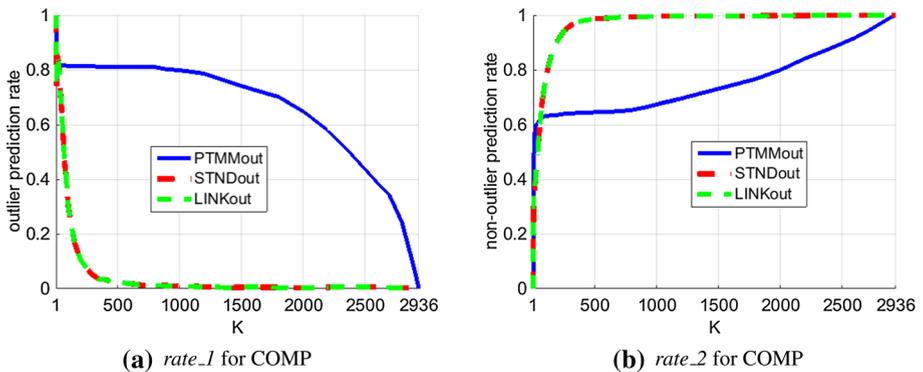


Fig. 13 *rate_1* and *rate_2* for COMP. Only PTMMout prevents sharp changes

additional update step. As mentioned L^+ is considerably slower, even compared to MPM (see Table 6).

To demonstrate the generality of our algorithm, we have performed experiments on datasets from other domains. Figure 9 shows the runtime of different methods on non-textual

datasets (IRIS, OLVT, and DGTS). We observe a consistent performance behavior: PTMM runs only slightly slower than STND and much faster than MPMM. Link is slightly slower than PTMM.

The usefulness of Minimax distances for K -NN classification has been investigated in previous studies, e.g. in Kim and Choi (2007) and Kim and Choi (2013). However, for the proof of concept, in Figs. 10 and 11, we illustrate the accuracy of different methods on the textual and non-textual datasets. Note that PTMM and MPMM yield the same accuracy scores. They differ only in the way they compute the Minimax nearest neighbors of v , thereby we have compared them only w.r.t runtime. Consistent to the previous study, we observe that Minimax measure is always the best choice, i.e. it yields the highest accuracy. We emphasize that this is the first study that demonstrates the effectiveness and performance of Minimax K -nearest neighbor classification beyond Euclidean distances (i.e. *cosine* similarity, a measure which is known to be more appropriate for textual data). As mentioned earlier, the previous studies (Kim and Choi 2007, 2013) always build the Minimax distances on squared Euclidean measures.

Scalability Finally, we study the scalability of different methods on larger datasets from two-moons data. As shown in Kim and Choi (2013) only L^+ and Minimax measures perform correctly on this data which has a non-linear geometric structure. Table 7 compares the runtime of these two methods for $K = 5$. PTMM performs significantly faster than the alternatives, e.g. it is three times faster than MPMM and 250 times faster than L^+ for 10,000 objects.

5.4 Outlier detection along with Minimax K -NN search

At the end, we experimentally investigate the ability of Algorithm 3 for detecting outliers while performing K -nearest neighbor search. We run our experiments on the textual datasets, i.e., on COMP, REC, SCI, and TALK. We add the unused documents of 20 news group collection to each dataset to play the role of outliers or the mismatching data source. This *ourlier dataset* consists of the 1664 documents of ‘alt.atheism’, ‘misc.forsale’, ‘soc.religion.christian’ categories. For each document, either from train or from outlier set, we compute its K nearest neighbors only from the documents of the train dataset and check whether it is predicted as an outlier or not.

We compare our algorithm (called here PTMMout), against the two others: STND and Link. These two algorithms do not directly provide outlier prediction. Thereby, in order to perform a fair comparison, we compute a Prim minimum spanning tree over the set containing v and the objects selected by the K -NN search, where v is the initial node. Then we check the same conditions to determine whether v is determined as an outlier. We call these two variants respectively STNDout and LINKout.

We obtain the rate of correct outlier prediction for the outlier dataset, called *rate_1*, by the ratio of the number of correct outlier reports to the size of this set. However, an algorithm might report an artificially high outlier number. Thereby, similar to *precision-recall* trade-off, we investigate the train dataset too and compute the ratio of the number of *non-outlier* reports to the total size of train dataset and call it *rate_2*. Finally, similar to F-measure, we compute the harmonic average of *rate_1* and *rate_2* and call it *o_rate*, i.e.

$$o_rate = 2 \cdot \frac{rate_1 \cdot rate_2}{rate_1 + rate_2}. \quad (12)$$

Figure 12 shows the o_rate scores for different datasets and for different values of K . For all datasets, we observe that PTMMout yields significantly a higher o_rate compared to the other algorithms. Moreover, the o_rate of PTMMout is considerably stable and constant over different K , i.e. the choice of optimal K is much less critical. For STNDout and LINKout, o_rate is maximal for only a fairly small range of K (e.g. between 30 and 50) and then it sharply drops down to zero. Figure 13 illustrates an in-depth analysis of $rate_1$ and $rate_2$ for the COMP dataset. Generally, as we increase K , $rate_1$ decreases but $rate_2$ increases. The reason is that as we involve more neighbors in $\mathcal{N}\mathcal{P}(v)$, then the chance of having an edge whose both sides are inside $\mathcal{N}\mathcal{P}(v)$ and its weight represents a Minimax distance increases. This case might particularly occur when $\mathcal{N}\mathcal{P}(v)$ contains objects from more than one class. Thereby, the probability that an object is labeled as a *non-outlier* increases. However, this transition is very sharp for STNDout and LINKout, such that for a relatively small K they label every object (either from train set or from outlier set) as a *non-outlier*. Moreover, as mentioned before, even for small K , STNDout and LINKout yield significantly lower o_rate than PTMMout.

6 Conclusion

We developed a framework to apply Minimax distances to any learning algorithm that works on numerical data, which takes into account both generality and efficiency. We studied both computing the pairwise Minimax distances for all pairs of objects and as well as computing the Minimax distances of all the objects to/from a fixed (test) object. For the *all-pair* case, we first employed the equivalence of Minimax distances over a graph and over a minimum spanning tree constructed on that, and developed an efficient algorithm to compute the pairwise Minimax distances on a tree. Then, we studied computing an embedding of the objects into a new vector space, such that their pairwise Minimax distance in the original data space equals to their squared Euclidean distance in the new space. We then extended our approach to the cases wherein there are multiple pairwise Minimax matrices. In the following, we studied computing Minimax distances from a fixed (test) object which can be used for instance in K -nearest neighbor search. Moreover, we investigated in detail the edges selected by the Minimax distances and thereby augmented the Minimax K -nearest neighbor search with the ability of detecting outlier objects. For each setting, we demonstrated the effectiveness of our framework via several experimental studies.

Acknowledgements Open access funding provided by Chalmers University of Technology. This work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation (No. 37200022). Parts of this work have been done at Xerox Research.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Aho, A. V., & Hopcroft, J. E. (1974). *The design and analysis of computer algorithms* (1st ed.). Boston, MA: Addison-Wesley Longman Publishing Co., Inc.
- Chang, H., & Yeung, D.-Y. (2008). Robust path-based spectral clustering. *Pattern Recognition*, *41*(1), 191–203.
- Chebotaev, P. (2011). A class of graph-geodetic distances generalizing the shortest-path and the resistance distances. *Discrete Applied Mathematics*, *159*(5), 295–302.
- Chehreghani, M. H. (2017). Efficient computation of pairwise minimax distance measures. In *2017 IEEE international conference on data mining, ICDM* (pp. 799–804). IEEE Computer Society.
- Chehreghani, M. H. (2020). Hierarchical correlation clustering and tree preserving embedding. CoRR, abs/2002.07756.
- Chehreghani, M. H. (2016). Adaptive trajectory analysis of replicator dynamics for data clustering. *Machine Learning*, *104*(2–3), 271–289.
- Cormen, T. H., Stein, C., Rivest, R. L., & Leiserson, C. E. (2001). *Introduction to algorithms* (2nd ed.). New York: McGraw-Hill Higher Education.
- Dijkstra, E. W. (1959). A note on two problems in connexion with graphs. *Numerische Mathematik*, *1*, 269–271.
- Dua, D., & Graff, C. (2019). *UCI machine learning repository* [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.
- Fiedler, M. (1998). Ultrametric sets in euclidean point spaces. *ELA. The Electronic Journal of Linear Algebra*, *3*, 23–30.
- Fischer, B., & Buhmann, J. M. (2003). Path-based clustering for grouping of smooth curves and texture segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *25*(4), 513–518.
- Fouss, F., Francoisse, K., Yen, L., Pirotte, A., & Saerens, M. (2012). An experimental investigation of kernels on graphs for collaborative recommendation and semisupervised classification. *Neural Networks*, *31*, 5372.
- Fouss, F., Pirotte, A., Renders, J.-M., & Saerens, M. (2007). Random-walk computation of similarities between nodes of a graph with application to collaborative recommendation. *IEEE Transactions on Knowledge and Data Engineering*, *19*(3), 355–369.
- Gabow, H. N., Galil, Z., Spencer, T., & Tarjan, R. E. (1986). Efficient algorithms for finding minimum spanning trees in undirected and directed graphs. *Combinatorica*, *6*(2), 109–122.
- Globerson, A., Chechik, G., Pereira, F., & Tishby, N. (2007). Euclidean embedding of co-occurrence data. *Journal of Machine Learning Research*, *8*, 2265–2295.
- Hofmann, T., Schölkopf, B., & Smola, A. J. (2008). Kernel methods in machine learning. *Annals of Statistics*, *36*(3), 1171–1220.
- Horn, R. A., & Johnson, C. R. (Eds.). (1990). *Matrix analysis*. Cambridge: Cambridge University Press.
- Hu, T. C. (1961). The maximum capacity route problem. *Operations Research*, *9*, 898–900.
- Hubert, L., & Arabie, P. (1985). Comparing partitions. *Journal of Classification*, *2*(1), 193–218.
- Khoshneshin, M., & Street, W. N. (2010). Collaborative filtering via euclidean embedding. In *Proceedings of the 2010 ACM conference on recommender systems, RecSys 2010, Barcelona, Spain, 26–30 September 2010* (pp. 87–94).
- Kim, K.-H., & Choi, S. (2007). Neighbor search with global geometry: A minimax message passing algorithm. In *ICML* (pp. 401–408).
- Kim, K.-H. & Choi, S. (2013). Walking on minimax paths for k-nn search. In *AAAI*.
- Kschischang, F. R., Frey, B. J., & Loeliger, H. A. (2006). Factor graphs and the sum-product algorithm. *IEEE Transactions on Information Theory*, *47*(2), 498–519.
- Luxburg, U. (2007). A tutorial on spectral clustering. *Statistics and Computing*, *17*(4), 395–416.
- Nadler, B., & Galun, M. (2007). Fundamental limitations of spectral clustering. *Advanced in Neural Information Processing Systems*, *19*, 1017–1024.
- Pavan, M., & Pelillo, M. (2007). Dominant sets and pairwise clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *29*(1), 167–172.
- Prim, R. C. (1957). Shortest connection networks and some generalizations. *The Bell Systems Technical Journal*, *36*(6), 1389–1401.
- Quarteroni, A., Sacco, R., & Saleri, F. (2007). Approximation of eigenvalues and eigenvectors. *Numerical Mathematics*, *37*, 183–244.
- Schoenberg, I. J. (1937). On certain metric spaces arising from euclidean spaces by a change of metric and their imbedding in hilbert space. *Annals of Mathematics*, *38*(4), 787–793.
- Shawe-Taylor, J., & Cristianini, N. (2004). *Kernel methods for pattern analysis*. Cambridge: Cambridge University Press.

- Tenenbaum, J. B., de Silva, V., & Langford, J. C. (2000). A global geometric framework for nonlinear dimensionality reduction. *Science*, *290*(5500), 2319.
- Torgerson, W. S. (1958). *Theory and methods of scaling*. Hoboken: Wiley.
- Varga, R. S., & Nabben R. (1993). On symmetric ultrametric matrices. *Numerical Linear Algebra*, 193–200.
- Veenman, C. J., Reinders, M. J. T., & Backer, E. (2002). A maximum variance cluster algorithm. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *24*, 1273–1280.
- Vinh, N. X., Epps, J., & Bailey, J. (2010). Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance. *Journal of Machine Learning Research*, *11*, 2837–2854.
- Weinberger, K. Q., & Saul, L. K. (2009). Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, *10*, 207–244.
- Yen, L., Saerens, M., Mantrach, A., & Shimbo, M. (2008). A family of dissimilarity measures between nodes generalizing both the shortest-path and the commute-time distances. In *KDD* (pp. 785–793).
- Young, G., & Householder, A. (1938). Discussion of a set of points in terms of their mutual distances. *Psychometrika*, *3*(1), 19–22.
- Zadeh, R., & Ben-David, S. (2009). A uniqueness theorem for clustering. In *UAI* (pp. 639–646).

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.