

On novelty detection for multi-class classification using non-linear metric learning

Samuel Rocha Silva^{a,*}, Thales Vieira^b, Dimas Martínez^c, Afonso Paiva^a

^a*ICMC, Universidade de São Paulo, São Carlos, Brazil*

^b*Instituto de Computação, Universidade Federal de Alagoas, Maceió, Brazil*

^c*Departamento de Matemática, Universidade Federal do Amazonas, Manaus, Brazil*

Abstract

Novelty detection is a binary task aimed at identifying whether a test sample is novel or unusual compared to a previously observed training set. A typical approach is to consider distance as a criterion to detect such novelties. However, most previous work does not focus on finding an optimum distance for each particular problem. In this paper, we propose to detect novelties by exploiting non-linear distances learned from multi-class training data. For this purpose, we adopt a kernelization technique jointly with the Large Margin Nearest Neighbor (LMNN) metric learning algorithm. The optimum distance tries to keep each known class's instances together while pushing instances from different known classes to remain reasonably distant. We propose a variant of the K-Nearest Neighbors (KNN) classifier that employs the learned distance to detect novelties. Besides, we use the learned distance to perform multi-class classification. We show quantitative and qualitative experiments conducted on synthetic and real data sets, revealing that the learned metrics are effective in improving novelty detection compared to other metrics. Our method also outperforms previous work regularly used for novelty detection.

Keywords: novelty detection, outlier detection, anomaly detection, metric learning, multi-class classification

*Corresponding author. Tel: +55-16-3373-9700

Email addresses: samuelrs@usp.br (Samuel Rocha Silva), thales@ic.ufal.br (Thales Vieira), dimas@ufam.edu.br (Dimas Martínez), apneto@icmc.usp.br (Afonso Paiva)

URL: <http://www.ic.ufal.br/professor/thales> (Thales Vieira)

1. Introduction

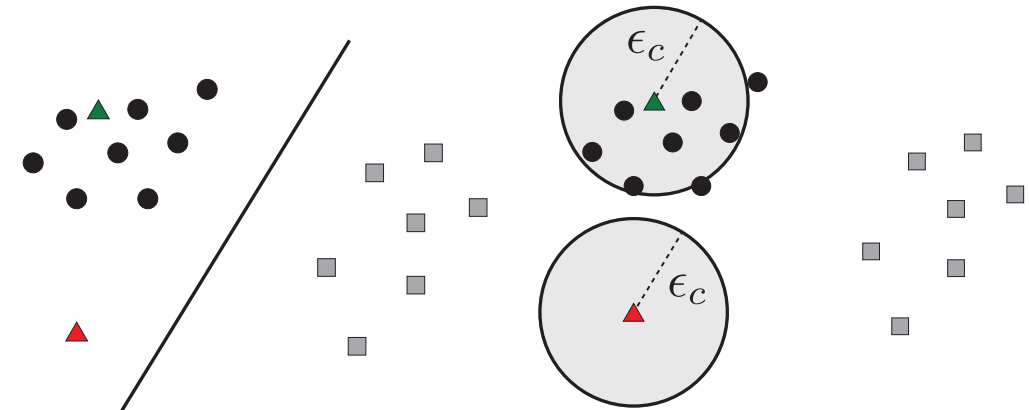
Novelty detection can be defined as the task of recognizing that test data differ in some respect from the data that are available during training (Pimentel et al., 2014). A large number of applications are covered by novelty detection methods including fraud detection (Jyothisna et al., 2011), medical diagnosis (Clifton et al., 2011), video surveillance (Diehl & Hampshire, 2002), mobile robotics (Sofman et al., 2011), among many others. It is worth mentioning that novelty detection is similar to anomaly and outlier detection and one-class classification, although it originated from different application domains. We refer the readers to the excellent surveys on novelty (Pimentel et al., 2014) and anomaly (Chandola et al., 2009) detection as well as one-class classification (Khan & Madden, 2010) for more information.

In general, novelty detection methods are applied when only data from a specific pattern is available, which is usually called the “normal” class, in contrast to novel data revealing a different pattern (“abnormal” class). In many real-world pattern recognition applications, however, many distinct patterns may be given during a training phase, and others may only appear over time (Faria et al., 2013).

As remarked by Bodesheim et al. (2013), it is not always possible to define a complete set of classes and acquire training instances from each of them. According to that authors, in multi-class novelty detection, one wants to detect whether a test sample is a novelty or if it belongs to one of many available classes, no matter to which class. For instance, training visual recognition systems require the trainer to provide images displaying examples of objects. However, the potential number of classes of real-world objects visible in an image is virtually infinite, and thus it is not possible to enumerate and collect training samples from all types of objects. Another critical problem where this situation occurs is gesture recognition. Many methods represent gestures as sequences of trainable classes of key poses, where only a few key poses are selected for training (Miranda et al., 2014b; Lv & Nevatia, 2007). However, during a gesture execution, some unimportant poses (novelties) will also occur in-between consecutive key poses. As a consequence, a key pose classifier must be capable of robustly identifying both: poses that are instances of a trained (normal) class; and unimportant poses that are not instances of any of the trained key pose classes, *i.e.* instances of abnormal classes.

36 It is worth emphasizing that the multi-class novelty detection problem
 37 is not well addressed by classifiers such as multi-class Support Vector Ma-
 38 chine (SVM) (Schölkopf & Smola, 2002), where the boundary of each class
 39 is estimated by only requiring that it properly separates the regions of the
 40 space representing the trained classes. Consequently, SVM may incorrectly
 41 classify an unimportant pose as an instance of a trained class, as shown in
 42 Fig. 1a. To overcome such issues, we aim to provide a multi-class classifica-
 43 tion solution that includes a novelty detection step to filter abnormal data.
 44

45 According to Pimentel et al. (2014), novelty detection methods can be
 46 classified as: (i) probabilistic, such as the Gaussian Mixture Models pro-
 47 posed in (Ilonen et al., 2006); (ii) reconstruction-based, in which we em-
 48 phasize the *Kernel PCA* (KPCA) approach (Hoffmann, 2007); (iii) domain-
 49 based, including the well-known one-class SVM (Schölkopf et al., 2001); (iv)



(a) SVM boundary: the hyperplane (here a line) that better separates both classes is chosen as the boundary between them. The green triangle is correctly classified as an instance of the class of circles. However, the red triangle is incorrectly classified as an instance of the same class of circles, although it is a novelty (or outlier).

(b) Distance-based novelty filter: by considering distances between an input instance and training instances, the red triangle can be recognized as an outlier and filtered, since none of the training instances are inside the circle of radius ϵ_c centered in the red triangle. In contrast, the green triangle can be recognized as an inlier of the black circles class.

Fig. 1: SVM classifier limitations in a simple binary classification problem and the advantages of employing a distance-based novelty filter: black circles and grey squares represent training instances of two classes, and triangles represent input instances to be classified.

50 information-theoretic (Keogh et al., 2007); and (v) distance-based. We focus
51 on the latter, more specifically on KNN approaches. Although a few novelty
52 detection methods are based on the distance to the K -nearest neighbors, as
53 in (Zhang & Wang, 2006), not much attention has been given to learning
54 more sophisticated metrics.

55 In the last decade, the multi-class novelty detection problem has received
56 special attention from researchers. In (Faria et al., 2013), an evaluation ap-
57 proach for multi-class data streams novelty detection problems was proposed.
58 Also focused on data streams, de Faria et al. (2016) proposed an algorithm
59 for novelty detection that addresses the problem as a multi-class task.

60 More related to our work, Bodesheim et al. (2013) presented a kernel null
61 space based discriminant analysis for novelty detection known as KNFST,
62 which still achieves state-of-the-art performance. To improve the scalability
63 of the KNFST method, an incremental version was later proposed by (Liu
64 et al., 2017), reducing computing time with similar accuracy.

65 We propose a multi-class novelty detection method that employs a non-
66 linear distance learned from data to both detect novelties and classify normal
67 samples from multi-class datasets. To the best of our knowledge, non-linear
68 distances have never been employed for novelty detection. The optimum
69 distance is found using a kernelized extension of the *Large Margin Near-*
70 *est Neighbor* (LMNN) algorithm (Weinberger et al., 2006), named *Kernel*
71 *Large Margin Nearest Neighbors* (KLMNN) (Chatpatanasiri et al., 2010)).
72 Such distance tries to keep instances of the same class nearby while push-
73 ing instances with different labels and novelties to remain reasonably dis-
74 tant. Consequently, it provides relevant information for novelty detection,
75 and also for multi-class classification. Firstly, novelty detection is performed
76 by a variation of the well-known KNN classifier, using the learned distance
77 and distance thresholds also learned from data. Then, normal samples are
78 classified by a different proposed variant of the KNN classifier that also em-
79 ploys the learned distance and thresholds. It is worth mentioning that our
80 solution is also suitable for anomaly and outlier detection.

81 We performed experiments to show that non-linear learned distances in-
82 deed outperform the use of other distances such as linear learned distances
83 (LMNN) and the Euclidean distance when adopted by the same classifier.
84 We also show results of comparisons revealing that our method outperforms
85 previous work on novelty detection.

86 **2. Non-linear metric learning**

87 In this section, we review metric learning concepts and techniques that
88 we employ in the proposed method. First, we introduce the Mahalanobis
89 distance and give some key insights on how it could be used for novelty
90 detection (Section 2.1). Then, we describe how to learn linear metrics from
91 training data (Section 2.2) using the LMNN algorithm (Weinberger et al.,
92 2006; Weinberger & Saul, 2009), in such a way that each instance is nearest
93 to instances of its class than to instances of other classes. Finally, we briefly
94 describe the KLMNN algorithm (Chatpatanasiri et al., 2010): a non-linear
95 version of the LMNN algorithm built over the KPCA (Section 2.3).

96 *2.1. Transforming the space with the Mahalanobis distance*

97 Given a set $\mathcal{X} = \{\mathbf{x}_i\} \subset \mathbb{R}^n$ of data points, the original Mahalanobis
98 distance is defined using the covariance matrix \mathbf{C} of \mathcal{X} . However, in the
99 metric learning literature, many methods aim to compute a distance given
100 in the form:

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

101 where \mathbf{M} is some positive semi-definite matrix found by optimization, instead
102 of \mathbf{C} . Distance functions from this class of functions can be adequately
103 represented by its matrix \mathbf{M} , and are usually called generalized Mahalanobis
104 distance.

105 A key insight for such distances is that matrix \mathbf{M} can be factorized by
106 using Cholesky decomposition as $\mathbf{M} = \mathbf{G}^\top \mathbf{G}$, where a linear transformation
107 is \mathbf{G} . Consequently, we have

$$d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{G}\mathbf{x}_i - \mathbf{G}\mathbf{x}_j\|_2^2, \quad (2)$$

108 which means that any generalized Mahalanobis distance is just the squared
109 Euclidean distance applied to linearly transformed data. Thus, the problem
110 of learning an optimal metric for a specific problem can be reduced to the
111 problem of optimally deforming the feature space \mathbb{R}^n , before applying the
112 squared Euclidean distance. It is worth mentioning that distinct features
113 may be concatenated to form the data points $\mathbf{x}_i \in \mathbb{R}^n$, including continuous,
114 binary, and categorical data as well. In particular, binary features may be
115 represented by an integer with two possible values (for instance, 0 and 1),
116 and categorical features may be one-hot encoded.

117 In this paper, we require such distances to be learned from data, satis-
118 fying specific constraints. In what follows, we describe such constraints and

119 appropriate metric learning techniques we adopted to solve the constrained
 120 optimization problem.

121 2.2. Large Margin Nearest Neighbor (LMNN)

122 The LMNN algorithm was originally presented with the goal of improv-
 123 ing k-NN classification accuracy (Weinberger et al., 2006). The algorithm
 124 receives as input a labeled dataset $\mathcal{P} = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_m, c_m)\}$ with $\mathbf{x}_i \in \mathcal{X}$
 125 and labels $c_i \in \mathcal{C}$, where \mathcal{C} is a set of classes.

126 The aim is to find a generalized Mahalanobis distance that tries to keep
 127 instances of the same class as nearest neighbors, while repelling instances
 128 from different classes (impostors). More specifically, the aim is to minimize

$$\mathcal{L}(\mathbf{M}) = \sum_{(i,j) \in \mathcal{S}} d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j) + \lambda \sum_{(i,j,k) \in \mathcal{R}} [1 + d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j) - d_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_k)]_+, \quad (3)$$

129 where $[\cdot]_+ = \max(\cdot, 0)$ and $\lambda \in \mathbb{R}$ is a multiplier associated with the penalty
 130 term. \mathcal{S} is the set of all pairs (i, j) where \mathbf{x}_j is one of the K nearest neighbors
 131 in the same class as \mathbf{x}_i and \mathcal{R} is the set of all tuples (i, j, k) such that $(i, j) \in \mathcal{S}$
 132 and \mathbf{x}_k is an instance from a different class. For more details, see (Weinberger
 133 & Saul, 2009). Moreover, we highlight three key insights:

- 134 1. According to the authors, KNN classification is improved when a learned
 135 distance function is employed (Weinberger et al., 2006). By attract-
 136 ing nearest neighbors from instances of the same class while repelling
 137 nearest neighbors from distinct classes, the learned linear transforma-
 138 tion \mathbf{G} tends to better separate clusters of data points according to
 139 their classes. This behavior can be observed in Fig. 2;
- 140 2. To the best of our knowledge, this strategy has never been experimented
 141 for novelty, anomaly, or outlier detection;
- 142 3. As already mentioned, the learned Mahalanobis distance is just a squared
 143 Euclidean distance calculated over data previously transformed by \mathbf{G} .
 144 However, a linear transformation of the feature space may not be suffi-
 145 cient to evaluate distances in a high-dimensional space correctly. In or-
 146 der to overcome this limitation, Chatpatanasiri et al. (2010) proposed
 147 KLMNN as a non-linear extension for the LMNN.

148 2.3. Kernel Large Margin Nearest Neighbor (KLMNN)

149 The KLMNN approach (Chatpatanasiri et al., 2010) relies on the *kernel*
 150 *trick* in KPCA derivation. The main advantage of this non-linear extension

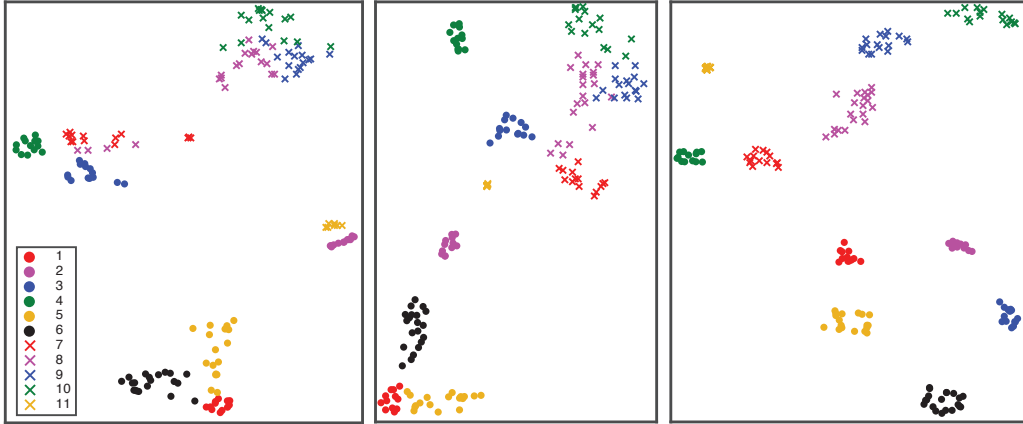


Fig. 2: Visualizing how LMNN and KLMNN improve data clusters: on the left, a t-SNE projection (Van Der Maaten, 2014) of a small training set composed of instances from 11 classes. After linearly transforming the data using LMNN (center), clusters are better separated from each other. Even better results, with more disjunct and spherical clusters, are achieved when the non-linear transformation computed using KLMNN is applied (right). Note that a multidimensional projection was applied because the dataset lies in \mathbb{R}^{17} .

151 is that it does not require the derivation of new mathematical formulas from
 152 the original LMNN formulation.

The first step of the framework is to (non-linearly) project the input dataset \mathcal{X} using the KPCA (Schölkopf et al., 1997). It is known that there exists a non-linear, possibly very-high-dimensional mapping $\Phi: \mathbb{R}^n \mapsto \mathbb{R}^N$ capable of transforming \mathcal{X} into a linearly separable set in a higher dimensional feature space \mathbb{R}^N (J. Mercer, 1909). By applying kernelization, the KPCA implicitly uses the unknown map Φ to diagonalize the covariance matrix of the transformed data, given by

$$\tilde{\mathbf{C}} = \frac{1}{m} \sum_{i=1}^m \Phi(\mathbf{x}_i) \Phi(\mathbf{x}_i)^\top.$$

153 Note that the number of eigenvectors of $\tilde{\mathbf{C}}$ equals to $|\mathcal{X}|$, which is in general
 154 much higher than n . Thus, it is expected that many eigenvalues of $\tilde{\mathbf{C}}$ will
 155 be very small, and consequently, a dimensionality reduction strategy may
 156 be applied. Following this idea, a transformed set $\tilde{\mathcal{X}}$ can be computed by
 157 projecting the data points over the most relevant eigenvectors of $\tilde{\mathcal{X}}$. From
 158 now on, we will denote the complete KPCA projection step by $\Pi: \mathcal{X} \rightarrow \tilde{\mathcal{X}}$. In

159 the following step, the original LMNN algorithm is applied to the transformed
 160 data to find the Mahalanobis matrix \mathbf{M} . Finally, from Equation (2), the non-
 161 linear learned distance is given by

$$\tilde{d}_{\mathbf{M}}(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{G}\Pi(\mathbf{x}_i) - \mathbf{G}\Pi(\mathbf{x}_j)\|_2^2. \quad (4)$$

162 It is worth mentioning that the parameters of the KPCA projection must
 163 be kept to calculate distances between unprojected points. The superiority of
 164 using the non-linear composition $\mathbf{G}\circ\Pi$ to improve data clusters configuration,
 165 in comparison with the simple linear transformation from LMNN, can be seen
 166 in Fig. 2. It is also worth mentioning that, we compute KPCA using the
 167 Gaussian kernel with a σ hyperparameter tuned according to the procedure
 168 described in Section 4.1. For more details about the kernelization, we refer
 169 the reader to the papers (Schölkopf et al., 1997; Chatpatanasiri et al., 2010).

170 3. Multi-class novelty detection and classification

171 In this section, we present a multi-class novelty classifier that may also
 172 be employed to provide multi-class classification of normal data. Here, we
 173 mainly consider novelties (or abnormal data) or instances from classes not
 174 included in the training set. If a test instance is classified as a novelty, it is
 175 filtered. Otherwise, it follows for multi-class classification.

176 3.1. Novelty classifier training

177 Following the supervised learning approach, the proposed novelty classi-
 178 fier is trained from a labeled dataset \mathcal{P} (defined in Section 2.2). The training
 179 phase is comprised of two steps: non-linear metric learning and distance
 180 threshold estimation.

181 3.1.1. Non-linear metric learning

182 In the first step, the KLMNN framework (Section 2.3) is employed to learn
 183 a non-linear distance function $\tilde{d}_{\mathbf{M}}$ from \mathcal{P} . From another perspective, a non-
 184 linear transformation of the feature space \mathbb{R}^n is learned better to reposition
 185 inliers (normal data) in well-defined clusters and isolate outliers (abnormal
 186 data).

187 Consequently, we expect inliers to be near the K nearest neighbors of
 188 their corresponding classes, and far from instances of other normal classes,
 189 as illustrated in Fig. 2 (right). Besides, our central hypothesis is that inliers

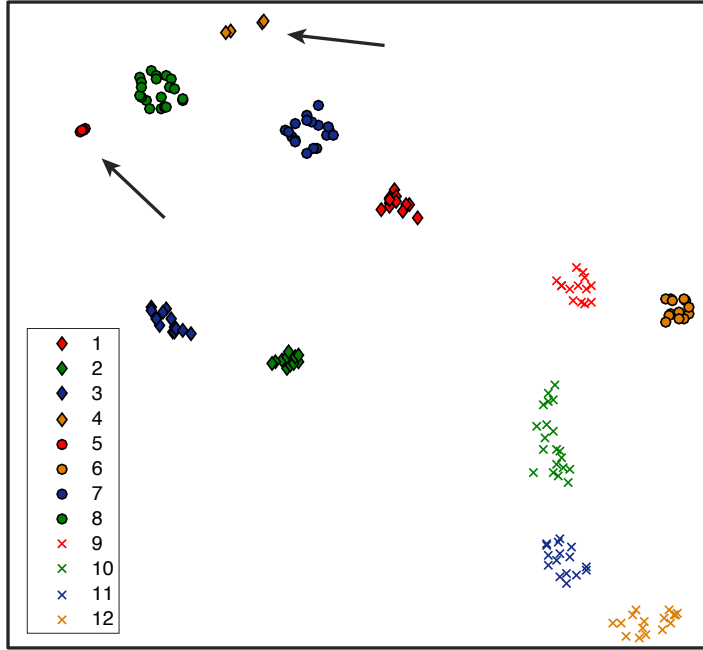


Fig. 3: Visualizing the behavior of abnormal classes after non-linearly transforming the space using instances from 10 trained classes, through t-SNE projections: instances from each inlier class are attracted to its representative cluster, while instances from two untrained classes (4 and 5) are isolated, as indicated by the arrows. In particular, instances of untrained class 4 are the only ones that are not well clustered, since classes 4 and 5 instances are not considered in the distance learning optimization (Equation (3)).

190 should become far from instances of abnormal classes (outliers) because it
 191 is expected that such outliers would lie outside the attracting region of each
 192 class. This idea is visually validated in the projection shown in Fig. 3, where
 193 outliers from two untrained (abnormal) classes are far from the regions of
 194 ten trained classes.

195 3.1.2. Distance threshold estimation

196 We propose a distance-based filter to detect and filter novelties composed
 197 of data from untrained classes and/or noise. To evaluate if a point \mathbf{x} is an
 198 inlier of a specific class c , we take into account the distance from \mathbf{x} to its
 199 K -nearest neighbors, considering only training instances from c .

Since classes may be represented by point sets with different density, we compute distance thresholds individually for each class. Considering the set

$\mathcal{X}^c = \{\mathbf{x}_i \mid (\mathbf{x}_i, c) \in \mathcal{P}\}$, *i.e.* the instances of class c in the training set \mathcal{P} . For each $\mathbf{x}_i^c \in \mathcal{X}^c$, we use $\tilde{d}_{\mathbf{M}}$ to find the distance d_i^c from \mathbf{x}_i^c to its K -nearest neighbors in \mathcal{X}^c . The distance threshold for class c is set to

$$\epsilon_c = \tau [\text{mean}(\mathcal{D}^c) + \text{std}(\mathcal{D}^c)] \quad \text{with} \quad \mathcal{D}^c = \{d_i^c \mid i = 1, \dots, |\mathcal{X}^c|\},$$

200 where $\text{mean}(\mathcal{D}^c)$ and $\text{std}(\mathcal{D}^c)$ are respectively the arithmetic mean and stan-
 201 dard deviation of \mathcal{D}^c , and $\tau \in \mathbb{R}$ is a tolerance to strengthen or soften the
 202 novelty classifier. We empirically found that $\tau = 1.2$ showed good results,
 203 although a cross-validation procedure could be applied for better tuning,
 204 as discussed in Section 4.1.

205 3.2. Novelty filtering

206 Let $\mathbf{x} \in \mathbb{R}^n$ be an input instance (from an unseen example) and its neigh-
 207 borhood $\mathcal{N}^c(\mathbf{x}) = \{\mathbf{x}_i \in \mathcal{X}^c \mid \tilde{d}_{\mathbf{M}}(\mathbf{x}, \mathbf{x}_i) < \epsilon_c\}$. An input instance \mathbf{x} is con-
 208 sidered to be an outlier of class c if $|\mathcal{N}^c(\mathbf{x})| < \kappa$, where $\kappa \leq K$ is an integer
 209 hyperparameter. Consequently, \mathbf{x} is filtered as a novelty if \mathbf{x} is an outlier for
 210 all trained classes. Only instances that pass this step will be considered for
 211 classification.

212 Fig. 1b depicts an example where a red instance would be filtered as a
 213 novelty since no training instance is inside its circle.

214 3.3. Inlier classification

215 We propose a modified KNN method to classify inlier instances, making
 216 use of the distance thresholds estimated in the training phase. We aim to
 217 take into account the density of each set \mathcal{X}^c for classification. For instance,
 218 we consider that compact clusters, such as the cluster representing class 11
 219 of Fig. 2 (right), should only influence classification if an instance is very
 220 near the cluster. On the opposite, the larger spread of class 8 data implies a
 221 larger influence region.

Given an input instance \mathbf{x} , let $\mathcal{K}(\mathbf{x}) \subset \mathcal{P}$ be the set of the K nearest *valid* neighbors of \mathbf{x} , defined as

$$\mathcal{K}(\mathbf{x}) = \{(\mathbf{x}_i, c_i) \in \mathcal{P} \mid \tilde{d}_{\mathbf{M}}(\mathbf{x}, \mathbf{x}_i) < \epsilon_c\}.$$

222 To find $\mathcal{K}(\mathbf{x})$, we recursively search for the K_i nearest neighbors of \mathbf{x} using
 223 a k-d tree data structure built over \mathcal{P} . We initiate the search with $K_1 = K$,
 224 and stop if at least K nearest neighbors are valid. Otherwise, the search
 225 continues for the K_{i+1} nearest neighbors with $K_{i+1} = 2K_i$, until the condition
 226 is satisfied. Thus, \mathbf{x} is classified as the mode of the labels of $\mathcal{K}(\mathbf{x})$.

227 4. Experiments

228 In this section, we describe experiments thoroughly performed on syn-
229 thetic and real datasets to validate, evaluate, and compare the proposed
230 methodology to well-established methods available in the literature. We pro-
231 vide experiments setup details in Section 4.1. Simulation studies on synthetic
232 data are described in Section 4.2, including computational performance eval-
233 uations. Then, the results of a comparison between our method and other
234 novelty detection approaches on real datasets are revealed in Section 4.3.
235 Finally, in Section 4.5, we described a simple proof-of-concept visual experi-
236 ment focused on gesture recognition applications.

237 4.1. Experiments Setup

238 We intended to answer the following research questions in our experi-
239 ments:

240 **Q.1** How does the proposed novelty detection method perform on linearly
241 and non-linearly separable data?

242 **Q.2** How does the proposed method accuracy behave when the number of
243 training instances increases?

244 **Q.3** Is the proposed method robust to the curse of dimensionality?

245 **Q.4** How does the proposed method compares to well established and state-
246 of-the-art novelty methods?

247 **Q.5** How does the proposed method training time increase when the number
248 of training instances and/or the number of features increase?

249 **Q.6** How does the proposed method perform on real noisy datasets?

250 To answer such questions, we adopted two specific protocol for cross-
251 validation and hyperparameters tuning, and two appropriate error metrics,
252 which we detail next. Then, we briefly describe the compared methods.

253 *Cross-validation protocols.* We aim to evaluate how novelty classifiers, in-
254 cluding ours and other compared methods, perform to discriminate data from
255 trained (normal) classes and data from classes that are not included in the
256 training set (abnormal). For this purpose, we propose two cross-validation
257 protocols: the first for simulation studies in synthetic datasets and the second
258 for experiments on real datasets.

259 **P.1** To conduct simulation studies on each synthetic dataset, we build a su-
260 pervised training set \mathcal{P}_t , representing the negative class (no novelties),
261 comprised of synthetically generated instances of each normal class and
262 including their corresponding labels. We define negative classes using
263 analytical models, which we detail in Section 4.2. It is worth emphasizing
264 that data from at least two classes is required to learn the non-linear
265 metrics. To tune the hyperparameters, random samples from an ab-
266 normal/positive class are synthetically generated to compose \mathcal{P}_u . We
267 then perform a uniform grid search with ten validation experiments per
268 hyperparameter configuration using data from \mathcal{P}_t and \mathcal{P}_u . The hyper-
269 parameter configuration with the best average performance in the ten
270 experiments is selected for the evaluation. Finally, a different test set
271 is composed of synthetic data uniformly sampled from the analytical
272 model (negative class) and the feature space (positive class). More de-
273 tails on the error metrics and hyperparameters tuning procedure follow
274 at the end of this section. The number of samples employed for each
275 experiment is detailed in Section 4.2.

276 **P.2** For real datasets, in each experiment, we exclude all instances of several
277 classes randomly selected from the set of classes. More specifically, let
278 \mathcal{C} be the set of classes of a dataset \mathcal{P} . We randomly sample a few
279 classes from \mathcal{C} to compose the trained class set \mathcal{C}_t , and the remain-
280 ing are included in the untrained class set \mathcal{C}_u . We denote by \mathcal{P}_u and
281 \mathcal{P}_t subsets of instances of \mathcal{P} belonging to \mathcal{C}_u and \mathcal{C}_t , respectively. For
282 cross-validation purposes, the training set is composed of a random
283 split of 80% of the instances in \mathcal{P}_t . The remaining instances of \mathcal{P}_t are
284 randomly included in the test set (20%). Besides, the test set receives
285 the same ratios of instances (20%) from \mathcal{P}_u . Each random selection
286 of classes is repeated ten times so that different subsets of classes are
287 considered normal/abnormal, and the random split for training/test is
288 performed. For instance, in a skeleton body poses dataset (described
289 in Section 4.3), we randomly exclude all instances of 9 key pose classes,
290 using only data from the remaining nine key poses for training. In this
291 manner, we can verify if poses from untrained classes are correctly fil-
292 tered as novelties. In practice, if an individual performs a body pose
293 which was not trained before (such as a pose similar to the instances of
294 one of the nine excluded key pose classes), the classifier should succeed
295 in filtering such pose as a novelty, instead of classifying it as one of the

296 trained key poses. To tune the hyperparameters, we performed a uni-
 297 form grid search with ten validation experiments per hyperparameter
 298 configuration. The hyperparameter configuration with the best aver-
 299 age performance in the ten experiments is selected, and the respective
 300 average score is used as a result of the final evaluation of the model.

301 *Error metrics.* To evaluate the accuracy of the novelty detection classi-
 302 fiers, we adopted two appropriate metrics for binary classification problems:
 303 the F_1 score and the Matthews correlation coefficient (MCC). As we describe
 304 next, both can be written in terms of TP (true positives); TN (true nega-
 305 tives), FP (false positive), and FN (false negatives). We consider abnormal
 306 data (novelties) to be a positive class and normal data to be the negative
 307 class. In both metrics, higher scores are better.

- F_1 score: a measure of accuracy of the classifier, defined as the har-
 monic mean of precision and recall, as follows

$$F_1 = \frac{TP}{TP + 0.5(FP + FN)}.$$

- *Matthews correlation coefficient* (MCC): a more informative measure
 of accuracy for binary classifier, which is generally better than F_1 score
 for novelty detection problems, since it takes into account the balance
 ratios of the four confusion matrix categories:

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}}.$$

308 In the following experiments, we consider the MCC as the reference error
 309 metric, although the F_1 score is also shown.

310 *Hyperparameters tuning details.* In all experiments, we tune hyperparame-
 311 ters of our approach, and also each compared methods to perform fair com-
 312 parisons. We applied uniform grid searches over the hyperparameters space
 313 of each method, considering the resulting MCC metric over validation data
 314 from the cross-validation approach. Our experiments revealed that such cal-
 315 ibration was crucial to achieve great results, as exemplified in Fig. 4: the
 316 non-linear classifier outperformed the other metrics only when the Gaussian
 317 kernel parameter (σ) was assigned to a value in a specific interval. For all

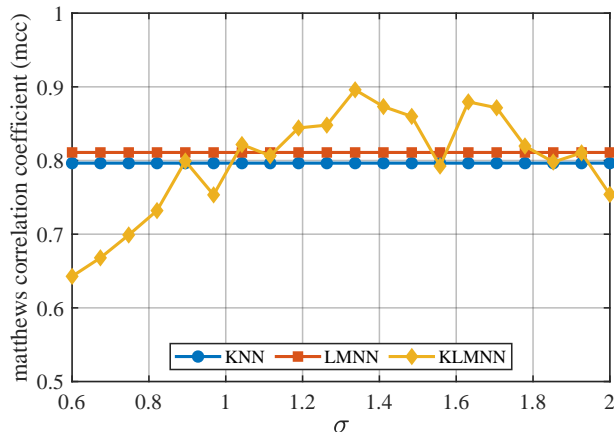


Fig. 4: Hyperparameter tuning on the real KGD dataset (see Section 4.3): outlier filtering achieved better results when our non-linear metric was employed compared to the Euclidean and linear metric. However, hyperparameter tuning was essential, since only when kernel parameter σ in the interval $[1.0, 1.8]$ were used for KPCA projection, resulting in non-linear distance outperformed the other distances (note that the other distances do not depend on σ).

318 KNN-based methods, we also performed cross-validation experiments to tune
 319 the hyperparameters K , κ , and τ for each dataset. For K and κ , we con-
 320 sidered all pairs of integers in the range $[1, 5]$ where $\kappa \leq K$. We evaluate
 321 50 equally spaced values in the range $[0.5, 3.0]$ for τ . Tables 2 and 4 reveal
 322 results of tuning experiments on a synthetic and on a real dataset, respec-
 323 tively.

324 *Compared methods.* We perform comparisons with well established previous
 325 work regularly used for novelty and anomaly detection: one-against-all multi-
 326 class SVM (Hsu & Lin, 2002) (MCSVM), in which an instance is classified as
 327 novelty if it does not belong to any trained class; one-class SVM (Schölkopf
 328 et al., 2001) (OSVM); KPCA for novelty detection (Hoffmann, 2007) (KP-
 329 CANOV); and Kernel Null Space Method (Bodesheim et al., 2013) (KNFST),
 330 which was recently enhanced by Liu et al. (2017) and still achieves state of
 331 the art accuracy. We also compare the usage of our proposed non-linear
 332 metrics (KLMNN) with linear (LMNN) and Euclidean (KNN) metrics. For
 333 all compared methods, we experimented with both the usual Gaussian and
 334 polynomial kernels. Regarding the KNFST, we adjust the kernel param-
 335 eter and a decision threshold that varies from zero to the minimum distance

336 between the trained classes’ target points. For the KPCANOV, we adjust
337 the kernel parameter and the number of relevant eigenvectors to the feature
338 space’s error reconstruction. For the OSVM algorithm, we only adjust the
339 kernel parameter. Already for MCSVM, both the kernel parameter and a
340 posterior probability threshold are adjusted. Although this posterior proba-
341 bility is not standard for SVM methods, it is an attractive way to estimate
342 decision boundaries for support vector machines. More details can be found
343 in the official documentation of the main libraries that implement SVM.

344 *4.2. Simulation studies on synthetic data*

345 To evaluate our novelty detection approach in different situations, we cre-
346 ated simulation studies on four datasets depicted in Figures 5 and 6: three
347 horizontal lines; three parabolas; three concentric circles; and four uniform
348 distributions. Note that two datasets are linearly separable (horizontal lines
349 and uniform distributions), and two are not linearly separable (parabola and
350 concentric circles). The uniform distribution simulations will be employed
351 to evaluate our novelty detection approach when the number of training ex-
352 amples or the number of features grows. For the process of hyperparameters
353 tuning and methods evaluation, we apply the cross-validation protocol **P.1**
354 (see Subsection 4.1).

355 In the first simulation, we created three classes with 25, 50, and 25 points
356 from the sampling of 3 horizontal lines. Similarly, we created the simulation
357 of parabolas and concentric circles. To assist the process of hyperparameters
358 tuning, we added a fourth class of random samples composed of 192, 187,
359 and 193 samples for the first, second, and third simulation, respectively.

360 During the evaluation of hyperparameters, we noticed that the kernel-
361 based methods presented better results with a polynomial kernel of type
362 $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + 1)^d$ than the traditional Gaussian kernel. On the other
363 hand, for OSVM and MCSVM, the Gaussian kernel was more advantageous.
364 To study the best hyperparameters K and κ for KNN, LMNN and KLMNN
365 methods, we chose to apply the protocol **P.1** individually for different combi-
366 nations of K and κ , so that the other hyperparameters are adjusted according
367 to each combination.

368 Table 1 summarizes the results, where we also show the value of the F_1
369 score. The best methods were KLMNN and KNFST, which achieved max-
370 imum MCC and F_1 scores in all simulation studies. KPCANOV performed
371 well for the simulation of horizontal lines and parabolas, while for the sim-
372 ulation of circles, the result was unsatisfactory. SVM-based methods were

Table 1: Comparison of the methods in 4 simulation studies (best results in bold).

simulation studies	train / validation		test		metrics	KNN-based methods			compared methods			
	normal	novelty	normal	novelty		KNN	LMNN	KLMNN	KNFST	OSVM	MCSVM	KPCANOV
horizontal lines	100	192	900	89 100	F_1	0.14	1.00	1.00	1.00	0.15	0.14	0.99
					MCC	0.25	1.00	1.00	1.00	0.26	0.26	0.99
parabolas	100	187	1 500	98 500	F_1	0.31	0.31	1.00	1.00	0.23	0.16	0.99
					MCC	0.30	0.30	1.00	1.00	0.32	0.24	0.99
concentric circles	100	193	3 300	197 970	F_1	0.16	0.16	1.00	1.00	0.10	0.17	0.25
					MCC	0.26	0.26	1.00	1.00	0.18	0.22	0.35

373 unsatisfactory in all simulations. These results can be attested visually by
 374 Fig. 5, which shows the decision regions for each method in the three simu-
 375 lations. Any point outside these regions is considered abnormal data.

376 For the simulation of horizontal lines, both linear (LMNN) and non-linear
 377 metric learning (KLMNN) were able to capture the exact distribution of
 378 the 3 lines. In this case, the MCC score was 1.0 for all tested K and κ
 379 combinations. On the other hand, for KNN the best result was only 0.25
 380 with hyperparameters $K = 2$ and $\kappa = 1$.

381 Meanwhile, in the 3 parabolas simulation (Fig. 5, *middle*), we observed
 382 that the linear metric learning was insufficient to capture the correct distri-
 383 bution of the data. In this case, the best result for KNN and LMNN was
 384 0.3 when $K = 1$ and $\kappa = 1$. In turn, with non-linear metric learning, it was
 385 possible to capture the exact distribution of the 3 parabolas, regardless of K
 386 and κ values. These results can be checked in more detail in Table 2, where
 387 κ is defined as a function of K .

Table 2: MCC of the KNN-based methods for different combinations of hyperparameters K and κ in the 3 parabolas simulation. The scores for values of κ less than 1 are omitted with symbol \emptyset . Best results in bold.

K	KNN				LMNN				KLMNN			
	κ				κ				κ			
	K	$K - 1$	$K - 2$	$K - 3$	K	$K - 1$	$K - 2$	$K - 3$	K	$K - 1$	$K - 2$	$K - 3$
1	0.30	\emptyset	\emptyset	\emptyset	0.30	\emptyset	\emptyset	\emptyset	1.00	\emptyset	\emptyset	\emptyset
2	0.14	0.29	\emptyset	\emptyset	0.07	0.29	\emptyset	\emptyset	1.00	1.00	\emptyset	\emptyset
3	0.12	0.14	0.16	\emptyset	0.12	0.16	0.29	\emptyset	1.00	1.00	1.00	x
4	0.10	0.13	0.17	0.17	0.11	0.13	0.16	0.17	1.00	1.00	1.00	1.00
5	0.10	0.12	0.13	0.15	0.10	0.12	0.13	0.14	1.00	1.00	1.00	1.00

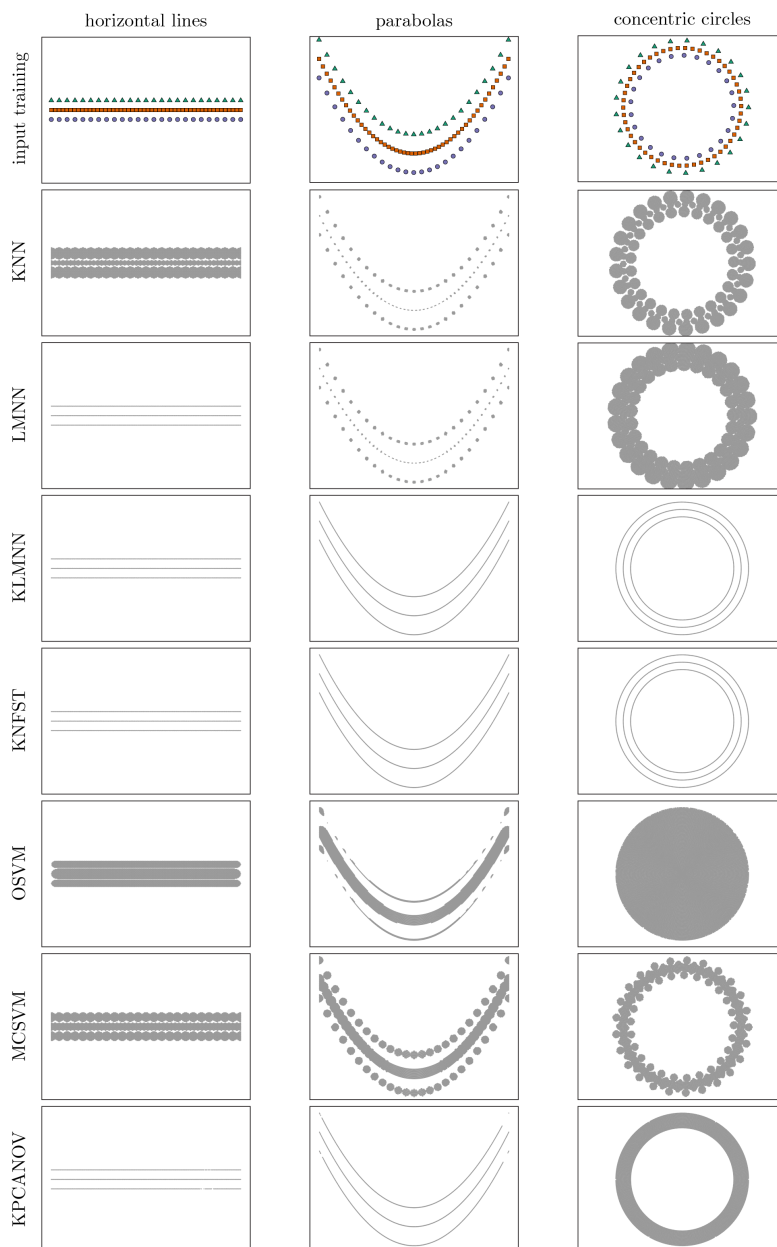


Fig. 5: Decision regions of all methods on simulation studies 1, 2 and 3. Every point belonging to the gray region is considered normal data, which indicates that they belong to a class of the training data. On the other hand, any point outside the region is considered novelty (abnormal class).

388 As in the previous case, in the simulation of concentric circles, only the
 389 KLMNN succeeded to correctly determine the circles. In this case, all tested
 390 combinations (K, κ) resulted in an MCC score of 1.0. In turn, with KNN
 391 and LMNN, the results were unsatisfactory, with a better score of only 0.26
 392 for both methods when $K = 2$ and $\kappa = 1$.

393 Fig. 6 shows another simulation study to evaluate the novelty detection
 394 performance of our KNN-based approach when we increase the number of
 395 training samples and the number of dimensions (features) as well. In this
 396 case, we create five classes with uniform probability distribution in 5 hy-
 397 percubes centered on the xy -plane of \mathbb{R}^n . For instance, in \mathbb{R}^2 classes are
 398 composed of random points in squares of side 0.9 centered on $(1.0, 0.5)$,
 399 $(-0.5, 1.0)$, $(-1.0, -0.5)$, $(0.5, -1.0)$ and $(0.0, 0.0)$. For the purpose of hy-
 400 perparameters tuning and methods evaluation, we select one class and con-
 401 sidered it as abnormal. In the first experiment, we fixed the number of
 402 dimensions at ten and varied the number of training samples by 100, 200,
 403 300, 400, 500, 750, 1000, 1500, and 2000. While in the second experiment,
 404 we fixed the number of training points and varied the number of dimensions
 405 from 2 to 20. Again, we adopted the protocol **P.1** to adjust the hyperparam-
 406 eters, and the overall performance of novelty detection was calculated in a
 407 test set with 50K samples (40K from trained classes and 10K from the abnor-
 408 mal class). For both KNN, LMNN, and KLMNN, the other hyperparameters
 409 were adjusted with $K = 3$ and $\kappa = 1$. Fig. 7 depicts the results.

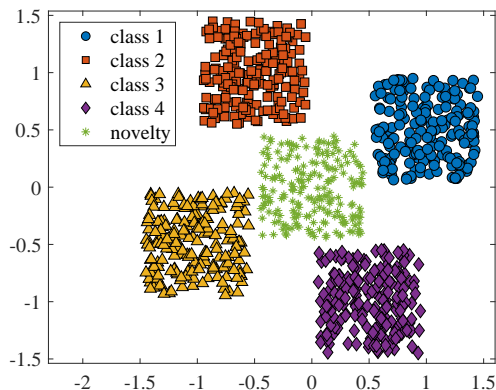


Fig. 6: Synthetic dataset with 5 hypercubes to evaluate the performance of the methods when the number of samples grows and also when the number of dimensions grows.

410 The first experiment (Fig. 7, *left*), when the number of training samples
 411 grows, our both LMNN and KLMNN improves the performance of KNN

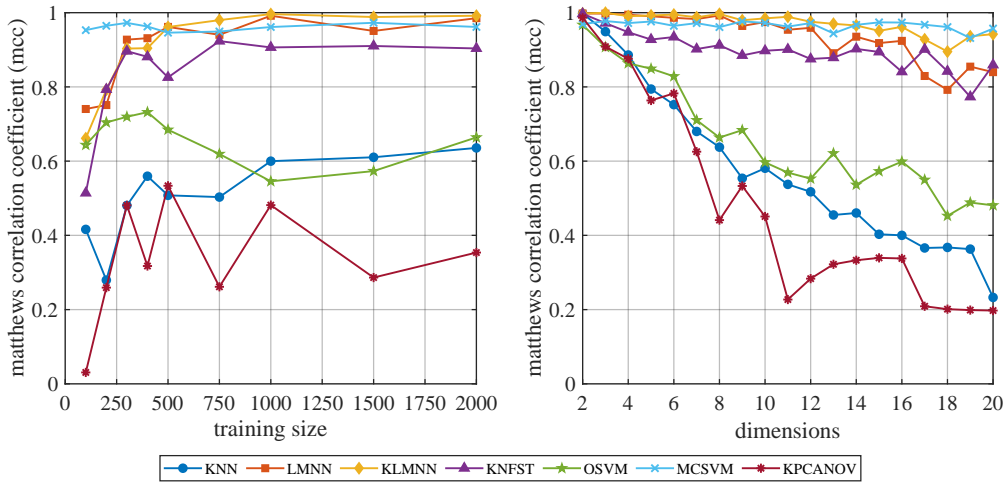


Fig. 7: Evaluation of methods with MCC when the number of training examples increases (*left*) and when the number of dimensions increases (*right*). All experiments were evaluated in a test set with 50K samples (40K normal and 10K novelty).

412 novelty detector, note that in the range 100 to 500 the MCC value it rises
 413 from approximately 0.66 to above 0.95 with KLMNN, reaching a value of
 414 0.99 with 1000 and 2000 training samples. However, our version of the KNN
 415 with Euclidean distance improves from 0.40 to just 0.55 in the same range
 416 of training data, reaching a maximum of 0.63 with 2000 training samples.
 417 For the other methods, MCSVM proved to be efficient in this experiment,
 418 reaching a score of 0.95 in almost all cases. This possibly occurs due to the
 419 simplicity of positioning and distribution of the classes, since the multi-class
 420 SVM calculates hyperplanes that best separate the classes from each other.
 421 OSVM is inefficient in the same task, achieving MCC values ranging around
 422 0.6, with a maximum value of 0.73. KNFST improves considerably in the
 423 range [100, 750], varying from 0.50 to 0.92, but it remained at the same level
 424 of 0.90 for experiments with more than 750 training samples. KPCANOV
 425 had the worst performance, with MCC of at most 0.53 remaining the score
 426 unchanged, even increasing the training samples.

427 The second experiment (Figure 7, *right*), when the number of features
 428 grows, both linear and non-linear metrics mitigate the impact of the curse of
 429 dimensionality on the KNN. It is enough to note that in the KNN, the MCC
 430 metric varies from 0.99 in dimension 2 to 0.23 in dimension 20, while LMNN
 431 and KLMNN varied from 0.99 to 0.84 and 0.95, respectively. Regarding the
 432 compared methods, we see that both OSVM and KPCANOV had the same

433 problem as KNN with Euclidean distance. In contrast, KNFST remained
434 relatively stable, with an MCC value of 0.85 in dimension 20. MCSVM is
435 also stable, with a score always around 0.95 due to positioning characteristics
436 and class distribution.

437 Therefore, the experiments above presented provide answers to the ques-
438 tions from **Q.1** to **Q.4** of our research (Section 4.1).

439 *Computational performance.* We use the simulation depicted by Fig. 6 to
440 measure the computational performance of our method and previous work.
441 We aim to evaluate how training/evaluation times grow when the number
442 of dimensions (features) and the number of training examples grows when
443 evaluated on a set of 50K test samples. The experiments were performed
444 on a single core 1.8GHz of an Intel Core i7-8550U processor and 8GB of
445 RAM, without using GPU processing. All models were trained using Matlab
446 R2018b on Windows 10. To measure training/evaluation time in terms of
447 the number of training examples, we fixed the dimension to 10. As shown
448 in the left chart of Fig. 8, both SVM-based approaches performed better
449 than the other methods. KNN-based methods reached satisfactory results,
450 although some additional time was required to learn the non-linear metrics
451 of our KLMNN. It is worth noting that the state-of-the-art KNFST method
452 performed similarly to KLMNN, in which case 2000 samples were trained in
453 less than a minute. Contrastingly, KPCANOV required substantially higher
454 training times and showed the worst growth pattern. The results of the ex-
455 periments varying the number of features are shown in the right chart of
456 Fig. 8. Here, we fixed the number of training examples to 800. Notably,
457 all methods showed almost constant training times when the number of di-
458 mensions varies. Specifically, in our KLMNN, the number of features only
459 influences kernel computations, which do not substantially contribute to the
460 total training time. Thus, answering the question **Q.5**.

461 4.3. Experiments on real datasets

462 We also experimented on four multi-class (more than two classes) datasets
463 with different characteristics. All datasets contain noisy data. Two ac-
464 cessible datasets were chosen from the UCI repository of machine learning
465 databases (Dheeru & Karra Taniskidou, 2017): *iris* and *glass*. The small iris
466 dataset contains three classes of 50 instances each, where each class refers
467 to a type of iris plant, and four attributes represent each instance. Due to
468 its simplicity, we chose this dataset to demonstrate that general multi-class

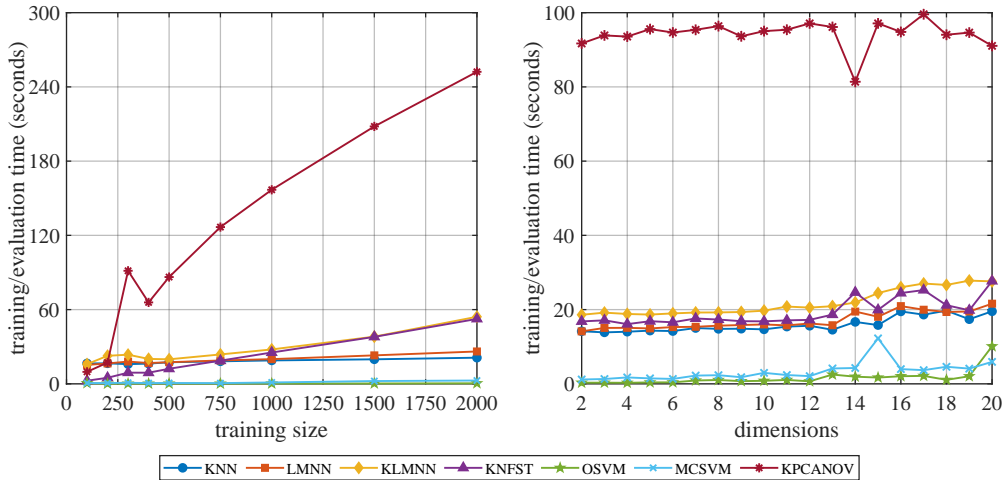


Fig. 8: Computational performance of the methods when the number of training examples grows (left) and when the number of dimensions grows (right). All experiments were evaluated on a test set with 50K samples (40K normal and 10K abnormal).

469 classifiers are not appropriate to handle data from untrained classes. The un-
 470 balanced glass is an extremely noisy dataset composed of instances of 6 types
 471 of glass (classes), where the number of instances of each class is, respectively:
 472 70, 76, 17, 13, 9 and 29, as depicted in glass. This dataset was adopted to
 473 investigate our method’s performance when dealing with low-quality data
 474 compared with the methods. The other experimented datasets are related to
 475 the gesture recognition problem, which strongly relies on multi-class novelty
 476 detection. The first one is a public gesture dataset for body key poses named
 477 *KGD* (Miranda et al., 2014a,b). This dataset is composed of instances from
 478 18 body key poses (classes), described by nine joint angles extracted from
 479 body skeletons, with approximately 11 instances per class. The second is a
 480 novel hand pose dataset representing signs of the Brazilian Sign Language
 481 (*libras*), and now made publicly available¹. The dataset is composed of 20
 482 signs (classes) represented by 14 joint angles of hand skeletons captured from
 483 a Leap Motion sensor (Leap Motion, 2018), with approximately 19 instances
 484 per class.

485 In Table 3, we summarize the results for all experimented methods and
 486 datasets, when the protocol **P.2** (Section 4.1) and optimal hyperparame-

¹<https://ic.ufal.br/professor/thales/leaplibras/>

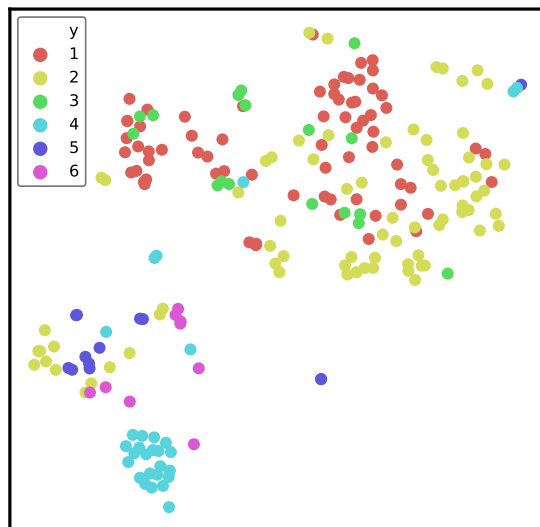


Fig. 9: t-SNE projection of the *glass* dataset: all classes are noisy and/or not well-clustered.

487 ters are employed on the test sets (for several random selections of nor-
 488 mal/abnormal classes). For instance, the K and κ calibration on the KGD
 489 dataset are revealed in Table 4 for all KNN-based methods.

490 Results on the iris dataset corroborate the limitations of MCSVM for
 491 novelty detection: by excluding one of the three classes from the training
 492 set, MCSVM fails to identify that instances from the excluded class are nov-
 493 elties, achieving the worst MCC result among all compared methods. This
 494 result may be visually explained by the sketch already shown in Fig. 1a: since
 495 only two classes were used for training, we claim that the resulting MCSVM
 496 hyperplane separates the trained classes, but does not define any untrained
 497 region as in the illustration shown in Fig. 1b for a distance-based novelty clas-
 498 sifier. Our approach outperformed all methods under both metrics, closely
 499 followed by its linear variant.

500 Except for the KGD dataset, the best results were achieved by the met-
 501 ric learning approaches KLMNN and LMNN. When both approaches are
 502 directly compared, KLMNN is only outperformed by the LMNN variant on
 503 the glass dataset under the F_1 score. This may be a consequence of severe
 504 noise found in that dataset, as shown in Fig. 9. In such scenarios, we claim
 505 that non-linear metric learning is more susceptible to degenerate the data.
 506 Overall, considering the MCC as the reference error metric, our approach
 507 outperformed the other methods in three out of the four datasets.

508 Finally, the study performed in this section provides the answer to the
 509 question **Q.6**.

Table 3: Comparison between our non-linear novelty classifier (KLMNN) and other methods on four real datasets. Here, **nc** represents the number of classes of the dataset, and **nuc** the number of excluded classes. The best results are highlighted in bold and the second best results underlined.

dataset	nc	nuc	acc. score	KNN-based methods			other methods			
				KNN	LMNN	KLMNN	KNFST	OSVM	MCSVM	KPCANOV
iris	3	1	F_1	0.94	<u>0.96</u>	0.96	0.91	0.94	0.91	0.94
			MCC	0.80	<u>0.88</u>	0.89	0.78	0.82	0.76	0.83
glass	6	3	F_1	0.66	0.74	0.68	0.65	0.72	0.62	<u>0.73</u>
			MCC	0.35	0.35	0.39	0.31	0.22	<u>0.36</u>	<u>0.35</u>
KGD	18	9	F_1	0.94	0.92	<u>0.95</u>	0.96	0.85	0.94	<u>0.95</u>
			MCC	0.87	0.84	<u>0.90</u>	0.92	0.71	0.89	<u>0.90</u>
libras	20	5	F_1	<u>0.93</u>	0.93	0.93	0.91	0.92	<u>0.93</u>	<u>0.93</u>
			MCC	0.69	0.69	0.70	0.66	0.67	0.68	<u>0.70</u>

Table 4: MCC of the KNN-based methods for different combinations of hyperparameters K and κ in the *KGD* dataset. The scores for values of κ less than 1 are omitted with symbol \emptyset . Best results in bold.

K	KNN				LMNN				KLMNN			
	κ				κ				κ			
	K	$K-1$	$K-2$	$K-3$	K	$K-1$	$K-2$	$K-3$	K	$K-1$	$K-2$	$K-3$
1	0.87	\emptyset	\emptyset	\emptyset	0.84	\emptyset	\emptyset	\emptyset	0.89	\emptyset	\emptyset	\emptyset
2	0.73	0.80	\emptyset	\emptyset	0.81	0.81	\emptyset	\emptyset	0.89	0.90	\emptyset	\emptyset
3	0.51	0.60	0.73	\emptyset	0.75	0.82	0.81	\emptyset	0.89	0.89	0.89	\emptyset
4	0.52	0.52	0.61	0.76	0.75	0.75	0.82	0.82	0.86	0.88	0.88	0.89
5	0.56	0.50	0.53	0.65	0.62	0.64	0.68	0.75	0.79	0.83	0.83	0.83

510 4.4. Multi-class classification results

511 We evaluated the non-linear multi-class KNN classifier described in Sec-
 512 tion 3.3 for a specific classification task. In this case, we compared only to
 513 the KNFST and MCSVM methods, since the other OSVM and KPCANOV
 514 were not developed for this purpose. To evaluate each method’s performance,
 515 we calculated the precision, recall, and F_1 score for each class individually.

516 Then, we calculated the well-known macro and weighted averages of each
 517 of these measures. We calculated these measures from the aggregate confu-
 518 sion matrix for experiments on real datasets, summarizing the results of 10
 519 validation experiments.

520 In Tables 5 and 6, we respectively show comparisons of our approach
 521 for multi-class classification on both a synthetic dataset (uniform distribu-
 522 tions with 1000 training examples in dimension 10); and a real data set
 523 (iris). In both tables, **pre** represent the precision, **rec** is the recall and **f1**
 524 is the F_1 score. In addition, **m. avg** is the macro average and **w. avg**
 525 is the weighted average measure. Both LMNN and KLMNN variations of our
 526 multi-class classifier performed better than the KNFST and MCSVM. Also,
 527 it is worth mentioning that learning metrics was important to considerably
 528 improved the classification performance of the knn classifier. The KNFST
 529 and MCSVM also performed well in this case, reaching average recalls of
 530 approximately 1.0, but with average precisions slightly lower. As for the
 531 iris dataset, we can see from Table 6 that our approach KLMNN outper-
 532 formed the compared methods, reaching 0.95 and 0.96 for macro average
 533 and weighted average F_1 score, respectively.

Table 5: Multi-class classification results on the uniform distributions from experiment of Fig. 6. Best results for macro average and weighted average F_1 score are in bold and underlined.

	KNN			LMNN			KLMNN			KNFST			MCSVM		
	pre	rec	f1	pre	rec	f1	pre	rec	f1	pre	rec	f1	pre	rec	f1
class 1	0.90	0.94	0.92	1.00	1.00	1.00	1.00	1.00	1.00	0.97	1.00	0.98	0.99	1.00	0.99
class 2	0.91	0.96	0.93	1.00	1.00	1.00	1.00	1.00	1.00	0.97	1.00	0.98	0.98	1.00	0.99
class 3	0.91	0.94	0.93	1.00	1.00	1.00	1.00	1.00	1.00	0.97	1.00	0.98	0.99	1.00	0.99
class 4	0.91	0.93	0.92	1.00	1.00	1.00	1.00	1.00	1.00	0.96	1.00	0.98	0.98	1.00	0.99
m. avg	0.91	0.94	0.93	1.00	1.00	<u>1.00</u>	1.00	1.00	1.00	0.97	1.00	0.98	0.98	1.00	0.99
w. avg	0.91	0.94	0.93	1.00	1.00	<u>1.00</u>	1.00	1.00	1.00	0.97	1.00	0.98	0.98	1.00	0.99

534 4.5. Proof-of-concept visual experiment for gestures

535 We finally performed a simple visual experiment on gestures represented
 536 by key pose classes, to demonstrate that gesture recognition can be improved
 537 by applying our non-linear novelty filtering step. In Fig. 10, we show a
 538 t-SNE projection of instances of key pose classes and poses of a gesture
 539 example with its respective dashed trajectory. In Fig. 10a, although the

Table 6: Multi-class classification results on the iris dataset. The first and second best scores for macro average and weighted average F_1 score are highlighted in bold and underline, respectively.

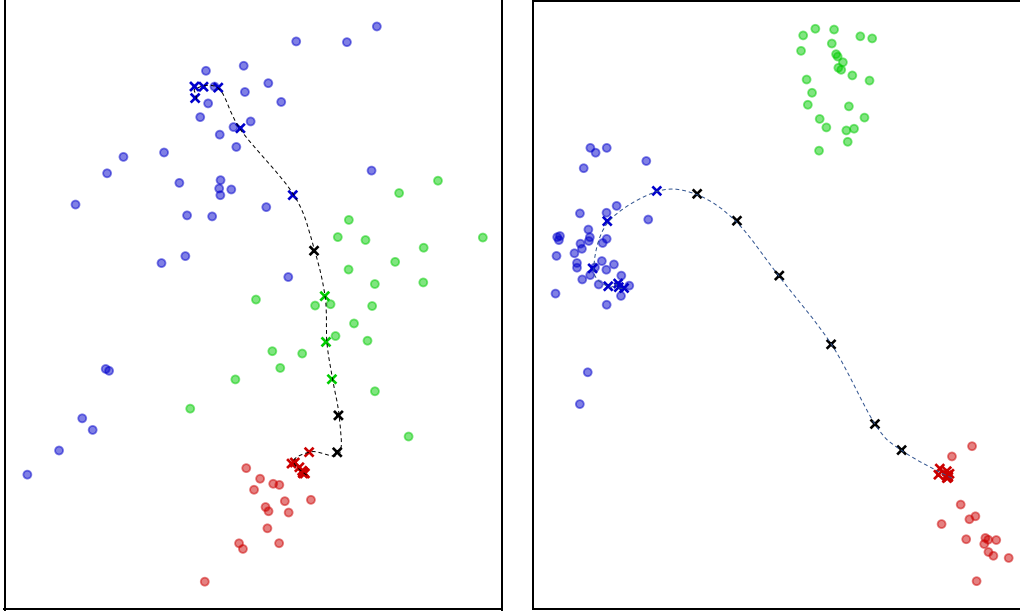
	KNN			LMNN			KLMNN			KNFST			MCSVM		
	pre	rec	f1	pre	rec	f1	pre	rec	f1	pre	rec	f1	pre	rec	f1
class 1	0.92	0.99	0.95	0.99	0.96	0.97	0.99	0.99	0.99	0.98	0.98	0.98	0.97	0.95	0.96
class 2	0.83	0.98	0.90	0.89	0.97	0.93	0.88	0.97	0.92	0.81	0.87	0.84	0.83	0.92	0.87
class 3	0.89	0.95	0.92	0.94	0.97	0.95	0.94	0.97	0.95	0.98	0.87	0.92	0.94	0.83	0.88
m. avg	0.88	0.97	0.92	0.94	0.97	<u>0.95</u>	0.93	0.97	0.95	0.92	0.90	0.91	0.92	0.90	0.91
w. avg	0.88	0.98	0.93	0.94	0.97	<u>0.95</u>	0.94	0.98	0.96	0.93	0.91	0.92	0.92	0.91	0.91

540 gesture trajectory intersects a region with several green instances, in the
541 high-dimensional feature space, they do not. Thus, the green crosses are the
542 result of incorrect pose classifications that were not detected as novelties. In
543 Fig. 10b, only the correct key poses are detected (red and blue). This gives
544 us evidence that our non-linear novelty filter may improve distance-based
545 gesture recognition methods.

546 5. Limitations and future work

547 We proposed a method for multi-class novelty detection that is also ca-
548 pable of performing multi-class classification. The method proposed here is
549 based on a non-linear distance learned from training data and may be ap-
550 plied for novelty, outlier or anomaly detection. Our main conclusions are:
551 1) multi-class SVMs are not appropriate for novelty detection and filtering
552 in situations where data from abnormal classes are given as input for classi-
553 fication, while our methodology succeeds; 2) in many situations, non-linear
554 distances learning (KLMNN) achieves better results when compared to linear
555 distances learning (LMNN) and the standard Euclidean distance; and 3) our
556 method outperforms previous work on novelty detection in most experiments.
557 In addition, our experiments revealed that our method scales well when both
558 the number of features and/or samples grows when compared to existing
559 novelty detection approaches, both in terms of accuracy and computational
560 performance.

561 Differently from one-class classification, such as OSVM and KPCANOV,
562 the proposed solution is limited to multi-class novelty detection problems
563 (Bodesheim et al., 2013), relying on training data from at least two normal
564 classes to learn a distance function. Nevertheless, there is still great potential



(a) Euclidean metric: a few poses of the gesture are incorrectly classified as the green class. (b) Non-linear metric: after transforming the space, the green class cluster is far from the trajectory, and only the correct classes red and blue are detected over the trajectory.

Fig. 10: t-SNE projection of instances of 3 key pose classes (dots), and the trajectory of a gesture (whose poses are represented by crosses) between two key pose classes (red and blue). The color of the crosses encodes pose classifications by our multi-class classifier (red, green or blue), or novelties (black).

565 for usage in a wide range of applications where only data from a subset of
 566 classes are given as input for training, as mentioned in Section 1.

567 As future work, we plan to apply the proposed solution for applications
 568 such as real-time pose and gesture recognition; and visual recognition sys-
 569 tems. Experimenting with different kernels for the KPCA projection may
 570 also be the topic of future research, as well as more sophisticated hyperpa-
 571 rameter tuning strategies.

572 **Acknowledgments**

573 We would like to thank the reviewers for their comments. The authors
574 would like to thank Coordenação de Aperfeiçoamento de Pessoal de Nível
575 Superior – Brazil (CAPES), National Council for Scientific and Technological
576 Development – Brazil (CNPq) fellowship #301642/2017-6, São Paulo
577 Research Foundation (FAPESP) under grant #2019/23215-9, and Alagoas
578 Research Foundation (FAPEAL) for partially financing this research.

579 **References**

- 580 Bodesheim, P., Freytag, A., Rodner, E., Kemmler, M., & Denzler, J. (2013).
581 Kernel null space methods for novelty detection. In *Proceedings of the*
582 *IEEE conference on computer vision and pattern recognition* (pp. 3374–
583 3381).
- 584 Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A
585 survey. *ACM computing surveys (CSUR)*, *41*, 15:1–15:58.
- 586 Chatpatanasiri, R., Korsrilabutr, T., Tangchanachaianan, P., & Kijirikul, B.
587 (2010). A new kernelization framework for mahalanobis distance learning
588 algorithms. *Neurocomputing*, *73*, 1570–1579.
- 589 Clifton, L., Clifton, D. A., Watkinson, P. J., & Tarassenko, L. (2011). Identifi-
590 cation of patient deterioration in vital-sign data using one-class support
591 vector machines. In *Computer Science and Information Systems (FedC-*
592 *SIS), 2011 Federated Conference on* (pp. 125–131).
- 593 Dheeru, D., & Karra Taniskidou, E. (2017). UCI machine learning repository.
594 URL: <http://archive.ics.uci.edu/ml>.
- 595 Diehl, C. P., & Hampshire, J. B. (2002). Real-time object classification and
596 novelty detection for collaborative video surveillance. In *Neural Networks,*
597 *2002. IJCNN'02. Proceedings of the 2002 International Joint Conference*
598 *on* (pp. 2620–2625). IEEE volume 3.
- 599 Faria, E. R., Gonçalves, I. J., Gama, J., & Carvalho, A. C. (2013). Evaluation
600 methodology for multiclass novelty detection algorithms. In *2013 Brazilian*
601 *Conference on Intelligent Systems* (pp. 19–25). IEEE.

- 602 de Faria, E. R., de Leon Ferreira, A. C. P., Gama, J. et al. (2016). Minas:
603 multiclass learning algorithm for novelty detection in data streams. *Data*
604 *mining and knowledge discovery*, *30*, 640–680.
- 605 Hoffmann, H. (2007). Kernel pca for novelty detection. *Pattern recognition*,
606 *40*, 863–874.
- 607 Hsu, C.-W., & Lin, C.-J. (2002). A comparison of methods for multiclass
608 support vector machines. *IEEE Trans. Neural Netw.*, *13*, 415–425.
- 609 Ilonen, J., Paalanen, P., Kamarainen, J.-K., & Kalviainen, H. (2006). Gaus-
610 sian mixture pdf in one-class classification: computing and utilizing confi-
611 dence values. In *ICPR* (pp. 577–580). IEEE volume 2.
- 612 J. Mercer, B. A. (1909). Functions of positive and negative type, and their
613 connection the theory of integral equations. *Phil. Trans. R. Soc. Lond. A*,
614 *209*, 415–446.
- 615 Jyothsna, V., Prasad, V. R., & Prasad, K. M. (2011). A review of anomaly
616 based intrusion detection systems. *Int. J. Comput. Appl.*, *28*, 26–35.
- 617 Keogh, E., Lin, J., Lee, S.-H., & Van Herle, H. (2007). Finding the most
618 unusual time series subsequence: algorithms and applications. *Knowl. Inf.*
619 *Syst.*, *11*, 1–27.
- 620 Khan, S. S., & Madden, M. G. (2010). A survey of recent trends in one
621 class classification. In *Artificial Intelligence and Cognitive Science* (pp.
622 188–197). Berlin, Heidelberg: Springer Berlin Heidelberg.
- 623 Leap Motion, I. (2018). Leap motion sdk. URL: [https://developer.
624 leapmotion.com/get-started/](https://developer.leapmotion.com/get-started/).
- 625 Liu, J., Lian, Z., Wang, Y., & Xiao, J. (2017). Incremental kernel null space
626 discriminant analysis for novelty detection. In *Proceedings of the IEEE*
627 *Conference on Computer Vision and Pattern Recognition* (pp. 792–800).
- 628 Lv, F., & Nevatia, R. (2007). Single view human action recognition using
629 key pose matching and Viterbi path searching. In *CVPR* (pp. 1–8).
- 630 Miranda, L., Vieira, T., Martinez, D., Lewiner, T., Vieira, A. W., & Campos,
631 M. F. M. (2014a). Keypose and Gesture Database. URL: [http://im.
632 ufal.br/professor/thales/gesturedb.html](http://im.ufal.br/professor/thales/gesturedb.html).

- 633 Miranda, L., Vieira, T., Martinez, D., Lewiner, T., Vieira, A. W., & Campos,
634 M. F. M. (2014b). Online gesture recognition from pose kernel learning
635 and decision forests. *Pat. Rec. Letters*, *39*, 65–73.
- 636 Pimentel, M. A., Clifton, D. A., Clifton, L., & Tarassenko, L. (2014). A
637 review of novelty detection. *Signal Processing*, *99*, 215–249.
- 638 Schölkopf, B., Platt, J. C., Shawe-Taylor, J., Smola, A. J., & Williamson,
639 R. C. (2001). Estimating the support of a high-dimensional distribution.
640 *Neural computation*, *13*, 1443–1471.
- 641 Schölkopf, B., Smola, A., & Müller, K.-R. (1997). Kernel principal compo-
642 nent analysis. In *International Conference on Artificial Neural Networks*
643 (pp. 583–588). Springer.
- 644 Schölkopf, B., & Smola, A. J. (2002). *Learning with Kernels*. MIT.
- 645 Sofman, B., Neuman, B., Stentz, A., & Bagnell, J. A. (2011). Anytime online
646 novelty and change detection for mobile robots. *Journal of Field Robotics*,
647 *28*, 589–618.
- 648 Van Der Maaten, L. (2014). Accelerating t-sne using tree-based algorithms.
649 *The Journal of Machine Learning Research*, *15*, 3221–3245.
- 650 Weinberger, K. Q., Blitzer, J., & Saul, L. K. (2006). Distance metric learning
651 for large margin nearest neighbor classification. In *Advances in neural*
652 *information processing systems* (pp. 1473–1480).
- 653 Weinberger, K. Q., & Saul, L. K. (2009). Distance metric learning for large
654 margin nearest neighbor classification. *Journal of Machine Learning Re-*
655 *search*, *10*, 207–244.
- 656 Zhang, J., & Wang, H. (2006). Detecting outlying subspaces for high-
657 dimensional data: the new task, algorithms, and performance. *Knowledge*
658 *and information systems*, *10*, 333–355.