

3.1 Motivations

This work focuses on defining a 'good' metric for classification of time series. The definition of a metric to compare samples is a fundamental issue in data analysis or machine learning. As seen in Chapter 2, temporal data may be compared based on one or several characteristics, called **modalities** (amplitude, behavior, frequency) and they might be subjected to delays. In some classification applications, the most discriminative characteristic between time series of different classes can be localized on a smaller part of the signal (scale). We believe that the definition of a temporal metric should consider at least these different aspects (modality, delay, scale) in order to improve the performance of a classifier. Fig. 3.1 illustrates a result obtained with our proposition. There is a significant improvement in classification performances by taking into account in the metric definition, several modalities (amplitude d_A , behavior d_B , frequency d_F) located at different scales (illustrated by black rectangles in the figure). The performance of the learned combined metric is compared with the ones of the standard metrics that take into account for each, only one modality on a global scale (involving all time series elements).

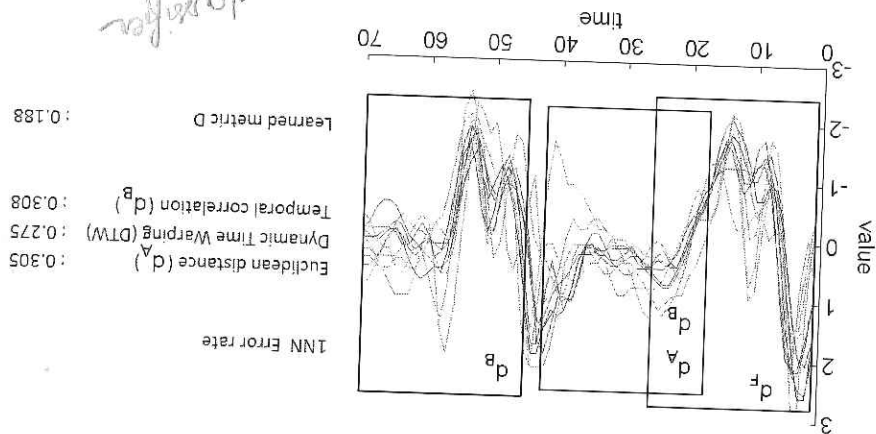


Figure 3.1: SonyAIBO dataset and error rate using a k NN ($k = 1$) with standard metrics (Euclidean distance, Dynamic Time Warping, temporal correlation) and a learned combined metric D . The figure shows the 4 major metrics involve in the combined metric D and their respective temporal scale (black rectangles).

Our aim is to take leverage from the metric learning framework [WS09b; BHS12] to learn a multi-modal and multi-scale temporal metric for time series nearest neighbors classification. Specifically, our objective is to learn from the data a linear or non linear function that combines several temporal modalities at several temporal scales, that satisfies metric properties (Section 2.2), and that generalizes the case of unimodal metrics at the global scale. ¹ Metric learning can be defined as learning, from the data and for a task, a pairwise function (i.e., a similarity, dissimilarity or a distance) that brings closer samples that are expected to be similar, and pushes far away those expected to be dissimilar. Such similarity and dissimilarity expectations, is inherently task- and application-dependent, generally given *a priori* and fixed during the learning process. Metric learning has become an active area of research in the

applying the k -NN classification:

$$(3.2) \quad D_L^2(\mathbf{x}_i, \mathbf{x}_j) = D^2(\mathbf{L}\mathbf{x}_i, \mathbf{L}\mathbf{x}_j)$$

$$(3.3) \quad D_L^2(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{L}(\mathbf{x}_i - \mathbf{x}_j)\|_2^2$$

Commonly, the squared distances can be expressed in terms of a square matrix:

$$(3.4) \quad D_L^2(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{L}^T \mathbf{L} (\mathbf{x}_i - \mathbf{x}_j)$$

Let $\mathbf{M} = \mathbf{L}^T \mathbf{L}$. It is proved that any matrix \mathbf{M} formed as below from a real-valued matrix \mathbf{L} is positive semidefinite (i.e., no negative eigenvalues) [WSS09a]. Using the matrix \mathbf{M} , squared distances can be expressed as:

$$(3.5) \quad D_M^2(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j)$$

The computation of the learned metric D_M can thus be seen as a two steps procedure: first, it computes a linear transformation of the samples \mathbf{x}_i given by the transformation \mathbf{L} ; second, it computes the Euclidean distance in the transformed space:

$$(3.6) \quad D_M^2(\mathbf{x}_i, \mathbf{x}_j) = D_L^2(\mathbf{L}\mathbf{x}_i, \mathbf{L}\mathbf{x}_j)$$

Learning the linear transformation \mathbf{L} is thus equivalent to learn the corresponding Mahalanobis metric D parametrized by \mathbf{M} . This equivalence leads to two different approaches to metric learning: we can either estimate the linear transformation \mathbf{L} , or estimate a positive semidefinite matrix \mathbf{M} . LMNN solution refers on the latter one.

Mathematically, the metric learning problem can be formalized as an optimization problem involving two terms for each sample \mathbf{x}_i : one term penalizes large distances between nearby inputs with the same label (pull), while the other term penalizes small distances between inputs with different labels (push). For all samples \mathbf{x}_i , this implies a minimization problem:

$$(3.7) \quad \left\{ \begin{array}{l} \arg \min_{\mathbf{M}, \xi} \sum_{i,j \sim i} D_M^2(\mathbf{x}_i, \mathbf{x}_j) + C \sum_{i,j \not\sim i} \xi_{ijl} \\ \text{s.t. } \forall j \sim i, l \neq i, \\ D_M^2(\mathbf{x}_i, \mathbf{x}_j) - D_M^2(\mathbf{x}_i, \mathbf{x}_l) \geq 1 - \xi_{ijl} \\ \xi_{ijl} \geq 0 \\ \mathbf{M} \succeq 0 \end{array} \right.$$

where ξ_{ijl} are slack variables, C is a trade-off between the push and pull term and $\mathbf{M} \succeq 0$ means that \mathbf{M} is a positive semidefinite matrix. Generally, the parameter C is tuned via cross validation and grid search (Section 1.1.2). Similarly to Support Vector Machine (SVM) approach, slack variables ξ_{ijl} are introduced to relax the optimization problem.

*Don pousser on a appelé cette méthode
qui est-ce qui on va en faire?*

In that space, the norm of a pairwise vector $\|x_{ij}\|$ refers to the proximity between the time series x_i and x_j . In particular, if $\|x_{ij}\| = 0$ then x_j is identical to x_i according to all metrics d_h .

3.3.2 Interpretation in the pairwise dissimilarity space

In this section, we give more detailed interpretations in the dissimilarity space. We recall that the norm of a pairwise vector is given by:

$$\|x_{ij}\| = \sqrt{\sum_{h=1}^p d_h(x_i, x_j)^2} \quad (3.10)$$

In the following, we denote the norm $\|x_{ij}\|$ as an initial distance in the dissimilarity space and call it D_0 . Any other initial metric could have been chosen. The norm of a pairwise vector x_{ij} can be interpreted as a proximity measure: the lower the norm of x_{ij} is, the closer are the time series x_i and x_j . Two pairwise vectors x_{ij} and x_{kl} that are on a same line that passes through the origin $x_{ii} = 0$ represent differences in the same proportions between their respective modalities (Fig. 3.4).

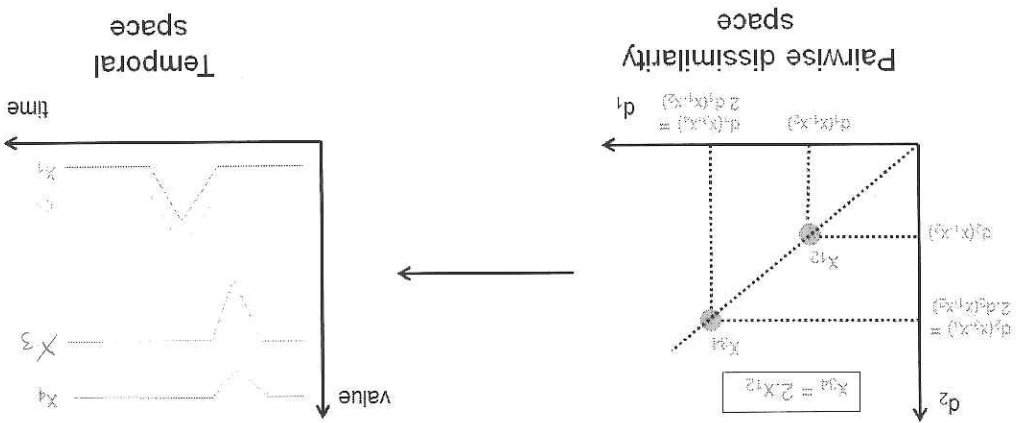


Figure 3.4: Example of interpretation of two pairwise vectors x_{12} and x_{34} on a same line passing through the origin in the pairwise dissimilarity space.

The Euclidean distance $\sqrt{\sum_{h=1}^p (d_h(x_i, x_j) - d_h(x_k, x_l))^2}$ between two pairwise vectors x_{ij}

and x_{kl} represents the similarity between the differences among the same modalities, in the same proportions. Note that if the Euclidean distance is close to 0 (x_{ij} and x_{kl} are close in the dissimilarity space), it doesn't mean that the time series x_i, x_j, x_k and x_l are similar. Fig 3.5 shows an example of two pairwise vectors x_{ij} and x_{kl} close together in the pairwise space. However, in the temporal space, the time series x_1 and x_3 are not similar for example. It means that x_i is as similar to x_j as x_k is to x_l , i.e., the distance D_0 between x_i and x_j is nearly the same than the distance D_0 between x_k and x_l .

l'important c'est de passer de l'un à l'autre

Je ne comprends pas bien l'intuit de ce passage

Je -> distance part 0?

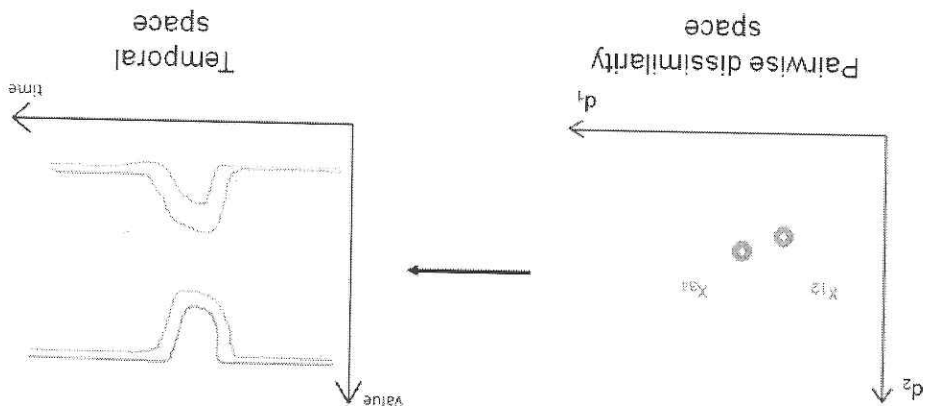


Figure 3.5: Example of two pairwise vectors x_{12} and x_{34} close in the pairwise dissimilarity space. However, the time series x_1 and x_3 are not similar in the temporal space.

3.3.3 Multi-scale description for time series

The multi-modal representation in the dissimilarity space can be enriched for time series by measuring each unimodal metric d_h at different scales. Note that the distance measures (amplitude-based d_A , frequential-based d_F , behavior-based d_B) in Eqs. 2.1, 2.4 and 2.6 implies systematically the total time series elements x_n and thus, restricts the distance measures to capture local temporal differences. In this work, we provide a multi-scale framework for time series comparison using a hierarchical structure. Many methods exist in the literature such as the sliding window [Keo+03] or the dichotomy [DCA11]. We detail here the latter one.

A multi-scale description can be obtained by repeatedly segmenting a time series expressed at a given temporal scale to induce its description at a more local level. Many approaches have been proposed assuming fixed either the number of the segments or their lengths [Fu11]. In this work, we consider a binary segmentation at each level. Let $I = [a; b]$ be a temporal interval of size $(b - a)$. The interval I is decomposed into two equal overlapped intervals I_L (left interval) and I_R (right interval). A parameter μ allows to overlap the two intervals I_L and I_R , covering discriminating subsequences in the central region of I (around $\frac{b+a}{2}$): $I = [a; b]; I_L = [a; a + \mu(b - a)]; I_R = [b - \mu(b - a); b]$. For $\mu = 0.6$, the overlap covers 10% of the size of the interval I . Then, the process is repeated on the intervals I_L and I_R . We obtain a set of intervals I_s illustrated in Fig. 3.6.

A multi-scale dissimilarity description between two time series is obtained by computing the usual time series metrics (d_A, d_B, d_F) on each of the resulting segments I_s . Note that for two time series x_i and x_j , the comparison between x_i and x_j is done on the same interval I_s . For a multi-scale amplitude-based measures $d_{I_s}^A$ is $\{d_{I_s1}^A, d_{I_s2}^A, \dots\}$ where $d_{I_s}^A$ is defined as:

$$d_{I_s}^A(x_i, x_j) = \sqrt{\sum_{t \in I_s} (x_{it} - x_{jt})^2} \quad (3.11)$$

The local behaviors- and frequential- based measures $d_{I_s}^B$ and $d_{I_s}^F$ are obtained similarly.

Then define a
multi-scale
on localised temporal segments

segmentation
as a temp

scale → temporal localisation

Our proposition is inspired from the LMNN framework where the optimization problem involves a pull term that penalizes large distances between sample of same labels ($Pull_i$). It can be interpreted as a regularization term on $Pull_i$. In LMNN, the push term penalizes small distances between samples of different labels ($Push_i$). It can be interpreted as a loss term on $Push_i$. To ensure a safety margin between similar and dissimilar samples, a constraint is added:

$$D^2(\mathbf{x}_i, \mathbf{x}_l) - D^2(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \xi_{ijl}.$$

Similarly, we formalize the M²TML problem as an optimization problem involving both a regularization term on D and the pull set $Pull_i$, denoted $R_{Pull}(D)$, and a loss term on ξ and the push set $Push_i$, denoted $L_{Push}(\xi)$. A set of constraints is added to control the push term in order to have a large margin between $Pull_i$ and $Push_i$:

$$\begin{aligned} \text{argmin}_{D, \xi} \{ & R_{Pull}(D) + L_{Push}(\xi) \} \\ \text{s.t. } & \forall i, j, l \in Push_i, l \in Push_i, \\ & D(\mathbf{x}_{il}) - D(\mathbf{x}_{ij}) \geq 1 - \xi_{ijl} \end{aligned} \quad (3.12)$$

Among the possibilities for the regularization term, we decide to choose to minimize the sum of the distances of the pull pairs. Among the possibilities for the loss term, we decide to choose to minimize the sum of the slack variables on the push pairs:

$$R_{Pull}(D) = \sum_{j \in Pull_i} D(\mathbf{x}_{ij}) \quad (3.13)$$

$$L_{Push}(\xi) = \sum_{j \in Pull_i, l \in Push_i} \xi_{ijl} \quad (3.14)$$

The M²TML problem for large margin k -NN classification can be written as the following optimization problem:

$$\begin{aligned} \text{argmin}_{D, \xi} \left\{ \sum_{j \in Pull_i} D(\mathbf{x}_{ij}) + C \sum_{j \in Pull_i, l \in Push_i} \xi_{ijl} \right\} \\ \text{s.t. } \forall i = 1, \dots, n, \forall j \in Pull_i, l \in Push_i, \\ D(\mathbf{x}_{il}) - D(\mathbf{x}_{ij}) \geq 1 - \xi_{ijl} \\ \xi_{ijl} \geq 0 \end{aligned} \quad (3.15)$$

where ξ_{ijl} are the slack variables and C , the trade-off between the pull (regularization) and push (loss) costs. In the next section, we detail different strategies to define the $Pull_i$ and $Push_i$ sets.

3.4.2 Push and pull set definition

To build the pairwise training set, we associate for each \mathbf{x}_i , two sets, $Pull_i$ and $Push_i$, where the two sets are chosen according to one of the following strategies, illustrated in Fig 3.7. Recall that the norm $D_0(\mathbf{x}_{ij}) = \|\mathbf{x}_{ij}\|_2$ is set as our initial distance D_0 .

1. k -NN vs impostors: for a given \mathbf{x}_i , the sets of pairs to pull and to push corresponds respectively to:

$$\forall i \in 1, \dots, n, \quad Pull_i = \{\mathbf{x}_{ij} / y_j = y_i, D_0(\mathbf{x}_{ij}) \text{ is among the } k\text{-lowest distance}\} \quad (3.16)$$

$$Push_i = \{\mathbf{x}_{il} / y_l \neq y_i, D_0(\mathbf{x}_{il}) \leq \max_{\mathbf{x}_{ij} \in Pull_i} D_0(\mathbf{x}_{ij})\} \quad (3.17)$$

2. k -NN vs all: for a given \mathbf{x}_i , the sets of pairs to pull and to push corresponds respectively to:

$$\forall i \in 1, \dots, n, \quad Pull_i = \{\mathbf{x}_{ij} / y_j = y_i, D_0(\mathbf{x}_{ij}) \text{ is among the } k\text{-lowest distance}\} \quad (3.18)$$

$$Push_i = \{\mathbf{x}_{il} / y_l \neq y_i\} \quad (3.19)$$

3. m -NN⁺ vs m -NN⁻: for a given \mathbf{x}_i , the pull and push sets are defined respectively as the set of the m -nearest neighbors of the same class ($y_j = y_i$), and the m -nearest neighbor of \mathbf{x}_i of a different class ($y_j \neq y_i$). More precisely, our proposition states: $m = \alpha \cdot k$ with $\alpha \geq 1$. Other propositions for m are possible:

$$\forall i \in 1, \dots, n, \quad Pull_i = \{\mathbf{x}_{ij} / y_j = y_i, D_0(\mathbf{x}_{ij}) \text{ is among the } m\text{-lowest distance}\} \quad (3.20)$$

$$Push_i = \{\mathbf{x}_{il} / \text{s.t. } y_l \neq y_i, D_0(\mathbf{x}_{il}) \text{ is among the } m\text{-lowest distance}\} \quad (3.21)$$

In the following, we denote m -NN⁺ = $\bigcup Pull_i$ and m -NN⁻ = $\bigcup Push_i$

Finally, let discuss about the similarities and differences between LMNN (Weinberger & Saul [WS09a]) and our M²TML proposition. In LMNN, the sets $Pull_i$ and $Push_i$ are defined according the k -NN vs impostors strategy (Eqs. 3.16 & 3.17) and may be unbalanced. The sets are defined and fixed during the optimization process according to the initial metric D_0 . In M²TML the sets $Pull_i$ and $Push_i$ are defined according the m -NN⁺ vs m -NN⁻ strategy (Eqs. 3.20 & 3.21) and are balanced. The sets are defined and fixed during the optimization process according to the initial metric D_0 , but the m -neighborhood is larger than the k -neighborhood. By considering a neighborhood larger than the k -neighborhood, we believe that the generalization properties of the learned metric D will be improved.

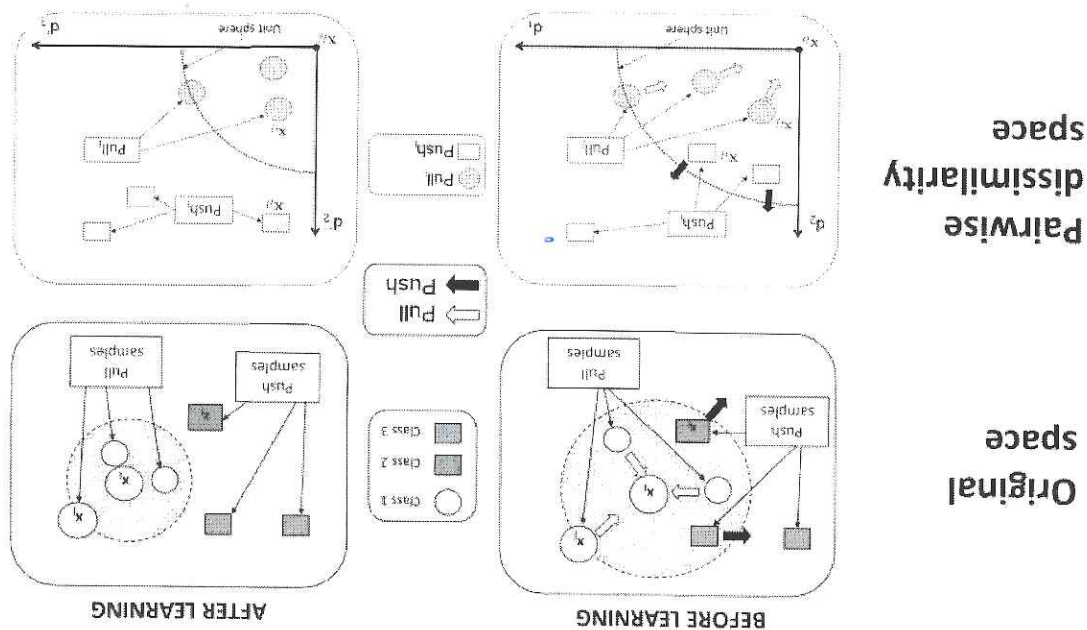


Figure 3.8: Metric learning problem in the original space (top) and the pairwise dissimilarity space (bottom) for a $k = 3$ neighborhood of x_i . Before learning (left), push samples x_i invade the targets perimeter x_j . In the dissimilarity pairwise space, this is equivalent to have push pairwise vectors x_{ij} with an initial distance d_0 lower than the distance of pull pairwise vectors x_{ij} . The aim of M^2TML is to learn a metric D to push x_{ij} (black arrow) and pull x_{ij} from the origin (white arrow).

- If $D(x_{ij}) < D(x_{ij})$, then the pairs x_{ij} is an impostor pair that invades the neighborhood of the target pairs x_{ij} . The slack variable $\xi_{ij} > 1$ will be penalized in the objective function.
- If $D(x_{ij}) \leq D(x_{ij})$ but $D(x_{ij}) > D(x_{ij}) + 1$, the pair x_{ij} is within the safety margin of the target pairs x_{ij} . The slack variable $\xi_{ij} \in [0, 1]$ will have a small penalization effect in the objective function.
- If $D(x_{ij}) > D(x_{ij}) + 1$, $\xi_{ij} = 0$ and the slack variable has no effect in the objective function.

In the following, we propose different regularizers for the pull term $R_{pull}(D)$. First, we use a linear regularization. Secondly, we use a quadratic regularization that enables to extend the approach to learn non-linear function for D by using the "kernel" trick. Thirdly, we formulate the problem as a SVM problem to solve a large margin problem between P_{pull_i} and P_{push_i} sets, and then, we define the combined metric D based on the SVM solution. Finally, we sum up the retained solution (SVM-based solution) and give the main steps of the algorithm.

3.7.2 Solution for the linearly separable Pull and Push sets

Let $\mathbf{x}_{i, test}$ be a new sample, $\mathbf{x}_{i, test} \in \mathcal{E}$ gives the proximity between \mathbf{x}_i and $\mathbf{x}_{i, test}$ based on the p multi-modal and multi-scale metrics d_h . We review in this section different interpretations in the dissimilarity space.

M²TML metric definition
 Given a test pair $\mathbf{x}_{i, test}$, the norm of the pair allows to estimate the proximity between \mathbf{x}_i and $\mathbf{x}_{i, test}$. In particular, for M²TML, two quantities are used to define the dissimilarity measure: the projected norm and the distance to the margin.

Let denote $\mathbf{P}^w(\mathbf{x}_{i, test})$, the orthogonal projection of $\mathbf{x}_{i, test}$ on the axis of direction \mathbf{w} :

$$\mathbf{P}^w(\mathbf{x}_{i, test}) = \frac{\langle \mathbf{w}, \mathbf{x}_{i, test} \rangle}{\|\mathbf{w}\|^2} \mathbf{w} = \frac{\mathbf{w}^T \mathbf{x}_{i, test}}{\|\mathbf{w}\|^2} \mathbf{w} \quad (3.47)$$

The projected norm $\|\mathbf{P}^w(\mathbf{x}_{i, test})\|$ of $\mathbf{x}_{i, test}$ on the direction \mathbf{w} limits the comparison of \mathbf{x}_i and $\mathbf{x}_{i, test}$ to the features separating pull and push sets (Fig. 3.9), it is defined as:

$$\|\mathbf{P}^w(\mathbf{x}_{i, test})\| = \frac{\|\mathbf{w}\|}{\|\mathbf{w}^T \mathbf{x}_{i, test}\|} \quad (3.48)$$

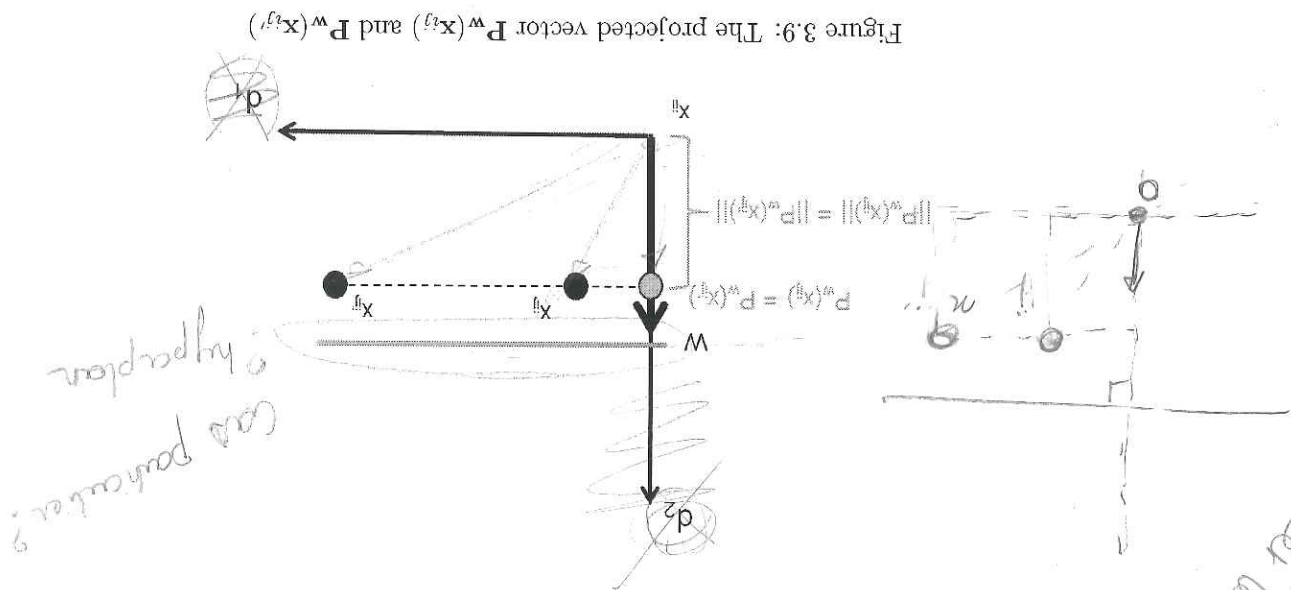


Figure 3.9: The projected vector $\mathbf{P}^w(\mathbf{x}_{ij})$ and $\mathbf{P}^w(\mathbf{x}_{ij'})$

Although the norm $\|\mathbf{P}^w(\mathbf{x}_{i, test})\|$ satisfies positivity, it doesn't guarantee lower distances for pull pairs than for push pairs as illustrated in Fig 3.10.

Note that the distance of the projection to the margin $\mathbf{w}^T \mathbf{P}^w(\mathbf{x}_{i, test}) + b$ gives the membership of the projected vector $\mathbf{P}^w(\mathbf{x}_{i, test})$ in the pull or push side. However, it can't used as a dissimilarity (non-positivity).

3.8 SVM-based solution and algorithm for M²TML

In this section, we review the main steps of the retained SVM solution. In particular, we detail two pre-processing steps needed to adapt the SVM framework to our metric learning problem that are the pairwise space normalization and the neighborhood scaling.

Pairwise space normalization. The scale between the p basic metrics d_h can be different. Thus, there is a need to scale the data within the pairwise space and ensure comparable ranges for the p basic metrics d_h . In our experiment, we use dissimilarity measures with values in $[0; +\infty]$. Therefore, we propose to Z-normalize their log distributions as explained in Section 2.5.2.

Neighborhood scaling. In real datasets, local neighborhoods may have very different scales as illustrated in Fig. 3.12. To make the pull neighborhood spreads comparable, we propose for each \mathbf{x}_i to scale each pairs \mathbf{x}_{ij} such that the L_2 norm (radius) of the farthest m -th nearest neighbor is 1:

$$\mathbf{x}_{ij}^{\text{norm}} = \left[\frac{d_1(\mathbf{x}_i, \mathbf{x}_j)}{r_i}, \dots, \frac{d_p(\mathbf{x}_i, \mathbf{x}_j)}{r_i} \right]^T \quad (3.52)$$

where r_i is the radius associated to \mathbf{x}_i corresponding to the maximum norm of its m -th nearest neighbor of same class in P_{pull_i} :

$$r_i = \max_{\mathbf{x}_{ij} \in P_{\text{pull}_i}} D_0(\mathbf{x}_{ij}) \quad (3.53)$$

For simplification purpose, we denote \mathbf{x}_{ij} as $\mathbf{x}_{ij}^{\text{norm}}$. Fig. 3.12 illustrates the effect of neighborhood scaling in the dissimilarity space.

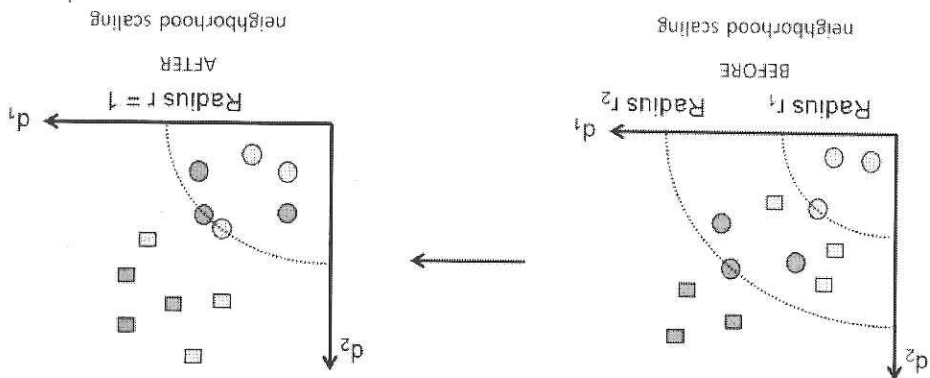


Figure 3.12: Effect of neighborhood scaling before (left) and after (right) on the neighborhood of two time series \mathbf{x}_1 (green) and \mathbf{x}_2 (red). Circle represent pairs P_{pull_i} and square represents pairs P_{push_i} for $m = 3$ neighbors. Before scaling, the problem is not linearly separable with a global SVM approach and the spread of each neighborhood are not comparable. After scaling, the target neighborhood becomes comparable and in this example, the problem becomes linearly separable.