

Metric Learning in dissimilarity space: formalization

Sommaire

3.1	Motivations	51
3.2	Dissimilarity space	54
3.2.1	Pairwise embedding	55
3.2.2	Interpretation in the dissimilarity space	55
3.3	Linear Programming (LP) formalization	56
3.4	Quadratic Programming (QP) formalization	59
3.5	Support Vector Machine (SVM) approximation	62
3.5.1	Motivations	62
3.5.2	Similarities and differences in the constraints	63
3.5.3	Similarities and differences in the objective function	64
3.5.4	Geometric interpretation	65
3.6	Conclusion of the chapter	66

In this chapter, we formalize the problem of Metric Learning in Dissimilarity space (MLD). We first motivate the problem of learning a metric that combines several metrics at different scales for a robust k -NN classifier. Secondly, we introduce the concept of dissimilarity space. Finally, we transpose the metric learning problem in the dissimilarity space and propose three possible formulations: Linear programming, Quadratic programming and SVM-based approximation.

3.1 Motivations

The definition of a metric to compare samples is a fundamental issue in data analysis, pattern recognition or machine learning. Contrary to static data, temporal data are more complex: they may be compared not only on their amplitudes but also on their dynamic, frequential spectrum or other inherent characteristics. For time series comparison, a large number of metrics have been proposed, most of them are designed to capture similitudes and differences based on one temporal modality. For amplitude-based comparison, measures cover variants

of Mahalanobis distance or the dynamic time warping (DTW) to cope with delays [BC94b]; [Rab89]; [SC78b]; [KL83]. Other propositions refer to temporal correlations or derivative dynamic time warping for behavior-based comparison [AT10b]; [RBK08]; [CCP06]; [KP01]; [DM09]. For frequential aspects, comparisons are mostly based on the Discret Fourier or Wavelet Transforms [SS12a]; [KST98]; [DV10]; [Zha+06]. A detailed review of the major metrics is proposed in [MV14]. In general, the most discriminant modality (amplitude, behavior, frequency, etc.) varies from a dataset to another.

In some applications, the most discriminative characteristic between time series of different classes can be localized on a smaller part of the signal. Involving totally or partially time series elements, rather than systematically the whole elements should be taken into consideration in the metric definition, a crucial key to localize discriminative features. In other applications, the combinations of several factors (modality, scale) can improve the performance of the classifier. Thus, there is a need for combined metrics.

Some works propose to combine several modalities through a priori models as in [DCDG10]; [DCA12]; [SB08]. Ideally, a combined metric should answer two scenarios depending on the datasets: 1) combines several modalities at several scales; 2) select one modality at one particular scale and thus, coming back to a uni-modal and uni-scale metric framework. Figs. 3.1 and 3.2 shows these two cases on two dataset examples used in classification of univariate time series. For SonyAIBO dataset (Fig. 3.1), by learning the metric, several modalities (frequential d_F , behavior d_B and amplitude d_A) at different locally temporal interval are combined together. Thanks to this combination, the error of the 1-NN classifier is decreased: global uni-modal metrics achieves 0.305 for d_A , 0.308 for d_B , 0.258 for d_F ; and the combined metric achieves a score of 0.188. For ECG200 dataset (Fig. 3.2), the learned metric mainly includes the global behavior component d_B and the error rate remains statistically the same. More detailed explanations will be given in Chapter 6.

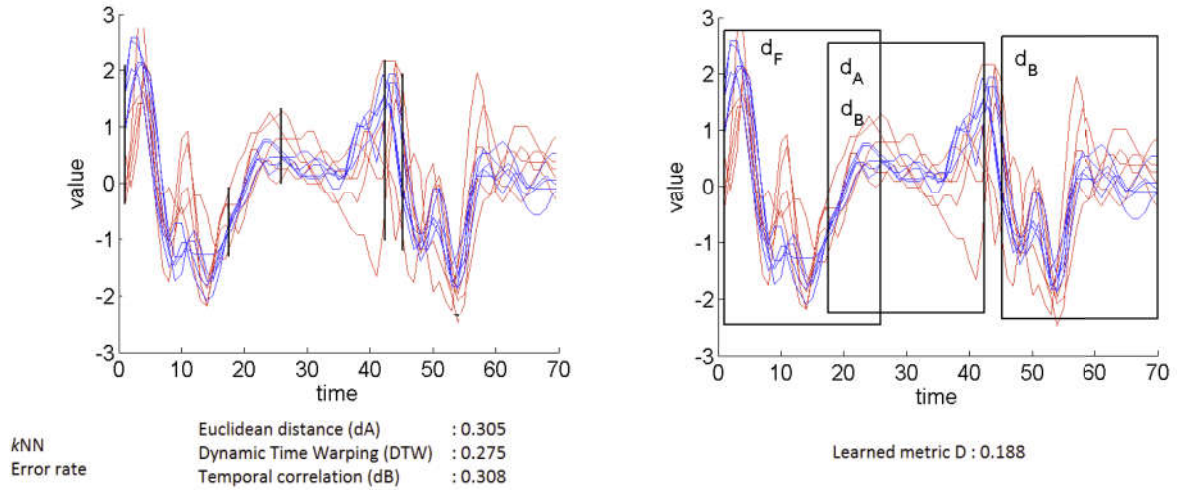


Figure 3.1: SonyAIBO dataset and error rate using a k NN with $k = 1$ with standard metrics (Euclidean distance, Dynamic Time Warping, temporal correlation) (left) and a learned combined metric (right). For the learned combined metric D , the figure shows the 4 major metrics involves in the combination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

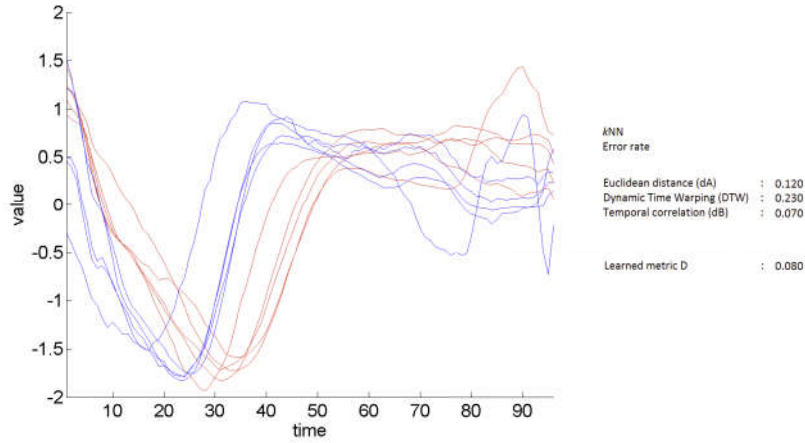


Figure 3.2: ECG200 dataset and error rate using a k NN with $k = 1$ with standard metrics (Euclidean distance, Dynamic Time Warping, temporal correlation) and a learned combined metric. For the learned combined metric D , the major discriminant feature is the behavior-based metric d_B (90% of D) computed on a global scale (including all time series elements).

Our aim is to take benefice of 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, and that satisfies metric properties (Section 2.2).

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) to make closer samples that are expected to be similar, and far away those expected to be dissimilar. Similar and dissimilar samples, are 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 last decades for various machine learning problems (supervised, semi-supervised, unsupervised, online learning) [??] and has received many interests in its theoretical background (generalization guarantees) [BHS13]. From the surge of recent research in metric learning, one can identify mainly two categories: the linear and non linear approaches. The former is the most popular, it defines the majority of the propositions, and focuses mainly on the Mahalanobis distance learning [WS09a]. The latter addresses non linear metric learning which aims to capture non linear structure in the data. In Kernel Principal Component Analysis (KPCA) [ZY10]; [Cha+10], the aim is to project the data into a non linear feature space and learn the metric in that projected space. In Support Vector Metric Learning (SVML) approach [XWC12], the Mahalanobis distance is learned jointly with the learning of the SVM model in order to minimize the validation error. In general, the optimization problems are more expensive to solve, and the methods tends to favor overfitting as the constraints are generally easier to satisfy in a nonlinear kernel space. A more detailed review is done in [BHS13].

Contrary to static data, metric learning for structured data (e.g. sequence, time series, trees, graphs, strings) remains less numerous. While for sequence data most of the works focus on string edit distance to learn the edit cost matrix [OS06]; [BHS12], metric learning for

time series is still in its infancy. Without being exhaustive, major recent proposals rely on weighted variants of dynamic time warping to learn alignments under phase or amplitude constraints [Rey11]; [JJO11]; [ZLL14], enlarging alignment learning framework to multiple temporal matching guided by both global and local discriminative features [FDCG13]. For the most of these propositions, temporal metric learning process is systematically: a) Uni-modal (amplitude-based), the divergence between aligned elements being either the Euclidean or the Mahalanobis distance and b) Uni-scale (global level) by involving the whole time series elements, which restricts its potential to capture local characteristics. Bellet & al. enlightened in [BHS13] perspectives for metric learning, especially, the learning of richer metrics that could take into account of the multi-modality within the data.

We propose in this work to learn a multi-modal and multi-scale temporal metric for a robust k -NN classifier. For this, the main idea is to embed time series into a dissimilarity space [PPD02]; [DP12] where a linear function combining several modalities at different temporal scales can be learned, driven jointly by a SVM and nearest neighbors metric learning framework [WS09b]. Thanks to the "Kernel trick", the proposed solution is extended to non-linear temporal metric learning context. A sparse and interpretable variant of the algorithm confirms its ability to localize finely discriminative modalities as well as their temporal scales. In the following, the term metric is used to reference both a distance or a dissimilarity measure.

In this chapter, we first present the concept of dissimilarity space. Then, we formalize the metric learning problem in the dissimilarity space and formalize the problem under three formulations: Linear Programming, Quadratic Programming, SVM approximation. Note that these formulations doesn't concern only time series and can be applied to any type of data. In the next chapter, we detailed our proposed solution to learn a multi-modal and multi-scale temporal metric for a robust k -NN classifier.

3.2 Dissimilarity space

Let $d_1, \dots, d_h, \dots, d_p$ be p given metrics that allow to compare samples. For instance, in Chapter 2, we have proposed three types of metrics for time series: amplitude-based d_A , behavior-based d_B and frequential-based d_F . Our objective is to learn a metric D that combines the p metrics in order to optimize the performance of a k -NN classifier. Formally:

$$D = f(d_1, \dots, d_p) \quad (3.1)$$

In this section, we first introduce the dissimilarity space. Then, we give some interpretations in the dissimilarity space.

3.2.1 Pairwise embedding

The computation of a metric d , and of course D , always takes into account a pair of samples $(\mathbf{x}_i, \mathbf{x}_j)$. We introduce a new space representation referred as the **dissimilarity space**. In this new space, illustrated in Fig. 4.1, a vector \mathbf{x}_{ij} represents a pair of time series $(\mathbf{x}_i, \mathbf{x}_j)$ described by the p unimodal metrics d_h : $\mathbf{x}_{ij} = [d_1(\mathbf{x}_i, \mathbf{x}_j), \dots, d_p(\mathbf{x}_i, \mathbf{x}_j)]^T$. We denote N the number of pairwise vectors \mathbf{x}_{ij} generated by this embedding.

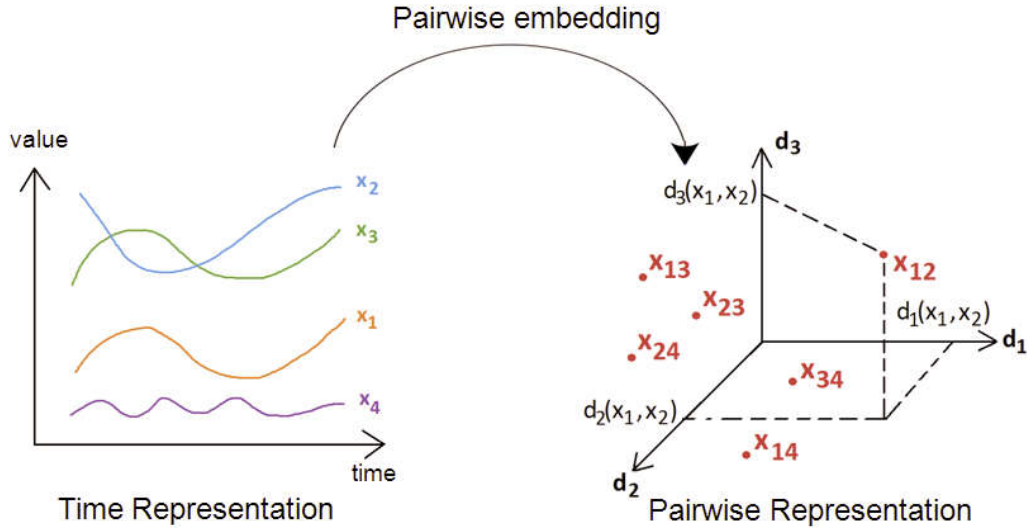


Figure 3.3: Example of embedding of time series \mathbf{x}_i from the temporal space (left) into the dissimilarity space (right). In this example, a pair of time series $(\mathbf{x}_1, \mathbf{x}_2)$ is projected into the dissimilarity space as a vector \mathbf{x}_{12} described by $p = 3$ basic metrics: $\mathbf{x}_{12} = [d_1(\mathbf{x}_1, \mathbf{x}_2), d_2(\mathbf{x}_1, \mathbf{x}_2), d_3(\mathbf{x}_1, \mathbf{x}_2)]^T$.

A combination function D of the metrics d_h can be seen as a function in this space. In the following, we propose first to use a linear combination of d_h : $D(\mathbf{x}_i, \mathbf{x}_j) = \sum_h w_h \cdot d_h(\mathbf{x}_i, \mathbf{x}_j)$. For simplification purpose, we denote $D(\mathbf{x}_i, \mathbf{x}_j) = D(\mathbf{x}_{ij})$ and the pairwise notation gives:

$$D(\mathbf{x}_i, \mathbf{x}_j) = D(\mathbf{x}_{ij}) = \mathbf{w}^T \mathbf{x}_{ij} \quad (3.2)$$

where \mathbf{w} is the vector of weights w_h : $\mathbf{w} = [w_1, \dots, w_p]^T$.

3.2.2 Interpretation in the dissimilarity space

The interpretation of the data in the dissimilarity space is particular since the dissimilarity space is not a standard Euclidean space. The interpretation in this space requires to be careful.

If $\mathbf{x}_{ij} = \mathbf{0}$ then \mathbf{x}_j is identical to \mathbf{x}_i according to all metrics d_h . The norm of the vector \mathbf{x}_{ij} can be interpreted as a proximity measure: the lower the norm of \mathbf{x}_{ij} is, the closer are the time series \mathbf{x}_i and \mathbf{x}_j . Nevertheless, if two pairwise vectors \mathbf{x}_{ij} and \mathbf{x}_{kl} has their norms closed, it doesn't mean that the time series \mathbf{x}_i , \mathbf{x}_j , \mathbf{x}_k and \mathbf{x}_l are similar. Fig 4.5 shows an example of

two pairwise vectors \mathbf{x}_{ij} and \mathbf{x}_{kl} that are close together in the dissimilarity space. However, in the temporal space, the time series \mathbf{x}_1 and \mathbf{x}_3 are not similar for example. It means that \mathbf{x}_i is as similar to \mathbf{x}_j as \mathbf{x}_k is to \mathbf{x}_l .

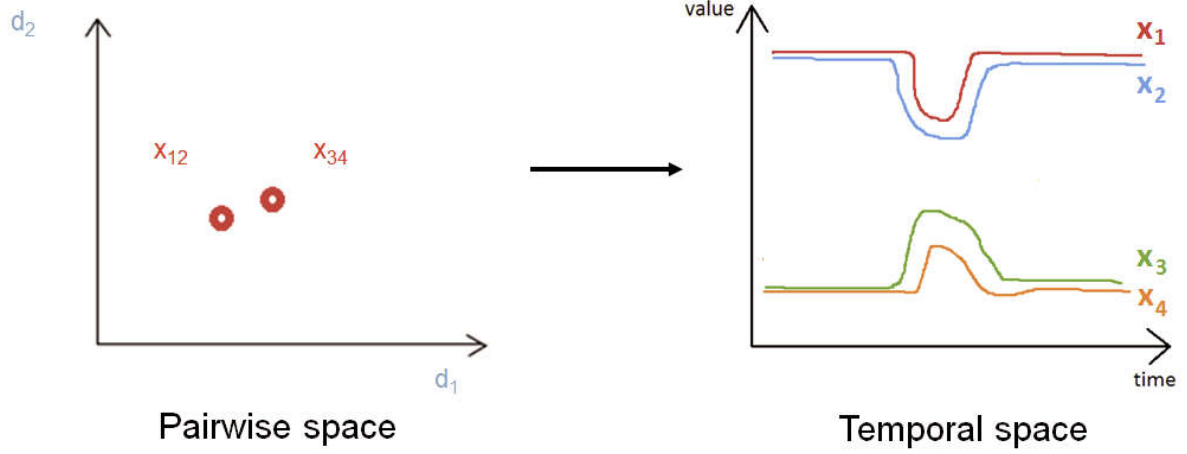


Figure 3.4: Example of two pairwise vectors \mathbf{x}_{12} and \mathbf{x}_{34} close in the dissimilarity space. However, the time series \mathbf{x}_1 and \mathbf{x}_3 are not similar in the temporal space.

A metric D that combines the p metrics d_1, \dots, d_p can be seen as a function of the dissimilarity space. It can be noticed that when the time series \mathbf{x}_i are embedded in the pairwise, the information of their original class y_i is lost. Any multi-class problem is transformed in the dissimilarity space as a binary classification problem.

In the next sections, we transpose the metric learning problem for large margin nearest neighbor classifier in the dissimilarity space. We propose three formulations: Linear programming, Quadratic programming and SVM-based approach.

3.3 Linear Programming (LP) formalization

Let $\{\mathbf{x}_i, y_i\}_{i=1}^n$ be a set of n time series labeled y_i . We embed the n time series in the dissimilarity space and define $\{\mathbf{x}_{ij}, y_{ij}\}$ as the set of N pairwise vectors \mathbf{x}_{ij} described by p metrics d_h and labeled $y_{ij} = +1$ if $y_j \neq y_i$ and -1 otherwise. Our objective is to define a metric D as a linear combination of the p metrics d_h (Eq. 4.2). In the dissimilarity space, the metric D should:

- **pull** to the origin the k nearest neighbors pairs \mathbf{x}_{ij} of same labels ($y_{ij} = -1$)
- **push** from the origin all the pairs \mathbf{x}_{il} of different classes ($y_{ij} = +1$)

Let \mathbf{X}_{tar} be a $p \times (k.n)$ matrix containing all targets \mathbf{x}_{ij} . Fig. 4.6 illustrates our idea to learn the metric D . For each time series \mathbf{x}_i , we build the set of target pairs \mathbf{x}_{ij} ($j \rightsquigarrow i$) and the set

of pairs \mathbf{x}_{il} of different class. Then, we optimize the weight vector \mathbf{w} so that the pairs \mathbf{x}_{ij} are pulled to the origin and the pairs \mathbf{x}_{il} are pushed from the origin.

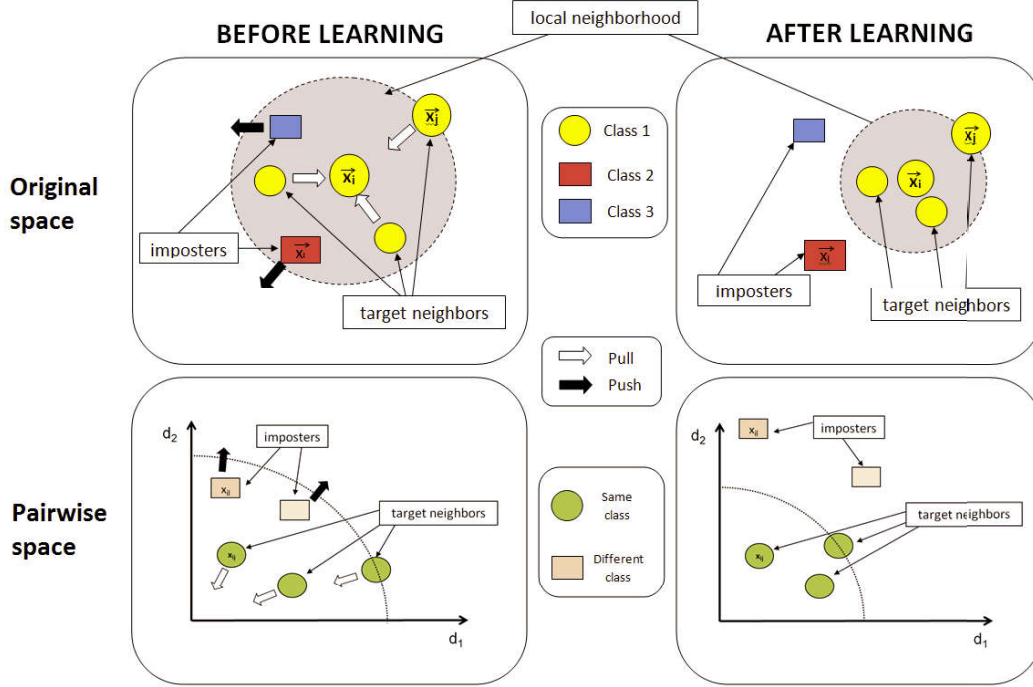


Figure 3.5: Geometric representation of the adaptation of metric learning problem from the original space (top) to the dissimilarity space (bottom) for a $k = 3$ target neighborhood of \mathbf{x}_i . Before learning (left), imposters \mathbf{x}_l invade the targets perimeter \mathbf{x}_j . In the dissimilarity space, this is equivalent to have pairwise vectors \mathbf{x}_{il} with a norm lower to some pairwise target \mathbf{x}_{ij} . The aim of metric learning is to push pairwise \mathbf{x}_{il} (black arrow) and pull pairwise \mathbf{x}_{ij} from the origin (white arrow).

Inspired from the Large Margin Nearest Neighbors (LMNN) framework proposed by Weinberger & Saul in Section 2.6.2, we transpose the metric learning problem into the dissimilarity space to learn a metric D that combines several unimodal metric d_h . In our problem, the optimal metric D is learned as the solution of a minimization problem, such that for each time series \mathbf{x}_i , it pulls its targets \mathbf{x}_j and pushes all the samples \mathbf{x}_l with a different label ($y_l \neq y_i$). The Metric Learning in Dissimilarity space (MLD) problem is formalized as:

$$\underset{D, \xi}{\operatorname{argmin}} \left\{ \underbrace{\sum_{i, j \rightsquigarrow i} D(\mathbf{x}_{ij})}_{\text{pull}} + C \underbrace{\sum_{i, j \rightsquigarrow i, l} \frac{1 + y_{il}}{2} \xi_{ijl}}_{\text{push}} \right\} \quad (3.3)$$

$$\text{s.t. } \forall j \rightsquigarrow i, l,$$

$$D(\mathbf{x}_{il}) - D(\mathbf{x}_{ij}) \geq 1 - \xi_{ijl} \quad (3.4)$$

$$\xi_{ijl} \geq 0 \quad (3.5)$$

where ξ_{ijl} are the slack variables and C , the trade-off between the pull and push costs. The proposed MLD differs from LMNN in which the push term in MLD considers all samples \mathbf{x}_l with a different label from \mathbf{x}_i , whereas in LMNN, only the imposters are taken into consideration (those whose invade the target perimeter). Intuitively, this due to the fact that we do not want that samples \mathbf{x}_l with a different class that were not at the beginning imposters, become imposters during the optimization process. By considering all the samples \mathbf{x}_l , we ensure that at each step of the optimization process, if a sample \mathbf{x}_l becomes an imposter, then it will violate the constraints in Eq. 4.7 and thus, its slack variables ξ_{ijl} will be penalized in the objective function (Eq. 4.5) :

- If $D(\mathbf{x}_{il}) < D(\mathbf{x}_{ij})$, then the pairs \mathbf{x}_{il} is an imposter pair that invades the neighborhood of the target pairs \mathbf{x}_{ij} . The slack variable $\xi_{ijl} > 1$ will be penalized in the objective function (Eq. 4.5).
- If $D(\mathbf{x}_{il}) \geq D(\mathbf{x}_{ij})$ but $D(\mathbf{x}_{il}) \leq D(\mathbf{x}_{ij}) + 1$, the pair \mathbf{x}_{il} is within the safety margin of the target pairs \mathbf{x}_{ij} . The slack variable $\xi_{ijl} \in [0; 1]$ will have a small penalization effect in the objective function (Eq. 4.5).
- If $D(\mathbf{x}_{il}) > D(\mathbf{x}_{ij}) + 1$, $\xi_{ijl} = 0$ and the slack variable has no effect in the objective function (Eq. 4.5).

By considering a linear combination of the unimodal distance d_h (Chapter 2): $D(\mathbf{x}_i, \mathbf{x}_j) = \sum_h w_h \cdot d_h(\mathbf{x}_i, \mathbf{x}_j)$, optimizing the metric D is equivalent to optimizing the weight vector \mathbf{w} . Eqs. 4.5 and 4.6 leads to the MLD primal formulation:

$$\underset{\mathbf{w}, \xi}{\operatorname{argmin}} \left\{ \underbrace{\|\mathbf{X}_{tar}^T \mathbf{w}\|}_{pull} + C \underbrace{\sum_{i,j \rightsquigarrow i,l} \frac{1 + y_{il}}{2} \xi_{ijl}}_{push} \right\} \quad (3.6)$$

$$\text{s.t. } \forall j \rightsquigarrow i, l,$$

$$\mathbf{w}^T (\mathbf{x}_{il} - \mathbf{x}_{ij}) \geq 1 - \xi_{ijl} \quad (3.7)$$

$$\xi_{ijl} \geq 0 \quad (3.8)$$

where $\|\mathbf{X}_{tar}^T \mathbf{w}\|$ denotes the norm of the vector $\mathbf{X}_{tar}^T \mathbf{w}$. Similarly to SVM, a L_1 or L_2 norm can be chosen. L_1 norm will privilege sparse solution of \mathbf{w} .

MLD can be seen as a large margin problem in the dissimilarity space and parallels can be done with SVM. The "pull" term acts as a regularizer which aims to minimize the norm of \mathbf{w} . Similarly to SVM, minimizing the norm of \mathbf{w} is equivalent to maximizing the margin $\frac{1}{\|\mathbf{w}\|_2}$ between target pairs \mathbf{x}_{ij} and pairs of different class \mathbf{x}_{il} .

3.4 Quadratic Programming (QP) formalization

The primal formulation of MLD (Eqs. 4.8, 4.9 and 4.10) supposed that the metric D is a linear combination of the metrics d_h . The primal formulation being similar to the one of SVM, it can be derived into its dual form to obtain non-linear solutions for D . For that, we consider in the objective function (Eq. 4.8), the square of the L_2 -norm on \mathbf{w} as the regularizer term, $\frac{1}{2} \|\mathbf{X}_{tar}^T \mathbf{w}\|_2^2$:

$$\underset{\mathbf{w}, \xi}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{X}_{tar}^T \mathbf{w}\|_2^2 + C \sum_{i,j \rightsquigarrow i,l} \frac{1 + y_{il}}{2} \xi_{ijl} \right\} \quad (3.9)$$

$$\text{s.t. } \forall j \rightsquigarrow i, l,$$

$$\mathbf{w}^T (\mathbf{x}_{il} - \mathbf{x}_{ij}) \geq 1 - \xi_{ijl} \quad (3.10)$$

$$\xi_{ijl} \geq 0 \quad (3.11)$$

This formulation can be reduced to the minimization of the following Lagrange function $L(\mathbf{w}, \xi, \boldsymbol{\alpha}, \mathbf{r})$, consisting of the sum of the objective function (Eq. 4.11) and the constraints (Eqs. 4.12 and 4.13) multiplied by their respective Lagrange multipliers $\boldsymbol{\alpha}$ and \mathbf{r} :

$$\begin{aligned} L(\mathbf{w}, \xi, \boldsymbol{\alpha}, \mathbf{r}) = & \frac{1}{2} \|\mathbf{X}_{tar}^T \mathbf{w}\|_2^2 + C \sum_{ijl} \frac{1 + y_{il}}{2} \xi_{ijl} - \sum_{ijl} r_{ijl} \xi_{ijl} \\ & - \sum_{ijl} \alpha_{ijl} (\mathbf{w}^T (\mathbf{x}_{il} - \mathbf{x}_{ij}) - 1 + \xi_{ijl}) \end{aligned} \quad (3.12)$$

where $\alpha_{ijl} \geq 0$ and $r_{ijl} \geq 0$ are the Lagrange multipliers. At the minimum value of $L(\mathbf{w}, \xi, \boldsymbol{\alpha}, \mathbf{r})$, we assume the derivatives with respect to \mathbf{w} and ξ_{ijl} are set to zero:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}} &= \mathbf{X}_{tar}^T \mathbf{X}_{tar} \mathbf{w} - \sum_{ijl} \alpha_{ijl} (\mathbf{x}_{il} - \mathbf{x}_{ij}) = 0 \\ \frac{\partial L}{\partial \xi_{ijl}} &= C - \alpha_{ijl} - r_{ijl} = 0 \end{aligned}$$

that leads to:

$$\mathbf{w} = (\mathbf{X}_{tar} \mathbf{X}_{tar}^T)^{-1} \sum_{ijl} \alpha_{ijl} (\mathbf{x}_{il} - \mathbf{x}_{ij}) \quad (3.13)$$

$$r_{ijl} = C - \alpha_{ijl} \quad (3.14)$$

Substituting Eq. 4.15 and 4.16 back into $L(\mathbf{w}, \xi, \boldsymbol{\alpha}, \mathbf{r})$ in Eq. 4.14, we get the MLD dual

formulation¹:

$$\operatorname{argmax}_{\alpha} \left\{ \sum_{ijl} \alpha_{ijl} - \frac{1}{2} \sum_{ijl} \sum_{i'j'l'} \alpha_{ijl} \alpha_{i'j'l'} (\mathbf{x}_{il} - \mathbf{x}_{ij})^T (\mathbf{X}_{tar} \mathbf{X}_{tar}^T)^{-1} (\mathbf{x}_{i'l'} - \mathbf{x}_{i'j'}) \right\} \quad (3.15)$$

s.t. $\forall i, j \rightsquigarrow i$ and l s.t. $y_{il} = +1$:

$$0 \leq \alpha_{ijl} \leq C \quad (3.16)$$

For any new pair of samples $\mathbf{x}_{i'}$ and $\mathbf{x}_{j'}$, the resulting metric D writes:

$$D(\mathbf{x}_{i'j'}) = \mathbf{w}^T \mathbf{x}_{i'j'} \quad (3.17)$$

$$D(\mathbf{x}_{i'j'}) = \sum_{ijl} \alpha_{ijl} (\mathbf{x}_{il} - \mathbf{x}_{ij})^T (\mathbf{X}_{tar} \mathbf{X}_{tar}^T)^{-1} \mathbf{x}_{i'j'} \quad (3.18)$$

with \mathbf{w} defined in Eq. 4.15. At the optimality, only the triplets $(\mathbf{x}_{il} - \mathbf{x}_{ij})$ with $\alpha_{ijl} > 0$ are considered as the support vectors. The direction \mathbf{w} of the metric D is lead by these triplets. All other triplets have $\alpha_{ijl} = 0$ (non-support vector), and the metric D is independent from this triplets. If we remove some of the non-support vectors, the metric D remains unaffected. From the viewpoint of optimization theory, we can also see this from the Karush-Kuhn-Tucker (KKT) conditions: the complete set of conditions which must be satisfied at the optimum of a constrained optimization problem. At the optimum, the Karush-Kuhn-Tucker (KKT) conditions apply, in particular:

$$\alpha_{ijl} (\mathbf{w}^T (\mathbf{x}_{il} - \mathbf{x}_{ij}) - 1 + \xi_{ijl}) = 0$$

from which we deduce that either $\mathbf{w}^T (\mathbf{x}_{il} - \mathbf{x}_{ij}) > 1$ and $\alpha_{ijl} = 0$ (the triplet $(\mathbf{x}_{il} - \mathbf{x}_{ij})$ is a non-support vector), or $\mathbf{w}^T (\mathbf{x}_{il} - \mathbf{x}_{ij}) = 1 - \xi_{ijl}$ and $\alpha_{ijl} > 0$ (the triplet is a support vector). Therefore, the learned metric D is a combination of scalar products between new pairs $\mathbf{x}_{i'j'}$ and a few number of triplets \mathbf{x}_{ijl} of the training set.

Extension to non-linear function of D

The above formula can extended to non-linear function for the metric D . The dual formulation in Eq. 4.17 only relies on the inner product $(\mathbf{x}_{i'l'} - \mathbf{x}_{i'j'})^T (\mathbf{X}_{tar} \mathbf{X}_{tar}^T)^{-1} (\mathbf{x}_{il} - \mathbf{x}_{ij})$. We can hence apply the kernel trick on Eqs. 4.19 and 4.20 to find non-linear solutions for D :

$$\begin{aligned} D(\mathbf{x}_{i'j'}) &= \mathbf{w}^T \phi(\mathbf{x}_{i'j'}) \\ D(\mathbf{x}_{i'j'}) &= \sum_{ijl} \alpha_{ijl} \phi(\mathbf{x}_{il} - \mathbf{x}_{ij}) \phi(\mathbf{x}_{i'j'}) \\ D(\mathbf{x}_{i'j'}) &= \sum_{ijl} \alpha_{ijl} K(\mathbf{x}_{il} - \mathbf{x}_{ij}; \mathbf{x}_{i'j'}) \end{aligned}$$

These equations suppose that the null vector $\mathbf{0}$ in the original space is transformed through the

¹complete details of the calculations in Appendix D

transformation ϕ into the null vector: $\phi(\mathbf{0}) = \mathbf{0}$ in the feature space. We recall that $D(\mathbf{x}_{ii} = \mathbf{0})$ is expected to be equal to zero (distinguishability property in Section 2.2). However, if the vectors \mathbf{x}_{ij} are projected in a feature space by a transformation ϕ , it doesn't guarantee that $\phi(\mathbf{0}) = \mathbf{0}$. Fig. 4.7 illustrates the idea for a polynomial kernel in which $\phi(\mathbf{0}) = [0, 0, 0, 1]^T$. Thus, the metric measure needs to be computed in the feature space relatively to the projection of $\phi(\mathbf{0})$. This is done by adding a term $\mathbf{w}^T \phi(\mathbf{0})$ to Eqs. 4.19 and 4.20:

$$D(\mathbf{x}_{i'j'}) = \mathbf{w}^T \phi(\mathbf{x}_{i'j'}) - \mathbf{w}^T \phi(\mathbf{0}) \quad (3.19)$$

$$D(\mathbf{x}_{i'j'}) = \sum_{ijl} \alpha_{ijl} \phi(\mathbf{x}_{il} - \mathbf{x}_{ij}) \phi(\mathbf{x}_{i'j'} - \mathbf{x}_{ij}) - \sum_{ijl} \alpha_{ijl} \phi(\mathbf{x}_{il} - \mathbf{x}_{ij}) \phi(\mathbf{0} - \mathbf{x}_{ij}) \quad (3.20)$$

$$D(\mathbf{x}_{i'j'}) = \sum_{ijl} \alpha_{ijl} K(\mathbf{x}_{il} - \mathbf{x}_{ij}; \mathbf{x}_{i'j'} - \mathbf{x}_{ij}) - \sum_{ijl} \alpha_{ijl} K(\mathbf{x}_{il} - \mathbf{x}_{ij}; \mathbf{0} - \mathbf{x}_{ij}) \quad (3.21)$$

where $\mathbf{0}$ denotes the null vector. The resulting metric D is made of two terms. The first one, $\mathbf{w}^T \phi(\mathbf{x}_{i'j'})$, is the metric measure for a new pair $\mathbf{x}_{i'j'}$. The second term, $\mathbf{w}^T \phi(\mathbf{0})$, adapts the metric measure relatively to the origin point.

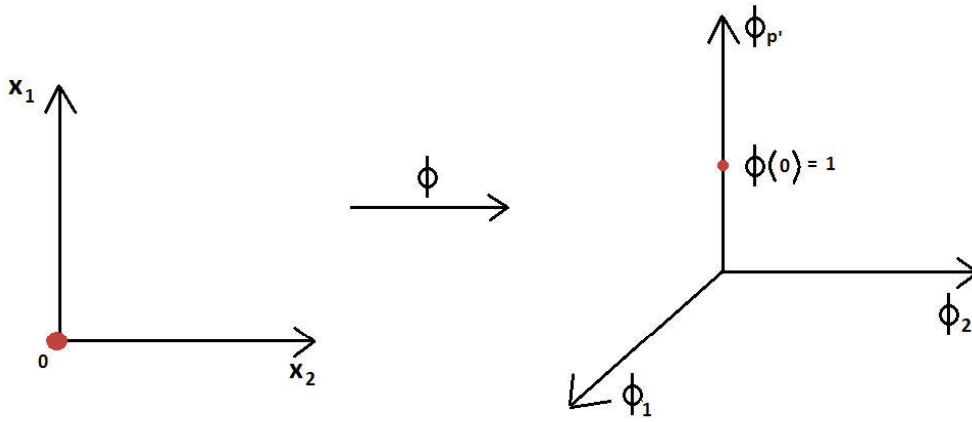


Figure 3.6: Illustration of samples in \mathbb{R}^2 . The transformation ϕ for a polynomial kernel $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + c)^d$ with $c = 1$ and $d = 2$ can be written explicitly: $\phi(\mathbf{x}_i) = [x_{i1}^2, x_{i2}^2, \sqrt{2}x_{i1}x_{i2}, 1]^T$. The origin point $\mathbf{x}_i = [0, 0]^T$ is projected in the Hilbert space as $\phi(\mathbf{x}_i = \mathbf{0}) = [0, 0, 0, 1]^T$.

However, to define proper metrics that respects the properties of metrics (Section 2.2), specific kernels must be used. Our work don't propose any solutions to this problem but open the field for new research on this topic.

3.5 Support Vector Machine (SVM) approximation

3.5.1 Motivations

Many parallels have been studied between Large Margin Nearest Neighbors (LMNN) and SVM (Section 2.6.3). Similarly, the proposed MLD approach can be linked to SVM: both are convex optimization problem based on a regularized and a loss term. SVM is a well known framework: its has been well implemented in many libraries (e.g., LIBLINEAR [FCH08] and LIBSVM [HCL08]), well studied for its generalization properties and extension to non-linear solutions.

Motivated by these advantages, we propose to solve the MLD problem by solving a similar SVM problem. Then, we can naturally extend MLD approach to find non-linear solutions for the metric D thanks to the 'kernel trick'. In the following, we show the similarities and the differences between LP/QP and SVM formulation.

For a time series \mathbf{x}_i , we define the set of pairs $\mathbf{X}_{pi} = \{(\mathbf{x}_{ij}, y_{ij}) \text{ s.t. } j \rightsquigarrow i \text{ or } y_{ij} = +1\}$. It corresponds for a time series \mathbf{x}_i to the set of pairs with target samples \mathbf{x}_j (k nearest samples of same labels $j \rightsquigarrow i$) or samples \mathbf{x}_l that has a different label from \mathbf{x}_i ($y_l \neq y_i$). Identity pairs \mathbf{x}_{ii} are not considered. We refer to $\mathbf{X}_p = \bigcup_i \mathbf{X}_{pi}$ and consider the following standard soft-margin weighted SVM problem on \mathbf{X}_p ²:

$$\begin{aligned} \underset{\mathbf{w}, b, \xi}{\operatorname{argmin}} \quad & \left\{ \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i,j, y_{ij}=-1} p_i^- \xi_{ij} + C \sum_{i,j, y_{ij}=+1} p_i^+ \xi_{ij} \right\} \\ \text{s.t. } & y_{ij}(\mathbf{w}^T \mathbf{x}_{ij} + b) \geq 1 - \xi_{ij} \end{aligned} \quad (3.22)$$

where p_i^- and p_i^+ are the weight factors for target pairs and pairs of different class.

We show in the following that solving the SVM problem in Eq. 4.24 for \mathbf{w} and b solves a similar MLD problem in Eq. 4.11 for $D(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{2}(\mathbf{w}^T \mathbf{x}_{ij} + b)$. If we set p_i^+ being the half of the number of targets of \mathbf{x}_i and p_i^- , the half of the number of time series L of a different class than \mathbf{x}_i :

$$p_i^+ = \frac{k}{2} = \sum_{j \rightsquigarrow i} \frac{1}{2} \quad (3.23)$$

$$p_i^- = \frac{L}{2} = \frac{1}{2} \sum_l \frac{1 + y_{il}}{2} \quad (3.24)$$

²the SVM formulation below divides the loss part into two terms similarly to asymmetric SVM

3.5.2 Similarities and differences in the constraints

First, we recall the SVM constraints in Eq. 4.24:

$$y_{ij}(\mathbf{w}^T \mathbf{x}_{ij} + b) \geq 1 - \xi_{ij}$$

These constraints can be split into two sets of constraints:

$$\begin{aligned} -(\mathbf{w}^T \mathbf{x}_{ij} + b) &\geq 1 - \xi_{ij} && \text{(same class: } y_{ij} = -1) \\ (\mathbf{w}^T \mathbf{x}_{il} + b) &\geq 1 - \xi_{il} && \text{(different classes: } y_{ij} = +1) \end{aligned}$$

By defining $D(\mathbf{x}_{ij}) = \frac{1}{2}(\mathbf{w}^T \mathbf{x}_{ij} + b)$, this leads to:

$$\begin{aligned} -D(\mathbf{x}_{ij}) &\geq \frac{1}{2} - \frac{\xi_{ij}}{2} \\ D(\mathbf{x}_{il}) &\geq \frac{1}{2} - \frac{\xi_{il}}{2} \end{aligned}$$

By summing each constraint two by two, this set of constraints implies the following set of constraints:

$$\left\{ \begin{aligned} &\bullet \forall i, j, k, l \text{ such that } y_{ij} = -1, \text{ and } y_{kl} = +1, i \neq j \text{ and } i \neq k : \\ &\quad D(\mathbf{x}_k, \mathbf{x}_l) - D(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \frac{\xi_{kl} + \xi_{ij}}{2} \\ &\bullet \forall i, j, l \text{ such that } y_{ij} = -1, \text{ and } y_{il} = +1, i \neq j : \\ &\quad D(\mathbf{x}_i, \mathbf{x}_l) - D(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \frac{\xi_{il} + \xi_{ij}}{2} \end{aligned} \right. \quad (3.25)$$

By defining $\xi_{ijl} = \frac{\xi_{ij} + \xi_{il}}{2}$, the second constraint in Eq. 4.27 from the SVM formulation is the same as the constraints in the MLD formulation in Eq. 4.12.

However, an additional set of constraints is present in the SVM formulation (first set of constraints in Eq. 4.27) and not in the proposed MLD. Geometrically, this can be interpreted as superposing the neighborhoods of all samples \mathbf{x}_i , making the union of all of their target sets \mathbf{X}_{pi} , and then pushing away all imposters \mathbf{x}_{il} from this resulting target set. This is therefore creating "artificial imposters" \mathbf{x}_{kl} that don't violate the local target space of sample \mathbf{x}_k , but are still considered as imposters because they invade the target of sample \mathbf{x}_i (because of the neighborhoods superposition) (Figure 4.8). This is more constraining in the SVM resolution for the resulting metric D especially if the neighborhoods have different spread.

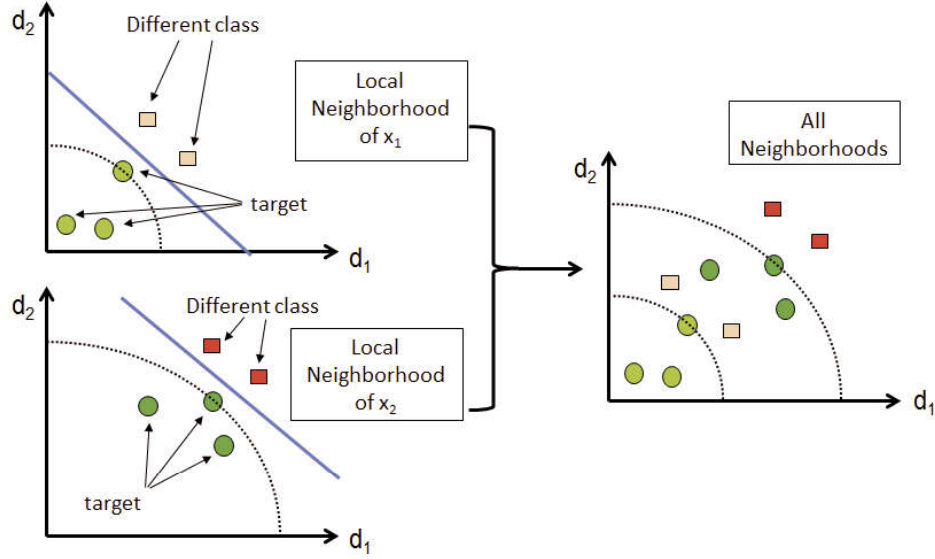


Figure 3.7: Geometric representation of the neighborhood of $k = 3$ for two time series \mathbf{x}_1 and \mathbf{x}_2 (left). For each neighborhood, time series of different class are represented by a square and the margin by a blue line. Taking each neighborhood separately, the problem is linearly separable (LP/QP formulation). By combining the two neighborhoods (SVM formulation), the problem is no more linearly separable and in this example, the time series of different class of \mathbf{x}_1 (orange square) are "artificial imposters" of \mathbf{x}_2 .

3.5.3 Similarities and differences in the objective function

Mathematically, from Eq. 4.25, we write:

$$\begin{aligned}
 \sum_{i,l,y_{il}=+1} p_i^+ \xi_{il} &= \sum_{il} p_i^+ \frac{1+y_{il}}{2} \xi_{il} \\
 &= \sum_{il} \left(\sum_{j \rightsquigarrow i} \frac{1}{2} \right) \frac{1+y_{il}}{2} \xi_{il} \\
 &= \frac{1}{2} \sum_{i,j \rightsquigarrow i,l} \frac{1+y_{il}}{2} \xi_{il}
 \end{aligned} \tag{3.26}$$

And from Eq. 4.26, we write:

$$\begin{aligned}
 \sum_{i,j,y_{ij}=-1} p_i^- \xi_{ij} &= \sum_{i,j \rightsquigarrow i} p_i^- \xi_{ij} \\
 &= \sum_{i,j \rightsquigarrow i} \left(\frac{1}{2} \sum_l \frac{1+y_{il}}{2} \right) \xi_{ij} \\
 &= \frac{1}{2} \sum_{i,j \rightsquigarrow i,l} \frac{1+y_{il}}{2} \xi_{ij}
 \end{aligned} \tag{3.27}$$

By replacing Eqs. 4.28 and 4.29 back into Eq. 4.24, the objective function becomes:

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i,j \rightsquigarrow i,l} \frac{1 + y_{il}}{2} \frac{\xi_{ij} + \xi_{il}}{2} \\ \min_{\mathbf{w}, \xi} \quad & \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w}}_{\text{Regularization}} + C \underbrace{\sum_{i,j \rightsquigarrow i,l} \frac{1 + y_{il}}{2} \xi_{ijl}}_{\text{Loss}} \end{aligned} \quad (3.28)$$

Even if the loss part (push cost) is the same for both objective functions, the regularization part (pull cost) is different. In the SVM formulation (Eq. 4.30), the regularization part tends to minimize the norm of \mathbf{w} whereas in MLD (Eq. 4.11), it tends to minimize the norm of \mathbf{w} after a linear transformation through \mathbf{X}_{tar} . This transformation can be interpreted as a Mahalanobis norm in the dissimilarity space with $\mathbf{M} = \mathbf{X}_{tar} \mathbf{X}_{tar}^T$. Nevertheless, both have the same objective: improve the conditioning of the problem by enforcing solutions with small norms. In practice, even with these differences, the SVM provides suitable solutions for our time series metric learning problem.

3.5.4 Geometric interpretation

Michèle pense que l'état, cette section est dure à comprendre. D'après Michèle, il faut 1) soit prendre + de place pour expliquer la signification géométrique 2) ou soit ne pas mettre cette partie car étant compliquée, cela pourrait nuire au lecteur. Qu'en penses tu Ahlame?

In this section, we give a geometric understanding of the differences between LP/QP resolution (left) and SVM-based resolution (right). Fig. 4.10 shows the Linear Programming (LP) and SVM resolutions of a k -NN problem with $k = 2$ neighborhoods.

For LP, the problem is solved for each neighborhood (blue and red) independently as shown in Fig. 4.9. We recall that LP/QP resolutions, support vectors are triplets of time series made of a target pair \mathbf{x}_{ij} and a pair of different classes \mathbf{x}_{il} (black arrows). Support vectors represent triplet which resulting distance $D(\mathbf{x}_{ij}, \mathbf{x}_{il})$ are the lowest. The optimization problem tends to maximize the margin between these triplets. The global solution (Fig. 4.10 (left)) is a compromise of all of the considered margins. In this case, the global margin is equal to one of the local margin. Note that the global LP solution is not always the same as the best local solution. For SVM-based resolution (Fig. 4.10 (right)), the problem involves all pairs and the margin is optimized so that pairs \mathbf{x}_{ij} and \mathbf{x}_{il} are globally separated.

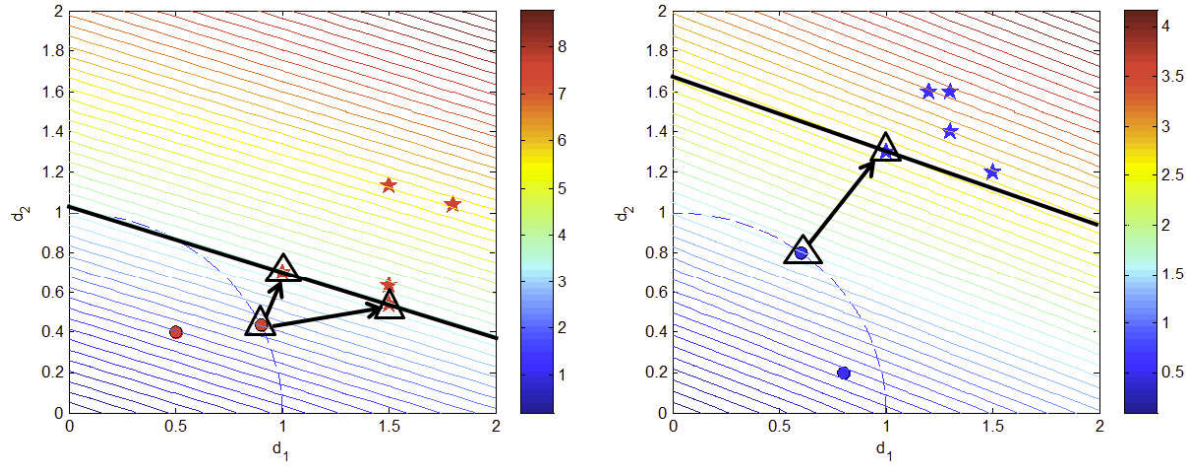


Figure 3.8: Solutions found by solving the LP problem for $k = 2$ neighborhood. Positive pairs (different classes) are indicated in stars and negative pairs (target pairs) are indicated in circle. Red and blue lines shows the margin when solving the problem for each neighborhood (red and blue points) separately. Support vector are indicated in black triangles: in the red neighborhood (left), 2 support vectors are retained and in the blue neighborhood (right), only one support vector is necessary.

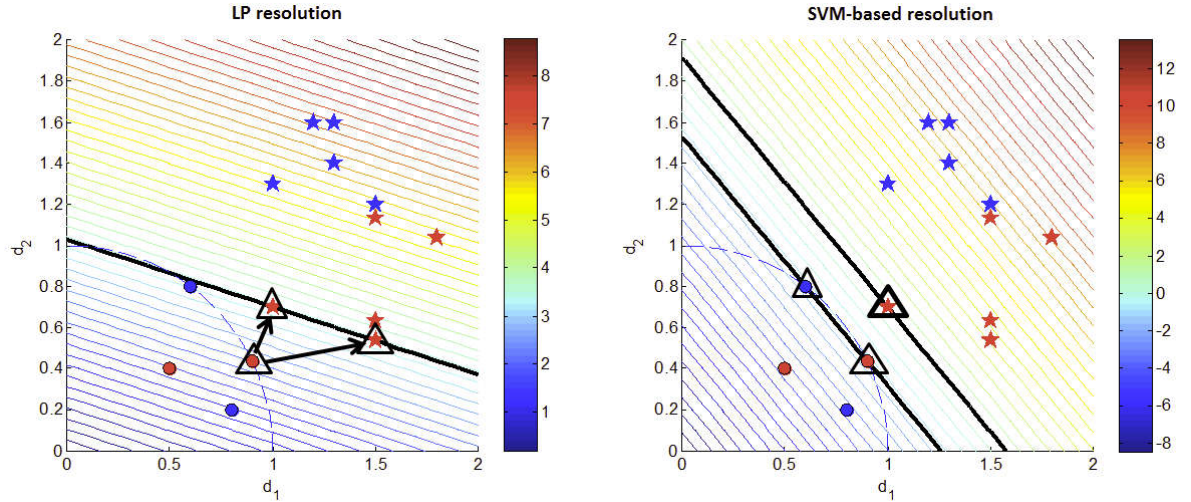


Figure 3.9: Solutions found by solving the LP problem (left) and the SVM problem (right). The global margin is indicated in black and the metric is represented in color levels. Support vectors made of triplets are indicated in black triangles. For the SVM, the black lines indicates the SVM canonical hyperplane where the support vector lies (black triangles).

3.6 Conclusion of the chapter

To learn a combined metric D from several unimodal metrics d_h that optimizes the k -NN performances, we first proposed a new space representation, the dissimilarity space where

each pair of time series is projected as a vector described the unimodal metrics. Then, we propose three formalizations of our metric learning problem: Linear Programming, Quadratic Programming, SVM-based approximation. Table 3.1 sums up the main pros and cons of each formulation.

	LP	QP	SVM-based
Linear	Yes	Yes	Yes
Non-linear extension	No	Yes	Yes
Exact/Approximation resolution	Exact	Exact	Approximation
Sparcity	Yes	Yes	Yes/No

Table 3.1: The different formalizations for Metric Learning in Dissimilarity space

In the following, we consider the SVM-based approximation because SVM framework is well known and well implementated. In the next chapter, we give the details of the steps of our proposed algorithm: Multi-modal and Multi-scale Time series Metric Learning (M^2TML).