

THÈSE

pour obtenir le grade de

DOCTEUR DE L'UNIVERSITÉ DE GRENOBLE

Spécialité : **Informatique et Mathématiques appliquées**

Arrêté ministériel : 7 août 2006

Présentée par

Cao Tri DO

Thèse dirigée par **Ahlame DOUZAL-CHOUAKRIA**,
codirigée par **Michèle ROMBAUT** et
co-encadré par **Sylvain MARIÉ**

préparée au sein du
Laboratoire d'Informatique de Grenoble (LIG)
dans l'école doctorale **Mathématiques, Sciences et
Technologies de l'Information, Informatique (MSTII)**

Metric Learning for Time Series Analysis

Thèse soutenue publiquement le **date de soutenance**,
devant le jury composé de:

Patrick GALLINARI

Laboratoire LIP6, Examineur

Stéphane CANU

Laboratoire LITIS, Rapporteur

Marc SEBBAN

Laboratoire LAHC, Rapporteur

Gustavo CAMPS-VALLS

Laboratoire IPL, Examineur

Ahlame DOUZAL-CHOUAKRIA

Laboratoire LIG, Directeur de thèse

Michèle ROMBAUT

Laboratoire GIPSA-Lab, Co-Directeur de thèse

Sylvain MARIÉ

Schneider Electric, Encadrant



UNIVERSITÉ DE GRENOBLE
ÉCOLE DOCTORALE MSTII
Description de complète de l'école doctorale

T H È S E

pour obtenir le titre de

docteur en sciences

de l'Université de Grenoble-Alpes

Mention : INFORMATIQUE ET MATHÉMATIQUES APPLIQUÉES

Présentée et soutenue par

Cao Tri DO

Metric Learning for Time Series Analysis

Thèse dirigée par Ahlame DOUZAL-CHOUAKRIA

préparée au Laboratoire d'Informatique de Grenoble (LIG)

soutenue le date de soutenance

Jury :

<i>Rapporteurs :</i>	Stéphane CANU	-	Laboratoire LITIS
	Marc SEBBAN	-	Laboratoire LAHC
<i>Directeur :</i>	Ahlame DOUZAL-CHOUAKRIA	-	Laboratoire LIG
<i>Co-Directeur :</i>	Michèle ROMBAUT	-	Laboratoire GIPSA-Lab
<i>Encadrant :</i>	Sylvain MARIÉ	-	Schneider Electric
<i>Examineur :</i>	Patrick GALLINARI	-	Laboratoire LIP6
<i>Examineur :</i>	Gustavo CAMPS-VALLS	-	Laboratoire IPL

Todo list

A compléter par Sylvain	1
ref sylvain?	12
références normalization	12
Comment [AD1]: dit d'enlever 1ère phrase mais Michèle dit de garder	13
Comment [AR2]: Ahlame trouve que ce n'est pas clair. A refaire	13
Comment [AD3]: Expliquer d'avantage	14
Comment [AD4]: Mettre dans les figures des + et - pour les classes	15
Partie non encore rédigée. A faire à la fin.	26
A refaire complètement à la fin.	26
Ajouter virtual sensors?	30
Comment [CTD5]: je préfère garder l'espace pour + de visibilité	32
figure à revoir	32
Comment [SMA6]: reformulation à revoir: let's try to	33
Figure: Ajouter la représentation fréquentielle des signaux	34
A faire à la fin, pas urgent	36
Comment [SMA7]: formulation totale à revoir	37
Comment [MR8]: Modifier figure. enlever 'one' et mettre la même échelle temporelle	37
Comment [AD9]: Ahlame pas fan des notations	37
voir si je mets la dtw avec fonction de coût cort dans les annexes	38
ref	40
références normalization	40
Compléter avec le mail de Sylvain. Partie peut être à mettre en annexe	42
A faire, avec papier PRL et papier Aurélien Bellet. Refaire cette intro avec les notes de Sylvain.	44

ref	46
Intro modifié pour tenir compte des chapitres précédents	50
Formulation du pb à valider	50
Correction d'anglais et précision des phrases	50
separate	50
Comment [SMA10]: Appuyer différence	51
Comment [MR11]: ajouter pairwise pour différencier	51
Comment [SMA12]: modif pour dire que le processus est général et pas que SVM . . .	51
Mettre la section LMNN ici?	51
titre à changer?	52
ref: dissimilarity space [PPD02]; [DP12]	52
Comment [SMA13]: label maintenant?	52
Comment [SMA14]: Proposition pour mettre en valeur que c'est une fonction	52
Comment [SMA15]: Basculer en chap 2?	53
ref	53
Sous section ajoutée, ok pour SMA et MR	53
figure enlevée ici car trop tôt	55
Comment [SMA16]: mettre $R(D, Pull)$?	55
Comment [SMA17]: Enlever les y_{il} , ils ne servent plus	56
Ajouter une remarque sur D_0 ?	56
Ajout ok pour SMA et MR	58
Voir si ce n'est pas redondant avec la partie interprétation	60
ref	76
Michèle trouve que ce serait bien de ranger les datasets par warp/non-warp et de classer des moins challenge au plus challenge pour aider la lecture.	77
Proposition de Sylvain	87
papier Springer?	91

Acknowledgements

I would like to thanks:

- my directors
- my GIPSA colleagues
- my AMA colleagues
- my Schneider colleagues
- my parents

Contents

Table des sigles et acronymes	xv
Introduction	1
1 Related work	5
1.1 Classification, Regression	5
1.2 Machine learning algorithms	13
1.3 Conclusion of the chapter	26
2 Time series metrics and metric learning	29
2.1 Definition of a time series	29
2.2 Properties and representation of a metric	32
2.3 Unimodal metrics for time series	33
2.4 Time series alignment and dynamic programming approach	37
2.5 Combined metrics for time series	39
2.6 Metric learning	43
2.7 Conclusion of the chapter	47
3 Multi-modal and Multi-scale Time series Metric Learning (M²TML)	49
3.1 Motivations	50
3.2 Multi-modal and multi-scale dissimilarity space	52
3.3 M ² TML general problem	55
3.4 Linear formalization for M ² TML	59
3.5 Quadratic formalization for M ² TML	61
3.6 SVM-based formalization for M ² TML	66

3.7	SVM-based solution and algorithm for M^2TML	70
3.8	Conclusion	71
4	Experiments	73
4.1	Description	73
4.2	Experimental protocol	76
4.3	Results and discussion	77
4.4	Conclusion of the chapter	83
	Conclusion	85
	List of articles	91
	A Detailed presentation of the datasets	93
	B Solver library	95
	C SVM library	97
	D QP resolution	99
	E SVM equivalency	101
	Bibliography	110

List of Figures

1.1	An example of overfitting in the case of classification. The objective is to separate blue points from red points. Black line shows a classifier f_1 with low complexity where as green line illustrates a classifier f_2 with high complexity. On training examples (blue and red points), the model f_2 separates all the classes perfectly but may lead to poor generalization on new unseen examples. Model f_1 is often preferred.	7
1.2	Example of a 2 dimensional grid search for parameters C and γ . It defines a grid where each cell of the grid contains a combination (C, γ) . Each combination is used to learn the model and is evaluated on the validation set.	7
1.3	v -fold Cross-validation for one combination of parameters. For each of v experiments, use $v - 1$ folds for training and a different fold for Testing, then the training error for this combination of parameter is the mean of all testing errors. This procedure is illustrated for $v = 4$	8
1.4	General framework for building a supervised (classification/regression) model. Example with 3 features and 2 classes ('Yes' and 'No').	9
1.5	Division of a dataset into 3 datasets: training, test and operational.	9
1.6	Example of k -NN classification. The test sample (green circle) is classified either to the first class (red stars) or to the second class (blue triangles). If $k = 3$ (solid line circle) it is assigned to the second class because there are 2 triangles and only 1 star inside the inner circle. If $k = 5$ (dashed line circle) it is assigned to the first class (3 stars vs. 2 triangles inside the outer circle).	13
1.7	Example of linear classifiers in a 2-dimensional classification problem. For a set of points of classes $+1$ and -1 that are linearly separable, there exists an infinite number of separating hyperplanes corresponding to $\mathbf{w}^T \mathbf{x} + b = 0$	15
1.8	The argument inside the decision function of a classifier is $\mathbf{w}^T \mathbf{x} + b$. The separating hyperplane corresponding to $\mathbf{w}^T \mathbf{x} + b = 0$ is shown as a line in this 2-dimensional plot. This hyperplane separates the two classes of data with points on one side labeled $y_i = +1$ ($\mathbf{w}^T \mathbf{x}_i + b \geq 0$) and points on the other side labeled $y_i = -1$ ($\mathbf{w}^T \mathbf{x}_i + b < 0$). Support vectors are circled in purple and lies on the hyperplanes $\mathbf{w}^T \mathbf{x} + b = +1$ and $\mathbf{w}^T \mathbf{x} + b = -1$	16
1.9	Hyperplane obtained after a dual resolution (full blue line). The 2 canonical hyperplanes (red lines) contain the support vectors whose $\alpha_i > 0$. Other points have their $\alpha_i = 0$ and the equation of the hyperplane is only affected by the support vectors.	19

1.10	Left: in two dimensions the two classes of data (-1 for cross and +1 for circle) are mixed together, and it is not possible to separate them by a line: the data is not linearly separable. Right: using a kernel, these two classes of data become separable by a hyperplane in feature space, which maps to the nonlinear boundary shown, back in input space. ¹	20
1.11	Illustration of the Gaussian kernel in the 1-dimensional input space for a small and large γ when \mathbf{x}_i is fixed and \mathbf{x}_j varies.	21
1.12	Geometric representation of SVM.	23
1.13	Example of several SVMs and how to interpret the weight vector \mathbf{w}	23
1.14	Illustration of SVM regression (left), showing the regression curve with the ϵ -insensitive "tube" (right). Samples \mathbf{x}_i above the ϵ -tube have $\xi_1 > 0$ and $\xi_1 = 0$, points below the ϵ -tube have $\xi_2 = 0$ and $\xi_2 > 0$, and points inside the ϵ -tube have $\xi = 0$	25
2.1	2 modes of representation	31
2.2	The Beveridge wheat price index is the average in nearly 50 places in various countries measured in successive years from 1500 to 1869 ²	31
2.3	Example of metric representation: (a) data points can be fixed and the distance sphere is shown. (b) sphere can be fixed and the data points and the axis are moving.	32
2.4	3 toy time series. Time series in blue and red are two sinusoidal signals. Time series in green is a random signal.	34
2.5	The signal from Fig. 2.4 and a signal \mathbf{x}_4 which is signal \mathbf{x}_1 and an added translation. Based on behavior comparison, \mathbf{x}_4 is the closest to \mathbf{x}_1	36
2.6	Example of a same sentence said by two different speakers. Time series are shifted, compressed and dilatated in the time.	37
2.7	Example of DTW grid between 2 time series \mathbf{x}_i and \mathbf{x}_j (top) and the signals before and after warping (bottom). On the DTW grid, the two signals can be represented on the left and bottom of the grid. The optimal path $\boldsymbol{\pi}^*$ is represented in green line and shows how to associate elements of \mathbf{x}_i to element of \mathbf{x}_j . Background show in grey scale the value of the considered metric (amplitude-based distance d_A in classical DTW)	39
2.8	Contour plot of the resulting combined metrics: D_{Lin} (1 st line), D_{Geom} (2 nd line) and D_{Sig} (3 rd line), for different values of α . For the three combined metrics, the first and second dimensions are respectively the amplitude-based metrics d_A and the behavior-based metric d_B	41

2.9	A nearly log-normal distribution, and its log transform ³	42
2.10	Pushed and pulled samples in the $k = 3$ target neighborhood of \mathbf{x}_i before (left) and after (right) learning. The pushed (vs. pulled) samples are indicated by a white (vs. black) arrows (Weinberger & Saul [WS09a]).	44
2.11	(a) Standard LMNN model view (b) LMNN model view under an SVM-like interpretation [Do+12]	46
2.12	(a) LMNN in a local SVM-like view (b) LMNN in an SVM metric learning view [Do+12]	47
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).	50
3.2	Example of embedding of time series \mathbf{x}_i from the temporal space (left) into the dissimilarity space (right) for $p = 3$ basic metrics.	52
3.3	Multi-scale decomposition	54
3.4	Example of two pairwise vectors \mathbf{x}_{12} and \mathbf{x}_{34} close in the pairwise space. However, the time series \mathbf{x}_1 and \mathbf{x}_3 are not similar in the temporal space.	54
3.5	Example of two pairwise vectors \mathbf{x}_{12} and \mathbf{x}_{34} close in the pairwise space. However, the time series \mathbf{x}_1 and \mathbf{x}_3 are not similar in the temporal space.	55
3.6	Example of different strategies to build $Pull_i$ and $Push_i$ sets for a $k = 2$ neighborhood.	57
3.7	Metric learning problem from the original space (top) to the dissimilarity space (bottom) for a $k = 3$ neighborhood of \mathbf{x}_i . Before learning (left), push samples \mathbf{x}_l invade the targets perimeter \mathbf{x}_j . In the pairwise space, this is equivalent to have push pairwise vectors \mathbf{x}_{il} with an initial distance D_0 lower than the distance of pull pairwise vectors \mathbf{x}_{ij} . The aim of metric learning is to learn a metric D to push \mathbf{x}_{il} (black arrow) and pull \mathbf{x}_{ij} from the origin (white arrow).	58
3.8	The projected vector $\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{ij})$ and $\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{ij'})$	67
3.9	Example of SVM solutions and of the resulting metric D defined by the norm of the projection on \mathbf{w} . Fig. (a) represents common expected configuration where pull pairs $Pull_i$ are situated in the same side as the origin $\mathbf{x}_{ii} = 0$. In Fig. (b), the vector $\mathbf{w} = [-1 - 1]^T$ indicates that push pairs $Push_i$ are on the side of the origin point. One problem arises in Fig. (b): distance of push pairs $D(\mathbf{x}_{il})$ is lower than the distance of pull pairs $D(\mathbf{x}_{ij})$	68

3.10	The behavior of the learned metric D ($p = 2; \lambda = 2.5$) with respect to common (a) and challenging (b) configurations of pull and push pairs.	69
3.11	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 pull pairs $Pull_i$ and square represents push pairs $Push_i$ for $m = 3$ neighbors. Before scaling, the problem is not linearly separable. The spread of each neighborhood are not comparable. After scaling, the target neighborhood becomes comparable and in this example, the problem becomes linearly separable between the circles and the squares.	71
4.1	Temporal representation of some datasets (SonyAIBO, ECG200, BME, UMD, FaceFour, PowerConsumption) considered in the experiments.	75
4.2	(a) Standard (Euclidean distance d_A and DTW) <i>vs.</i> M^2TML (D and $D_{\mathcal{H}}$) metrics. (b) Best Uni-modal (DTW and d_{B-DTW}) <i>vs.</i> M^2TML (D and $D_{\mathcal{H}}$) metrics.	79
4.3	M^2TML feature weights for 4 datasets.	80
4.4	Temporal representation of the top M^2TML feature weights for 4 datasets.	82
4.5	MDS visualization of the d_{B-DTW} (Fig. a) and D (Fig. b) dissimilarities for FaceFour	82
4.6	(a) Two classes of time series from the Sony AIBO accelerometer. (b) The and-shapelets from the walk cycle on carpet. (c) The Sony AIBO Robot. ⁴	87
4.7	Example of discretization by binning a continuous label y into $Q = 4$ equal-length intervals. Each interval is associated to a unique class label. In this example, the class label for each interval is equal to the mean in each interval.	88
4.8	Border effect problems. In this example, \mathbf{x}_2 and \mathbf{x}_3 have closer value labels y_2 and y_3 than \mathbf{x}_3 and \mathbf{x}_4 . However, with the discretization \mathbf{x}_2 and \mathbf{x}_3 don't belong to the same class and thus are consider as not similar.	89
4.9	Example of pairwise label definition using an ϵ -tube (red lines) around the time series \mathbf{x}_i (circled in blue). For, time series \mathbf{x}_j that falls into the tube, the pairwise label is $y_{ij} = -1$ (similar) and outside of the tube, $y_{ij} = +1$ (not similar).	89

- E.1 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}_2 103

List of Tables

1.1	Confusion matrix for a 2-class problem.	10
3.1	The different formalizations for M^2TML	72
4.1	Dataset description giving the number of classes (Nb. Class), the number of time series for the training (Nb. Train) and the testing (Nb. Test) sets, and the length of each time series (TS length).	74
4.2	Considered metric in the experiments	75
4.3	Parameter ranges	76
4.4	1-NN test error rates for standard, <i>a priori</i> combined and M^2TML measures. . .	77
4.5	Top 5 multi-modal and multi-scale features involved in D	81

Table of Acronyms

LIG	<i>Laboratoire d'Informatique de Grenoble</i>
AMA	<i>Apprentissage, Méthode et Algorithme</i>
GIPSA-Lab	<i>Grenoble Images Parole Signal Automatique Laboratoire</i>
AGPiG	<i>Architecture, Géométrie, Perception, Images, Gestes</i>
A4S	<i>Analytics for Solutions</i>
k-NN	<i>k-nearest neighbors</i>
SVM	<i>Support Vector Machines</i>
SVR	<i>Support Vector Regression</i>
d_E	<i>Euclidean distance</i>
$corr$	<i>Pearson correlation</i>
$cort$	<i>Temporal correlation</i>
dtw	<i>Dynamic Time Warping</i>
IoT	<i>Internet of Things</i>
Acc	<i>Classification accuracy</i>
Err	<i>Classification error rate</i>
MAE	<i>Mean Absolute Error</i>
RMSE	<i>Root Mean Square Error</i>
FAQ	Frequently Asked Questions

Introduction

Motivation

The work of this PhD is in the context of a CIFRE¹ thesis with Schneider Electric and two public research laboratories, the LIG² and the GIPSA-lab³. Within Schneider Electric, the PhD took place in the Analytics for Solutions (A4S) team, part of the Technology and Strategy entity whose main missions are à compléter par Sylvain. Among the wide activities of the A4S team, in the context of physical system modeling (*e.g.*, building, sensor network, Internet of Things), two topics are at least studied: modeling by physical models (white/grey box) and modeling by machine learning algorithms (black-box). With the increase of the amount of data and sensors that collect data, modeling accurately systems through *a priori* equations (white/grey box) for some prediction tasks has become more and more difficult. Within the vast amount of applications in Schneider Electric, some applications in particular involve temporal data, *e.g.*, forecasting the energy consumption in a building, virtual sensors, fault detection. More generally, Schneider Electric, like many other companies and other diverse application domains (medicine, marketing, meteorology, etc.) has taken a growing interest these last decades in machine learning problems (classification, regression, clustering) that involves time series of one or several dimensions, of different sampling, of two or more classes, etc. A time series can be seen in signal processing and in control theory as the response of a dynamic system. Contrary to static data, time series are more challenging in the sense that the temporal aspect (*i.e.*, order of appearance of the observations) is an additional key information.

A compléter
par Sylvain

Problem statement and contributions

In this work, we focus on classification problems of monovariate time series (data of 1 dimension) with a fixed sampling rate and of same lengths. Among the wide variety of algorithms that exist in machine learning, some approaches (*e.g.*, *k*-Nearest Neighbors (*k*-NN)) classify samples using a concept of neighborhood based on the comparison between samples. In general, the concept of 'near' and 'far' between samples is expressed through a distance measure. Time series can be compared based not only on their amplitudes like static data but also on other characteristics or modalities such that their dynamic or frequency components. Many metrics for time series have been proposed in the literature such that the euclidean distance [Din+08], the temporal correlation [FDcG13], the Fourier-based distance [SS12a]. A detailed review of the major metrics is proposed in [MV14]. In general, the existing metrics involve

¹Conventions Industrielles de Formation par la REcherche

²Laboratoire d'Informatique de Grenoble

³Grenoble Images Parole Signal Automatique

one modality at the global scale (*i.e.*, implying systematically all the time series observations). We believe also that the multi-scale aspect of time series, not present in static data, could enrich the definition of the existing metrics.

In this work, our objective is to learn a combined multi-modal and multi-scale time series metric for a robust k -NN classifier. The main contributions of the PhD are:

- The definition of a new space representation: the dissimilarity space which embeds a pair of time series into a vector described by basic temporal metrics.
- The definition of basic temporal metrics that involves one modality at one specific scale.
- The learning of a multi-modal and multi-scale temporal metric for a large margin k -NN classifier of univariate time series.
- The definition of the general problem of learning a combined metric as a metric learning problem using the dissimilarity representation.
- The proposition of a framework based on Support Vector Machine (SVM) and a linear and non-linear solution to define the combined metric that satisfies at least the properties of a dissimilarity measure.
- The comparison of the proposed approach with standard metrics on a large number of public datasets.
- The analysis of the proposed approach to extract the discriminative features that are involved in the definition of the learned combined metric.

Organization of the manuscript

The first part makes a review of existing methods in machine learning and metrics for time series. The first chapter presents classical approaches in machine learning. In particular, we recall the general principle, framework and focus on two standard machine learning algorithms: the k -Nearest Neighbors (k -NN) and the Support Vector Machine (SVM) approach. In the second chapter, we review some basic terminology for time series and recall three types of metrics proposed at least for time series: amplitude-, behavior- and frequential-based.

The second part of the manuscript propose a multi-modal and multi-scale metric learning (M^2TML) method. In the third chapter, we first review the concept of metric learning for static data and focus on a framework of metric learning for nearest neighbors classification proposed by Weinberger & Saul [WS09b]. Secondly, we present a new space representation, the pairwise dissimilarity space based on a multi-modal and multi-scale time series description and their corresponding basic metrics. Then, we formalize the general M^2TML optimization problem using the pairwise dissimilarity space representation. From the general formalization, we propose at least three different formalizations. The first and second proposition involve different regularizers, allowing to learn an *a priori* linear or non-linear form of the combined

metric. The third proposition presents a framework based on SVM and a solution to build the combined metric, in the linear and non-linear context, satisfying at least the properties of a dissimilarity measure. Finally, Chapter 4 presents the experiments conducted on a wide range of 30 public and challenging datasets, and discusses the results obtained.

Notations

\mathbf{x}_i	a vector sample / a time series
y_i	a label (discrete or continuous)
\hat{y}_i	a predicted label (discrete or continuous)
$X = \{\mathbf{x}_i, y_i\}_{i=1}^n$	a training set of $n \in \mathbb{N}$ labeled time series
$X_{Test} = \{\mathbf{x}_j, y_j\}_{j=1}^m$	a test set of $m \in \mathbb{N}$ labeled time series
d_A	Amplitude-based distance
d_B	Behavior-based distance
d_F	Frequential-based distance
d_E	Euclidean distance
L_q	Minkovski q-norm
$\ \mathbf{x}\ _q$	q-norm of the vector \mathbf{x}
$cort$	Temporal correlation
D	Linear learned combined distance
$D_{\mathcal{H}}$	Non-linear learned combined distance
κ	Kernel function
$\phi(\mathbf{x}_i)$	embedding function from the original space to the Hilbert space
\mathbf{x}_{ij}	a pair of time series \mathbf{x}_i and \mathbf{x}_j
y_{ij}	the pairwise label of \mathbf{x}_{ij}
t	time stamp/index with $t = 1, \dots, T$
q	length of the time series (supposed fixed)
f	frequential index
F	length of the Fourier transform
ξ	Slack variable
p	number of metric measure considered in the metric learning process
r	order of the temporal correlation
k	number of nearest neighbors
C	Hyper-parameter of the SVM (trade-off)
α	Parameter to control the size of the neighborhood
λ	Parameter to control the push term
\mathbf{w}	weight vector
v	number of folds in cross-validation

Related work

Contents

1.1	Classification, Regression	5
1.1.1	Machine learning principle	5
1.1.2	Model selection in supervised learning	6
1.1.3	Model evaluation	10
1.1.4	Data normalization	12
1.2	Machine learning algorithms	13
1.2.1	k -Nearest Neighbors (k -NN) classifier	13
1.2.2	Support Vector Machine (SVM) algorithm	14
1.2.3	Other classification algorithms	26
1.3	Conclusion of the chapter	26

In this chapter, we recall some concepts of machine learning. First, we review the principles, the learning framework and the evaluation protocol in supervised learning. Then, we present the algorithms used in our work: k -Nearest Neighbors (k -NN) and Support Vector Machines (SVM).

1.1 Classification, Regression

In this section, we review some terminology used in machine learning. First, we recall the principle of machine learning. Then, we detail how to design a framework for supervised learning. After that, we present model evaluation. Finally, we review data normalization.

1.1.1 Machine learning principle

The idea of machine learning (also known as pattern learning or pattern recognition) is to imitate with algorithms executed on computers, the ability of living beings to learn from examples. For instance, to teach a child how to read letters, we show him during a training phase, labeled examples of letters ('A', 'B', 'C', etc.) written in different styles and fonts. We don't give him a complete and analytic description of the topology of the characters but

labeled examples. Then, during a testing phase, we want the child to be able to recognize and to label correctly the letters that have been seen during the training, and also to generalize to new instances [G. 06].

Let $X = \{\mathbf{x}_i, y_i\}_{i=1}^n$ be a training set of n vector samples $\mathbf{x}_i \in \mathbb{R}^p$ and y_i their corresponding labels. The aim of supervised machine learning is to learn a relationship (model) f between the samples \mathbf{x}_i and their labels y_i based on examples [Bis06]; [G. 06]; [OE73]. After the training phase based on labeled examples (\mathbf{x}_i, y_i) , the model f has to be able to generalize on the testing phase, *i.e.*, to give a correct prediction \hat{y}_j for new instances \mathbf{x}_j that haven't been seen during the training.

When y_i are class labels (*e.g.*, class 'A', 'B', 'C' in the case of child's reading), learning the model f is a classification problem; when y_i is a continuous value (*e.g.*, the energy consumption in a building), learning f is a regression problem. For both problems, when a part of the labels y_i are known and an other part of y_i is unknown during training, learning f is a semi-supervised problem [Zhu07]. Note that when the labels y_i are totally unknown, learning f refers to a clustering problem (unsupervised learning) [JMF99]; [CHY96], out of the scope of this work.

1.1.2 Model selection in supervised learning

A key objective of supervised learning algorithms is to build models f with good generalization capabilities, *i.e.*, models f that correctly predict the labels y_j of new unknown samples \mathbf{x}_j . There exists two types of errors committed by a classification or regression model f : training error and generalization error. **Training error** is the error on the training set and **generalization error** is the error on the testing set. A good supervised model f must not only fit the training data X well, it must also accurately classify records it has never seen before (test set X_{Test}). In other words, a good model f must have low training error as well as low generalization error. This is important because a model that fits the training data too much can have a poorer generalization error than a model with a higher training error. Such a situation is known as model overfitting (Fig. 1.1).

In most cases, learning algorithms require to tune some hyper-parameters. A first approach could consist in trying all the possible combinations of hyper-parameters values and keep the one with the lowest training error. However, as discussed above, the model with the lowest training error is not always the one with the best generalization error. To avoid overfitting, the training set can be divided into 2 sets: a learning and a validation set. Suppose we have two hyper-parameters to tune: C and γ . We make a grid search for each combination (C, γ) of the hyper-parameters, that is in this case a 2-dimensional grid (Fig. 1.2). For each combination (a cell of the grid), the model is learnt on the learning set and evaluated on the validation set. At the end, the model with the lowest error on the validation set is retained. This process is referred as the model selection.

An alternative is cross-validation with v folds, illustrated in Fig. 1.3. In this approach, we partition the training data into v equal-sized subsets. The objective is to evaluate the error for each combination of hyper-parameters. For each run, one fold is chosen for validation,

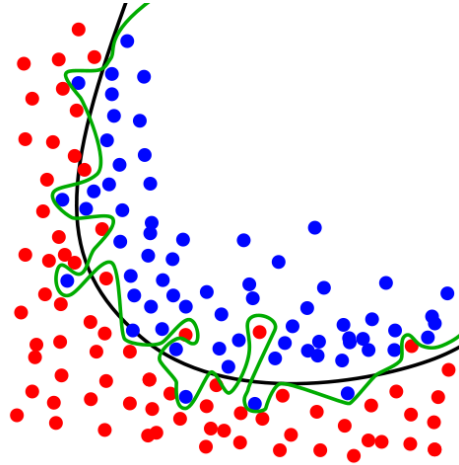


Figure 1.1: An example of overfitting in the case of classification. The objective is to separate blue points from red points. Black line shows a classifier f_1 with low complexity where as green line illustrates a classifier f_2 with high complexity. On training examples (blue and red points), the model f_2 separates all the classes perfectly but may lead to poor generalization on new unseen examples. Model f_1 is often preferred.

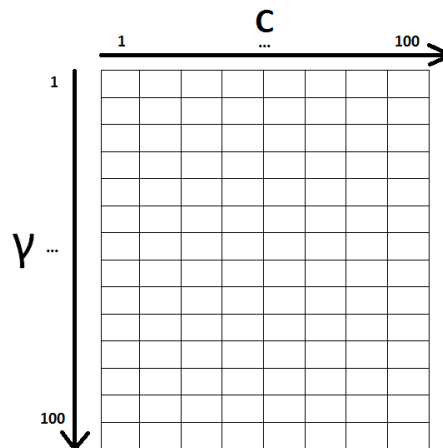


Figure 1.2: Example of a 2 dimensional grid search for parameters C and γ . It defines a grid where each cell of the grid contains a combination (C, γ) . Each combination is used to learn the model and is evaluated on the validation set.

while the $v - 1$ remaining folds are used as the learning set. We repeat the process for each fold, thus v times. Each fold gives one validation error and thus we obtain v errors. The total error for the current combination of hyper-parameters is obtained by summing up the errors for all v folds. When $v = n$, the size of training set, this approach is called leave-one-out or Jackknife. Each test set contains only one sample. The advantage is that as much data as possible are used for training. Moreover, the validation sets are exclusive and they cover the entire data set. The drawback is that it is computationally expensive to repeat the procedure n times. Furthermore, since each validation set contains only one record, the variance of the

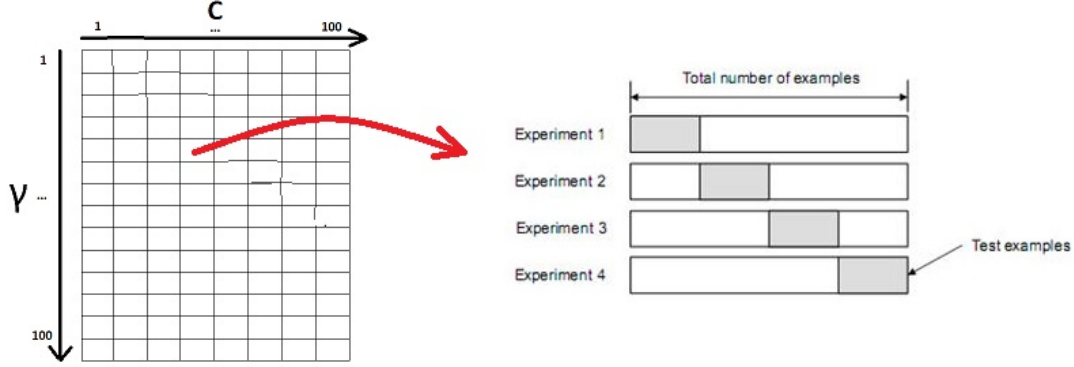


Figure 1.3: v -fold Cross-validation for one combination of parameters. For each of v experiments, use $v - 1$ folds for training and a different fold for Testing, then the training error for this combination of parameter is the mean of all testing errors. This procedure is illustrated for $v = 4$.

estimated performance metric is usually high. This procedure is often used when n , the size training set, is small. There exists other methods such as sub-sampling or bootstraps [OE73]; [G. 06]. We only use cross-validation in our experiments.

To sum up, Fig. 1.4 shows a general approach for solving machine learning problems. In general, a dataset can be divided into 3 sub-datasets (illustrated in Fig. 1.5):

- A **training set** $X = \{\mathbf{x}_i, y_i\}_{i=1}^n$, which consists of n samples \mathbf{x}_i whose labels y_i are known. The training set is used to build the supervised model f . When the learning algorithm requires some hyper-parameters to be tuned, there exists a risk of model overfitting. To avoid such risks, the training set X has to be divided into two subsets :
 - A **learning set** which is used to build the supervised model f for each value of the hyper-parameters.
 - A **validation set** which is used to evaluate the supervised model f for each value of the hyper-parameter. The model f with the lowest error on the validation set is kept, thus ensuring that it has the best generalization abilities.
- A **test set** $X_{Test} = \{\mathbf{x}_j, y_j\}_{j=1}^m$, which consists of m samples \mathbf{x}_j whose labels y_j are also known but are not used during the training step. The model f is applied to predict the label \hat{y}_j of samples \mathbf{x}_j to evaluate the performance of the learnt model by comparing \hat{y}_j and y_j .
- An **operational set** $X_{op} = \{\mathbf{x}_l, y_l\}_{l=1}^L$, which consists of L samples \mathbf{x}_l whose labels y_l are totally unknown. The operational set is in general a new dataset on which the learnt algorithm is applied.

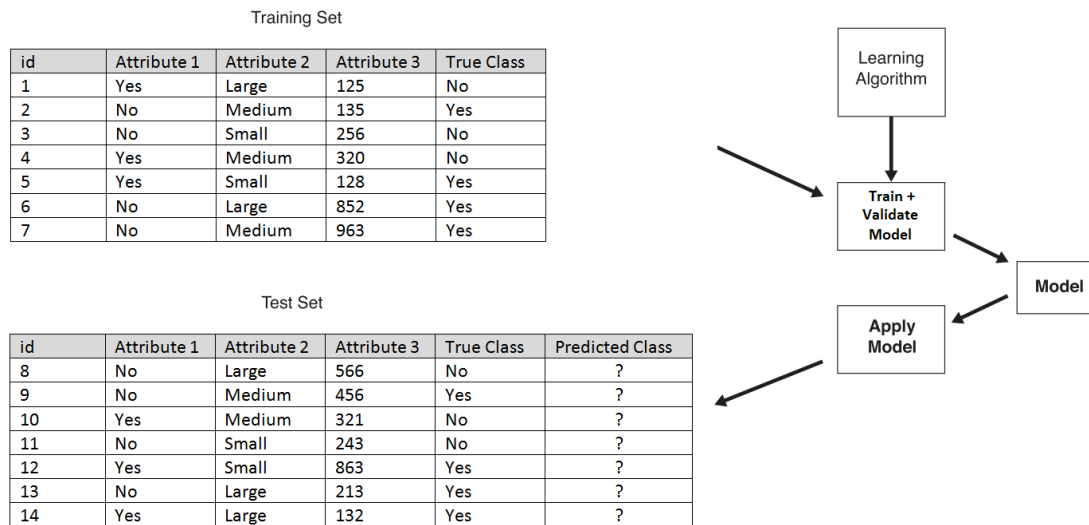


Figure 1.4: General framework for building a supervised (classification/regression) model. Example with 3 features and 2 classes ('Yes' and 'No').

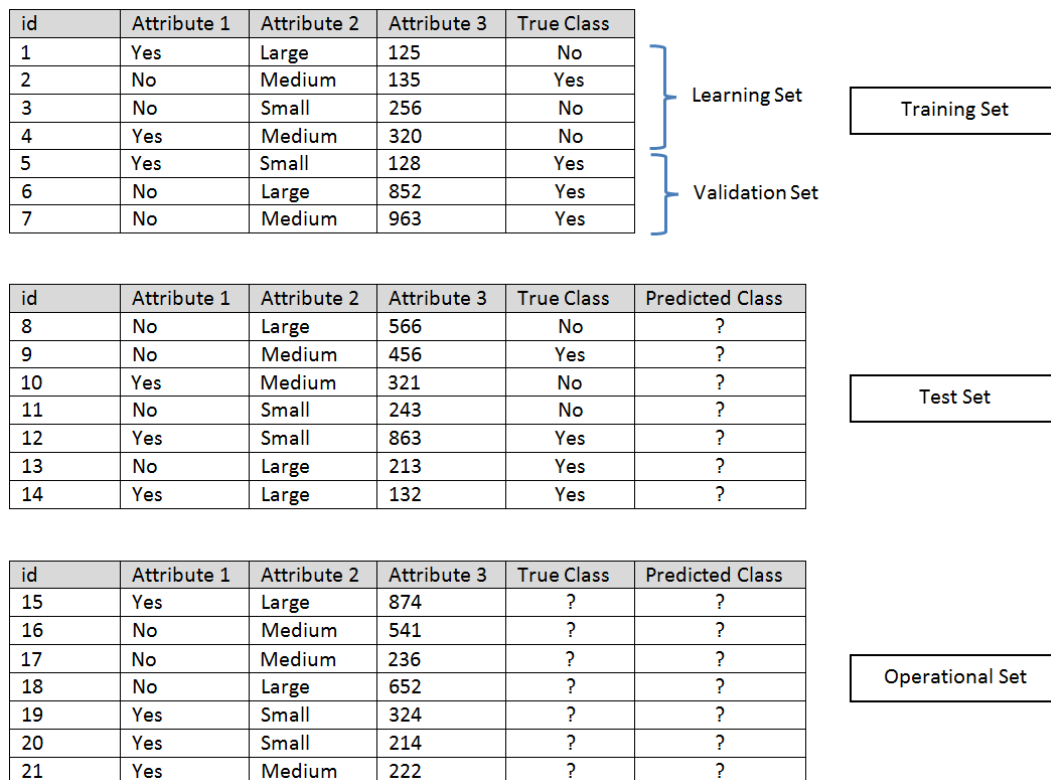


Figure 1.5: Division of a dataset into 3 datasets: training, test and operational.

1.1.3 Model evaluation

As seen in the previous section, model selection is inherently based on the ability to quantify its error on the training and validation sets. In this section, we recall how this error is computed for classification and regression problems.

1.1.3.a Classification evaluation

The performance of a classification model is based on the counts of test samples \mathbf{x}_j correctly and incorrectly predicted by the model f . These counts are tabulated in a table called the confusion matrix. Table 1.1 illustrates the concept for a binary classification problem. Each cell g_{ij} of the table stands for the number of samples from class i predicted to be of class j . Based on this matrix, the number of correct predictions made by the model is $\sum_{i=1}^C g_{ii}$, where C is the number of classes.

		Predicted class	
		Class = 1	Class = 0
Actual Class	Class = 1	g_{11}	g_{10}
	Class = 0	g_{01}	g_{00}

Table 1.1: Confusion matrix for a 2-class problem.

For binary classification problems, g_{11} is the number of true positives, g_{10} is the number of false negatives, g_{01} is the number of false positives and g_{00} is the number of true negatives.

To summarize the information, it is generally more convenient to use performance metrics such as the classification accuracy (Acc) or error rate (Err). This allows to compare several models with a single number. Note that $Err = 1 - Acc$.

$$Acc = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} = \frac{\sum_{i=1}^C g_{ii}}{\sum_{i,j=1}^C g_{ij}} \quad (1.1)$$

$$Err = \frac{\text{Number of wrong predictions}}{\text{Total number of predictions}} = \frac{\sum_{i,j=1, i \neq j}^C g_{ij}}{\sum_{i,j=1}^C g_{ij}} \quad (1.2)$$

Using these performance metrics allows to compare the performance of different classifiers f . It allows to determine in particular whether one learning algorithm outperforms another on a particular learning task on a given test dataset X_{Test} . However, depending on the size of the test dataset, the difference in error rate Err between two classifiers may not be statistically significant. Snedecor & Cochran proposed in 1989 a statistical test based on measuring the

difference between two learning algorithms [Coc77]. It has been used by many researchers [Die97]; [DHB95].

Let consider 2 classifiers f_A and f_B . We test these classifiers on the test set X_{Test} and denote p_A and p_B their respective error rates. The intuition of this statistical test is that when algorithm A classifies an example \mathbf{x}_j from the test set X_{Test} , the probability of misclassification is p_A . Thus, the number m_A (resp. m_B) of misclassification of m test examples made by classifier f_A (resp. f_B) is a binomial random variable with mean mp_A and variance $p_A(1 - p_A)m$. The binomial distribution can be approximated by a normal distribution when m has a reasonable value (Law of large numbers). The difference between two independent normally distributed random variables is also normally distributed with a mean $m(p_A - p_B)$. Thus, the quantity $m_A - m_B$ is a normally distributed random variable. Under the null hypothesis (the two algorithm should have the same error rate), this will have a mean of zero and a standard error se of:

$$se = \sqrt{\frac{2p(1-p)}{m}} \quad (1.3)$$

where $p = \frac{p_A + p_B}{2}$ is the average of the two error probabilities. From this analysis, we obtain the statistic:

$$z = \frac{p_A - p_B}{\sqrt{2p(1-p)/m}} \quad (1.4)$$

which has (approximatively) a standard normal distribution. We can reject the null hypothesis if $|z| > Z_{0.975} = 1.96$ (for a 2-sided test with probability of incorrectly rejecting the null hypothesis of 0.05).

1.1.3.b Regression evaluation

As the concept of classes is restricted to classification problems, the performance of a regression model f is based on metrics that measure the difference between the predicted label \hat{y}_j and the known label y_j . The Mean Absolute Error function (MAE) computes the mean absolute error, a risk metric corresponding to the expected value of the absolute error loss or L_1 -norm loss.

$$MAE = \frac{1}{m} \sum_{j=1}^m |\hat{y}_j - y_j| \quad (1.5)$$

A commonly used performance metrics is the Root Mean Squared Error function ($RMSE$) that computes the root of the mean square error, a risk metric corresponding to the expected value of the squared (quadratic) error loss.

$$RMSE = \sqrt{\frac{1}{m} \sum_{j=1}^m (\hat{y}_j - y_j)^2} \quad (1.6)$$

Many works rely on the R^2 function, the coefficient of determination. It provides a measure of how well future samples are likely to be predicted by the model ¹. It can also be interpreted as a measure of how the model f is better than a constant model .

ref syl-
vain?

$$R^2 = 1 - \frac{\sum_{j=1}^m (\hat{y}_j - y_j)^2}{\sum_{j=1}^m (\bar{y} - y_j)^2} \quad (1.7)$$

where $\bar{y} = \sum_{j=1}^m y_j$ is the mean over the known labels y_j .

1.1.4 Data normalization

Real dataset are often subjected to noise or uneven scaling. Before applying any learning protocol, it is often necessary to pre-process the data: data scaling, data filtering (*e.g.*, denoising), outlier removal, etc. We focus on data normalization (data scaling) in our work.

Part 2 of Sarle's Neural Networks FAQ (1997) ² explains the importance of data normalization for neural networks but they can be applied to any learning algorithms. The main advantage of normalization is to avoid attributes in greater numeric ranges to dominate those in smaller numeric ranges. Another advantage is to avoid numerical difficulties during the calculation. For example, in the case of Support Vector Machine (svm), because kernel values usually depend on the inner products of feature vectors, *i.e.* the linear kernel and the polynomial kernel, large attribute values might cause numerical problems [HCL08].

In most cases, it is recommended to scale each attribute to the range $[-1; +1]$ or $[0; 1]$. Many normalization methods have been proposed such as Min/Max normalization, Z-normalization or normalization of the log distribution . Let $X = \{\mathbf{x}_i, y_i\}_{i=1}^n$ be a training set, \mathbf{x}_i being a sample described by p features x_1, \dots, x_p . We define μ_j and σ_j as the mean and the standard deviation of a variable x_j , applying the Z-normalized variable x_j^{norm} is given by:

références
normal-
ization

$$x_j^{norm} = \frac{x_j - \mu_j}{\sigma_j} \quad (1.8)$$

Finally, we recall some precautions to the practitioner in the learning protocol, experimented by Hsu & al. in the context of svm [HCL08]. First, training and testing data must be scaled using the same method. Second, training and testing data must not be scaled separately. Third, the whole dataset must not be scaled together at the same time. These often leads to poorer results. A proper way to do normalization is to scale the training data, store the parameters of the normalization (*e.g.* μ_i and σ_i for Z-normalization), then apply the same normalization to the testing data.

¹http://scikit-learn.org/stable/modules/model_evaluation.html

²<http://www.faqs.org/faqs/ai-faq/neural-nets/>

1.2 Machine learning algorithms

Many algorithms have been proposed in the context of supervised learning, such as Deep Neural Networks, Decision Trees or Relevance Vector Machines (RVM). Our proposition relies on both results from Support Vector Machines (SVM) in the context of k -Nearest Neighbors (k -NN) classification. We limit the section to present these two algorithms.

Comment [AD1]: dit d'enlever lère phrase mais Michèle dit de garder

1.2.1 k -Nearest Neighbors (k -NN) classifier

A simple approach to classify samples is to consider that "close" samples have a great probability to belong to the same class. Given a test sample \mathbf{x}_j , one can use the class y_i of its nearest neighbor \mathbf{x}_i in the training set in order to predict its labels: $\hat{y}_j = y_i$.

Comment [AR2]: Ahlan trouve que ce n'est pas clair. A refaire

More generally, we can consider the k nearest neighbors of \mathbf{x}_j . The class y_j of a test sample \mathbf{x}_j is assigned with a voting scheme among them, *i.e.*, using the majority of the class of nearest neighbors. This algorithm is referred as the k -Nearest Neighbors algorithm (k -NN) [SJ89]; [CH67]. Fig. 1.6 illustrates the concept for a neighborhood of $k = 3$ and $k = 5$.

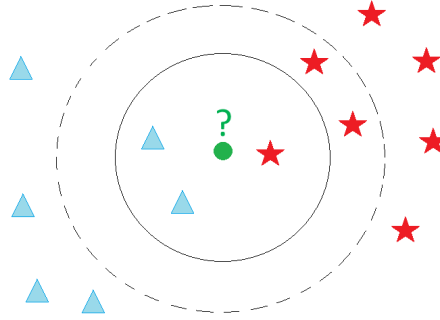


Figure 1.6: Example of k -NN classification. The test sample (green circle) is classified either to the first class (red stars) or to the second class (blue triangles). If $k = 3$ (solid line circle) it is assigned to the second class because there are 2 triangles and only 1 star inside the inner circle. If $k = 5$ (dashed line circle) it is assigned to the first class (3 stars vs. 2 triangles inside the outer circle).

In the k -NN algorithm, the notion of "closeness" between samples \mathbf{x}_i is based on the computation of a metric ³ D . For static data, frequently used metrics are the Euclidean distance, the Minkowski distance or the Mahalanobis distance. Considering a training set X of n samples, solving the 1-NN classification problem is equivalent to solve the following optimization problem: for a new sample \mathbf{x}_j , $\forall i \in \{1 \dots n\}$,

$$y_j = y_{i^*} \quad (1.9)$$

where $i^* = \underset{i \in \{1 \dots n\}}{\operatorname{argmin}} D(\mathbf{x}_i, \mathbf{x}_j)$.

³A clarification of the terms metric, distance, dissimilarity, etc. will be given in Chapter 2. For now, we refer all of them as metrics.

The k -NN algorithm can be extended to estimate continuous labels (regression problems). The procedure is similar. The label y_j is defined as :

$$y_j = \frac{1}{k} \sum_{i=1}^k y_i \quad (1.10)$$

where i corresponds to the index of the k -nearest neighbors [Alt92]. There exists other variants of the k -NN algorithms. In a weighed k -NN, the approach consists in weighting the k -NN decision by assigning to each neighbor \mathbf{x}_i from an unknown sample \mathbf{x}_j , a weight w_i defined as a function of the distance $D(\mathbf{x}_i, \mathbf{x}_j)$ [Dud76]. To cope with uncertainty or imprecision in the labeling of the training data \mathbf{x}_i , other authors propose in a fuzzy k -NN to determine the membership degree in each class of an unseen sample \mathbf{x}_j by combining the memberships of its neighbors [KGG85]. Denoeux propose a framework based on Dempster-Shafer theory where the k -NN rule takes into account the non-representativity of training data, the weighting rule and uncertainty in the labeling [Den95].

**Comment
[AD3]:** Expliquer d'avantage

Despite its simplicity, the k -NN algorithm has been shown to be successful on time series classification problems [BMP02]; [Xi+06]; [Din+08]. However, the k -NN algorithm presents some disadvantages, mainly due to its computational complexity, both in memory space (storage of the training samples \mathbf{x}_i) and time (search of the neighbors) [OE73]. Suppose that we have n labeled training samples in p dimensions, and find the closest neighbors to a test sample \mathbf{x}_j ($k = 1$). In the most simple approach, we look at each stored samples \mathbf{x}_i ($i = 1 \dots n$) one by one, calculate its distance to \mathbf{x}_j ($D(\mathbf{x}_i, \mathbf{x}_j)$) and retain the index of the current closest one. For the standard Euclidean distance, each metric computation is $O(p)$ and thus the search is $O(pn)$. Finally, note that using standard metrics (such as the Euclidean distance) in the k -NN relies on all p dimensions in its computation and thus assumes that all dimensions have the same effect on the metric. This assumption may be wrong and can impact the classification performances.

1.2.2 Support Vector Machine (SVM) algorithm

Support Vector Machine (SVM) is a classification method introduced in 1992 by Boser, Guyon, and Vapnik [BGV92]; [CV95] to solve at first linearly separable problems. The SVM classifier has demonstrated high accuracy, ability to deal with high-dimensional data, good generalization properties and interpretation for various applications from recognizing handwritten digits, to face identification, text categorization, bioinformatics and database marketing [Wan02]; [YL99]; [HHP01]; [SSB03]; [CY11]. SVMs belong to the category of kernel methods, algorithms that depends on the data only through dot-products [SS13]. It allows thus to solve non-linear problems. This section gives a brief overview of the mathematical key points and interpretation of the method. For more informations, the reader can consult [SS13]; [CY11]; [CV95].

We first present an intuition of maximum margin concept. We give the primal formulation of the SVM optimization problem. Then, by transforming the latter formulation into its dual

form, the kernel trick can be applied to learn non-linear classifiers. Finally, we detail how we can interpret the obtained coefficients and how SVMs can be extended for regression problems.

1.2.2.a Intuition

Let $\{\mathbf{x}_i, y_i\}_{i=1}^n$ be a set of n samples $\mathbf{x}_i \in \mathbb{R}^p$ and their labels $y_i = \pm 1$ (2 class-problem). The objective is to learn a hyperplane, whose equation is $\mathbf{w}^T \mathbf{x} + b = 0$, that can separate samples of class +1 from the ones of class -1. When the problem is linearly separable such as in Fig. 1.7, there exists an infinite number of valid hyperplanes.

Comment [AD4]: Mettre dans les figures des + et - pour les classes

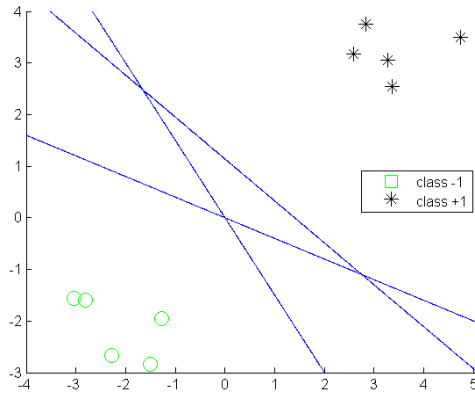


Figure 1.7: Example of linear classifiers in a 2-dimensional classification problem. For a set of points of classes +1 and -1 that are linearly separable, there exists an infinite number of separating hyperplanes corresponding to $\mathbf{w}^T \mathbf{x} + b = 0$.

Vapnik & al. [CV95] propose to choose the separating hyperplane that maximizes the margin, *e.g.* the hyperplane that leaves as much distance as possible between the hyperplane and the closest samples \mathbf{x}_i of each class, called the support vectors. This distance is equal to $\frac{1}{\|\mathbf{w}\|_2}$. We denote $\|\mathbf{w}\|_2$, the L_2 -norm of the vector \mathbf{w} and $\|\mathbf{w}\|_1$ the L_1 -norm of \mathbf{w} :

$$\|\mathbf{w}\|_2 = \sqrt{\mathbf{w}^T \mathbf{w}} = \sqrt{\sum_{h=1}^p w_h^2} \quad (1.11)$$

$$\|\mathbf{w}\|_1 = \sum_{h=1}^p |w_h| \quad (1.12)$$

where $\mathbf{w} = [w_1, \dots, w_p]^T$ denotes the weight vector.

The two hyperplanes passing through the support vectors of each class are referred as the canonical hyperplanes, and the region between the canonical hyperplanes is called the margin band (Fig. 1.8). From this, for a binary classification problem, to classify a new sample \mathbf{x}_j , the decision function is:

$$f(\mathbf{x}_j) = \text{sign}(\mathbf{w}^T \mathbf{x}_j + b) \quad (1.13)$$

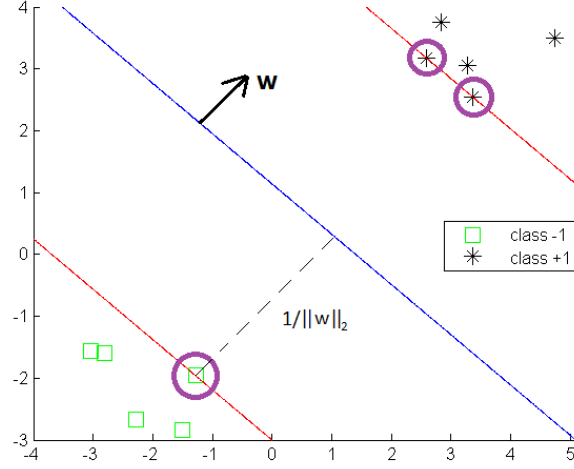


Figure 1.8: The argument inside the decision function of a classifier is $\mathbf{w}^T \mathbf{x} + b$. The separating hyperplane corresponding to $\mathbf{w}^T \mathbf{x} + b = 0$ is shown as a line in this 2-dimensional plot. This hyperplane separates the two classes of data with points on one side labeled $y_i = +1$ ($\mathbf{w}^T \mathbf{x}_i + b \geq 0$) and points on the other side labeled $y_i = -1$ ($\mathbf{w}^T \mathbf{x}_i + b < 0$). Support vectors are circled in purple and lies on the hyperplanes $\mathbf{w}^T \mathbf{x} + b = +1$ and $\mathbf{w}^T \mathbf{x} + b = -1$

1.2.2.b Primal formulation

Finding \mathbf{w} and b by maximizing the margin $\frac{1}{\|\mathbf{w}\|_2}$ is equivalent to minimizing the norm of \mathbf{w} such that all samples from the training set are correctly classified:

$$\underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad (1.14)$$

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad (1.15)$$

This is a constrained optimization problem in which we minimize an objective function (Eq. 1.14) subject to constraints (Eq. 1.15). This formulation is referred as the primal hard margin problem. When the problem is not linearly separable, slack variables $\xi_i \geq 0$ are introduced to relax the optimization problem:

$$\underset{\mathbf{w}, b, \xi}{\operatorname{argmin}} \left(\overbrace{\frac{1}{2} \|\mathbf{w}\|_2^2}^{\text{Regularization}} + C \overbrace{\sum_{i=1}^n \xi_i}^{\text{Loss}} \right) \quad (1.16)$$

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad (1.17)$$

$$\xi_i \geq 0 \quad (1.18)$$

where $C > 0$ is a penalty hyper-parameter.

This formulation is referred as the primal soft margin problem. It is a quadratic programming optimization problem subjected to constraints. Thus, it is a convex problem: any local solutions is a global solution. The objective function in Eq. 1.16 is made of two terms. The first one, the regularization term, penalizes the complexity of the model and thus, controls the ability of the algorithm to generalize on new samples. The second one, the loss term, is an adaptation term to the data. The hyper-parameter C is a trade-off between the regularization and the loss term. When C tends to $+\infty$, all the slack variables ξ_i have to be equal to zero for the loss not to be infinite. The problem is thus equivalent to the primal hard margin problem. The hyper-parameter C is learnt during the training phase (Section 1.1.2).

1.2.2.c Dual formulation

From the primal formulation, it is possible to have an equivalent dual form. This latter formulation allows samples \mathbf{x}_i to appear in the optimization problem through dot-products only. The kernel trick can be applied to extend the methods to learn non-linear classifiers.

First, to simplify the calculation development, let consider the hard margin formulation in Eqs. 1.16, 1.17 and 1.18. As a constrained optimization problem, the formulation is equivalent to the minimization of a Lagrange function $L(\mathbf{w}, b)$, consisting of the sum of the objective function and the n constraints multiplied by their respective Lagrange multipliers $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_n]^T$:

$$\underset{\boldsymbol{\alpha}}{\operatorname{argmax}} \left(L(\mathbf{w}, b) = \frac{1}{2}(\mathbf{w}^T \mathbf{w}) - \sum_{i=1}^n \alpha_i (y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1) \right) \quad (1.19)$$

$$\text{s.t. } \forall i = 1 \dots n :$$

$$\alpha_i \geq 0 \quad (1.20)$$

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0 \quad (1.21)$$

$$\alpha_i (y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1) = 0 \quad (1.22)$$

where $\alpha_i \geq 0$ are the Lagrange multipliers. In optimization theory, Eq. 1.20, 1.21 and 1.22 are called the Karush-Kuhn-Tucker (KKT) conditions. It corresponds to the set of conditions which must be satisfied at the optimum of a constrained optimization problem. The KKT conditions will play an important role in the interpretation of SVM in Section 1.2.2.e.

At the maximum value of $L(\mathbf{w}, b)$, we assume the derivatives with respect to b and \mathbf{w} are set to zero:

$$\begin{aligned} \frac{\partial L}{\partial b} &= - \sum_{i=1}^n \alpha_i y_i = 0 \\ \frac{\partial L}{\partial \mathbf{w}} &= \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0 \end{aligned}$$

that leads to:

$$\sum_{i=1}^n \alpha_i y_i = 0 \quad (1.23)$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad (1.24)$$

By substituting \mathbf{w} into $L(\mathbf{w}, b)$ in Eq. 1.19, we obtain the dual formulation (*Wolfe dual*):

$$\underset{\alpha}{\operatorname{argmax}} \left(\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \right) \quad (1.25)$$

s.t. $\forall i = 1 \dots n :$

$$\sum_{i=1}^n \alpha_i y_i = 0 \quad (1.26)$$

$$\alpha_i \geq 0 \quad (1.27)$$

The dual objective in Eq. 1.25 is quadratic in the parameters α_i . Adding the constraints in Eqs. 1.26 and 1.27, it is a constrained quadratic programming optimization problem (QP). Note that while the primal formulation is minimization, the equivalent dual formulation is maximization. It can be shown that the objective functions of both formulations reach the same value when the solution is found [CY11].

In the same spirit, considering the soft margin primal problem, it can be shown that it leads to the same formulation [CY11] (Eqs. 1.25 and 1.26), except that the Lagrange multipliers α_i are upper bounded by the trade-off C in the soft margin formulation:

$$0 \leq \alpha_i \leq C \quad (1.28)$$

The constraints in Eq. 1.28 are called the Box constraints [CY11]. From the optimal value of α_i , denoted α_i^* , it is possible to compute the weight vector \mathbf{w}^* and the bias b^* at the optimality:

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i \quad (1.29)$$

$$b^* = \sum_{i=1}^n (\mathbf{w}^{*T} \mathbf{x}_i - y_i) \quad (1.30)$$

At the optimality point, Eq. 1.22 leads $\alpha_i^* = 0$ for all datapoints that are well classified and that are not on the margin. Hence, only a few number of datapoints have $\alpha_i^* > 0$ as shown as in Fig. 1.9. These samples are the vector supports. All other datapoints have $\alpha_i^* = 0$, and the decision function is independent of them. Thus, the representation is sparse.

From Eqs. 1.13 & 1.29, to classify a new sample \mathbf{x}_j , the decision function for a binary

classification problem is:

$$f(\mathbf{x}_j) = \text{sign}\left(\sum_{i=1}^n \alpha_i^* y_i (\mathbf{x}_i^T \mathbf{x}_j) + b^*\right) \quad (1.31)$$

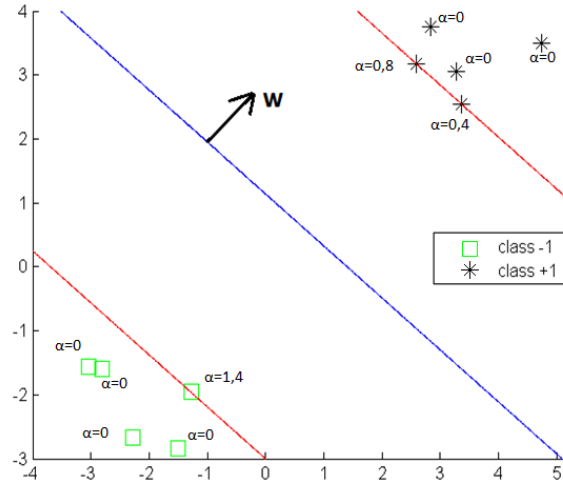


Figure 1.9: Hyperplane obtained after a dual resolution (full blue line). The 2 canonical hyperplanes (red lines) contain the support vectors whose $\alpha_i > 0$. Other points have their $\alpha_i = 0$ and the equation of the hyperplane is only affected by the support vectors.

1.2.2.d Kernel trick

The concept of kernels was introduced by Aizerman & al. in 1964 to design potential functions in the context of pattern recognition [ABR64]. The idea was re-introduced in 1992 by Boser & al. for Support Vector Machine (SVM) and has been received a great number of improvements and extensions to symbolic objects such as text or graphs [BGV92].

From the dual objective in Eq. 1.25, we note that the samples \mathbf{x}_i are only involved in a dot-product. Therefore, if we map these samples \mathbf{x}_i into a higher dimensional hyperspace, called the feature space, we only need to know the dot product in the feature space:

$$(\mathbf{x}_i \cdot \mathbf{x}_j) \rightarrow \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) \quad (1.32)$$

where Φ is the mapping function.

The intuition behind using such mapping is that for many datasets, it is not possible to find a hyperplane that can separate the two classes in the input space if the problem is not linearly separable. However, by applying a transformation Φ , data might become linearly separable in a higher dimensional space. Fig. 1.10 illustrates the idea: in the original 2-dimensional space (left), the two classes can't be separated by a line. However, with a third dimension such that

the +1 (circle) labeled points are moved forward and the -1 (cross) labeled moved back the two classes become separable.

In most of the case, the mapping function Φ does not need to be known since we only need the dot product $\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$. Therefore, we can use any kernel function κ such that: $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$. We call Gram matrix G , the matrix containing all $\kappa(\mathbf{x}_i, \mathbf{x}_j)$:

$$G = (\kappa(\mathbf{x}_i, \mathbf{x}_j))_{1 \leq i, j \leq n} = \begin{pmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \dots & \kappa(\mathbf{x}_1, \mathbf{x}_n) \\ \dots & & \dots \\ \kappa(\mathbf{x}_n, \mathbf{x}_1) & \dots & \kappa(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

Defining a kernel has to follow rules. One of these rules specifies that the kernel function has to define a proper inner product in the feature space. Mathematically, the Gram matrix has to be semi-definite positive (Mercer's theorem) [SS13]. These restricted feature spaces, containing an inner product are called Hilbert spaces.

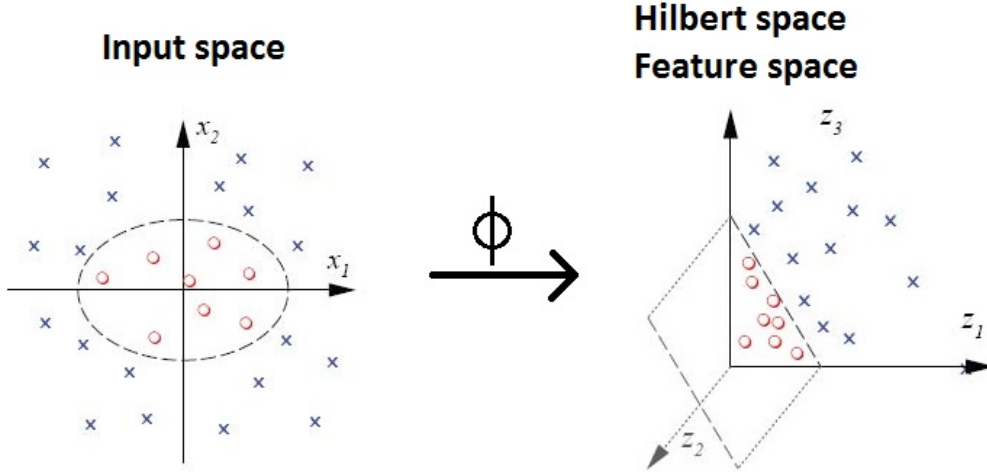


Figure 1.10: Left: in two dimensions the two classes of data (-1 for cross and +1 for circle) are mixed together, and it is not possible to separate them by a line: the data is not linearly separable. Right: using a kernel, these two classes of data become separable by a hyperplane in feature space, which maps to the nonlinear boundary shown, back in input space.⁴

Many kernels have been proposed in the literature such as the polynomial, sigmoid, exponential or wavelet kernels [SS13]. The most popular ones that we will use in our work are respectively the Linear and the Gaussian (or Radial Basis Function (RBF)) kernels:

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j \quad (1.33)$$

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_j - \mathbf{x}_i\|_2^2}{2\sigma^2}\right) = \exp(-\gamma\|\mathbf{x}_j - \mathbf{x}_i\|_2^2) \quad (1.34)$$

where $\gamma = \frac{1}{2\sigma^2}$ is the parameter of the Gaussian kernel and $\|\mathbf{x}_j - \mathbf{x}_i\|_2$ is the Euclidean distance

⁴source: <http://users.sussex.ac.uk/~christ/crs/ml/lec08a.html>

between \mathbf{x}_i and \mathbf{x}_j . Note that the Linear kernel is the identity transformation. In practice, for large scale problem (when the number of dimensions p is high), using a Linear kernel is sufficient [FCH08].

The Gaussian kernel computed between a sample \mathbf{x}_j and a support vector \mathbf{x}_i is an exponentially decaying function in the input space. The maximum value of the kernel ($\kappa(\mathbf{x}_i, \mathbf{x}_j)=1$) is attained at the support vector (when $\mathbf{x}_i = \mathbf{x}_j$). Then, the value of the kernel decreases uniformly in all directions around the support vector, with distance and ranges between zero and one. It can thus be interpreted as a similarity measure. Geometrically speaking, it leads to hyper-spherical contours of the kernel function as shown in Fig. 1.11⁵. The parameter γ controls the decreasing speed of the sphere. In practice, this parameter is learnt during the training phase.

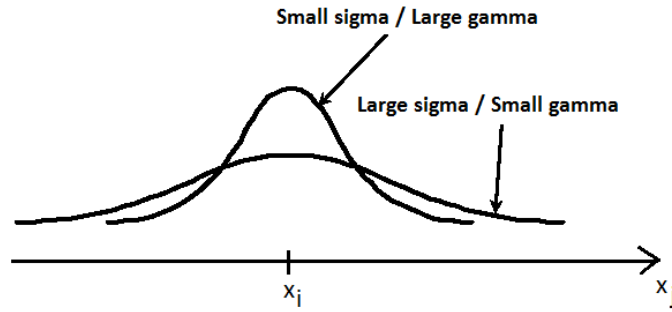


Figure 1.11: Illustration of the Gaussian kernel in the 1-dimensional input space for a small and large γ when \mathbf{x}_i is fixed and \mathbf{x}_j varies.

By applying the kernel trick to the soft margin formulation in Eqs. 1.25, 1.26 and 1.28, the following optimization problem allows to learn non-linear classifiers:

$$\underset{\alpha}{\operatorname{argmax}} \left(\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \right) \quad (1.35)$$

$$\text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad (1.36)$$

$$0 \leq \alpha_i \leq C \quad (1.37)$$

The decision function f becomes:

$$f(\mathbf{x}_j) = \operatorname{sign} \left(\sum_{i=1}^n \alpha_i^* y_i \kappa(\mathbf{x}_i, \mathbf{x}_j) + b^* \right) \quad (1.38)$$

Let n_{SV} be the number of support vectors ($n_{SV} \leq n$). To recover b^* , we recall that for support

⁵<https://www.quora.com/Support-Vector-Machines/What-is-the-intuition-behind-Gaussian-kernel-in-SVM>

vectors \mathbf{x}_i :

$$y_j \left(\sum_{i=1}^{n_{SV}} \alpha_i^* y_i \kappa(\mathbf{x}_i, \mathbf{x}_j) + b^* \right) = 1 \quad (1.39)$$

From this, we can solve b^* using an arbitrarily chosen support vector \mathbf{x}_i :

$$b^* = \frac{1}{y_j} - \sum_{i=1}^{n_{SV}} \alpha_i^* y_i \kappa(\mathbf{x}_i, \mathbf{x}_j) \quad (1.40)$$

Note that in this case, we can't recover the weight vector \mathbf{w}^* but it is not useful here for the decision function.

1.2.2.e Interpretation

Interpretation in the primal

We recall that \mathbf{x}_i is a sample in p dimensions: x_1, \dots, x_p . Geometrically, the vector \mathbf{w} represents the direction of the hyperplane and points towards the direction of positive decision function $f(\mathbf{x}) \geq 0$ (Fig. 1.12). The absolute value of the bias $|b|$ is equal to the distance of the hyperplane to the origin point $\mathbf{x} = \mathbf{0}$ ⁶ if the norm of the vector \mathbf{w} is equal to 1. In the soft margin problem, the slack variables ξ_i can be interpreted as follows:

- $\xi = 0$ implies that \mathbf{x}_i is correctly classified and is either on the margin or on the correct side of the margin.
- $0 < \xi \leq 1$ implies that \mathbf{x}_i lies inside the margin, but on the correct side of the decision boundary
- $\xi \geq 1$ implies that \mathbf{x}_i lies on the wrong side of the decision boundary and is misclassified.

In the primal, the weight vector $\mathbf{w} = [w_1, \dots, w_p]^T$ contains as many elements as there are dimensions in the dataset, *i.e.*, $\mathbf{w} \in \mathbb{R}^p$. The magnitude of each element in \mathbf{w} denotes the importance of the corresponding variable for the classification problem. If the element of \mathbf{w} for some variable is 0, these variables are not used for the classification problem.

In order to visualize the above interpretation of the weight vector \mathbf{w} , let us examine several hyperplanes $\mathbf{w}^T \mathbf{x} + b = 0$ shown in Fig. 1.13 with $p = 2$. Fig. 1.13(a) shows a hyperplane where elements of \mathbf{w} are the same for both variables x_1 and x_2 . The interpretation is that both variables contribute equally for classification of objects into positive and negative. Fig. 1.13(b) shows a hyperplane where the element of \mathbf{w} for x_1 is 1, while that for x_2 is 0. This is interpreted as that x_1 is important but x_2 is not. An opposite example is shown in Fig. 1.13(c) where x_2 is considered to be important but x_1 is not. Finally, Fig. 1.13(d) provides a 3-dimensional example ($p = 3$) where an element of \mathbf{w} for x_3 is 0 and all other elements are equal to 1. The interpretation is that x_1 and x_2 are important but x_3 is not.

⁶ $\mathbf{0}$ stands for the null vector: $\mathbf{0} = [0, \dots, 0]^T$

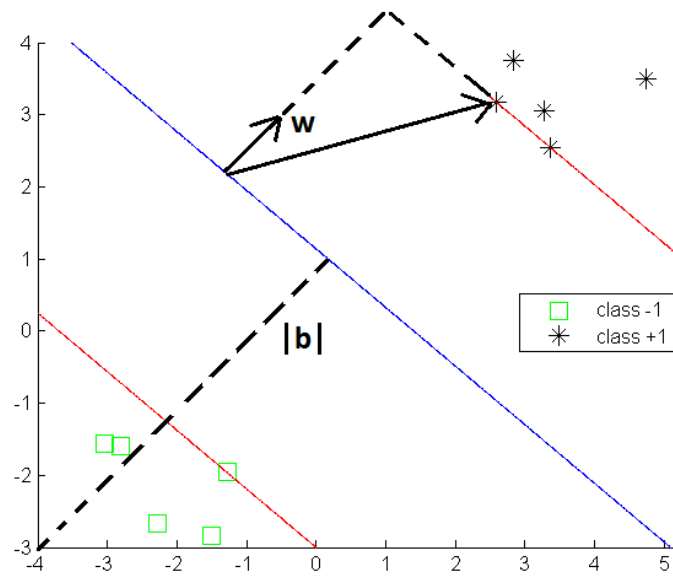
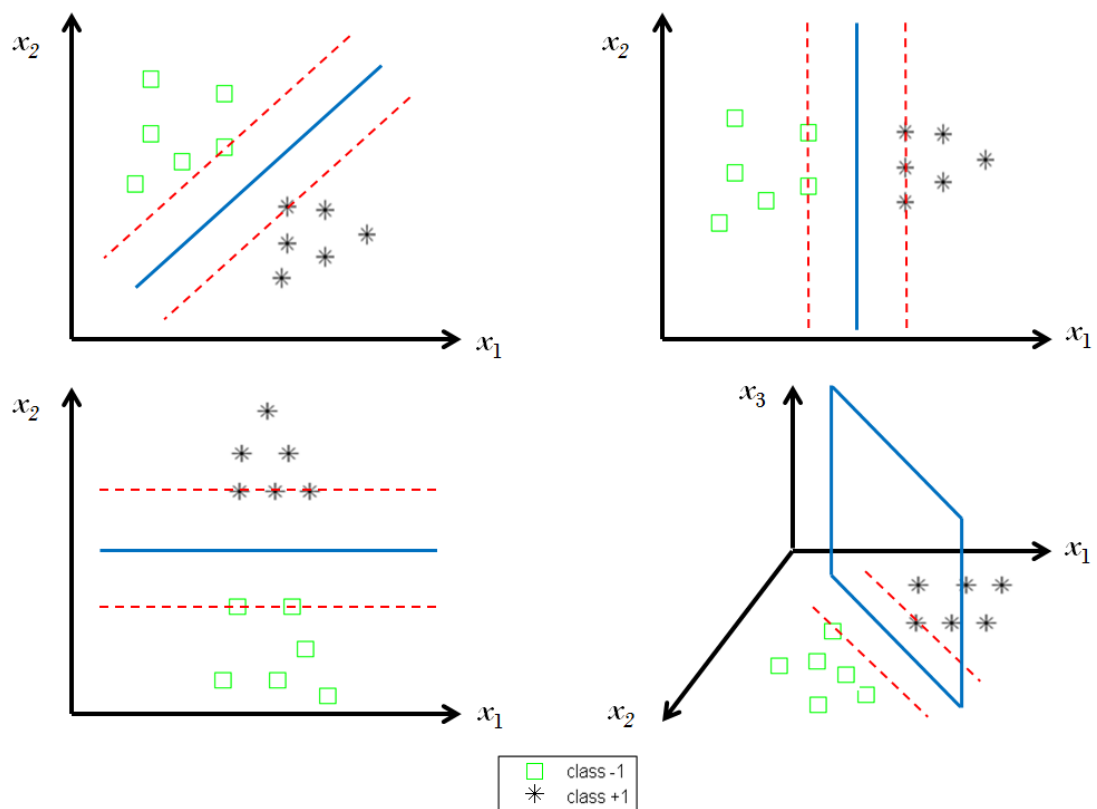


Figure 1.12: Geometric representation of SVM.

Figure 1.13: Example of several SVMs and how to interpret the weight vector \mathbf{w}

Another way to interpret how much a variable contributes in the vector \mathbf{w} is to express the contribution in percentage: the ratio $\frac{w_j}{\|\mathbf{w}\|_2} \cdot 100$ defines the percentage of contribution for each variable x_j in the SVM model. The interpretation is only valid if the variables x_j of the time series are normalized before learning the SVM model, they evolve in the same range.

Interpretation in the dual

As a constrained optimization, the dual form satisfies the Karush-Kuhn-Tucker (KKT) conditions (Eqs. 1.20, 1.21 and 1.22). We recall Eq. 1.22:

$$\alpha_i(y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1) = 0$$

From this, for every datapoint \mathbf{x}_i , either $\alpha_i^* = 0$ or $y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$. Any datapoint with $\alpha_i^* = 0$ do not appear in the sum of the decision function f in Eq. 1.31 or 1.38. Hence, they play no role for the classification decision of a new sample \mathbf{x}_j . The others \mathbf{x}_i such that $\alpha_i^* > 0$ corresponds to the support vectors. Looking at the distribution of α_i^* allows also to have either a better understanding of the datasets, or either to detect outliers. The higher the coefficient α_i^* for a sample \mathbf{x}_i is, the more the sample \mathbf{x}_i impacts on the decision function f . However, an unusually high value of α_i^* among the samples can lead to two interpretations: either this point is a critical point to the decision, or this point is an outlier. In the soft margin formulation, by constraining α_i^* to be inferior to C (Box constraints) the effect of outliers can be reduced and controlled.

1.2.2.f Variants of SVM

From the primal formulation of SVM (Eqs. 1.14 & 1.15), some works investigate the effect of modifications in the regularization and loss term [HCL08].

First, the two common regularizers are $\|\mathbf{w}\|_1$ and $\|\mathbf{w}\|_2$. The former is referred to as L_1 -Regularizer while the latter is L_2 -Regularizer. L_1 -Regularizer is used to obtain sparser models than L_2 -Regularizer, *i.e.*, the vector \mathbf{w} will contain many elements w_i that will equal to zero. Thus, it can be used for variable selection.

Secondly, the two common loss functions ξ_i are $\max(1 - y_i \mathbf{w}^T \mathbf{x}_i, 0)$ and $[\max(1 - y_i \mathbf{w}^T \mathbf{x}_i, 0)]^2$. The former is referred to as L_1 -Loss and the latter is L_2 -Loss function. L_2 -loss function will penalize more slack variables ξ_i during training and would be more sensitive to outliers. Theoretically, it should lead to less error in training and poorer generalization in most of the case [HCL08]. In general, L_1 -Loss is preferred.

1.2.2.g Extensions of SVM

SVM has received many interests in recent years. Many extensions has been developed such as ν -SVM, asymmetric soft margin SVM or multiclass SVM [KU02]; [CS01]. One interesting extension is the extension of Support Vector Machine to regression problems, also called

Support Vector Regression (svr). The objective is to find a linear regression model $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$. To preserve the property of sparseness, the idea is to consider an ϵ -insensitive error function. It gives zero error if the absolute difference between the prediction $f(\mathbf{x}_i)$ and the target y_i is less than ϵ where $\epsilon > 0$ penalize samples that are outside of a ϵ -tube as shown as in Fig. 1.14.

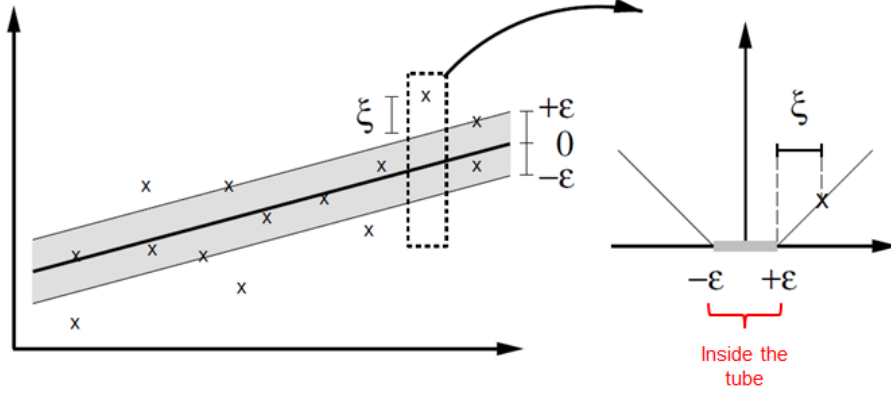


Figure 1.14: Illustration of SVM regression (left), showing the regression curve with the ϵ -insensitive "tube" (right). Samples \mathbf{x}_i above the ϵ -tube have $\xi_1 > 0$ and $\xi_1 = 0$, points below the ϵ -tube have $\xi_2 = 0$ and $\xi_2 > 0$, and points inside the ϵ -tube have $\xi = 0$.

The ϵ -insensitive error function E_ϵ is defined by:

$$E_\epsilon(f(\mathbf{x}_i) - y_i) = \begin{cases} 0 & \text{if } |f(\mathbf{x}_i) - y_i| < \epsilon \\ |f(\mathbf{x}_i) - y_i| - \epsilon & \text{otherwise} \end{cases} \quad (1.41)$$

The soft margin optimization problem in its primal form is formalized as:

$$\underset{\mathbf{w}, b}{\operatorname{argmin}} \left(\overbrace{\frac{1}{2} \|\mathbf{w}\|_2^2}^{\text{Regularization}} + C \overbrace{\sum_{i=1}^n (\xi_{i_1} + \xi_{i_2})}^{\text{Loss}} \right) \quad (1.42)$$

s.t. $\forall i = 1 \dots n :$

$$y_i - (\mathbf{w}^T \mathbf{x}_i + b) \geq \epsilon - \xi_{i_1} \quad (1.43)$$

$$(\mathbf{w}^T \mathbf{x}_i + b) - y_i \geq \epsilon - \xi_{i_2} \quad (1.44)$$

$$\xi_{i_1} \geq 0 \quad (1.45)$$

$$\xi_{i_2} \geq 0 \quad (1.46)$$

The slack variables are divided into 2 kind of slacks variables, one for samples above the decision function f (ξ_{i_1}), and one for samples under the decision function f (ξ_{i_2}). As for SVM,

it is possible to have a dual formulation:

$$\operatorname{argmax}_{\boldsymbol{\alpha}} \left(\sum_{i=1}^n y_i (\alpha_{i_1} - \alpha_{i_2}) - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_{i_1} - \alpha_{i_2}) (\alpha_{j_1} - \alpha_{j_2}) (\mathbf{x}_i \cdot \mathbf{x}_j) \right) \quad (1.47)$$

s.t. $\forall i = 1 \dots n :$

$$\sum_{i=1}^n \alpha_{i_1} = \sum_{i=1}^n \alpha_{i_2} \quad (1.48)$$

$$0 \leq \alpha_{i_1} \leq C \quad (1.49)$$

$$0 \leq \alpha_{i_2} \leq C \quad (1.50)$$

As in SVM, we obtain three possible regression functions for a new sample \mathbf{x}_j , respectively in its primal, dual, and non-linear form:

$$f(\mathbf{x}_j) = \mathbf{w}^T \mathbf{x}_j + b \quad (1.51)$$

$$f(\mathbf{x}_j) = \sum_{i=1}^n (\alpha_{i_1}^* - \alpha_{i_2}^*) (\mathbf{x}_i \cdot \mathbf{x}_j) + b \quad (1.52)$$

$$f(\mathbf{x}_j) = \sum_{i=1}^n (\alpha_{i_1}^* - \alpha_{i_2}^*) K(\mathbf{x}_i, \mathbf{x}_j) + b \quad (1.53)$$

More informations about the calculation development can be found in [Bis06].

1.2.3 Other classification algorithms

Partie non encore rédigée. A faire à la fin.

- Positionner les travaux par rapport aux autres méthodes d'apprentissage supervisé
- S'intéresser au Deep neural network (à la mode en ce moment)
- RVM, Decision Tree,
- Ne pas trop développer
- Dans notre cas, on ne s'intéressera pas à ce type d'algorithmes (type deep learning) car il ne repose pas sur une notion de distance et les features qui sont trouvés ne sont pas interprétables

1.3 Conclusion of the chapter

A refaire complètement à la fin.

This chapter reviews the different steps in a machine learning framework: data normalization, model selection and model evaluation. We focus on two machine learning algorithms used in our proposition: the k -Nearest Neighbors (k -NN) and the Support Vector Machine (SVM). In the following, we consider the k -NN as our classifier. The SVM will be used in our work for its large margin concept.

Our objective being the learning of a metric that optimizes the performances of the k -NN classifier, we review in the next section some metrics proposed for time series as well as metric learning concept for static data.

Time series metrics and metric learning

Contents

2.1	Definition of a time series	29
2.2	Properties and representation of a metric	32
2.3	Unimodal metrics for time series	33
2.3.1	Amplitude-based metrics	33
2.3.2	Frequential-based metrics	34
2.3.3	Behavior-based metrics	35
2.3.4	Other metrics and Kernels for time series	36
2.4	Time series alignment and dynamic programming approach	37
2.5	Combined metrics for time series	39
2.5.1	Combination functions	39
2.5.2	Impact of normalization	40
2.6	Metric learning	43
2.6.1	Review on metric learning work	43
2.6.2	Large Margin Nearest Neighbors (LMNN)	44
2.6.3	Parallels between LMNN and SVM	46
2.7	Conclusion of the chapter	47

In this chapter, we first present the definition of time series. Then, we recall the general properties of a metric and introduce some metrics proposed for time series. In particular, we focus on amplitude-based, behavior-based and frequential-based metrics. As real time series are subject to varying size and delays, we recall the concept of alignment and dynamic programming. Then, we present some proposed combined metrics for time series. Finally, we review the concept of metric learning.

2.1 Definition of a time series

Time series are data that can be frequently found in various emerging applications such as sensor networks, smart buildings, social media networks or Internet of Things (IoT) [Naj+12];

Ajouter
virtual
sensors?

[Ngu+12]; [YG08]. They are involved in many learning problems such as recognizing a human movement in a video, detecting a particular operating mode, etc. [PAN+08]; [Ram+08]. In **clustering** problems, one would like to organize similar time series together into homogeneous groups. In **classification** problems, the aim is to assign time series to one of several predefined categories (e.g., different types of defaults in a machine). In **regression** problems, the objective is to predict a continuous value from observed time series (e.g., forecasting the measurement of a power meter from pressure and temperature sensors). Due to their temporal and structured nature, time series constitute complex data to be analyzed by classic machine learning approaches.

For physical systems, a time series of duration T can be seen as a signal, sampled at a frequency f_e , in a temporal window $[0; T]$. From a mathematical perspective, a time series of length Q is a collection of a finite number of observations made sequentially at discrete time instants $t = 1, \dots, Q$. Note that $Q = T f_e$.

Let $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iQ})$ be a univariate time series of length Q . Each observation x_{it} is bounded (i.e., the infinity is not a valid value: $x_{it} \neq \pm\infty$). The time series \mathbf{x}_i is said to be univariate if the collection of observations x_{it} ($t = 1, \dots, Q$) comes from the observation of one variable (i.e., the temperature measured by one sensor). When simultaneous observation of p variables (several sensors such as the temperature, the pressure, etc.) are made at the same time, the time series is said to be multivariate. From this, one possible representation would be:

$$\mathbf{x}_i = (x_{i1,1}, \dots, x_{i1,p}, x_{i2,1}, \dots, x_{i2,p}, \dots, x_{iQ,1}, \dots, x_{iQ,p})$$

where $\mathbf{x}_{it,j}$ is the observation of the time series \mathbf{x}_i at the time instant t along the variable x_j . An other possible representation could consider a multivariate time series \mathbf{x}_i as the union of multiple univariate time series. In this case:

$$\mathbf{x}_i = (\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,p}) = (x_{i1,1}, \dots, x_{iQ,1}, x_{i1,2}, \dots, x_{iQ,2}, \dots, x_{iQ,p})$$

where $\mathbf{x}_{i,j} = (x_{i1,j}, \dots, x_{iQ,j})$. For simplification purpose, we only consider univariate time series in the following.

Some authors propose to extract representative features from time series. Fig. 2.2 illustrates a model for time series proposed by Chatfield in [Cha04]. It states that a time series can be decomposed into 3 components: a trend, a cycle (periodic component) and a residual (irregular variations).

According to Chatfield, most time series exhibit either or both a long term change in the mean (trend) and a periodic (cyclic) component. The trend can be linear, quadratic, etc. The cyclic component is a variation at a fixed period of time (seasonality) such as for example the seasonal variation of temperature. In practice, the 3 features (trend, cycle, residuals) are rarely sufficient for the classification or regression of real time series.

¹This time series can be downloaded from <http://www.york.ac.uk/depts/maths/data/ts/ts04.dat>

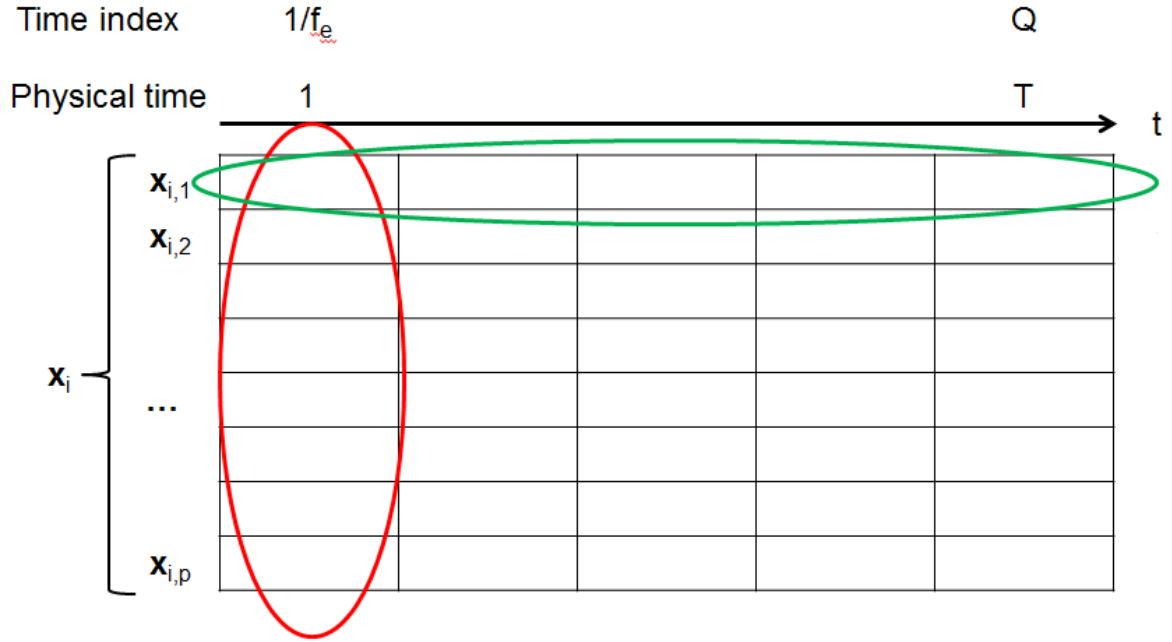
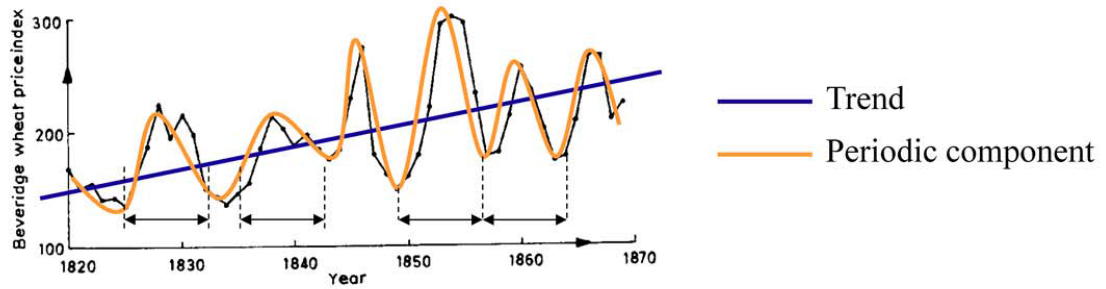


Figure 2.1: 2 modes of representation

Figure 2.2: The Beveridge wheat price index is the average in nearly 50 places in various countries measured in successive years from 1500 to 1869¹.

Other authors made the hypothesis of time independency between the observations x_{it} . They consider time series as a static vector data and use classic machine learning algorithms [Lia+12]; [CT01]; [HWZ13]; [HHK12]. Our work focuses on classification and regression problems, and on time series comparison through metrics.

2.2 Properties and representation of a metric

A mapping $D : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}^+$ over a vector space \mathbb{R}^p is called a metric or a distance if for all vectors $\forall \mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_l \in \mathbb{R}^p$, it satisfies the properties [DD09]:

Comment [CTD5]: je préfère garder l'espace pour + de visibilité

1. $D(\mathbf{x}_i, \mathbf{x}_j) \geq 0$ (positivity)
2. $D(\mathbf{x}_i, \mathbf{x}_j) = D(\mathbf{x}_j, \mathbf{x}_i)$ (symmetry)
3. $D(\mathbf{x}_i, \mathbf{x}_j) = 0 \Leftrightarrow \mathbf{x}_i = \mathbf{x}_j$ (distinguishability)
4. $D(\mathbf{x}_i, \mathbf{x}_j) + D(\mathbf{x}_j, \mathbf{x}_l) \geq D(\mathbf{x}_i, \mathbf{x}_l)$ (triangular inequality)

A mapping D that satisfies at least properties 1, 2 and 3 in its implication, called reflexivity ($\mathbf{x}_i = \mathbf{x}_j \Rightarrow D(\mathbf{x}_i, \mathbf{x}_j) = 0$) is called a dissimilarity, and the one that satisfies at least properties 1, 2, 4 a pseudo-metric. A metric, a dissimilarity and a pseudo metric can be both interpreted as a measure of how "different" two samples are. Note that they can be used for comparisons of several samples: if a time series \mathbf{x}_i is expected to be closer to \mathbf{x}_j than to \mathbf{x}_l , then $D(\mathbf{x}_i, \mathbf{x}_j) \leq D(\mathbf{x}_i, \mathbf{x}_l)$. On the contrary, the mapping is called a similarity S when the sample \mathbf{x}_i is expected to be closer to \mathbf{x}_j than to \mathbf{x}_l and then $S(\mathbf{x}_i, \mathbf{x}_j) \geq S(\mathbf{x}_i, \mathbf{x}_l)$. To simplify the discussion in the following, we refer to pseudo-metric and dissimilarity as metrics, pointing out the distinction only when necessary.

Metric can be represented in two ways (Fig. 2.3). First, data points (samples \mathbf{x}_i) can be fixed and the distance sphere is shown. Secondly, the distance sphere can be fixed and the data points and the axis are moving.

figure à revoir

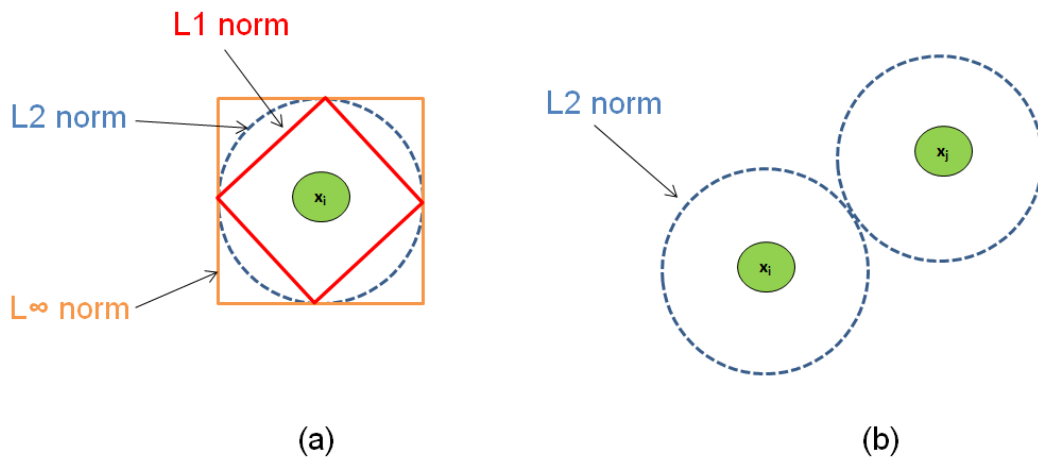


Figure 2.3: Example of metric representation: (a) data points can be fixed and the distance sphere is shown. (b) sphere can be fixed and the data points and the axis are moving.

2.3 Unimodal metrics for time series

Defining and evaluating metrics for time series has become an active area of research for a wide variety of problems in machine learning [Din+08]; [Naj+12]. In the following, we suppose that time series have the same duration T and have been regularly sampled at the frequency f_e . Therefore, they have the same length $Q = Tf_e$. Let $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iQ})$ and $\mathbf{x}_j = (x_{j1}, x_{j2}, \dots, x_{jQ})$ be two univariate time series of length Q .

A large number of distance measures have been proposed in the literature [MV14]. Contrary to static data, time series may exhibit modalities and specificities due to their temporal nature (e.g., value, shape, frequency, delay, temporal locality). In this section, we review 3 categories of time series metrics used in our work: amplitude-based, frequential-based and behavior-based.

2.3.1 Amplitude-based metrics

The most usual comparison measures are amplitude-based metrics, where time series are compared in the temporal domain on their amplitudes regardless of their behaviors or frequential characteristics. Among these metrics, there are the commonly used Euclidean distance that compares elements observed at the same time [Din+08]:

$$d_E(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{t=1}^Q (x_{it} - x_{jt})^2} \quad (2.1)$$

Note that the Euclidean distance is a particular case of the Minkowski L_p norm ($p = 2$). An other amplitude-based metric is the Mahalanobis distance [PL12], defined as a dissimilarity measure weighted by a matrix \mathbf{M} :

$$d_M(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M}^{-1} (\mathbf{x}_i - \mathbf{x}_j) \quad (2.2)$$

If the covariance matrix \mathbf{M} is the identity matrix, the Mahalanobis distance is equal to the Euclidean distance. If the covariance matrix \mathbf{M} is diagonal, then the resulting distance measure is called a normalized Euclidean distance:

$$d_M(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{t=1}^Q \frac{(x_{it} - x_{jt})^2}{m_t}} \quad (2.3)$$

where m_t is the variance of the x_{it} and x_{jt} over the sample set. Note that this is equivalent to normalize each feature: $x'_{il} = x_{il}/\sqrt{m_l}$ and use the Euclidean distance on the normalized features. In the following of the work, we consider the standard Euclidean distance d_E as the amplitude-based distance d_A .

In the example of Fig. 2.4, the aim is to determine which time series (\mathbf{x}_2 or \mathbf{x}_3) is the

Comment [SMA6]: reformulation à revoir: let's try to

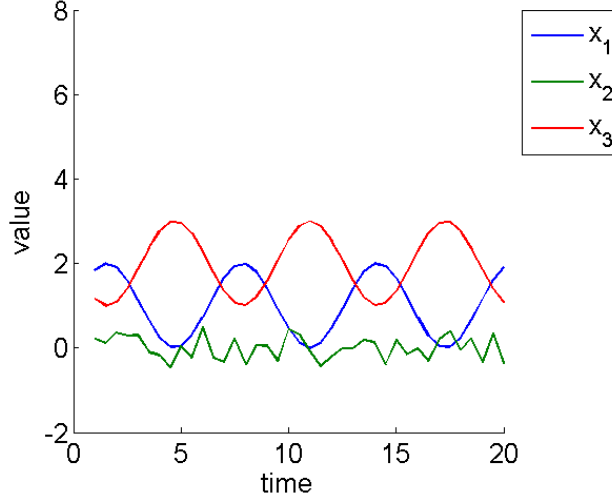


Figure 2.4: 3 toy time series. Time series in blue and red are two sinusoidal signals. Time series in green is a random signal.

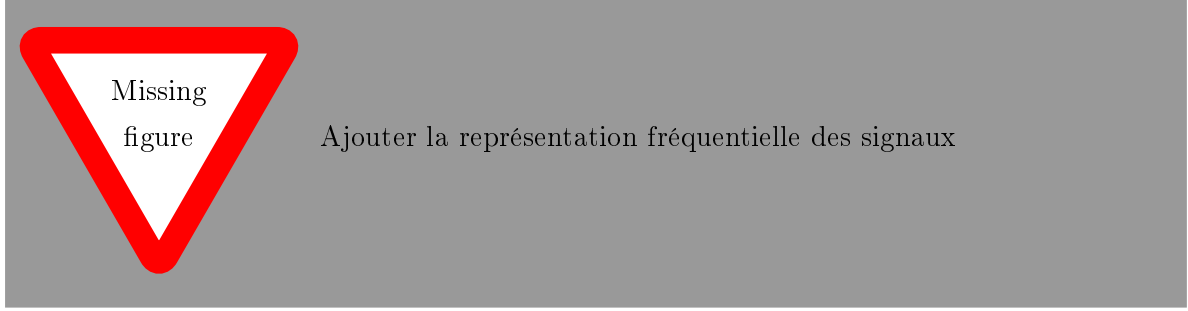
closest to \mathbf{x}_1 . The amplitude-based distance d_A states that \mathbf{x}_2 is closer to \mathbf{x}_1 than \mathbf{x}_3 since $d_A(\mathbf{x}_1, \mathbf{x}_2) = 7.8816 < d_A(\mathbf{x}_1, \mathbf{x}_3) = 31.2250$.

2.3.2 Frequential-based metrics

The second category, commonly used in signal processing, relies on comparing time series based on their frequential properties (e.g. Fourier Transform, Wavelet, Mel-Frequency Cepstral Coefficients [SS12b]; [TC98]; [BM67]). In our work, we limit the frequential comparison to Discrete Fourier Transform [Lhe+11], but other frequential properties can be used as well. Thus, for time series comparison, first the time series \mathbf{x}_i are transformed into their Fourier representation $\tilde{\mathbf{x}}_i = [\tilde{x}_{i1}, \dots, \tilde{x}_{iF}]$, with \tilde{x}_{if} the complex component at frequential index f . The Euclidean distance is then used between their respective complex number modules $|\tilde{x}_{if}|$:

$$d_F(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{f=1}^F (|\tilde{x}_{if}| - |\tilde{x}_{jf}|)^2} \quad (2.4)$$

In the example of Fig. 2.4, the frequential-based distance d_F states that the time series \mathbf{x}_3 is closer to \mathbf{x}_1 than \mathbf{x}_2 since $d_F(\mathbf{x}_1, \mathbf{x}_3) = 0.8519 < d_F(\mathbf{x}_1, \mathbf{x}_2) = 0.9250$. This can be illustrated in the Frequency domain (Fig. ??)



2.3.3 Behavior-based metrics

The third category of metrics aims to compare time series based on their shape or behavior despite the range of their amplitudes. By time series of similar behavior, it is generally intended that for all temporal window $[t, t']$, they increase or decrease simultaneously with the same growth rate. On the contrary, they are said of opposite behavior if for all $[t, t']$, if one time series increases, the other one decreases and (vise-versa) with the same growth rate in absolute value. Finally, time series are considered of different behaviors if they are not similar, nor opposite. Many applications refer to the Pearson correlation [AT10]; [Ben+09] for behavior comparison. A generalization of the Pearson correlation is introduced in [DCA11]:

$$cort_r(\mathbf{x}_i, \mathbf{x}_j) = \frac{\sum_{t,t'=1}^Q (x_{it} - x_{it'})(x_{jt} - x_{jt'})}{\sqrt{\sum_{t,t'=1}^Q (x_{it} - x_{it'})^2} \sqrt{\sum_{t,t'=1}^Q (x_{jt} - x_{jt'})^2}} \quad (2.5)$$

where $|t - t'| \leq r$, $r \in [1, \dots, Q - 1]$. The parameter r can be tuned or fixed a priori. It measures the importance of noise in data. For non-noisy data, low orders r is generally sufficient. For noisy data, the practitioner can either use de-noising data techniques (Kalman or Wiener filtering [Kal60]; [Wie42]), or fix a high order r .

The temporal correlation $cort$ computes the sum of growth rate between \mathbf{x}_i and \mathbf{x}_j between all pairs of values observed at $[t, t']$ for $t' \leq t + r$ (r -order differences). The value $cort_r(\mathbf{x}_i, \mathbf{x}_j) = +1$ means that \mathbf{x}_i and \mathbf{x}_j have similar behavior, i.e, there exists some constant c such that $\mathbf{x}_i = \mathbf{x}_j + c$. The value $cort_r(\mathbf{x}_i, \mathbf{x}_j) = -1$ means that \mathbf{x}_i and \mathbf{x}_j have opposite behavior, i.e, there exists some constant c such that $\mathbf{x}_i = -\mathbf{x}_j + c$. Finally, $cort_r(\mathbf{x}_i, \mathbf{x}_j) = 0$ expresses that their growth rates are stochastically linearly independent (different behaviors).

When $r = Q - 1$, it leads to the Pearson correlation. As $cort_r$ is a similarity measure, it can be transformed into a dissimilarity measure:

$$d_B(\mathbf{x}_i, \mathbf{x}_j) = \frac{1 - cort_r(\mathbf{x}_i, \mathbf{x}_j)}{2} \quad (2.6)$$

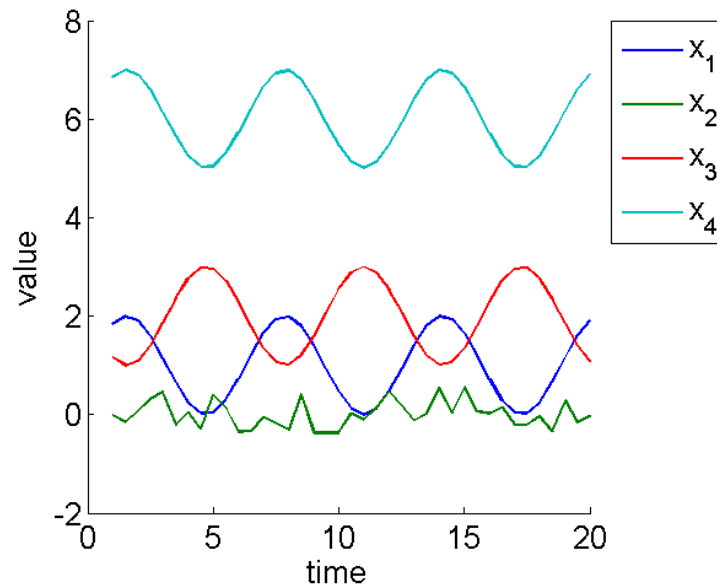


Figure 2.5: The signal from Fig. 2.4 and a signal \mathbf{x}_4 which is signal \mathbf{x}_1 and an added translation. Based on behavior comparison, \mathbf{x}_4 is the closest to \mathbf{x}_1 .

Considering Fig. 2.5, we obtain:

$$d_B(\mathbf{x}_1, \mathbf{x}_2) = 0.477$$

$$d_B(\mathbf{x}_1, \mathbf{x}_3) = 1$$

$$d_B(\mathbf{x}_1, \mathbf{x}_4) = 0$$

2.3.4 Other metrics and Kernels for time series

A faire à la fin, pas urgent

- Il existe dans la littérature de nombreuses autres métriques pour les séries temporelles (laisser la porte ouverte).
- Certaines métriques sont utilisées dans le domaine temporelle
- D'autres métriques sont utilisés dans d'autres représentations (Wavelet, etc.)
- Certaines combinent la représentation temporelles et fréquentielles (Représentation spectrogramme en temps-fréquence)
- Se baser sur l'article "TSclust : An R Package for Time Series Clustering".
- Fermer le cadre : dans la suite de notre travail, on ne va pas les utiliser mais elles pourront être intégrées dans le framework qui suivra au chapitre suivant

2.4 Time series alignment and dynamic programming approach

In some applications, time series needs to be compared at different time t (i.e. energy data [Naj+12]) whereas in others, comparing time series on the same time t is essential (i.e. gene expression [DCN07]). When time series are asynchronous (i.e. varying delays or dynamic changes), they must be aligned before any analysis process. The asynchronous effects can be of various natures: time shifting (phase shift in signal processing), time compression or time dilatation. For example, in the case of voice recognition (Fig. 2.6), it is straightforward that a same sentence said by two different speakers will produce different time series: one speaker may speak faster than the other; one speaker may take more time on some vowels, etc.

Comment [SMA7]: formulation totale à revoir

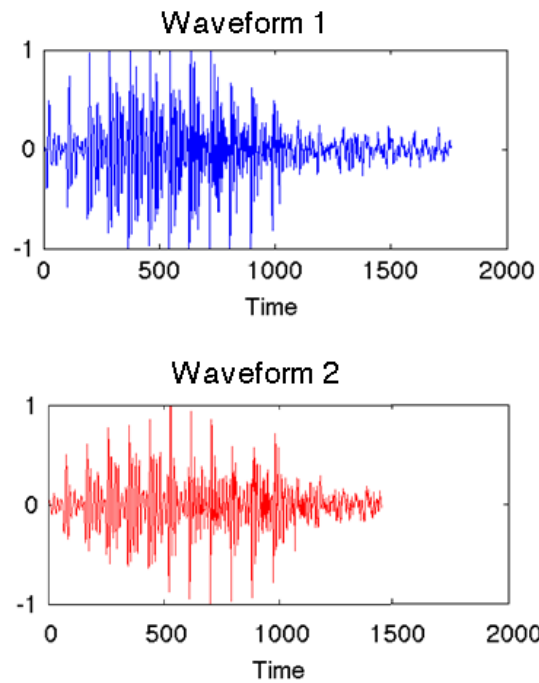


Figure 2.6: Example of a same sentence said by two different speakers. Time series are shifted, compressed and dilatated in the time.

To cope with delays and dynamic changes, dynamic programming approach has been introduced [BC94]. An alignment π of length $|\pi_{ij}| = m$ between two time series \mathbf{x}_i and \mathbf{x}_j of length Q is defined as the set of m ($Q \leq m \leq 2Q - 1$) couples of aligned elements of \mathbf{x}_i to m elements of \mathbf{x}_j :

$$\pi_{ij} = ((\pi_i(1), \pi_j(1)), (\pi_i(2), \pi_j(2)), \dots, (\pi_i(m), \pi_j(m))) \quad (2.7)$$

where the applications π_i and π_j defined from $\{1, \dots, m\}$ to $\{1, \dots, Q\}$ obey the following

Comment [MR8]: Modifier figure. enlever 'one' et mettre la même échelle temporelle

Comment [AD9]: Ahlan pas fan des notations

boundary monotonicity conditions:

$$1 = \pi_i(1) \leq \pi_i(2) \leq \dots \leq \pi_i(m) = Q \quad (2.8)$$

$$1 = \pi_j(1) \leq \pi_j(2) \leq \dots \leq \pi_j(m) = Q \quad (2.9)$$

$$\forall l \in \{1, \dots, m\},$$

$$\pi_i(l+1) \leq \pi_i(l) + 1 \quad (2.10)$$

$$\text{and} \quad \pi_j(l+1) \leq \pi_j(l) + 1 \quad (2.11)$$

$$\text{and} \quad (\pi_i(l+1) - \pi_i(l)) - (\pi_j(l+1) - \pi_j(l)) \geq 1. \quad (2.12)$$

In the following, we denote $\boldsymbol{\pi} = \pi_{ij}$ to simplify the notation. Intuitively, an alignment $\boldsymbol{\pi}$ defines a way to associate elements of two time series. Alignments can be described by paths in the $Q \times Q$ grid that crosses the elements of \mathbf{x}_i and \mathbf{x}_j (Fig. 2.7). We denote $\boldsymbol{\pi}$ a valid alignment and A_{ij} , the set of all possible alignments between \mathbf{x}_i and \mathbf{x}_j ($\boldsymbol{\pi} \in A$). To find the best alignment $\boldsymbol{\pi}^*$ between two time series \mathbf{x}_i and \mathbf{x}_j , the Dynamic Time Warping (DTW) algorithm has been proposed [KR04]; [SC].

DTW requires to choose a cost function φ to be optimised, such as a dissimilarity function (d_A, d_B, d_F , etc.). Standard DTW uses the Euclidean distance d_A (Eq. 2.1) as the cost function [BC94]. The warp path $\boldsymbol{\pi}$ is optimized for the chosen cost function φ :

$$\boldsymbol{\pi}^* = \operatorname{argmin}_{\boldsymbol{\pi} \in A_{ij}} \frac{1}{|\boldsymbol{\pi}|} \sum_{(t,t') \in \boldsymbol{\pi}} \varphi(x_{it}, x_{jt'}) \quad (2.13)$$

When the cost function φ is a similarity measure, the optimization involves maximization instead of minimization. When other constraints are applied on $\boldsymbol{\pi}$, Eq. (2.13) leads to other variants of DTW (Sakoe-Shiba [SC78], Itakura parallelogram [RJ93]). Finally, the warped signals $\mathbf{x}_{i,\boldsymbol{\pi}^*}$ and $\mathbf{x}_{j,\boldsymbol{\pi}^*}$ are defined as:

$$\mathbf{x}_{i,\boldsymbol{\pi}^*} = (x_{i\pi_i(1)}, \dots, x_{i\pi_i(m)}) \quad (2.14)$$

$$\mathbf{x}_{j,\boldsymbol{\pi}^*} = (x_{j\pi_j(1)}, \dots, x_{j\pi_j(m)}) \quad (2.15)$$

Once an optimal alignment $\boldsymbol{\pi}^*$ has been found, and whatever cost function φ have been chosen to find it, the metric presented in Section 2.3 (amplitude-based d_A , behavior-based d_B , frequential-based d_F) can be then computed on the warped signals $\mathbf{x}_{i,\boldsymbol{\pi}^*}$ and $\mathbf{x}_{j,\boldsymbol{\pi}^*}$. In the following, we suppose that the best alignment $\boldsymbol{\pi}^*$ is found. For simplification purpose, we refer to $\mathbf{x}_{i,\boldsymbol{\pi}^*}$ and $\mathbf{x}_{j,\boldsymbol{\pi}^*}$ as \mathbf{x}_i and \mathbf{x}_j .

voir si je mets la dtw avec fonction de coût cort dans les annexes

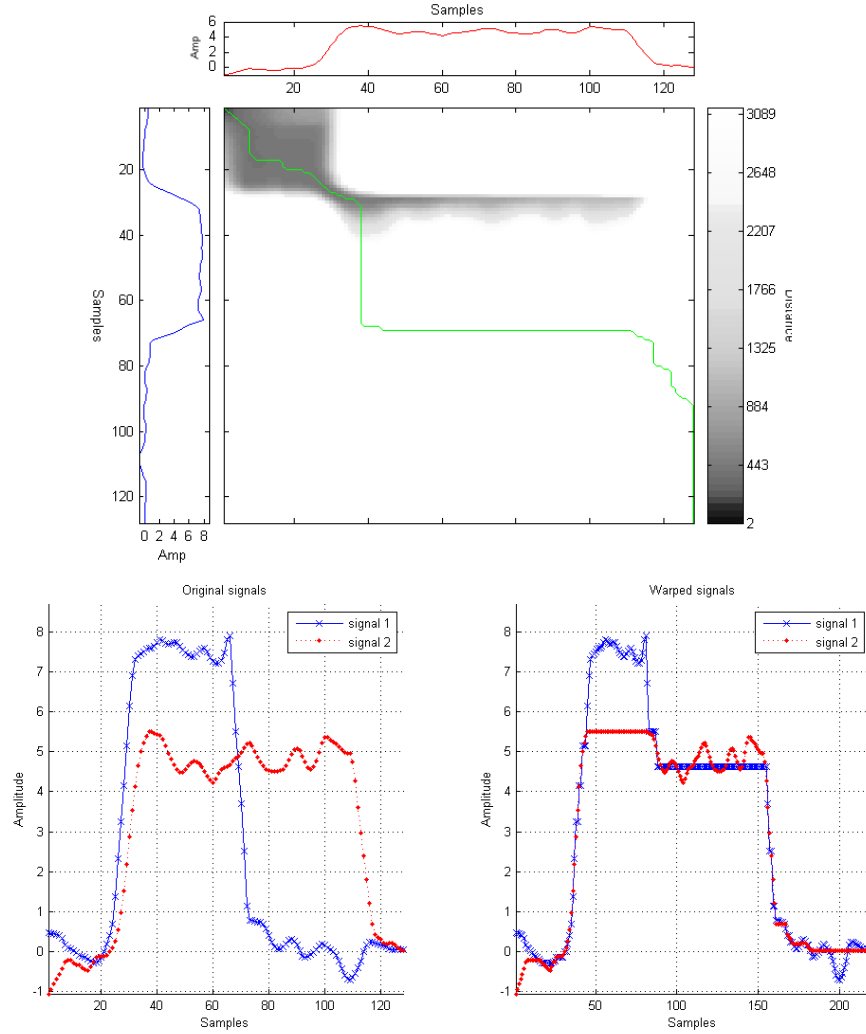


Figure 2.7: Example of DTW grid between 2 time series \mathbf{x}_i and \mathbf{x}_j (top) and the signals before and after warping (bottom). On the DTW grid, the two signals can be represented on the left and bottom of the grid. The optimal path π^* is represented in green line and shows how to associate elements of \mathbf{x}_i to element of \mathbf{x}_j . Background show in grey scale the value of the considered metric (amplitude-based distance d_A in classical DTW)

2.5 Combined metrics for time series

2.5.1 Combination functions

In most classification problems, it is not known a priori if time series of a same class exhibits same characteristics based on their amplitude, behavior or frequential components alone. In some cases, several components (amplitude, behavior and/or frequential) may be implied.

A first technic considers a classifier for each p metric and combines the decision of the p

ref

resulting classifiers. This method is referred as post-fusion, not considered in our work. Other propositions show the benefit of involving both behavior and amplitude components through a combination function. They combine the unimodal metrics together to obtain a single metric used after that in a classifier. This is called pre-fusion. The most basic combination functions that we could use combine two unimodal metrics through linear and geometric functions. For example, with d_A and d_B , we obtain:

$$D_{Lin}(\mathbf{x}_i, \mathbf{x}_j) = \alpha d_B(\mathbf{x}_i, \mathbf{x}_j) + (1 - \alpha) d_A(\mathbf{x}_i, \mathbf{x}_j) \quad (2.16)$$

$$D_{Geom}(\mathbf{x}_i, \mathbf{x}_j) = (d_B(\mathbf{x}_i, \mathbf{x}_j))^\alpha (d_A(\mathbf{x}_i, \mathbf{x}_j))^{1-\alpha} \quad (2.17)$$

where $\alpha \in [0; 1]$ defines the trade-off between the amplitude d_A and the behavior d_B components, and is thus application dependent. For example, in classification problems, we could learn it through a grid search procedure. Without being restrictive, these combinations can be extended to take into account more unimodal metrics.

More specific work on d_A and $cort$ propose to combine the two components through a sigmoid combination function [DCA11]:

$$D_{Sig}(\mathbf{x}_i, \mathbf{x}_j) = \frac{2d_A(\mathbf{x}_i, \mathbf{x}_j)}{1 + \exp(\alpha cort_r(\mathbf{x}_i, \mathbf{x}_j))} = \frac{2d_A(\mathbf{x}_i, \mathbf{x}_j)}{1 + \exp(\alpha(1 - 2d_B(\mathbf{x}_i, \mathbf{x}_j)))} \quad (2.18)$$

where α is a parameter that defines the compromise between behavior and amplitude components.

Fig.2.8 illustrates the value of the resulting combined metrics (D_{Lin} , D_{Geom} and D_{Sig}) in 2-dimensional space using contour plots for different values of the trade-off α . For small value of α ($\alpha = 0$), the three metrics only include d_A . For high value of α ($\alpha = 1$), D_{Lin} and D_{Geom} only include d_B . For $\alpha = 0.6$ and for small values of d_B , D_{Sig} mostly includes d_B while for large value of d_B , D_{Sig} mostly includes d_A . Note that these combinations are fixed and defined independently from the analysis task at hand. Moreover, in the case of D_{Sig} , the component $cort_r$ can be seen as a penalizing factor of d_A : it doesn't represent a real compromise between value and behavior components. Finally, one could extend D_{Lin} and D_{Geom} by adding metrics, but that would imply to add parameters. The grid search to find the best parameters would thus become time consuming.

2.5.2 Impact of normalization

références
normal-
ization

In most cases, it is recommended to scale each attribute to the range $[-1; +1]$ or $[0; 1]$. Many normalization methods have been proposed such as Min/Max normalization, Z-normalization or normalization of the log distribution. Let $X = \{\mathbf{x}_i, y_i\}_{i=1}^n$ be a training set, \mathbf{x}_i being a sample described by p features x_1, \dots, x_p . We define μ_j and σ_j as the mean and the standard deviation of a variable x_j , applying the Z-normalized variable x_j^{norm} is given by:

$$x_j^{norm} = \frac{x_j - \mu_j}{\sigma_j} \quad (2.19)$$

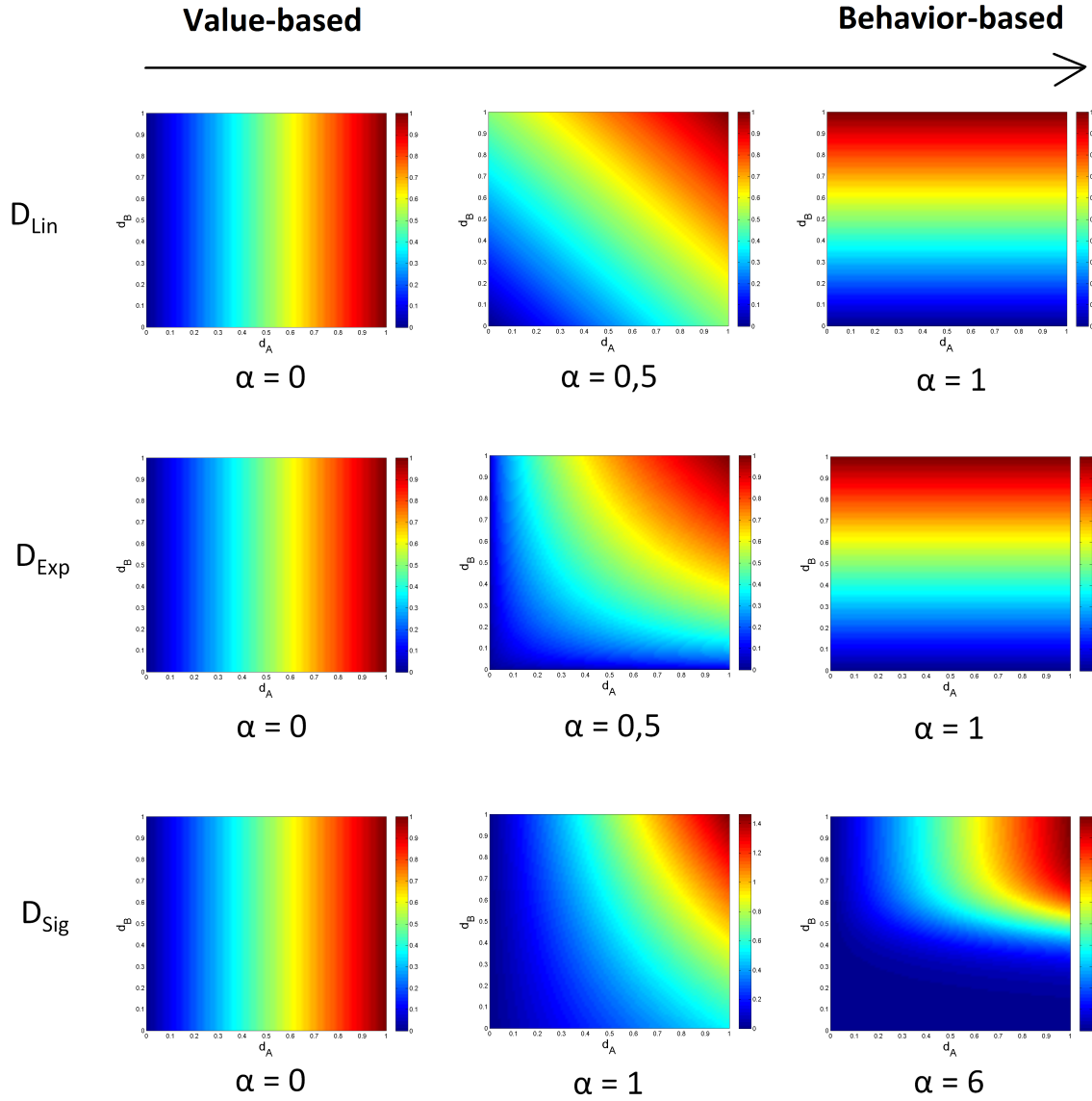


Figure 2.8: Contour plot of the resulting combined metrics: D_{Lin} (1st line), D_{Geom} (2nd line) and D_{Sig} (3rd line), for different values of α . For the three combined metrics, the first and second dimensions are respectively the amplitude-based metrics d_A and the behavior-based metric d_B .

Note that the underlying assumption supposes that the variable x_j is normally distributed: data evolves between $[-\infty; +\infty]$ and are coming from a Gaussian process. In some cases, the data are skewed such as monetary amounts or incomes. These data are sometimes log-normally distributed, *e.g.*, the log of the data is normally distributed (Fig. 2.9). The underlying idea is to take the log of the data (x_j^{ln}) to restore the symmetry, and then, to apply a Z-normalization

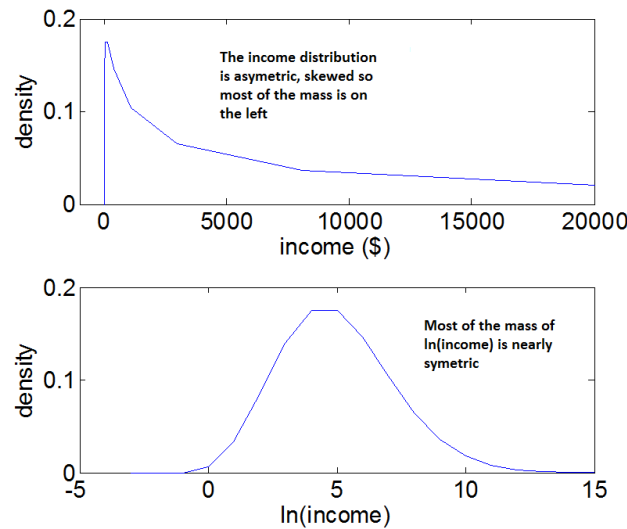


Figure 2.9: A nearly log-normal distribution, and its log transform ²

of this transformation:

$$x_j^{ln} = \ln(x_j); \quad (2.20)$$

$$x_j^{ln,norm} = \frac{x_j^{ln} - \mu_j^{ln}}{\sigma_j^{ln}} \quad (2.21)$$

$$x_j^{norm} = \exp(x_j^{ln,norm}) \quad (2.22)$$

where \ln denotes the Natural Logarithm function, μ_j^{ln} and σ_j^{ln} the mean and the standard deviation of a variable x_j^{ln} .

Compléter avec le mail de Sylvain. Partie peut être à mettre en annexe

Avec Cao on a discuté de la normalisation log hier et du coup j'ai regardé comment la justifier un peu mieux. Voici une intuition "avec les mains" montrant que l'on peut approximer la distribution de la distance euclidienne DE avec une loi log-normale:

Version simple : 1 feature, 1. On considère pour simplifier des données à 1 dimension (1 feature pour les vecteurs d'exemples). Soit x cette feature. On considère qu'elle suit une loi normale de variance (1/2) et de moyenne μ .

2. la différence entre deux exemples x et y , $d = (x-y) = (x + (-y))$, suit donc une loi normale de moyenne zéro et de variance 1. Explication: * inverse de lois normale $y =$ loi normale avec inverse de la moyenne et même variance ; * puis somme de deux lois normales (x) et $(-y)$

²source: <http://www.r-statistics.com/2013/05/log-transformations-for-skewed-and-wide-distributions-from-practical-data-science-with-r/>

= loi normale de moyenne égale à la somme des moyennes et variance égale à la somme des variances Voir: https://en.wikipedia.org/wiki/Sum_of_normally_distributed_random_variables

3. Le carré de la distance euclidienne entre deux exemples x et y est $DE^2 = (x - y)^2 = d^2$. Puisque d suit une loi normale de variance 1, DE^2 suit une loi du χ^2 non centrée d'ordre $k=1$. Voir: https://fr.wikipedia.org/wiki/Loi_du_%CF%87%C2%B2_non_centr%C3%A9e

4. Certaines personnes ont prouvé que la loi du χ^2 peut être approximée par une loi log-normale : http://www.rennes.supelec.fr/ren/perso/cmoy/papers/URSI_4pages.pdf (il faudrait regarder un peu mieux si il y a des hypothèses là-dedans)

5. puisque $\log(DE) = \log(\sqrt{DE^2}) = 1/2 * \log(DE^2)$, alors $\log(DE)$ suit une loi normale (car la loi d'une variable normale multipliée par une constante est normale).

=> DE peut donc être approximée par une loi log-normale

Extension à DE sur plusieurs features avec covariance diagonale égale à $(1/2)I$ 1. on considère qu'elles sont tirées d'une distribution jointe normale multivariée (une gaussienne en n dimension, avec covariance identité* $1/2$) 2. fonctionne encore 3. Fonctionne encore : "la loi du χ^2 non centrée est le carré de la norme d'un vecteur aléatoire de loi $N(\text{mean}, \text{covariance}=\text{Identité})$. Simplement cela augmente l'ordre k de la distribution, k devient le nombre de features.

Extension à (co-)variances non $(1/2)$ mon intuition est que si une feature a une variance beaucoup plus importante que les autres elle va "tuer" la variabilité des autres et donc ça va faire comme si il n'y avait qu'une seule feature. Du coup je pense que l'on retombe sur une loi du χ^2 , mais avec un ordre k plus réduit ($k < \text{nb de features}$)

<https://mail.google.com/mail/u/0/#search/sylvain+normalisation/14fb15fd972e75bc>

2.6 Metric learning

In this section, we first review metric learning concepts. Then, we focus on the framework proposed by Weinberger & Saul for Large Margin Nearest Neighbor (LMNN) classification [WS09a].

2.6.1 Review on metric learning work

In the case of static data, many work have demonstrated that k -NN classification performances depends highly on the considered metric and can be improved by learning an appropriate metric [She+02]; [Gol+04]; [CHL05]. Metric Learning can be defined as a process that aims to learn a distance from labeled examples by making closer samples that are expected to be

similar, and far away those expected to be dissimilar.

A faire, avec papier PRL et papier Aurélien Bellet. Refaire cette intro avec les notes de Sylvain.

2.6.2 Large Margin Nearest Neighbors (LMNN)

Let $\mathbf{X} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ be a set of N static vector samples, $\mathbf{x}_i \in \mathbb{R}^p$, p being the number of descriptive features and y_i the class labels. Weinberger & Saul proposed in [WS09a] an approach to learn a metric D for a large margin k -NN in the case of static data.

Large Margin Nearest Neighbor (LMNN) approach is based on two intuitions: according to the learnt metric D , first, each training sample \mathbf{x}_i should have the same label y_i as its k nearest neighbors; second, training samples with different labels should be widely separated. For this, the concept of **target** and **imposters** for each training sample \mathbf{x}_i is introduced. The training sample \mathbf{x}_i is referred as a **center point**. Given a metric D , target neighbors of \mathbf{x}_i , noted $j \rightsquigarrow i$, are the k closest \mathbf{x}_j of the same class ($y_j = y_i$), while imposters of \mathbf{x}_i , denoted, $l \nrightarrow i$, are the \mathbf{x}_l of different class ($y_l \neq y_i$) that invade the perimeter defined by the farthest targets of \mathbf{x}_i . Mathematically, for a sample \mathbf{x}_i , an imposter \mathbf{x}_l is defined by an inequality related to the targets \mathbf{x}_j : $\forall l, \exists j \in j \rightsquigarrow i /$

$$D(\mathbf{x}_i, \mathbf{x}_l) \leq D(\mathbf{x}_i, \mathbf{x}_j) + 1 \quad (2.23)$$

Geometrically, an imposter \mathbf{x}_l is a sample that invades the target neighborhood plus one unit margin as illustrated in Fig. 2.10. The target neighborhood is defined with respect to an initial metric. Without prior knowledge, L2-norm is often used. Metric Learning by LMNN aims at minimizing the number of impostors invading the target neighborhood. By adding a margin of safety of one, the model is ensured to be robust to small amounts of noise in the training sample (large margin). The learned metric D pulls the targets \mathbf{x}_j and pushes the imposters \mathbf{x}_l as shown in Fig. 2.10.

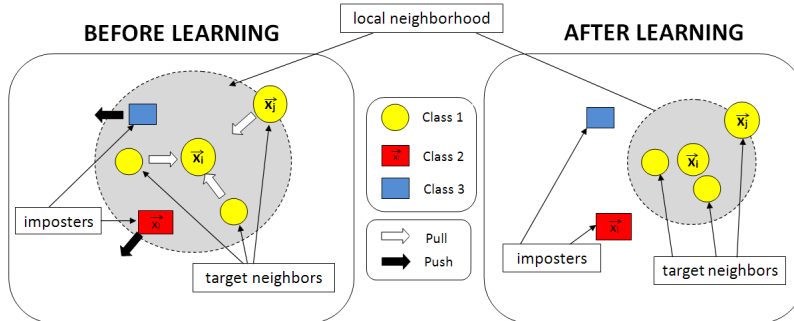


Figure 2.10: Pushed and pulled samples in the $k = 3$ target neighborhood of \mathbf{x}_i before (left) and after (right) learning. The pushed (vs. pulled) samples are indicated by a white (vs. black) arrows (Weinberger & Saul [WS09a]).

LMNN approach learns a Mahalanobis distance D for a robust k -NN. We recall that the k -NN decision rule will correctly classify a sample if the majority of its k nearest neighbors share the same label (Section 1.2.1). The objective of LMNN is to increase the number of samples with this property by learning a linear transformation \mathbf{L} of the input space ($\mathbf{x}_i = \mathbf{L}\mathbf{x}_i$) before applying the k -NN classification:

$$D(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{L}(\mathbf{x}_i - \mathbf{x}_j)\|_2^2 \quad (2.24)$$

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

$$\mathbf{M} = \mathbf{L}'\mathbf{L} \quad (2.25)$$

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) [WS09a]. Using the matrix \mathbf{M} , squared distances can be expressed as:

$$D_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)\mathbf{M}(\mathbf{x}_i - \mathbf{x}_j) \quad (2.26)$$

The computation of the learned metric $D_{\mathbf{M}}$ 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:

$$D_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_j) = D^2(\mathbf{L}\mathbf{x}_i, \mathbf{L}\mathbf{x}_j) \quad (2.27)$$

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, it can be formalized as an optimization problem involving two competing 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:

$$\begin{aligned} \underset{\mathbf{M}, \xi}{\operatorname{argmin}} \quad & \underbrace{\sum_{i, j \rightsquigarrow i} D_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_j)}_{\text{pull}} + C \underbrace{\sum_{i, j \rightsquigarrow i, l \nrightarrow i} \frac{1 + y_{il}}{2} \cdot \xi_{ijl}}_{\text{push}} \\ \text{s.t. } \quad & \forall j \rightsquigarrow i, l \nrightarrow i, \\ & D_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_l) - D_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \xi_{ijl} \\ & \xi_{ijl} \geq 0 \\ & \mathbf{M} \succeq 0 \end{aligned} \quad (2.28)$$

where C is a trade-off between the push and pull term and $y_{il} = -1$ if $y_i = y_l$ (same class) and $+1$ otherwise (different classes). Generally, the parameter C is tuned via cross validation and grid search. Similarly to Support Vector Machine (SVM) approach, slack variables ξ_{ijl} are introduced to relax the optimization problem.

2.6.3 Parallels between LMNN and SVM

Many connections can be made between LMNN and SVM: both are convex optimization problem based on a regularized and a loss term. In particular, Do & al. investigate this relationship and have shown that SVM can be formulated as a metric learning problem [Do+12]. The Mahalanobis distance \mathbf{M} learned by LMNN can be expressed as a quadratic mapping ϕ . For a center point \mathbf{x}_i , for any sample \mathbf{x} , we have [Do+12]:

$$D_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}) = D^2(\mathbf{L}\mathbf{x}_i, \mathbf{L}\mathbf{x}) \quad (2.29)$$

$$D_{\mathbf{M}}^2(\mathbf{x}_i, \mathbf{x}) = \mathbf{w}_i^T \phi(\mathbf{x}) + b_i \quad (2.30)$$

where \mathbf{w}_i and b_i are the coefficient of the hyperplane H_i in the quadratic space ϕ .

Do & al. show that LMNN can be seen as a set of local SVM classifiers in the quadratic space induced by ϕ . For each center point \mathbf{x}_i , LMNN tries in its objective function to have its target neighbors \mathbf{x}_j to have small value $\mathbf{w}_i^T \phi(\mathbf{x}_j) + b_i$, i.e. be at the small distance from the hyperplane H_i . Minimizing the target neighbor distances from the hyperplane H_i makes the distance between support vectors and H_i small. Fig. 2.11 gives the equivalent point of view from the original space (Fig. 2.11(a)) into the quadratic space (Fig. 2.11(b)). The circle \mathbf{C}_i with the center $\mathbf{L}\mathbf{x}_i$ in Fig. 2.11(a) corresponds to the hyperplane H_i in Fig. 2.11(b).

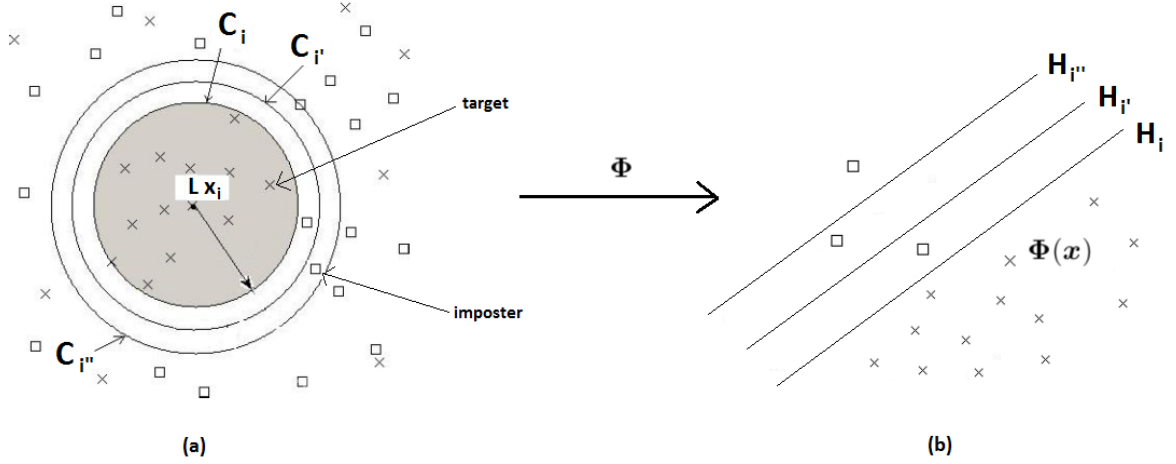


Figure 2.11: (a) Standard LMNN model view (b) LMNN model view under an SVM-like interpretation [Do+12]

Geometrically, SVM margin is defined globally with respect to a hyperplane, while LMNN margin is defined locally with respect to a center point \mathbf{x}_i . Fig. 2.12(a) illustrates the different local linear models in the quadratic space. The optimization process of LMNN combines the different local SVM hyperplane by bringing each point $\phi(\mathbf{x}_i)$ around a consensus hyperplane H .

From these connections, some authors extends the LMNN approach to work in non-linear feature spaces by using the “kernel trick”. Finally, note that LMNN differs from SVM in which LMNN requires no modification for multiclass problems.

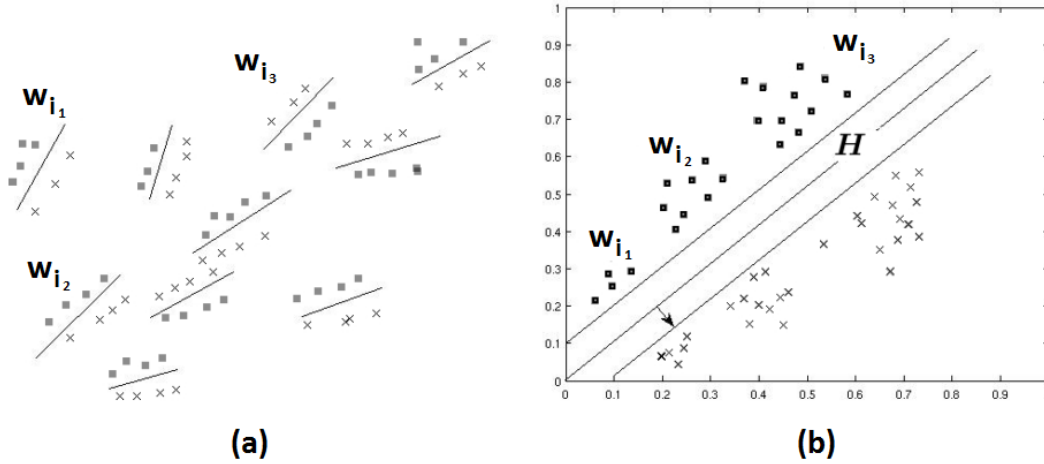


Figure 2.12: (a) LMNN in a local SVM-like view (b) LMNN in an SVM metric learning view [Do+12]

2.7 Conclusion of the chapter

To cope with modalities inherent to time series (amplitude, behavior, frequency, etc.), we review in this chapter several unimodal metrics for time series, in particular, the Euclidean distance d_A , the Temporal correlation d_B or the Fourier-based distance d_F . In practice, real time series may be subjected to delays and need to be re-aligned before any analysis task. For that, the Dynamic Time Warping (DTW) algorithm is used in practice. However, these metrics (d_A, d_B, d_F) only include one modality. In general, several modalities may be implied and authors proposed to combine temporal metrics together. They mainly combine the Euclidean distance d_A and the Temporal correlation d_B .

As k -NN performances is impacted by the choice of the metric, other work propose in the case of static data to learn the metric in order to optimize the k -NN classification. In the following, we extend this framework to learn a combined metric for a large margin k -NN classifier of time series.

Multi-modal and Multi-scale Time series Metric Learning (M^2TML)

Contents

3.1	Motivations	50
3.2	Multi-modal and multi-scale dissimilarity space	52
3.2.1	Pairwise embedding	52
3.2.2	Multi-scale description for time series	53
3.2.3	Interpretation in the dissimilarity space	53
3.3	M^2TML general problem	55
3.3.1	General formalization for M^2TML	55
3.3.2	Push and pull set definition	56
3.3.3	Interpretation in the dissimilarity space	58
3.4	Linear formalization for M^2TML	59
3.5	Quadratic formalization for M^2TML	61
3.5.1	Primal and dual formalization	61
3.5.2	Non-linear combined metric	63
3.5.3	Link between SVM and the quadratic formalization	64
3.6	SVM-based formalization for M^2TML	66
3.6.1	Support Vector Machine (SVM) resolution	66
3.6.2	Linearly separable Pull and Push sets	67
3.6.3	Non-linearly separable Pull and Push sets	69
3.7	SVM-based solution and algorithm for M^2TML	70
3.8	Conclusion	71

In this chapter, we first motivate the problem of learning a multi-modal and multi-scale temporal metric for nearest neighbors classification. Secondly, we introduce the concept of dissimilarity space. Thirdly, we formalize the general problem of learning a combined metric. Then, we propose three different formalizations (Linear, Quadratic and SVM-based), each involving a different regularization term. We give an interpretation of the solution and study the properties of the obtained metric. Finally, we give the algorithm.

3.1 Motivations

Intro
modi-
fi   pour
tenir
compte
des
chapitres
pr  c  -
dents

Formulation
du pb   
valider

This work focuses on designing 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 can be compared based on one or several characteristics, called modalities (amplitude, behavior, frequency), they may be subjected to delays, and 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 , frequential 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).

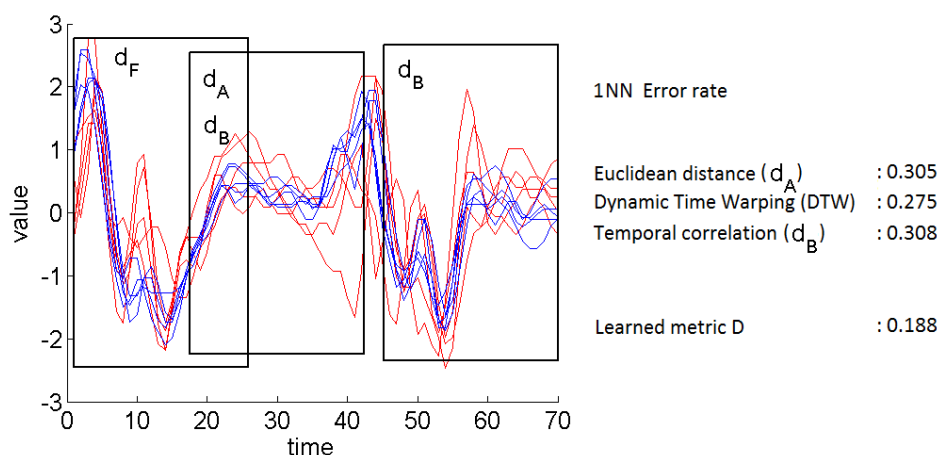


Figure 3.1: SonyAIBO dataset and error rate using a k NN ($k = 1$) with standard metrics (Euclidean distance, Dynamic Time Warping, temporal correlation) and a learned combined metric D . The figure shows the 4 major metrics involve in the combined metric D and their respective temporal scale (black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

Correction
d'anglais
et
pr  cision
des
phrases

separate

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. 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

last decade 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 at capturing non linear structure in the data, *e.g.*, Kernel Principal Component Analysis (KPCA) and Support Vector Metric Learning (SVML). In both cases, the metric is directly learned in the original space (*i.e.*, space described by the features of the samples). In KPCA, the aim is to project the data into a non linear feature space and learn the metric in that projected space [ZY10]; [Cha+10]. In SVML, the Mahalanobis distance is learned jointly with the learning of the SVM model in order to minimize the validation error [XWC12]. In general, the optimization problems in non linear approaches is more expensive to solve than in linear approaches, and the methods tends to favor overfitting as the constraints are generally easier to satisfy in a nonlinear kernel space. A more detailed review on metric learning is done in [BHS13].

Comment [SMA10]: Ajouter la différence

Contrary to static data, metric learning for structured data (*e.g.* sequence, time series, trees, graphs, strings) is less frequent. 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 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), involving all time series elements at once, which restricts its potential to capture local characteristics. We believe that perspectives for metric learning, in the case of time series, should include multi-modal and multi-scale aspects.

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 pairwise dissimilarity space where a linear function combining several modalities at different temporal scales can be learned, driven by a large margin optimization process inspired from the 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 proposed metrics 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.

Comment [MR11]: ajouter la notion de pairwise pour différencier

Comment [SMA12]: modifier pour dire que le processus est général et pas que SVM

In this chapter, we first present the concept of pairwise dissimilarity space. Then, we formalize the general problem of learning a combined metric for a robust k -NN as the learning a function in the dissimilarity space. From the general formalization, we propose three formalizations (Linear, Quadratic and SVM-based), give an interpretation of the solutions and study the properties of the learned metrics. Finally, we give the algorithm. Note that these formalizations don't concern only time series and can be applied to learn a combined metric on any type of data.

Mettre la section LMNN ici?

3.2 Multi-modal and multi-scale dissimilarity space

titre à
changer?

In this section, we first present the concept of dissimilarity space for multi-modal metrics. Then, in the case of time series, we enrich this representation with a multi-scale description.

ref: dis-
similarity
space
[PPD02];
[DP12]

3.2.1 Pairwise embedding

Comment
[SMA13]: la
bel main-
tenant?

Let $\{\mathbf{x}_i, y_i\}_{i=1}^n$ be a set of n time series $\mathbf{x}_i = [x_{i1}, \dots, x_{iq}] \in \mathbb{R}^q$ labeled y_i . Let d_1, \dots, d_p be p given metrics that allow to compare samples \mathbf{x}_i . As discussed in Chapter 2, three naturally modalities are involved for time series comparison: 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 basic temporal metrics for a robust k -NN classifier.

The computation of a metric d , and 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**. We note φ an embedding function that maps each pair of time series $(\mathbf{x}_i, \mathbf{x}_j)$ to a vector \mathbf{x}_{ij} in a dissimilarity space $\mathcal{E} = \mathbb{R}^p$ whose dimensions are d_1, \dots, d_p (Fig. 3.2):

$$\begin{aligned} \varphi : \mathbb{R}^q \times \mathbb{R}^q &\rightarrow \mathcal{E} = \mathbb{R}^p \\ (\mathbf{x}_i, \mathbf{x}_j) &\rightarrow \mathbf{x}_{ij} = [d_1(\mathbf{x}_i, \mathbf{x}_j), \dots, d_p(\mathbf{x}_i, \mathbf{x}_j)]^T \end{aligned} \quad (3.1)$$

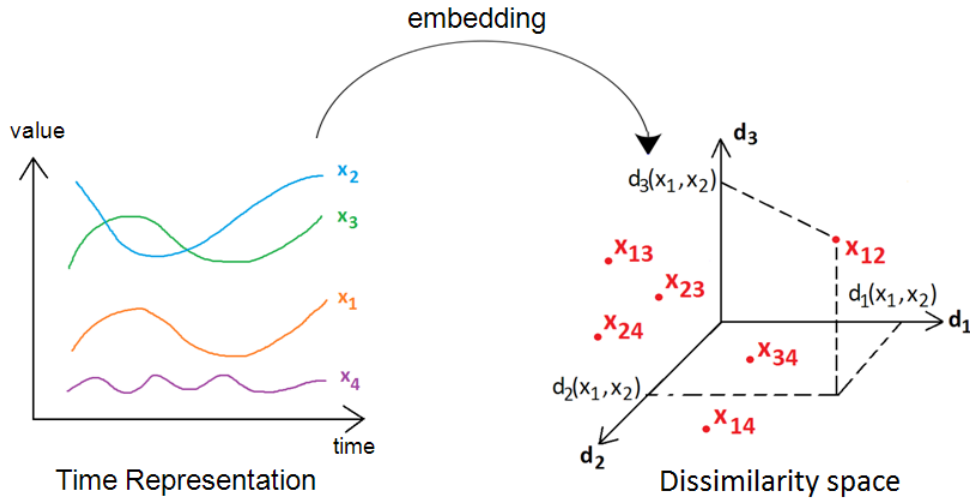


Figure 3.2: Example of embedding of time series \mathbf{x}_i from the temporal space (left) into the dissimilarity space (right) for $p = 3$ basic metrics.

Comment
[SMA14]: Propo-
sition
pour
mettre en
valeur que
c'est une
fonction

A metric D that combines the p metrics d_1, \dots, d_p can be seen as a function of the dissimilarity space: Propo-

$$\begin{aligned} D : \mathbb{R}^p &\rightarrow \mathbb{R} \\ \mathbf{x}_{ij} &\rightarrow D(\mathbf{x}_{ij}) = f(d_1(\mathbf{x}_i, \mathbf{x}_j), \dots, d_p(\mathbf{x}_i, \mathbf{x}_j)) \end{aligned} \quad (3.2)$$

In that space, the norm of a pairwise vector $\|\mathbf{x}_{ij}\|$ refers to the proximity between the time series \mathbf{x}_i and \mathbf{x}_j . In particular, if $\|\mathbf{x}_{ij}\| = 0$ then \mathbf{x}_j is identical to \mathbf{x}_i according to all metrics d_h . Note that in this space, the proximity between pairs can't be interpreted.

3.2.2 Multi-scale description for time series

As illustrated in Fig. 3.1, 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_{it} and thus, restricts the distance measures to capture local temporal differences. In our 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 or the dichotomy . We detail here the latter one.

Comment
[SMA15]: Ba
culer en
chap 2?

ref

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. In our 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 α allows to overlap the two intervals I_L and I_R , covering discriminating subsequences in the central region of I (around $\frac{b+a}{2}$): $I = [a; b]; I_L = [a; a + \alpha(b - a)]; I_R = [b - \alpha(b - a); b]$. For $\alpha = 0.6$, 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. 3.3.

A multi-scale dissimilarity description between two samples 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 \mathbf{x}_i and \mathbf{x}_j , the comparison between \mathbf{x}_i and \mathbf{x}_j is done on the same interval I_s . For a multi-scale amplitude-based comparison based on binary segmentation, the set of involved amplitude-based measures $d_A^{I_s}$ is $\{d_A^{I_1}, d_A^{I_2}, \dots\}$ where $d_A^{I_s}$ is defined as:

$$d_A^{I_s}(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{t \in I_s} (x_{it} - x_{jt})^2} \quad (3.3)$$

The local behaviors- and frequential- based measures $d_B^{I_s}$ and $d_F^{I_s}$ are obtained similarly.

In the following, for simplification purpose, we consider d_1, \dots, d_p as the set of multi-modal and multi-scale metrics.

3.2.3 Interpretation in the dissimilarity space

In this section, we give more detailed interpretations in the dissimilarity space. We recall that

Sous
section
ajoutée,
ok pour
SMA et
MR

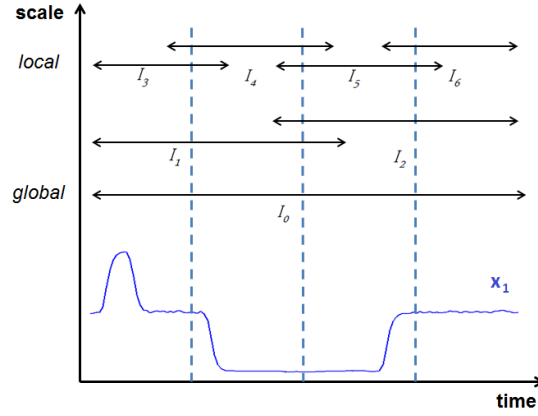


Figure 3.3: Multi-scale decomposition

the norm of a pairwise vector is given by:

$$\|\mathbf{x}_{ij}\| = \sum_{h=1}^p d_h(\mathbf{x}_i, \mathbf{x}_j) \quad (3.4)$$

In the following, we denote the norm as an initial distance in the dissimilarity space and call it D_0 . The norm of a pairwise 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 . Two pairwise vectors \mathbf{x}_{ij} and \mathbf{x}_{kl} that are on a same line that passes through the origin $\mathbf{x}_{ii} = \mathbf{0}$ represent differences in the the same proportions between their respective modalities (Fig. ??).

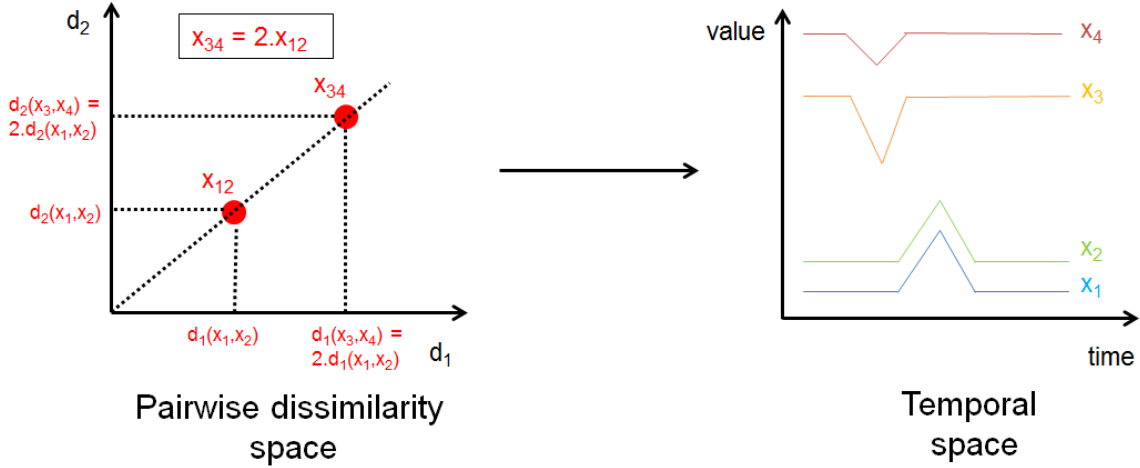


Figure 3.4: Example of two pairwise vectors \mathbf{x}_{12} and \mathbf{x}_{34} close in the pairwise space. However, the time series \mathbf{x}_1 and \mathbf{x}_3 are not similar in the temporal space.

The euclidean distance $\sqrt{\sum_{h=1}^p (d_h(\mathbf{x}_i, \mathbf{x}_j) - d_h(\mathbf{x}_k, \mathbf{x}_l))^2}$ between two pairwise vectors \mathbf{x}_{ij} and \mathbf{x}_{kl} represent the similarity between the differences among the same modality, in the

same proportion. Note that if the euclidean distance is close to 0 (\mathbf{x}_{ij} and \mathbf{x}_{kl} are close in the dissimilarity space) it doesn't mean that the time series \mathbf{x}_i , \mathbf{x}_j , \mathbf{x}_k and \mathbf{x}_l are similar. Fig 3.5 shows an example of two pairwise vectors \mathbf{x}_{ij} and \mathbf{x}_{kl} close together in the pairwise 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 , *i.e.*, the distance D_0 between \mathbf{x}_i and \mathbf{x}_j is roughly the same than the distance D_0 between \mathbf{x}_k and \mathbf{x}_l .

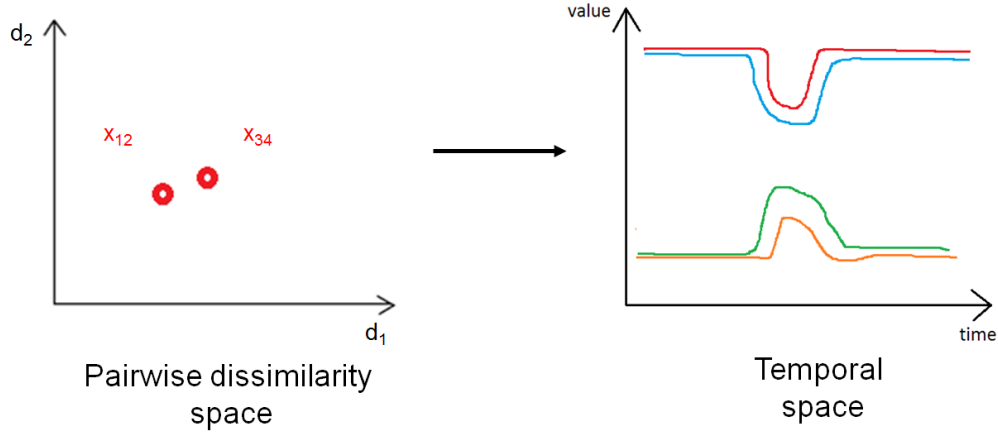


Figure 3.5: Example of two pairwise vectors \mathbf{x}_{12} and \mathbf{x}_{34} close in the pairwise space. However, the time series \mathbf{x}_1 and \mathbf{x}_3 are not similar in the temporal space.

3.3 M²TML general problem

In this section, we propose to define the Multimodal and Multiscale Time series Metric Learning (M²TML) problem in the initial space as a general problem of learning a function in the dissimilarity space. First, we give the intuition and formalize the general optimization problem. Secondly, we propose different strategies to define the neighborhood.

3.3.1 General formalization for M²TML

Our objective is to learn a dissimilarity $D = f(d_1, \dots, d_p)$ in \mathcal{E} , the embedding space, that combines the p dissimilarities d_1, \dots, d_p for a robust k -NN classifier. The function f can be linear or non-linear and must satisfy at least the properties of a dissimilarity, *i.e.*, positivity ($D(\mathbf{x}_{ij}) \geq 0$), reflexivity ($D(\mathbf{x}_{ii}) = 0 \forall i$) and symmetry ($D(\mathbf{x}_{ij}) = D(\mathbf{x}_{ji}) \forall i, j$) (Section 2.2). To simplify the discussion in the following, we refer to dissimilarity as metrics. The proposition is based on two standard intuitions in metric learning, *i.e.*, for each time series \mathbf{x}_i , the metric D should bring closer the time series \mathbf{x}_j of the same class ($y_j = y_i$) while pushing the time series \mathbf{x}_l of different classes ($y_l \neq y_i$). These two sets are called respectively $Pull_i$ and $Push_i$. Formally, the metric learning problem can be written as an optimization problem that involves both a regularization term $R(Pull)$ and a loss term $L(Push)$ under constraints that controls the push term:

figure
enlevée
ici car
trop tôt

Comment
[SMA16]: me
tre
 $R(D, Pull)$?

$$\begin{aligned} & \underset{D, \xi}{\operatorname{argmin}} [R(Pull) + L(Push)] \\ & \text{s.t. } constraints \end{aligned} \quad (3.5)$$

The problem of learning a combined metric D for large margin k -NN classification can be written as the following optimization problem:

$$\begin{aligned} & \underset{D, \xi}{\operatorname{argmin}} \left\{ \underbrace{\sum_{j \in Pull_i} D(\mathbf{x}_{ij})}_{pull} + C \underbrace{\sum_{\substack{j \in Pull_i \\ l \in Push_i}} \frac{1 - y_{il}}{2} \xi_{ijl}}_{push} \right\} \\ & \text{s.t. } \forall i = 1, \dots, n, \forall j \in Pull_i, l \in Push_i, \\ & \quad D(\mathbf{x}_{il}) - D(\mathbf{x}_{ij}) \geq 1 - \xi_{ijl} \\ & \quad \xi_{ijl} \geq 0 \end{aligned} \quad (3.6)$$

Comment [SMA17]: Enlever les y_{il} , ils ne servent plus

where $y_{il} = -1$ if $y_i \neq y_l$ and $+1$ otherwise, ξ_{ijl} are the slack variables and C , the trade-off between the pull and push costs. In the next section, we detailed different strategies to define the $Pull_i$ and $Push_i$ sets.

3.3.2 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 3.6:

Ajouter une remarque sur 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:

$$\forall i \in 1, \dots, n, \quad Pull_i = \{\mathbf{x}_{ij} / y_j = y_i, \|\mathbf{x}_{ij}\|_2 \text{ is among the } k\text{-lowest norms}\} \quad (3.7)$$

$$Push_i = \{\mathbf{x}_{il} / y_l \neq y_i, \|\mathbf{x}_{il}\|_2 \leq \max_{\mathbf{x}_{ij} \in Pull_i} \|\mathbf{x}_{ij}\|_2\} \quad (3.8)$$

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

$$\forall i \in 1, \dots, n, \quad Pull_i = \{\mathbf{x}_{ij} / y_j = y_i, \|\mathbf{x}_{ij}\|_2 \text{ is among the } k\text{-lowest norms}\} \quad (3.9)$$

$$Push_i = \{\mathbf{x}_{il} / y_l \neq y_i\} \quad (3.10)$$

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: $m = \alpha.k$ with $\alpha \geq 1$. Other propositions for m are possible:

$$\forall i \in 1, \dots, n, \quad Pull_i = \{\mathbf{x}_{ij} / y_j = y_i, \|\mathbf{x}_{ij}\|_2 \text{ is among the } m\text{-lowest norms}\} \quad (3.11)$$

$$Push_i = \{\mathbf{x}_{il} / \text{s.t. } y_l \neq y_i, \|\mathbf{x}_{il}\|_2 \text{ is among the } m\text{-lowest norms}\} \quad (3.12)$$

In the following, we denote $m\text{-NN}^+ = \bigcup_i Pull_i$ and $m\text{-NN}^- = \bigcup_i Push_i$

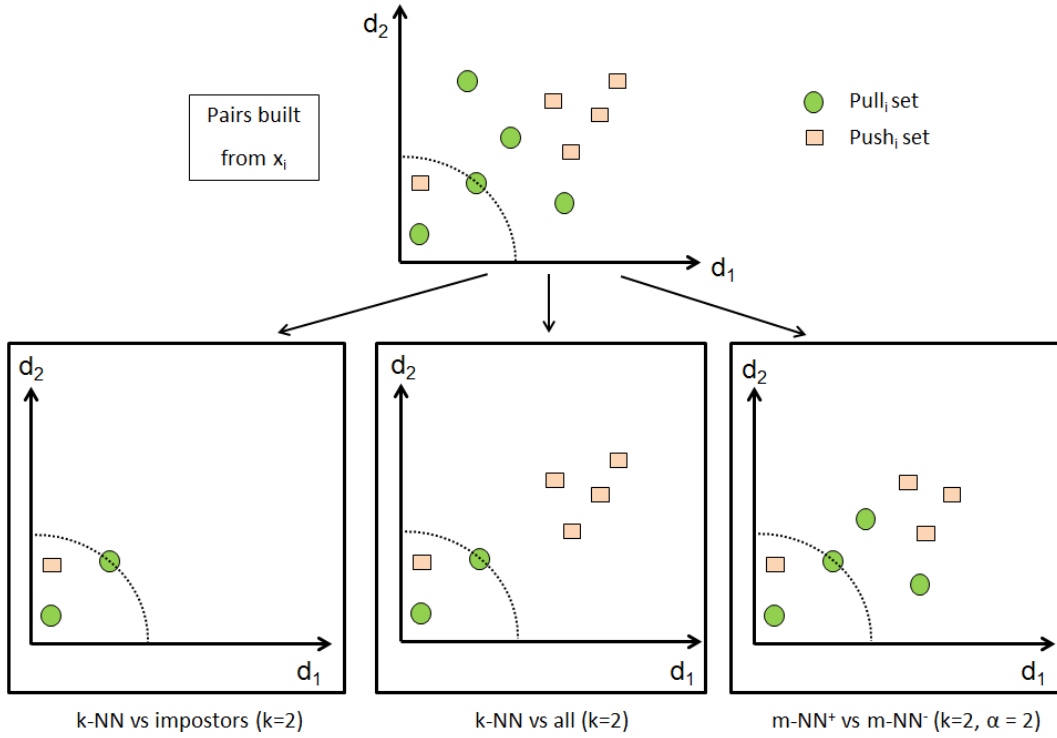


Figure 3.6: Example of different strategies to build $Pull_i$ and $Push_i$ sets for a $k = 2$ neighborhood.

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. 3.7 & 3.8) and may be unbalanced. The sets are defined and fixed during the optimization process according to an initial metric (Euclidean distance). In M²TML the sets $Pull_i$ and $Push_i$ are defined according the **k -NN vs impostors** strategy (Eqs. 3.11 & 3.12) and are balanced. The sets are defined and fixed during the optimization process according to an initial metric (Euclidean distance), but the m -neighborhood is larger than the k -neighborhood. By considering a neighborhood larger than the k -neighborhood (**k -NN vs impostors** strategy), we believe that the generalization properties of the learnt metric D would be improved.

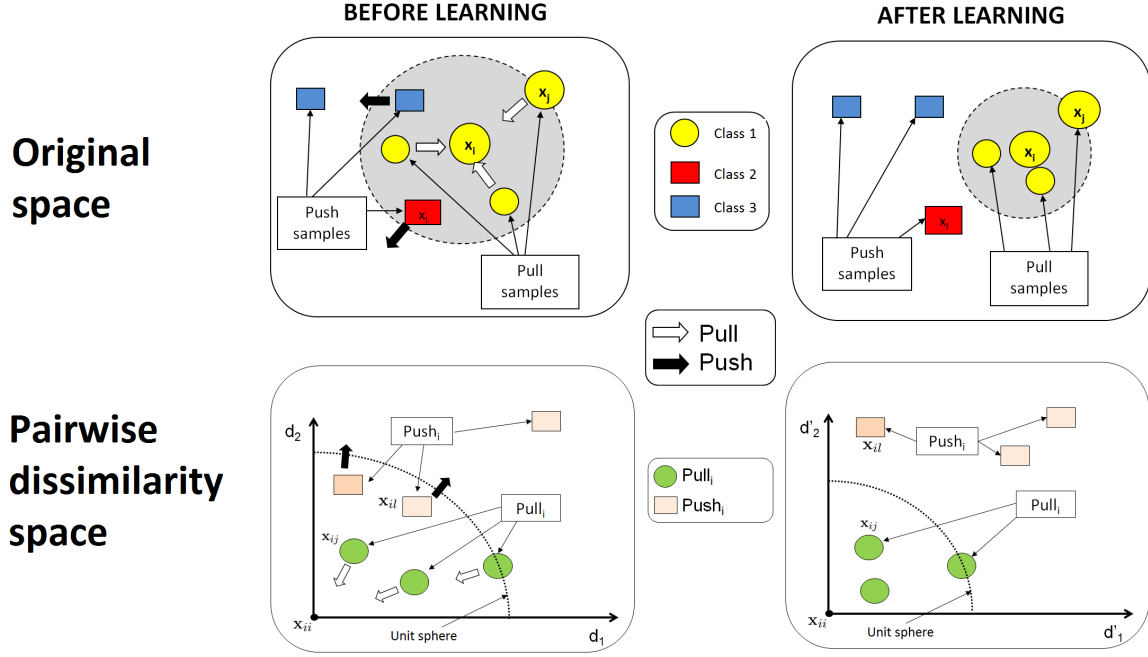


Figure 3.7: Metric learning problem from the original space (top) to the dissimilarity space (bottom) for a $k = 3$ neighborhood of \mathbf{x}_i . Before learning (left), push samples \mathbf{x}_l invade the targets perimeter \mathbf{x}_j . In the pairwise space, this is equivalent to have push pairwise vectors \mathbf{x}_{il} with an initial distance D_0 lower than the distance of pull pairwise vectors \mathbf{x}_{ij} . The aim of metric learning is to learn a metric D to push \mathbf{x}_{il} (black arrow) and pull \mathbf{x}_{ij} from the origin (white arrow).

3.3.3 Interpretation in the dissimilarity space

In this section, we give more detailed interpretations of the M²TML problem in the dissimilarity space. Our objective is to learn a metric D as a linear or non-linear combination of the p unimodal metrics d_1, \dots, d_p . The metric D can be seen as a function of the dissimilarity space that should:

- **pull** to the origin $\mathbf{x}_{ii} = \mathbf{0}$ the pairs \mathbf{x}_{ij} of $Pull_i$
- **push** away from the origin all the pairs \mathbf{x}_{il} of $Push_i$

Fig. 3.7 illustrates the idea in the original space and in the dissimilarity space: first, we build the sets $Pull_i$ and $Push_i$ according to an initial metric D_0 ; secondly, we optimize the metric D so that the pairs $Pull_i$ are pulled to the origin and the pairs $Push_i$ are pushed away from the origin.

Note that by considering a larger neighborhood, we ensure that pairs $Push_i$ doesn't invade the perimeter defined by pairs $Pull_i$. Similarly to the interpretation of slack variables in SVM,

Ajout
ok pour
SMA et
MR

if a push pair invade the perimeter defined by pairs $Pull_i$, then in Eq. 3.6, it will violate the constraints and the slack variables ξ_{ijl} will be penalized in the objective function:

- 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.
- 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.
- If $D(\mathbf{x}_{il}) > D(\mathbf{x}_{ij}) + 1$, $\xi_{ijl} = 0$ and the slack variable has no effect in the objective function.

In the following, we propose different regularizers for the pull term $R(Pull)$. First, we use a linear regularization. Secondly, we use a quadratic regularization that enables to extend the method 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 $Pull_i$ and $Push_i$ sets, and then, induce a combined metric D . Finally, we sum up the retained solution (SVM-based solution) and give the main steps of the algorithm.

3.4 Linear formalization for M²TML

In this section, we define the problem of learning a combined metric D as a linear combination in the dissimilarity space using a linear regularizer. First, we give the optimization problem. Then, we discuss on the properties of the learnt metric D .

Let $\mathbf{X} = \{\mathbf{x}_{ij}, y_{ij}\}_{i,j=1}^n$ be a set of pairwise vectors $\mathbf{x}_{ij} = [d_1(\mathbf{x}_i, \mathbf{x}_j), \dots, d_p(\mathbf{x}_i, \mathbf{x}_j)]^T$ described by p metrics d_1, \dots, d_p and labeled $y_{ij} = -1$ if $y_i \neq y_j$ and $+1$ otherwise. We consider a linear combination of the p metrics:

$$D(\mathbf{x}_{ij}) = \mathbf{w}^T \mathbf{x}_{ij} = \sum_{h=1}^p w_h \cdot d_h(\mathbf{x}_i, \mathbf{x}_j) \quad (3.13)$$

where $\mathbf{w} = [w_1, \dots, w_p]^T$ is the vector of weights w_h . From Eq. 3.6, learning a linear combined metric D can be formalized as follow:

$$\begin{aligned}
& \underset{\mathbf{w}, \xi}{\operatorname{argmin}} \left\{ \underbrace{\sum_{j \in \overset{i}{Pull}_i} \mathbf{w}^T \mathbf{x}_{ij}}_{pull} + C \underbrace{\sum_{\substack{j \in \overset{i}{Pull}_i \\ l \in \overset{i}{Push}_i}} \frac{1 - y_{il}}{2} \xi_{ijl}}_{push} \right\} \\
& \text{s.t. } \forall i = 1, \dots, n, \forall j \in \overset{i}{Pull}_i, l \in \overset{i}{Push}_i, \\
& \quad \mathbf{w}^T (\mathbf{x}_{il} - \mathbf{x}_{ij}) \geq 1 - \xi_{ijl} \\
& \quad \xi_{ijl} \geq 0 \\
& \quad \forall h = 1, \dots, p, \quad w_h \geq 0
\end{aligned} \tag{3.14}$$

where ξ_{ijl} are the slack variables, C the trade-off between the pull and push costs, and $\overset{i}{Pull}_i$ and $\overset{i}{Push}_i$ are defined in Eqs. 3.11 & 3.12. Similarly to SVM, note that the slack variables ξ_{ijl} can be interpreted. In particular, the pairs \mathbf{x}_{ij} and \mathbf{x}_{il} that violate the constraints ($\mathbf{w}^T \mathbf{x}_{il} < \mathbf{w}^T \mathbf{x}_{ij}$) will be penalized in the objective function. They corresponds to push pairs \mathbf{x}_{il} that invades the neighborhood of the pull pairs \mathbf{x}_{ij} .

The problem is very similar to a C -SVM classification problem. When C is infinite, we have a "strict" problem: the solver will try to find a direction \mathbf{w} in the dissimilarity space \mathcal{E} for which all $\xi_{ijl} = 0$, that means that only pull samples should be in the close neighborhood of each \mathbf{x}_i . Let denote \mathbf{x}_{ij}^* and \mathbf{x}_{il}^* , the vectors for which $\xi_{ijl} = 0$. In that case, if a solution is found, the margin $\min_{i,j,l} (\|\mathbf{x}_{il}^* - \mathbf{x}_{ij}^*\|_2)$ can be derived from the tightest constraint, for which equality holds:

$$\begin{aligned}
& \mathbf{w}^T (\mathbf{x}_{il}^* - \mathbf{x}_{ij}^*) = 1 \\
& \mathbf{w}^T (\mathbf{x}_{il}^* - \mathbf{x}_{ij}^*)^T (\mathbf{x}_{il}^* - \mathbf{x}_{ij}^*) \mathbf{w} = 1 \\
& \|\mathbf{x}_{il}^* - \mathbf{x}_{ij}^*\|_2^2 = \frac{1}{\|\mathbf{w}\|_2^2} \\
& \|\mathbf{x}_{il}^* - \mathbf{x}_{ij}^*\|_2 = \frac{1}{\|\mathbf{w}\|_2}
\end{aligned}$$

Concerning the properties of D , positivity is ensured with the constraints $w_h \geq 0$ and because d_1, \dots, d_p are dissimilarity measures ($d_h \geq 0$). As the metric D is defined as a linear combination of dissimilarity measures d_1, \dots, d_p , it can be shown that symmetry and reflexivity is verified.

Voir si ce n'est pas redondant avec la partie interprétation

3.5 Quadratic formalization for M²TML

In this section, we define the problem of learning D as a linear or non-linear combination in the dissimilarity space using a quadratic regularizer. First, we give the optimization problem and its dual formulation form involving only dot products. Then, we discuss on the properties of the learnt metric D . Finally, we study a link between SVM and the quadratic formalization.

3.5.1 Primal and dual formalization

The formulation in Eq. 3.14 suppose that the metric D is a linear combination of the metrics d_h . The linear formalization being similar to the one of SVM, it can be derived into its dual form to extend the method to find non-linear solutions for D . For that, we propose to change the linear regularizer $R(Pull)$ in the objective function of Eq. 3.14 into a quadratic regularizer:

$$\begin{aligned} \underset{\mathbf{w}, \xi}{\operatorname{argmin}} \left\{ R(Pull) + C \sum_{\substack{j \in \dot{Pull}_i \\ l \in Push_i}} \frac{1 - y_{il}}{2} \xi_{ijl} \right\} \\ \text{s.t. } \forall i = 1, \dots, n, \forall j \in Pull_i, l \in Push_i, \\ \mathbf{w}^T(\mathbf{x}_{il} - \mathbf{x}_{ij}) \geq 1 - \xi_{ijl} \\ \xi_{ijl} \geq 0 \end{aligned} \quad (3.15)$$

Two solutions for $R(Pull)$ are studied:

$$1. \quad R(Pull) = \frac{1}{2} \sum_{h=1}^p \sum_{j \in \dot{Pull}_i} (w_h d_h(\mathbf{x}_{ij}))^2 \quad (3.16)$$

$$2. \quad R(Pull) = \frac{1}{2} \sum_{h=1}^p \left(\sum_{j \in \dot{Pull}_i} w_h d_h(\mathbf{x}_{ij}) \right)^2 = \frac{1}{2} m.n \sum_{h=1}^p (w_h \bar{d}_h)^2 \quad (3.17)$$

where $\bar{d}_h = \frac{1}{mn} \sum_{j \in \dot{Pull}_i} d_h(\mathbf{x}_{ij})$ denotes the mean of the distances $d_h(\mathbf{x}_{ij})$ for each metric d_h .

Let $\bar{\mathbf{x}} = [\bar{d}_1, \dots, \bar{d}_p]^T$ be a vector of size p containing the mean of the metrics $\bar{d}_1, \dots, \bar{d}_p$.

Other regularizations are possible. We focus on these two propositions that can be reduced to the following formula:

$$R(Pull) = \frac{1}{2} \mathbf{w}^T \mathbf{M} \mathbf{w} \quad (3.18)$$

where \mathbf{M} denotes respectively the following matrix for each regularizer:

$$1. \quad \mathbf{M} = \text{Diag}(\mathbf{X}_{pull}^T \mathbf{X}_{pull}) = \begin{bmatrix} \sum_{j \in \text{Pull}_i} d_1^2(\mathbf{x}_{ij}) & & 0 \\ & \ddots & \\ 0 & & \sum_{j \in \text{Pull}_i} d_p^2(\mathbf{x}_{ij}) \end{bmatrix} \quad (3.19)$$

$$2. \quad \mathbf{M} = \text{Diag}(\bar{\mathbf{x}}) \text{Diag}(\bar{\mathbf{x}}) = \begin{bmatrix} \bar{d}_1^2 & & 0 \\ & \ddots & \\ 0 & & \bar{d}_p^2 \end{bmatrix} \quad (3.20)$$

where $\mathbf{X}_{pull} = \bigcup_i \text{Pull}_i$ be a $(m.n) \times p$ matrix containing the vector $\mathbf{x}_{ij} \in \text{Pull}_i$.

From this, the optimization problem can be written using a quadratic regularization for the pull term:

$$\begin{aligned} & \underset{\mathbf{w}, \xi}{\text{argmin}} \left\{ \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{M} \mathbf{w}}_{\text{pull}} + C \underbrace{\sum_{\substack{j \in \text{Pull}_i \\ l \in \text{Push}_i}} \frac{1 - y_{il}}{2} \xi_{ijl}}_{\text{push}} \right\} \\ & \text{s.t. } \forall i = 1, \dots, n, \forall j \in \text{Pull}_i, l \in \text{Push}_i, \\ & \quad \mathbf{w}^T (\mathbf{x}_{il} - \mathbf{x}_{ij}) \geq 1 - \xi_{ijl} \\ & \quad \xi_{ijl} \geq 0 \end{aligned} \quad (3.21)$$

Similarly to SVM, 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 and the constraints multiplied by their respective Lagrange multipliers $\boldsymbol{\alpha}$ and \mathbf{r} :

$$\begin{aligned} L(\mathbf{w}, \xi, \boldsymbol{\alpha}, \mathbf{r}) = & \frac{1}{2} \mathbf{w}^T \mathbf{M} \mathbf{w} + C \sum_{ijl} \frac{1 - y_{il}}{2} \xi_{ijl} - \sum_{ijl} r_{ijl} \xi_{ijl} \\ & - \sum_{ijl} \alpha_{ijl} (\mathbf{w}^T (\mathbf{x}_{il} - \mathbf{x}_{ij}) - 1 + \xi_{ijl}) \end{aligned} \quad (3.22)$$

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:

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

The matrix \mathbf{M} being diagonal in both case (Eqs. 3.19 & 3.20), it is thus inversible. The equations leads to:

$$\mathbf{w} = \mathbf{M}^{-1} \sum_{ijl} \alpha_{ijl} (\mathbf{x}_{il} - \mathbf{x}_{ij}) \quad (3.23)$$

$$r_{ijl} = C - \alpha_{ijl} \quad (3.24)$$

Substituting Eq. 3.23 and 3.24 back into $L(\mathbf{w}, \xi, \boldsymbol{\alpha}, \mathbf{r})$ in Eq. 3.22, we get the dual formula¹:

$$\begin{aligned} & \underset{\boldsymbol{\alpha}}{\operatorname{argmax}} \left\{ \sum_{ijl} \alpha_{ijl} - \frac{1}{2} \sum_{ijl} \sum_{i'j'l'} \alpha_{ijl} \alpha_{i'j'l'} (\mathbf{x}_{il} - \mathbf{x}_{ij})^T \mathbf{M}^{-1} (\mathbf{x}_{i'l'} - \mathbf{x}_{i'j'}) \right\} \\ & \text{s.t. } \forall i = 1, \dots, n, \forall j \in \text{Pull}_i, l \in \text{Push}_i, \\ & 0 \leq \alpha_{ijl} \leq C \end{aligned} \quad (3.25)$$

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

$$D(\mathbf{x}_{i'j'}) = \underbrace{\sum_{ijl} \alpha_{ijl} (\mathbf{x}_{il} - \mathbf{x}_{ij})^T \mathbf{M}^{-1} \mathbf{x}_{i'j'}}_{\mathbf{w}^T} \quad (3.26)$$

$$D(\mathbf{x}_{i'j'}) = \underbrace{\sum_{ijl} \alpha_{ijl} \mathbf{x}_{il}^T \mathbf{M}^{-1} \mathbf{x}_{i'j'}}_{\text{similarity of } \mathbf{x}_{i'j'} \text{ to Push set}} - \underbrace{\sum_{ijl} \alpha_{ijl} \mathbf{x}_{ij}^T \mathbf{M}^{-1} \mathbf{x}_{i'j'}}_{\text{similarity of } \mathbf{x}_{i'j'} \text{ to Pull set}} \quad (3.27)$$

3.5.2 Non-linear combined metric

The above formula can extended to non-linear function for the metric D . The inner product in Eq. 3.27 can be easily kernelized using the "kernel" trick:

$$\begin{aligned} \mathbf{x}_{il}^T \mathbf{M}^{-1} \mathbf{x}_{i'j'} &= \mathbf{x}_{il}^T \mathbf{M}^{-\frac{1}{2}} \mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{i'j'} \\ &= \left(\mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{il} \right)^T \left(\mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{i'j'} \right) \\ &= \langle \mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{il}; \mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{i'j'} \rangle \\ &= \kappa(\mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{il}; \mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{i'j'}) \end{aligned}$$

¹complete details of the calculations in Appendix D

As \mathbf{M}^{-1} is a diagonal matrix, it is inversible and can be written $\mathbf{M}^{-1} = \mathbf{M}^{-\frac{1}{2}}\mathbf{M}^{-\frac{1}{2}}$. For each regularization, we give below the matrix $\mathbf{M}^{-\frac{1}{2}}$:

$$1. \quad \mathbf{M}^{-\frac{1}{2}} = \begin{bmatrix} \frac{1}{\sum_{j \in Pull_i} d_1^2(\mathbf{x}_{ij})} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\sum_{j \in Pull_i} d_p^2(\mathbf{x}_{ij})} \end{bmatrix} \quad (3.28)$$

$$2. \quad \mathbf{M}^{-\frac{1}{2}} = \begin{bmatrix} \frac{1}{\bar{d}_1^2} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\bar{d}_p^2} \end{bmatrix} \quad (3.29)$$

The matrix $\mathbf{M}^{-\frac{1}{2}}$ can be interpreted in the first regularization proposition as a normalization by the variance of the distance for each metric d_h . In the second regularization, it can be interpreted as a normalization by the mean of the distance for each metric d_h .

By replacing the inner product by a kernel back into Eq. 3.27, we obtain:

$$D(\mathbf{x}_{i'j'}) = \overbrace{\sum_{ijl} \alpha_{ijl} \kappa(\mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{il}; \mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{i'j'})}^{\text{similarity of } \mathbf{x}_{i'j'} \text{ to } Push \text{ set}} - \overbrace{\sum_{ijl} \alpha_{ijl} \kappa(\mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{ij}; \mathbf{M}^{-\frac{1}{2}} \mathbf{x}_{i'j'})}^{\text{similarity of } \mathbf{x}_{i'j'} \text{ to } Pull \text{ set}} \quad (3.30)$$

Let's give some interpretation and discussion about the properties of D . Similarly to SVM, from Eq. 3.26, at the optimality, only the triplets $(\mathbf{x}_{il} - \mathbf{x}_{ij})$ with $\alpha_{ijl} > 0$ are considered as the support vectors and the computation of the metric D depends only on these support vectors. Note that in this case, there exists two categories of support vectors (Eqs. 3.27 & 3.30): the vectors \mathbf{x}_{il} from the push set $Push_i$ and the vectors \mathbf{x}_{ij} from the pull set $Pull_i$ which $\alpha_{ijl} > 0$. The resulting metric D can be interpreted as the difference involving two similarity terms: a new pair $\mathbf{x}_{i'j'}$ is as dissimilar as its similarity to the *Push* set is high while its similarity to the *Pull* set is low. Inversely, the pair $\mathbf{x}_{i'j'}$ is as similar as its similarity to the *Push* set is low while its similarity to the *Pull* set is high.

Concerning the properties of the metric D , it can be shown that symmetry and reflexivity is ensured. However, as D is a difference of two similarity terms, positivity is not always ensured, *e.g.*, the similarity of the pull term is greater than the similarity of the push term. The resulting metric D is not thus a dissimilarity measure.

3.5.3 Link between SVM and the quadratic formalization

Many parallels have been studied between Large Margin Nearest Neighbors (LMNN) and SVM (Section 2.6.3). 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. Similarly, we can make a link between the quadratic formalization and a SVM problem where the form of the metric D is defined *a priori*.

We show in Appendix E that solving the SVM problem in Eq. 3.31 for \mathbf{w} and b solves a similar problem with a quadratic regularization in Eq. 3.21 for $D(\mathbf{x}_{ij}) = -\frac{1}{2}(\mathbf{w}^T \mathbf{x}_{ij} + b)$.

$$\begin{aligned} & \underset{\mathbf{w}, b, \xi}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{\substack{j \in Pull_i \text{ or} \\ j \in Push_i}} p_i \xi_{ij} \right\} \\ & \text{s.t. } \forall i, j \in Pull_i \text{ or } j \in Push_i : \\ & y_{ij}(\mathbf{w}^T \mathbf{x}_{ij} + b) \geq 1 - \xi_{ij} \end{aligned} \quad (3.31)$$

where p_i is a weight factor for each slack variable ξ_{ij} (in classical SVM, $p_i = 1$). The loss part in the SVM formulation can be split into 2 terms involving the sets $Pull_i$ and $Push_i$:

$$\begin{aligned} & \underset{\mathbf{w}, b, \xi}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{j \in Pull_i} p_i^+ \xi_{ij} + C \sum_{l \in Push_i} p_i^- \xi_{il} \right\} \\ & \text{s.t. } : \\ & \forall i, j \in Pull_i : y_{ij}(\mathbf{w}^T \mathbf{x}_{ij} + b) \geq 1 - \xi_{ij} \\ & \forall i, l \in Push_i : y_{il}(\mathbf{w}^T \mathbf{x}_{il} + b) \geq 1 - \xi_{il} \end{aligned} \quad (3.32)$$

where p_i^+ and p_i^- are the weight factors for pull pairs $Pull_i$ and push pairs $Push_i$. To obtain an equivalence, we set p_i^- as the half of the number pairs in $Pull_i$ and p_i^+ as the half of the number of time series L in $Push_i$:

$$p_i^- = \frac{k}{2} = \sum_{j \in Pull_i} \frac{1}{2} \quad (3.33)$$

$$p_i^+ = \frac{L}{2} = \frac{1}{2} \sum_l \frac{1 - y_{il}}{2} \quad (3.34)$$

In particular, let's underline the main similarities and differences:

- Both problems share a same set of constraints between triplets:

$$\forall i, j \in Pull_i, l \in Push_i : D(\mathbf{x}_{il}) - D(\mathbf{x}_{ij}) \geq 1 - \xi_{ijl}$$

- The SVM problem includes an additional set of constraints that is not present in the quadratic formalization. SVM takes into account pull pairs \mathbf{x}_{ij} and push pairs \mathbf{x}_{kl} that

doesn't belong to the same neighborhood:

$$\forall i, j \in Pull_i, k, l \in Push_k, i \neq k : D(\mathbf{x}_{kl}) - D(\mathbf{x}_{ij}) \geq 1 - \frac{\xi_{kl} + \xi_{ij}}{2}$$

- The SVM problem includes in the loss term additional slack variables ξ_{ijl} that are not present in the quadratic formalization because of the additional set of constraints. It is not only a push term.
- The two problems involves different regularized term: in the quadratic formalization, the regularizer involves a pull action (Eqs. 3.16 & 3.17), not present in SVM.
- Both problems suppose at first a linear combination for D .

Concerning the properties of the metric D , positivity is not ensured in the primal and dual formulation. Symmetry and reflexivity is ensured.

3.6 SVM-based formalization for M²TML

In this section, we present a solution based on SVM where the form of the metric D is not known *a priori*. We formulate the problem as a SVM problem to solve a large margin problem between $Pull_i$ and $Push_i$ sets, and then, induce a combined metric D for the obtained SVM solution. Thanks to the SVM framework, the proposition can be naturally extended to learn both, linear or non-linear function for the metric D .

3.6.1 Support Vector Machine (SVM) resolution

Let $\{\mathbf{x}_{ij}; y_{ij} = \pm 1\}$, $\mathbf{x}_{ij} \in Pull_i \cup Push_i$ be the training set, with $y_{ij} = -1$ for $\mathbf{x}_{ij} \in Push_i$ and $+1$ for $\mathbf{x}_{ij} \in Pull_i$. For a maximum margin between the sets $Pull_i$ and $Push_i$, the problem is formalized in the dissimilarity space \mathcal{E} :

$$\begin{aligned} & \underset{\mathbf{w}, b, \xi}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i,j} \xi_{ij} \\ & \text{s.t. } y_{ij}(\mathbf{w}^T \mathbf{x}_{ij} + b) \geq 1 - \xi_{ij} \\ & \xi_{ij} \geq 0 \end{aligned} \tag{3.35}$$

In the linear case, a L_1 regularization in Eq. 3.35 leads to a sparse and interpretable \mathbf{w} that uncovers the modalities, periods and scales that differentiate best pull from push pairs for a robust nearest neighbors classification. Note that in practice, the local neighborhoods for each sample \mathbf{x}_i can have very different scales. Thanks to the unit radii normalization \mathbf{x}_{ij}/r_i , where r_i denotes the norm of the m -th neighbors in $Pull_i$, the SVM ensures a global large margin solution involving equally local neighborhood constraints (*i.e.* local margins).

Note that any multi-class problem is transformed in the dissimilarity space as a binary classification problem.

3.6.2 Linearly separable Pull and Push sets

Let \mathbf{x}_{test} be a new sample, $\mathbf{x}_{i,test} \in \mathcal{E}$ gives the proximity between \mathbf{x}_i and \mathbf{x}_{test} based on the p multi-modal and multi-scale metrics d_h . We denote $\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{i,test})$ the orthogonal projection of $\mathbf{x}_{i,test}$ on the axis of direction \mathbf{w} and $\|\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{i,test})\|$ its norm. 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}_{test} . In particular, for M²TML, two quantities are used to define the dissimilarity measure: the projected norm and the distance to the margin.

The projected norm $\|\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{i,test})\|$ of $\mathbf{x}_{i,test}$ on the direction \mathbf{w} limits the comparison of \mathbf{x}_i and \mathbf{x}_{test} to the features separating pull and push sets (Fig. 3.8), it is defined as:

$$\|\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{i,test})\| = \frac{\|\mathbf{w}^T \mathbf{x}_{i,test}\|}{\|\mathbf{w}\|} \quad (3.36)$$

with:

$$\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{i,test}) = \frac{\langle \mathbf{w}, \mathbf{x}_{i,test} \rangle}{\|\mathbf{w}\|^2} \mathbf{w} = \frac{\mathbf{w}^T \mathbf{x}_{i,test}}{\|\mathbf{w}\|^2} \mathbf{w} \quad (3.37)$$

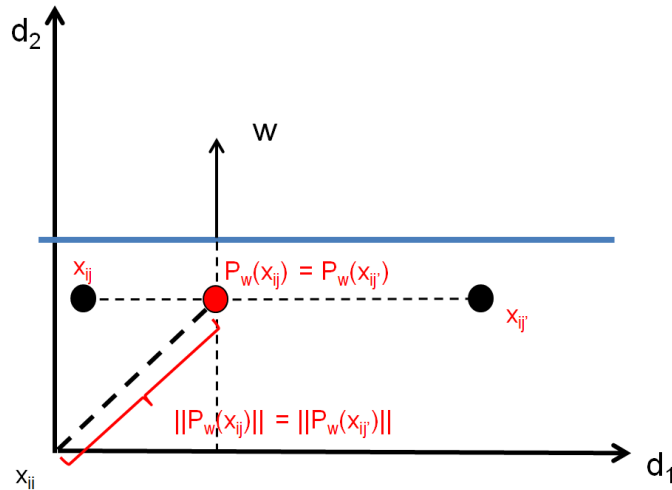


Figure 3.8: The projected vector $\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{ij})$ and $\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{ij'})$

Although the norm $\|\mathbf{P}_{\mathbf{w}}(\mathbf{x}_{i,test})\|$ satisfies positivity, it doesn't guarantee lower distances for

pull pairs than for push pairs as illustrated in Fig 3.9.

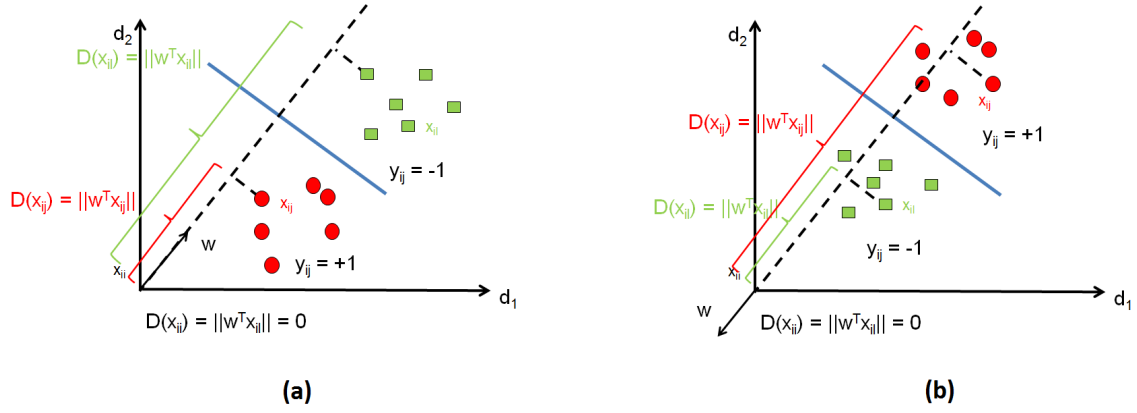


Figure 3.9: Example of SVM solutions and of the resulting metric D defined by the norm of the projection on w . Fig. (a) represents common expected configuration where pull pairs $Pull_i$ are situated in the same side as the origin $x_{ii} = 0$. In Fig. (b), the vector $w = [-1 -1]^T$ indicates that push pairs $Push_i$ are on the side of the origin point. One problem arises in Fig. (b): distance of push pairs $D(x_{il})$ is lower than the distance of pull pairs $D(x_{ij})$.

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

We propose to add an exponential term to operate a "push" on push pairs based on their distances to the separator hyperplane, that leads to the dissimilarity measure D of required properties:

$$D(x_{i,test}) = \|P_w(x_{i,test})\| \cdot \exp(\lambda[-(w^T P_w(x_{i,test}) + b)]_+) \quad \lambda > 0 \quad (3.38)$$

where λ controls the "push" term and $w^T P_w(x_{i,test}) + b$ defines the distance between the orthogonal projected vector and the separator hyperplane; $[t]_+ = \max(0; t)$ being the positive operator. Note that, for a pair lying into the pull side ($y_{ij} = +1$), $[-(w^T P_w(x_{i,test}) + b)]_+ = 0$, the exponential term is vanished (i.e. no "pull" action) and the dissimilarity leads to the norm term. For a pair situated in the push side ($y_{ij} = -1$), the norm is expanded by the push term, all the more the distance to the hyperplane is high.

Fig. 3.10, illustrates for $p = 2$ the behavior of the learned dissimilarity according to two extreme cases. The first one (Fig. 3.10-a), represents common expected configuration where pairs $Pull_i$ are situated in the same side as the origin. The dissimilarity increases proportionally to the norm in the pull side, then exponentially on the push side. Although the expansion operated in the push side is dispensable in that case, it doesn't affect nearest neighbors classification. Fig. 3.10-b, shows a challenging configuration where pairs $Push_i$ are situated in the same side as the origin. The dissimilarity behaves proportionally to the norm on the pull side, and increases exponentially from the hyperplane until an abrupt decrease induced by a norm near 0. Note that the region under the abrupt decrease mainly uncovers

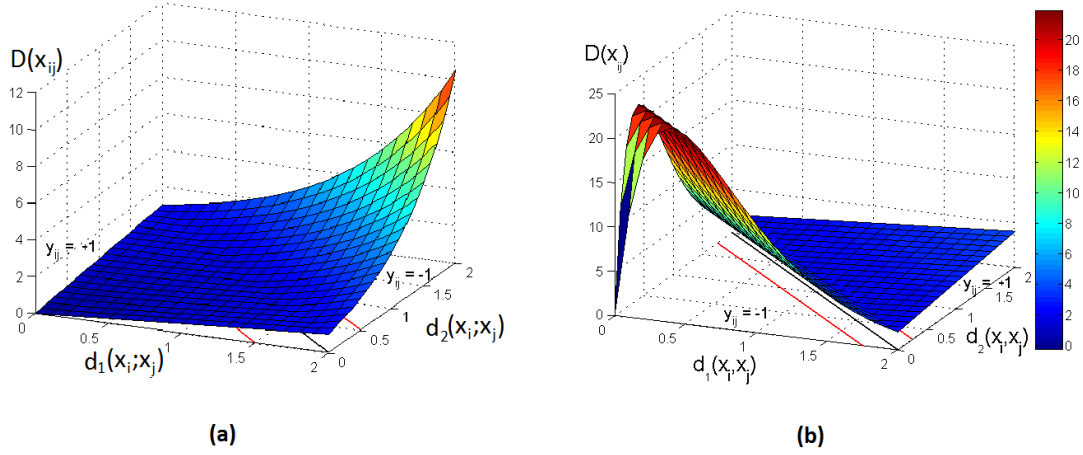


Figure 3.10: The behavior of the learned metric D ($p = 2$; $\lambda = 2.5$) with respect to common (a) and challenging (b) configurations of pull and push pairs.

false pairs $Push_i$, *i.e.*, pairs of norm zero labeled differently.

3.6.3 Non-linearly separable Pull and Push sets

The above solution holds true for any kernel κ and allows to extend the dissimilarity D given in Eq. 3.38 to non linearly separable pull and push pairs. Let κ be a kernel defined in the dissimilarity space \mathcal{E} and the related Hilbert space (feature space) \mathcal{H} . For a non linear combination function of the metrics $d_h, h = 1, \dots, p$ in \mathcal{E} , we define the dissimilarity measure $D_{\mathcal{H}}$ in the feature space \mathcal{H} as:

$$D_{\mathcal{H}}(\mathbf{x}_{i,test}) = |(\|\mathbf{P}_{\mathbf{w}}(\Phi(\mathbf{x}_{i,test}))\| - \|\mathbf{P}_{\mathbf{w}}(\Phi(\mathbf{0}))\|)| \cdot \exp \left(\lambda \left[- \left(\sum_{ij} y_{ij} \alpha_{ij} \kappa(\mathbf{x}_{ij}, \mathbf{x}_{i,test}) + b \right) \right]_+ \right) \quad \lambda > 0 \quad (3.39)$$

with $\Phi(\mathbf{w})$ the image of \mathbf{w} into the feature space \mathcal{H} . Based on Eq. 3.37, substituting Eq. 3.23 back into \mathbf{w} , the inner product gives $\langle \mathbf{w}; \Phi(\mathbf{x}_{i,test}) \rangle = \sum_{ij} y_{ij} \alpha_{ij} \kappa(\mathbf{x}_{ij}, \mathbf{x}_{i,test})$ and the norm of \mathbf{w} gives $\|\mathbf{w}\| = \sqrt{\sum_{ijkl} \alpha_{ij} \alpha_{kl} y_{ij} y_{kl} \kappa(\mathbf{x}_{ij}, \mathbf{x}_{kl})}$. Replacing back into Eq. 3.36, the norm of the orthogonal projection of $\Phi(\mathbf{x}_{i,test})$ on $\Phi(\mathbf{w})$ gives:

$$\|\mathbf{P}_{\mathbf{w}}(\Phi(\mathbf{x}_{i,test}))\| = \frac{\sum_{ij} y_{ij} \alpha_{ij} \kappa(\mathbf{x}_{ij}, \mathbf{x}_{i,test})}{\sqrt{\sum_{ijkl} \alpha_{ij} \alpha_{kl} y_{ij} y_{kl} \kappa(\mathbf{x}_{ij}, \mathbf{x}_{kl})}} \quad (3.40)$$

Note that as $\Phi(\mathbf{0})$ doesn't meet the origin in the feature space \mathcal{H} , the norms in Eq. 3.39 are centered with respect to $\Phi(\mathbf{0})$ to ensure the reflexivity property. It is easy to show that both

D and $D_{\mathcal{H}}$ ensure the properties of a dissimilarity (positivity, reflexivity, symmetry).

Note that the framework to define the metric D and $D_{\mathcal{H}}$ can also be used in the linear and quadratic formalization. However, the obtained solution for D and $D_{\mathcal{H}}$ can be far away from the original form of D that was optimized in the optimization problem.

3.7 SVM-based solution and algorithm for M²TML

In this section, we review the main steps of this latter 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 1.1.4.

Neighborhood scaling. In real datasets, local neighborhoods can have very different scales as illustrated in Fig. E.1. 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:

$$\mathbf{x}_{ij}^{norm} = \left[\frac{d_1(\mathbf{x}_i, \mathbf{x}_j)}{r_i}, \dots, \frac{d_p(\mathbf{x}_i, \mathbf{x}_j)}{r_i} \right]^T \quad (3.41)$$

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 $Pull_i$:

$$r_i = \max_{\mathbf{x}_{ij} \in Pull_i} \|\mathbf{x}_{ij}\|_2 \quad (3.42)$$

For simplification purpose, we denote \mathbf{x}_{ij} as \mathbf{x}_{ij}^{norm} . Fig. 3.11 illustrates the effect of neighborhood scaling in the dissimilarity space.

Finally, Algorithm 1 summarizes the main steps to learn a multi-modal and multi-scale metric D for a robust nearest neighbors classification. Algorithm 2 details the steps to classify a new sample \mathbf{x}_{test} using the learned metric D .

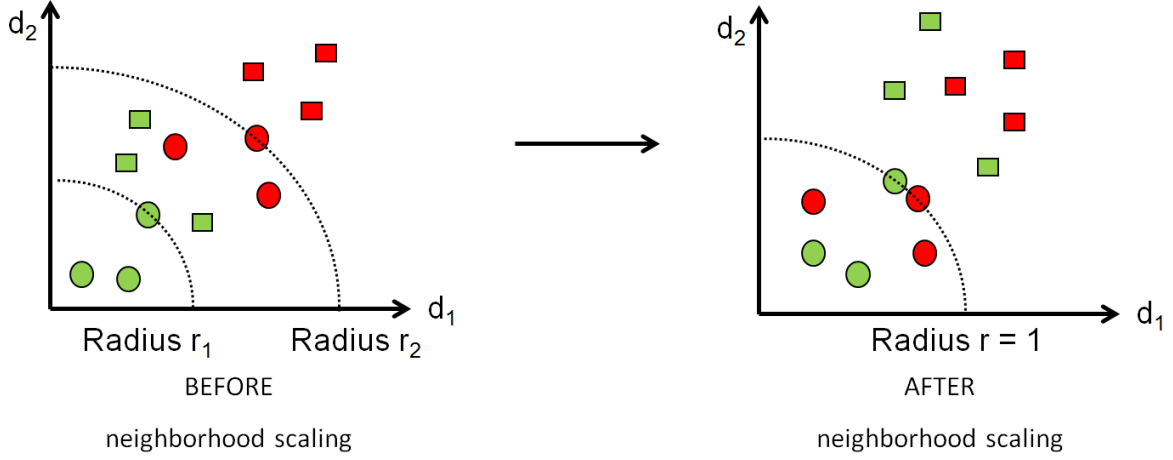


Figure 3.11: Effect of neighborhood scaling before (left) and after (right) on the neighborhood of two time series \mathbf{x}_1 (green) and \mathbf{x}_2 (red). Circle represent pull pairs $Pull_i$ and square represents push pairs $Push_i$ for $m = 3$ neighbors. Before scaling, the problem is not linearly separable. The spread of each neighborhood are not comparable. After scaling, the target neighborhood becomes comparable and in this example, the problem becomes linearly separable between the circles and the squares.

Algorithm 1 Multi-modal and Multi-scale Temporal Metric Learning (M^2TML) for k -NN classification

- 1: Input: $\{\mathbf{x}_i, y_i\}_{i=1}^n$ n labeled time series
 d_1, \dots, d_p metrics as described in Eqs. 2.1, 2.4, 2.6, 3.3
a kernel κ
 - 2: Output: the learned dissimilarity D or $D_{\mathcal{H}}$ depending of κ
 - 3: *Dissimilarity embedding*
Embed pairs $(\mathbf{x}_i, \mathbf{x}_j)$ $i, j \in 1, \dots, n$ into \mathcal{E} as described in Eq. 3.1 and normalize d_h s
 - 4: *Build $Pull_i$ and $Push_i$ sets*
Build the sets of pairs $Pull_i$ and $Push_i$ as described in Eq. 3.7 & 3.8 and scale the radii to 1 (Eq. 3.41).
 - 5: *SVM learning*
Train a SVM for a large margin classifier between $Pull_i$ and $Push_i$ sets (Eq. 3.35)
 - 6: *Dissimilarity definition*
Consider Eq. 3.38 (resp. Eq. 3.39) to define D (resp. $D_{\mathcal{H}}$) a linear (resp. non linear) combination function of the metrics d_h s.
-

3.8 Conclusion

To learn a multi-modal and multi-scale temporal combined metric, we propose in this chapter to embed time series into a dissimilarity space. The multi-modal and multi-scale metric learning (M^2TML) problem can be formalized as a problem of learning a function in the dissimilarity space, that ensures the properties of a dissimilarity. We formulate the M^2TML problem into a general optimization problem involving a pull and push term. Choosing a m -neighborhood,

Algorithm 2 k -NN classification using the learned metric D or $D_{\mathcal{H}}$

- 1: Input: $\{\mathbf{x}_i, y_i\}_{i=1}^n$ n labeled time series
 $\{\mathbf{x}_{test}, y_{test}\}$ a labeled time series to test
 d_1, \dots, d_p metrics as described in Eqs. 2.1, 2.4, 2.6, 3.3
the learned dissimilarity D or $D_{\mathcal{H}}$ depending of the kernel κ
 - 2: Output: Predicted label \hat{y}_{test}
 - 3: *Dissimilarity embedding*
Embed pairs $(\mathbf{x}_i, \mathbf{x}_{test})$ $i \in 1, \dots, n$ into \mathcal{E} as described in Eq. 3.1 and normalize d_h s using the same normalization parameters in Algorithm 1
 - 4: *Combined metric computation*
Consider Eq. 3.38 (resp. Eq. 3.39) to compute $D(\mathbf{x}_i, \mathbf{x}_{test})$ (resp. $D_{\mathcal{H}}(\mathbf{x}_i, \mathbf{x}_{test})$) a linear (resp. non linear) combination function of the metrics $d_h(\mathbf{x}_i, \mathbf{x}_{test})$.
 - 5: *Classification*
Consider the k lowest dissimilarities $D(\mathbf{x}_i, \mathbf{x}_{test})$ (resp. $D_{\mathcal{H}}(\mathbf{x}_i, \mathbf{x}_{test})$). Extract the labels y_i of the considered \mathbf{x}_i and make a vote scheme to predict the label \hat{y}_{test} of \mathbf{x}_{test}
-

greater than the k -neighborhood allows to generalize better the learnt metric. From the general formalization, we propose three different formalizations (Linear, Quadratic, SVM-based). Table 3.1 sums up the characteristics of each formalization and the induced dissimilarity.

	Linear formalization	Quadratic formalization	SVM-based formalization
D	Linear	Linear/Non-linear	Linear/Non-linear
Sparcity	Yes	No	Yes/No
Dissimilarity properties	Yes	No (non-positivity)	Yes

Table 3.1: The different formalizations for M²TML

The adaptation of SVM in the dissimilarity space to learn the multi-modal and multi-scale metric D have brought us to propose a pre-processing step before solving the problem such as the neighborhood scaling, and a post-processing step such as defining the metric D .

As we have defined all functions components of our algorithms (learning, testing), we test our proposed algorithms M²TML in the next chapter on large public datasets.

Experiments

Contents

4.1	Description	73
4.2	Experimental protocol	76
4.3	Results and discussion	77
4.3.1	Results	77
4.3.2	Comparison of the classification performances on the test set	78
4.3.3	Analysis of the discriminative features	79
4.3.4	Effect on the neighborhood before and after learning	81
4.4	Conclusion of the chapter	83

In this chapter, we evaluate the efficiency of the proposed M^2TML algorithm on public datasets for classification problems of univariate time series. First, we describe the datasets. Secondly, we detail the experimental protocol. Finally, we present and discuss the obtained results.

4.1 Description

The efficiency of the learned multi-modal and multi-scale dissimilarities D and $D_{\mathcal{H}}$ is evaluated through a 1–NN classification on 30 public datasets¹[Keo+11]. The 1–NN classifier is used to make the results comparable with the results of the UCR time series data mining archive. Time series comes from several fields (simulated data, medical data, electrical data, etc.), are from variable lengths (from small ($q = 24$) to long lengths ($q = 1882$)) and the number of classes to discriminate evolves between 1 and 37 classes. Note that some of the datasets have a small number of time series in the training set ($n < 30$) and others have a large number of time series in the training set ($n > 100$). The results using standard metrics (Euclidean distance, Dynamic time warping) show both easy and challenging classifications problems, the latter being opened for improvements. The datasets encompass time series that involve global or local temporal differences, require or not time warping, with linearly or non linearly separable neighborhoods, discussed later. Table 4.1 gives a description of the datasets considered in

¹PowerCons: <https://archive.ics.uci.edu/ml/datasets/Individual+household+electric+power+consumption>, BME and UMD: <http://ama.liglab.fr/~douzal/tools.html>.

the experiments and Fig. 4.1 gives the temporal representation for some of the datasets. Note that for some datasets (*e.g.*, SonyAIBO, ECG200, FaceFour, PowerConsumption), it is visually difficult to discriminate the classes using one modality (value, behavior, frequential).

Dataset	Nb. Class	Nb. Train	Nb. Test	TS length
1 ItalyPowerD	2	67	1029	24
2 CinCECGtorso	4	40	1380	1639
3 ECG200	2	100	100	96
4 SonyAIBOII	2	27	953	65
5 Coffee	2	28	28	286
6 ECG5Days	2	23	861	136
7 SonyAIBO	2	20	601	70
8 Adiac	37	390	391	176
9 Beef	5	30	30	470
10 Trace	4	100	100	275
11 CBF	3	30	900	128
12 CC	6	300	300	60
13 UMDsmooth	3	360	1440	150
14 BMEsmooth	3	300	1500	128
15 DiatomSizeReduc	4	16	306	345
16 Symbols	6	25	995	398
17 GunPoint	2	50	150	150
18 FacesUCR	14	200	2050	131
19 TwoLeadECG	2	23	1139	82
20 MoteStrain	2	20	1252	84
21 Lighting2	2	60	61	637
22 OliveOil	4	30	30	570
23 FISH	7	175	175	463
24 FaceFour	4	24	88	350
25 SwedishLeaf	15	500	625	128
26 MedicalImages	10	381	760	99
27 Lighting7	7	70	73	319
28 PowerCons	2	73	292	144
29 OSULeaf	6	200	242	427
30 InlineSkate	7	100	550	1882

Table 4.1: Dataset description giving the number of classes (Nb. Class), the number of time series for the training (Nb. Train) and the testing (Nb. Test) sets, and the length of each time series (TS length).

The results of the learned metrics D and $D_{\mathcal{H}}$ are compared to those of three *a priori* combined metrics D_{Lin} , D_{Geom} , D_{Sig} (Eqs. 2.16, 2.17, 2.18) and five alternative uni-modal metrics covering:

1. The standard Euclidean distance d_A (Eq. 2.1) and Dynamic time warping DTW
2. The behavior-based measures d_B (Eq. 2.6) and d_{B-DTW} its counterpart for asynchronous time series, that is d_B is evaluated once time series synchronized by dynamic programming
3. The frequential-based metric d_F (Eq. 2.4).

Table 4.2 recalls briefly the considered metrics in the experiments.

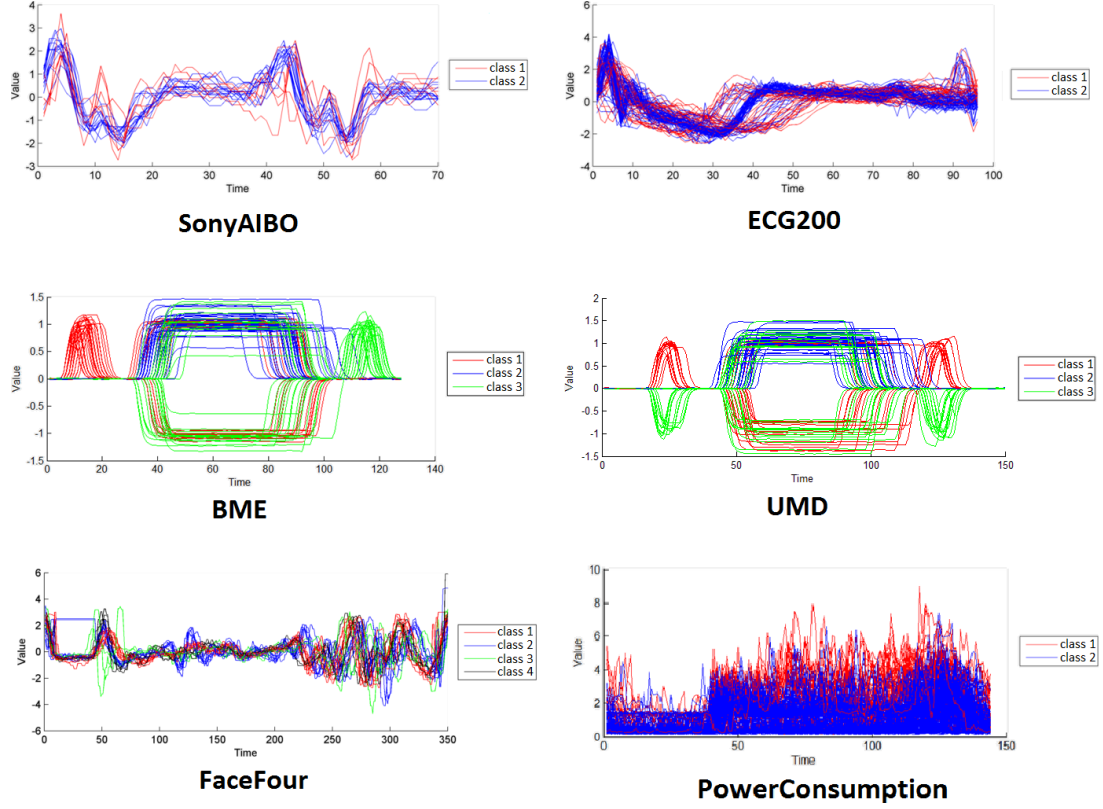


Figure 4.1: Temporal representation of some datasets (SonyAIBO, ECG200, BME, UMD, FaceFour, PowerConsumption) considered in the experiments.

Symbol	Name	Equation reference	Description
d_A	Value-based dissimilarity	Eq. 2.1	Euclidean distance
d_B	Behavior-based dissimilarity	Eq. 2.6	Behavior metric based on cort
DTW	Dynamic time warping	Eqs. 2.13 & 2.1	Euclidean distance after alignment
d_{B-DTW}	Behavior-based aligned dissimilarity	Eqs. 2.13 & 2.6	Behavior metric based on cort after alignment
d_F	Frequential-based dissimilarity	Eq. 2.4	Frequential based on Fourier transform
D_{Lin}	Linear combined metric	Eq. 2.16	Combines d_A and d_B (resp. DTW and d_{B-DTW})
D_{Geom}	Geometric combined metric	Eq. 2.17	Combines d_A and d_B (resp. DTW and d_{B-DTW})
D_{Sig}	Sigmoid combined metric	Eq. 2.18	Combines d_A and d_B (resp. DTW and d_{B-DTW})
D	Linear learned metric	Eq. 3.38	M ² TML linear combined metric
$D_{\mathcal{H}}$	Non-linear learned metric	Eq. 3.39	M ² TML non-linear combined metric with a Gaussian kernel

Table 4.2: Considered metric in the experiments

The *a priori* combined metrics (D_{Lin} , D_{Geom} , D_{Sig}) relies, for synchronous (resp. asynchronous), on 2 log-normalized dissimilarities d_A , d_B (resp. DTW, d_{B-DTW}). The alternative metrics and the *a priori* combined metrics are evaluated as usual by involving the all time series elements (*i.e.* at the global scale). For D and $D_{\mathcal{H}}$, we consider a 21-dimensional embedding space \mathcal{E} that relies, for synchronous (resp. asynchronous) data, on 3 log-normalized dissimilarities d_A^s , d_B^s (resp. DTW^s, d_{B-DTW}^s), and d_F^s , at 7 temporal granularities $s \in \{0, \dots, 6\}$ obtained by binary segmentation, described in Section 3.2.

4.2 Experimental protocol

The different metrics can be split into two categories. For those without parameters to tune (d_A , DTW), the 1-NN classifier is applied directly on the test set. For those that require to tune parameters (d_B , d_{B-DTW} , D_{Lin} , D_{Geom} , D_{Sig} , D , $D_{\mathcal{H}}$), we recall briefly the grid search and cross-validation procedure (Section 1.1.2). When a learning algorithm requires to tune some parameters, to avoid overfitting, the training set can be divided into two sets: a learning and a validation set. The model is learnt for each combination of parameters (grid search) on the learning set and evaluated on the validation set. The model with the lowest error on the validation set is retained. An other alternative is cross-validation, which partitions the training set into v folds, performs the learning on one subset, and validates on the $v - 1$ other subsets. To take into account of variability within the data, multiple rounds of cross-validation are performed using different partitions, and the validation results are averaged over the rounds. Note that for unbalanced datasets in classification problems, it is recommended to use stratified sampling. Table 4.3 resumes the parameter ranges for each metric. We recall that the parameters retained are those that:

- **First**, minimize the average classification error on the validation set.
- **Secondly**, in the case of multiple solutions leading to equal performances, the most discriminant one is retained (*i.e.*, making closer pull pairs and far away push pairs). Precisely, it minimizes the ratio $\frac{d_{intra}}{d_{inter}}$ where d_{intra} and d_{inter} stands respectively to the mean of all intraclass and interclass distances.

As D and $D_{\mathcal{H}}$ involves several parameters to be tuned, we detail hereafter the procedure. The combined metrics D and $D_{\mathcal{H}}$ (κ as the Gaussian kernel) are learned respectively under L_1 and L_2 regularization, using LIBLINEAR and LIBSVM libraries [FCH08]; [HCL08]. The parameters are estimated on a validation set by line/grid search. A cross-validation and stratified sampling for unbalanced datasets are used. Particularly, for each couple (r, λ) $r \in \{1, 4, 10\}$ and $\lambda \in \{0, 10, 30\}$, the pairwise SVM parameters (C, α, γ) are learned by grid search as indicated in Table 4.3. The temporal order r for the behavior-based metrics d_B is noise-dependent, typically 1 is retained for noise-free data. The parameter λ corresponds to the strength of the 'push' term; precisely, if no, moderate or strong 'push' is required during the training process, a λ value of 0, 10 and 30 is learned, respectively.

Dissimilarity	Parameter	Ranges	Description
d_B	r	$\{1, 2, 3, \dots, q - 1\}$	Order of behavior-based metric
$D_{Lin}, D_{Geom}, D_{Sig}$	α	$\{0, 0.1, \dots, 1\}$	Trade-off between value and behavior components
$D, D_{\mathcal{H}}$	λ	$\{0, 10, 30\}$	Strength of the 'push' term
$D, D_{\mathcal{H}}$	C	$\{10^{-3}, 0.5, 1, 5, 10, 20, 30, \dots, 150\}$	Parameter of svm
$D, D_{\mathcal{H}}$	α	$\{1, 2, 3\}$	Size of the $m = \alpha.k$ neighborhood
$D_{\mathcal{H}}$	γ	$\{10^{-3}, 10^{-2}, \dots, 10^3\}$	Parameter of the Gaussian kernel

Table 4.3: Parameter ranges

4.3 Results and discussion

In this section, we present first a summary table of the quantitative results obtained in the experiment. Secondly, we present an analysis of the performances of the different metrics. Finally, we present the ability of our proposed approach M²TML to extract discriminative features.

4.3.1 Results

Table 4.4 reports the 1-NN classification test errors when based on uni-modal metrics (first 5 columns), on three *a priori* combined metrics (D_{Lin} , D_{Geom} , D_{Sig}) and on D and $D_{\mathcal{H}}$. The results for each dataset that are statistically and significantly better than the rest are indicated in bold (z-test at 5% risk detailed in Section 1.1.3.a). The last column 'WARP' indicates the synchronous (✓) or asynchronous (×) data type.

Dataset	Alternative uni-modal metrics					A priori combinations			M ² TML		WARP
	d_A	d_B	d_F	DTW	d_{B-DTW}	D_{Lin}	D_{Geom}	D_{Sig}	$D(\lambda^*)$	$D_{\mathcal{H}}(\lambda^*)$	WARP
1 ItalyPowerD	0.045	0.028	0.078	0.050	0.055	0.028	0.028	0.030	0.028 (30)	0.037 (10)	×
2 CinCECGtorso	0.103	0.367	0.167	0.349	0.367	0.094	0.094	0.093	0.092 (0)	0.079 (0)	×
3 ECG200	0.120	0.070	0.160	0.230	0.190	0.070	0.070	0.070	0.080 (0)	0.080 (0)	×
4 SonyAIBOII	0.141	0.142	0.128	0.169	0.194	0.142	0.142	0.144	0.155 (0)	0.131 (0)	×
5 Coffee	0.250	0.000	0.357	0.179	0.143	0.000	0.000	0.071	0.000 (0)	0.000 (10)	×
6 ECG5Days	0.203	0.153	0.006	0.232	0.236	0.203	0.203	0.203	0.007 (10)	0.024 (0)	×
7 SonyAIBO	0.305	0.308	0.258	0.275	0.343	0.308	0.308	0.293	0.188 (0)	0.165 (30)	×
8 Adiac	0.389	0.297	0.261	0.396	0.338	0.373	0.363	0.402	0.358 (0)	0.361 (0)	×
9 Beef	0.467	0.300	0.500	0.500	0.500	0.367	0.267	0.467	0.033 (0)	0.267 (0)	×
10 Trace	0.240	0.240	0.140	0.000	0.000	0.000	0.000	0.000	0.000 (0)	0.010 (0)	✓
11 CBF	0.148	0.140	0.382	0.003	0.000	0.000	0.000	0.000	0.031 (30)	0.003 (0)	✓
12 CC	0.120	0.113	0.383	0.007	0.027	0.007	0.007	0.007	0.003 (0)	0.007 (0)	✓
13 UMDsmooth	0.023	0.017	0.094	0.011	0.016				0.000 (0)	0.003 (0)	✓
14 BMEsmooth	0.029	0.016	0.090	0.044	0.045				0.000 (0)	0.247 (0)	✓
15 DiatomSizeR	0.065	0.076	0.069	0.033	0.029	0.033	0.033	0.042	0.026 (0)	0.029 (0)	✓
16 Symbols	0.101	0.111	0.080	0.050	0.043	0.051	0.050	0.052	0.034 (30)	0.046 (30)	✓
17 GunPoint	0.087	0.113	0.027	0.093	0.027	0.027	0.027	0.040	0.020 (10)	0.040 (10)	✓
18 FacesUCR	0.231	0.227	0.175	0.095	0.102	0.098	0.098	0.099	0.068 (10)	0.059 (0)	✓
19 TwoLeadECG	0.253	0.153	0.103	0.096	0.008	0.005	0.005	0.018	0.006 (0)	0.016 (10)	✓
20 MoteStrain	0.121	0.263	0.278	0.165	0.171	0.260	0.248	0.188	0.185 (0)	0.153 (10)	✓
21 Lighting2	0.246	0.246	0.148	0.131	0.213	0.131	0.131	0.131	0.148 (0)	0.131 (0)	✓
22 OliveOil	0.133	0.133	0.167	0.200	0.100	0.133	0.133	0.133	0.167 (0)	0.100 (10)	✓
23 FISH	0.217	0.149	0.229	0.166	0.137	0.109	0.137	0.126	0.149 (0)	0.240 (0)	✓
24 FaceFour	0.216	0.216	0.239	0.170	0.136	0.170	0.170	0.170	0.000 (0)	0.034 (0)	✓
25 SwedishLeaf	0.211	0.186	0.146	0.208	0.109	0.115	0.110	0.125	0.110 (0)	0.114 (0)	✓
26 MedicalImages	0.316	0.313	0.345	0.263	0.290	0.263	0.263	0.263	0.237 (0)	0.236 (10)	✓
27 Lighting7	0.425	0.411	0.316	0.274	0.288	0.342	0.356	0.342	0.397 (0)	0.233 (0)	✓
28 PowerCons	0.366	0.445	0.315	0.397	0.401	0.401	0.401	0.401	0.318 (0)	0.308 (0)	✓
29 OSULeaf	0.484	0.475	0.426	0.409	0.265	0.264	0.264	0.322	0.380 (0)	0.376 (0)	✓
30 InlineSkate	0.658	0.658	0.675	0.616	0.623	0.605	0.605	0.602	0.733 (10)	0.625 (0)	✓

Table 4.4: 1-NN test error rates for standard, *a priori* combined and M²TML measures.

Michèle
trouve
que ce
serait
bien de
ranger
les
datasets
par
warp/non-
warp et
de classer
des
moins
challenge
au plus
challenge
pour
aider la
lecture.

4.3.2 Comparison of the classification performances on the test set

From Table 4.4, we can see first that the 1-NN classification reaches the best results in:

1. Less than one-third of the data when based on unimodal metrics d_A , d_B or d_F
2. Slightly more than one-third for unimodal metrics DTW and d_{B-DTW}
3. Two-thirds (19 - 20 times on 30) when based on *a priori* combined metrics D_{Lin} , D_{Geom} and D_{Sig}
4. More than two-thirds (23 times on 30) when based on learned metrics D or $D_{\mathcal{H}}$.

Particularly, note that for nearly all datasets for which an uni-modal metric succeeds, the M^2TML metrics succeed similarly or lead to equivalent results. However, for several challenging datasets (*e.g.* FaceFour, Beef, FaceUCR, SonyAIBO, BME), M^2TML realizes drastic improvements, to the best of our knowledge never achieved before for these challenging public data. For instance, the impressive scores of 3% obtained for Beef against an error rate varying from 30% to 50% for alternative metrics, and of 0% obtained for FaceFour v.s. 13% to 23% for alternative metrics. Finally, D and $D_{\mathcal{H}}$ are all the more outperforming if only compared to the standard metrics d_A (the Euclidean distance) and DTW.

If we compare the *a priori* combined metrics (D_{Lin} , D_{Geom} , D_{Sig}) based on only the unimodal metrics involved in the combination (either d_A and d_B or DTW and d_{B-DTW}), we observe that *a priori* combined metrics achieved on two-third of the data with an equivalent or better score. Compared to the learnt metrics (D , $D_{\mathcal{H}}$), the results are globally similar except for 8 datasets where the learned metrics perform better (FaceFour, Beef, ECG5Days, FaceUCR, SonyAIBO, PowerCons, BME, UMD) and one where the *a priori* combined metrics perform better (OSULeaf). Note that the combined metric D_{Sig} is limited to two components and can't be easily extend to other metrics in its combination. D_{Lin} and D_{Geom} could be easily extended and a proposition could be:

$$D_{Lin}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{h=1}^p \alpha_h d_h(\mathbf{x}_i, \mathbf{x}_j) \quad (4.1)$$

$$D_{Geom}(\mathbf{x}_i, \mathbf{x}_j) = \prod_{h=1}^p \alpha_h d_h(\mathbf{x}_i, \mathbf{x}_j) \quad (4.2)$$

However, by considering p metrics d_h the resulting models requires to optimize p parameters. The grid search to find the best parameters α_h can become time consuming.

In the second part, we perform a graphical analysis for a global comparison on the whole datasets. In Fig. 4.2-a, each dataset is projected according to, on the x-axis its best error rate obtained for D and $D_{\mathcal{H}}$, and on y-axis its best performance w.r.t the standard metrics d_A and DTW. In Fig. 4.2-b, the y-axis is related to the best error rate w.r.t DTW and d_{B-DTW} ,

the two most performant uni-modal metrics. For both plots we can note that the datasets are principally projected above the first bisector, indicating higher error rates mostly obtained for alternative metrics than for M^2TML . For the less challenging datasets (near the origin of each graph), although almost projected near the bisector denoting equal performances for the compared metrics, M^2TML still bring improvements with projections clearly positioned above the bisector. Finally, from Fig. 4.2 (b) we can see that M^2TML metrics perform significantly lower than d_{B-DTW} on OSUleaf, while InlineSkate dataset remains challenging for all studied metrics.

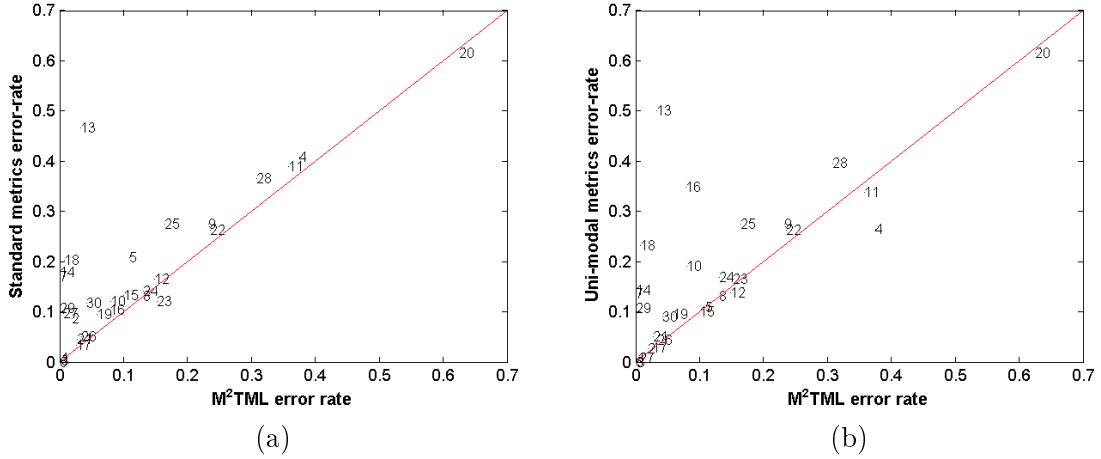


Figure 4.2: (a) Standard (Euclidean distance d_A and DTW) vs. M^2TML (D and $D_{\mathcal{H}}$) metrics. (b) Best Uni-modal (DTW and d_{B-DTW}) vs. M^2TML (D and $D_{\mathcal{H}}$) metrics.

4.3.3 Analysis of the discriminative features

For the learnt metric D , thanks to the L_1 regularization, the learned SVM reveals the features that most differentiate pull from push pairs. We recall that the weight for each feature can be analyzed through the weight vector \mathbf{w} obtained by learning the SVM classifier. Table 4.5 shows the sparse, muti-modal and multi-scale potential of M^2TML approach. It gives for each dataset, the weights of the top five 'discriminative' features that contribute to the definition of D . For instance, for FaceFour D reaches an error of 0% by combining, in the order of importance, the behavior d_{B-DTW} , frequential d_F and amplitude DTW modalities, at the global (I^0) and local (I^4 , I^5 , I^2) scales. For Beef, besides the impressive error rate of 3%, the learned model is very sparse as D involves only the behavior modality based on the segment I^3 (d_B^3). Similarly for Coffee, the obtained 0% involves only the behavior and frequential modalities at several scales. Note that if we look at only the most discriminative feature (1st column), the M^2TML method helps to localize discriminative modality and a specific temporal scale that could not be easily guessed *a priori* (e.g., Lightning7: behavior modality on the segment I_6 (d_{B-DTW}^6), OliveOil: frequential modality on the segment I_5 (d_F^5), TwoLeadECG: behavior modality on the segment I_4 (d_{B-DTW}^4)).

In Fig. 4.3, we plot the weights of all features for SonyAIBO, Beef, CincECGtorso and FaceFour cases as an example. It illustrates both the sparsity of the M^2TML approach (Beef, CincECGtorso and FaceFour) and the ability of the algorithm to combine all the features into the metric D (SonyAIBO). In particular, the approach is able to either select one single feature (Beef) or combine several selected features (CinCECGTorso, FaceFour). Fig. 4.4 illustrates the temporal locations of the most discriminative features for these datasets. Note that from looking at the temporal representation, it is not easy to determine *a priori* which modality (value, behavior, frequential) and at which temporal scale is the most discriminative feature to separate the classes.

In summary, we can emphasize that for almost all datasets, the definition of D involves no more than five features (the most contributive ones), that assesses not only the model's sparsity but also the representativeness of the revealed features.

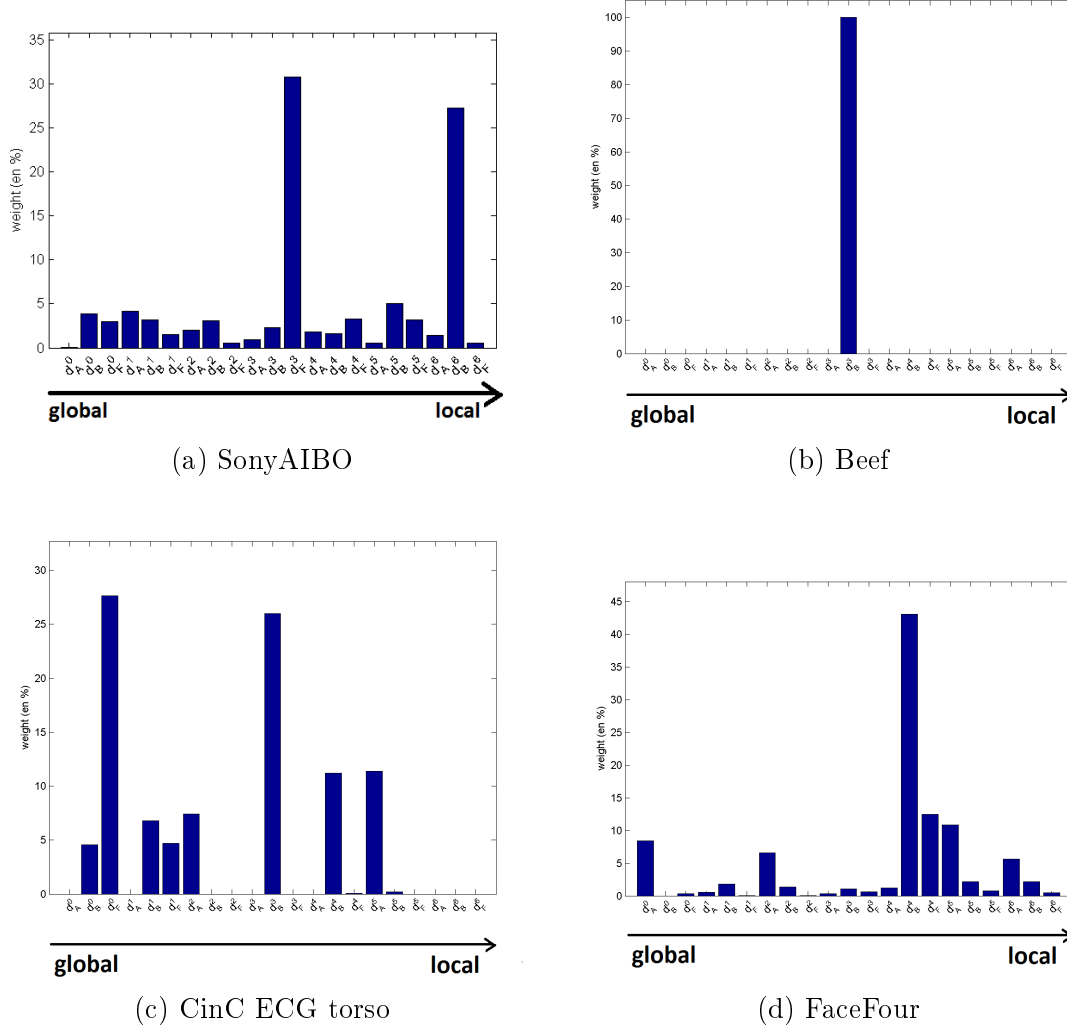


Figure 4.3: M^2TML feature weights for 4 datasets.

Dataset	Feature weights (%)				
CC	DTW ² (56.3%)	d_F^0 (18.8%)	d_F^4 (5.9%)	d_F^1 (4%)	d_{B-DTW}^0 (3.2%)
GunPoint	d_{B-DTW}^0 (42.1%)	DTW ⁵ (10.9%)	DTW ¹ (10.2%)	d_{B-DTW}^6 (8.4%)	d_F^4 (8.1%)
CBF	DTW ⁴ (56.5%)	d_F^3 (43.4%)	d_{B-DTW}^1 (0.2%)	DTW ⁰ (0%)	d_{B-DTW}^0 (0%)
OSULeaf	DTW ² (23.7%)	d_F^0 (19.5%)	d_{B-DTW}^0 (14.6%)	d_{B-DTW}^2 (9.4%)	DTW ¹ (9%)
SwedishLeaf	d_F^0 (21.5%)	DTW ⁰ (15.9%)	d_{B-DTW}^0 (15.2%)	DTW ⁶ (11.5%)	d_{B-DTW}^1 (6.1%)
Trace	DTW ⁰ (58.3%)	DTW ⁶ (6.9%)	d_{B-DTW}^0 (5.8%)	DTW ² (5.6%)	DTW ⁵ (5.5%)
FaceFour	d_{B-DTW}^4 (44.5%)	d_F^4 (12.7%)	DTW ⁵ (11.1%)	DTW ⁰ (8.3%)	DTW ² (6.4%)
Lighting2	d_{B-DTW}^0 (30.4%)	DTW ⁶ (18.7%)	d_F^1 (16.5%)	d_{B-DTW}^6 (13.4%)	DTW ⁰ (10.7%)
Lighting7	d_{B-DTW}^6 (87.4%)	d_F^6 (8.6%)	d_{B-DTW}^5 (4%)	-	-
ECG200	d_B^0 (89.6%)	d_B^2 (2.4%)	d_A^3 (2.3%)	d_B^1 (2.2%)	d_B^4 (2%)
Adiac	d_F^0 (79.2%)	d_F^4 (13.8%)	d_A^4 (3.5%)	d_F^5 (1.7%)	d_B^5 (1.2%)
FISH	d_{B-DTW}^5 (17.9%)	d_F^0 (10.5%)	d_{B-DTW}^6 (9.9%)	d_{B-DTW}^4 (8.3%)	d_{B-DTW}^3 (7.8%)
Beef	d_B^3 (100%)	-	-	-	-
Coffee	d_B^2 (22.4%)	d_F^4 (20.1%)	d_B^6 (14.6%)	d_B^0 (8.1%)	d_F^5 (7%)
OliveOil	d_F^5 (97%)	d_{B-DTW}^2 (3%)	-	-	-
CinCECGtorso	d_F^0 (38.4%)	d_A^5 (13.1%)	d_B^4 (11.5%)	d_F^1 (11.2%)	d_A^2 (9.8%)
DiatomSizeR	d_F^5 (39.1%)	d_F^0 (36%)	d_{B-DTW}^4 (24.9%)	-	-
ECG5Days	d_B^5 (59.5%)	d_B^6 (32.3%)	d_A^4 (3.9%)	d_B^2 (3.1%)	d_B^4 (1.2%)
FacesUCR	d_F^2 (21.5%)	d_{B-DTW}^0 (19.5%)	d_F^4 (16.7%)	DTW ⁰ (12.6%)	d_{B-DTW}^2 (8.6%)
InlineSkate	d_F^4 (42.5%)	DTW ⁵ (22.8%)	DTW ⁴ (17.6%)	DTW ² (6.7%)	d_{B-DTW}^6 (5.9%)
ItalyPowerD	d_B^6 (68.7%)	d_B^0 (25.9%)	d_B^3 (5.2%)	d_B^4 (0.2%)	-
MedicalImages	d_{B-DTW}^1 (53.3%)	d_F^3 (12.9%)	d_{B-DTW}^2 (10.7%)	d_{B-DTW}^3 (10.1%)	d_{B-DTW}^0 (3.8%)
MoteStrain	d_{B-DTW}^5 (93.2%)	d_{B-DTW}^6 (6.8%)	-	-	-
SonyAIBOII	d_B^3 (100%)	-	-	-	-
SonyAIBO	d_F^3 (30.8%)	d_B^6 (27.3%)	d_B^5 (5%)	d_A^1 (4.1%)	d_B^0 (3.9%)
Symbols	d_{B-DTW}^0 (45.6%)	d_{B-DTW}^6 (35.3%)	d_{B-DTW}^5 (19%)	DTW ⁰ (0.1%)	-
TwoLeadECG	d_{B-DTW}^4 (60%)	d_F^1 (12%)	DTW ⁴ (11.4%)	d_{B-DTW}^6 (7.6%)	d_{B-DTW}^1 (4.2%)
PowerCons	d_F^0 (26.1%)	DTW ⁰ (20.3%)	d_F^1 (19.3%)	d_{B-DTW}^0 (6.1%)	d_F^2 (5.1%)
BME	d_{B-DTW}^0 (75.2%)	d_F^4 (15.5%)	d_{B-DTW}^2 (5.8%)	d_{B-DTW}^1 (1.9%)	d_F^1 (0.7%)
UMD	d_{B-DTW}^0 (99.8%)	d_{B-DTW}^5 (0.2%)	-	-	-

Table 4.5: Top 5 multi-modal and multi-scale features involved in D

4.3.4 Effect on the neighborhood before and after learning

In the last part, we compare the global effect of the alternative and M^2TML metrics on the 1-NN neighborhood distribution and class discrimination. For that, a MultiDimensional Scaling (MDS)² is used to visualize the distribution of samples according to their pairwise dissimilarities. Briefly, MDS is a method of visualizing the proximity between samples in a dataset. A

²Matlab function: mdscale for metrics and non metrics

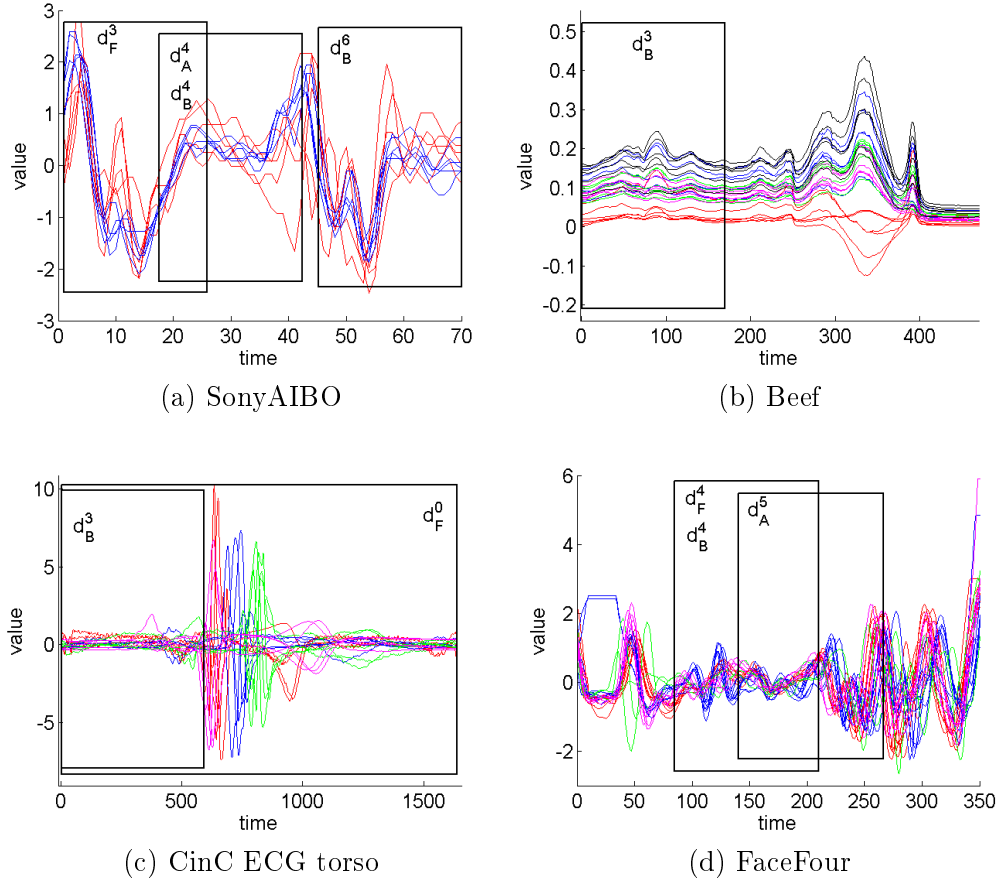


Figure 4.4: Temporal representation of the top M^2TML feature weights for 4 datasets.

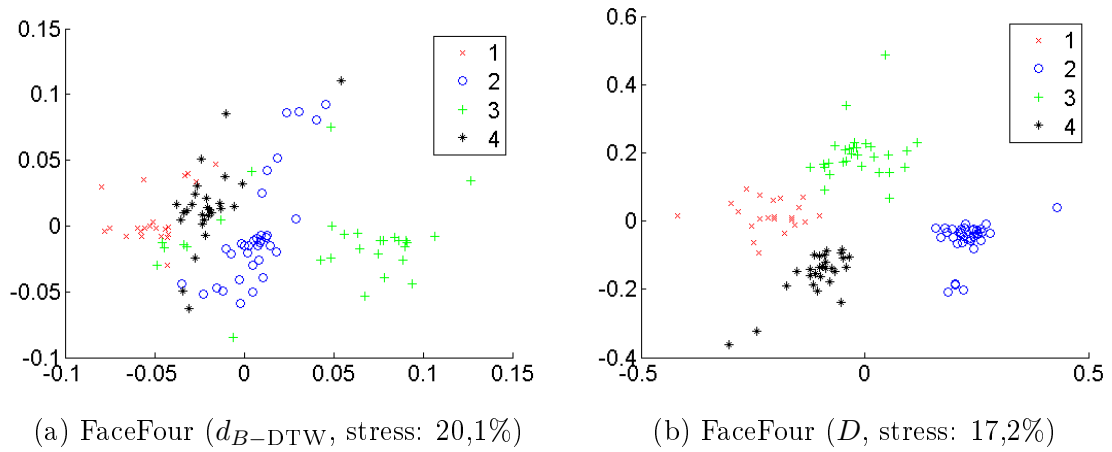


Figure 4.5: MDS visualization of the d_{B-DTW} (Fig. a) and D (Fig. b) dissimilarities for FaceFour

MDS algorithm aims to place each sample in P -dimensional space (in general, $P = 2$ or 3) such that the between-object distances are preserved as well as possible. Classical MDS takes an input matrix giving dissimilarities between pairs of samples and outputs a coordinate for each sample whose configuration minimizes a loss function called stress. Note that the MDS representation has no link with the dissimilarity space representation whose dimensions are basic temporal metrics.

For FaceFour, Fig. 4.5 shows the first obtained plans and their corresponding stresses, the classes being indicated in different symbols and colors. We can see distinctly the effect of the learned D that leads to more compact and more isolated classes with robust neighborhoods for 1-NN classification (*i.e.* closer pull pairs and far away push pairs) than the best alternative metric d_{B-DTW} that shows more overlapping classes and heterogeneous neighborhoods.

4.4 Conclusion of the chapter

The large conducted experiments and the impressive performances obtained attest the efficiency of the learned M^2TML metrics for time series nearest neighbors classification. Finally, let us underline the merit of the M^2TML solution, that not only leads to equivalent or better performances from the standard metrics (Euclidean distance, Dynamic time warping), but also provides a comprehensive and fine-grained information about which modalities are mostly discriminant, how they should be combined and precisely at which temporal granularity.

Conclusion and perspectives

Conclusion

In this PhD, we have proposed a framework to learn a combined multi-modal and multi-scale temporal metric (M^2TML) for a robust k -NN classifier. It is based on a new space representation, the dissimilarity space, where pairs of time series are embedded as vectors described by different basic temporal metrics. A metric combining the basic metrics can be seen as a function of the dissimilarity space, learned by using a large margin optimization process (SVM) inspired from the nearest neighbors metric learning framework. The obtained metric satisfies the properties of a dissimilarity, leads to good performances on a large number of public datasets, and gives an interpretable solution that allows to analyze the modalities and scales that are the most discriminant.

Temporal data may be compared based on various characteristics or modalities. They may be compared not only on their amplitudes like static data, but also on other modalities such that their behavior, frequency, etc. Some authors propose to combine several modalities but are either, limited to two modalities or in case of multiple modalities (more than 2), the number of parameters to optimize for the classifier may become time consuming. In general, state of the art approaches compared the time series by involving all observations, restricting the potential of comparison measures (metrics) to capture local differences. Real time series can be also subjected to delays. We believe that all of these considerations (modality, scale, delays) should be taken into consideration in the definition of a metric in order to improve the performance of the classifier.

In the second part of this work, we propose a general formalization of the problem of learning a combined multimodal and multiscale temporal metric (M^2TML) for a robust k -NN. Based on a pairwise dissimilarity representation of the pairs of time series, the metric learning can be reduced to the learning of a linear or non linear function of the dissimilarity space that satisfies the properties of a dissimilarity (positivity, reflexivity, symmetry). Inspired from metric learning work, the problem is formalized as an optimization problem involving a regularization and loss term which aims to pull samples that are expected to be similar and push away samples that should be dissimilar. By considering a linear combination of the basic metrics, changing the regularization term leads the general formalization to a linear and quadratic formalization. The latter allows to extend to the learning of non-linear functions thanks to the "kernel" trick. However, the methods can lead to functions that doesn't meet the properties of a dissimilarity (non-positivity). Secondly, we formulate the problem as a SVM problem which aims to separate pull and push samples, then we define a metric that satisfy the required properties.

The efficiency of the proposed SVM-based solution has been tested in the case of classification of univariate time series, on a wide variety of datasets coming from various fields

(simulated data, medicine, power consumption, etc.), diverse sizes of training and testing, various number of classes, etc. The M²TML solution achieves not only, either equivalent or better performances compared to the standard global metrics (euclidean distance, dynamic time warping, temporal correlation, Fourier distance), but it also provides a sparse and interpretable solution that allows to give a comprehensive analysis of the most discriminative modalities and their respective temporal granularity that may not be always intuitive *a priori*.

Perspectives

Extension to other modalities, multivariate problems and other type of data

In this work, we focus on three basic temporal metrics (euclidean distance, temporal correlation, Fourier-based distance). Montero & Vilar propose in [MV14] a review on a wide number of metrics dedicated to time series. For remaining challenging datasets in our experiments, it could be interesting to integrate other basic temporal metrics in our framework to see the obtained results.

The framework can be easily extend to multivariate problem. For each dimension, we consider the set of multimodal and multiscale description. Then, we consider the union over the dimensions as our new pairwise dissimilarity description d_1, \dots, d_p .

The proposed framework has been tested in the case of time series data but is more general. It can be applied to any other type of data (strings, graphs, images) to learn a combined metric. Deza & Deza makes a detailed review of metrics for various domains in [DD09].

Multiscale description

A second improvement is about the multiscale description. Our multiscale approach is based on a binary segmentation using a dichotomy process. Other solutions could be proposed to localize for finely event of interest. For example, in the case of the dataset SonyAIBO, the discriminative temporal location of the signal is known *a priori* (Fig. 4.6). With the actual multiscale description, it is not possible to extract exactly the two red patterns of interest. A solution based on a sliding window of variable lengths could be used to locate precisely these patterns.

Learning of local metrics

Some authors [WS09b]; [WWK12] suggest that in some datasets, global linear metric learning approach are not sufficient to improve the accuracy of k -NN classification. However, since the discriminatory power of the input features might vary between different neighborhoods,

³source: <http://www.cs.ucr.edu/~eamonn/LogicalShapelet.pdf>

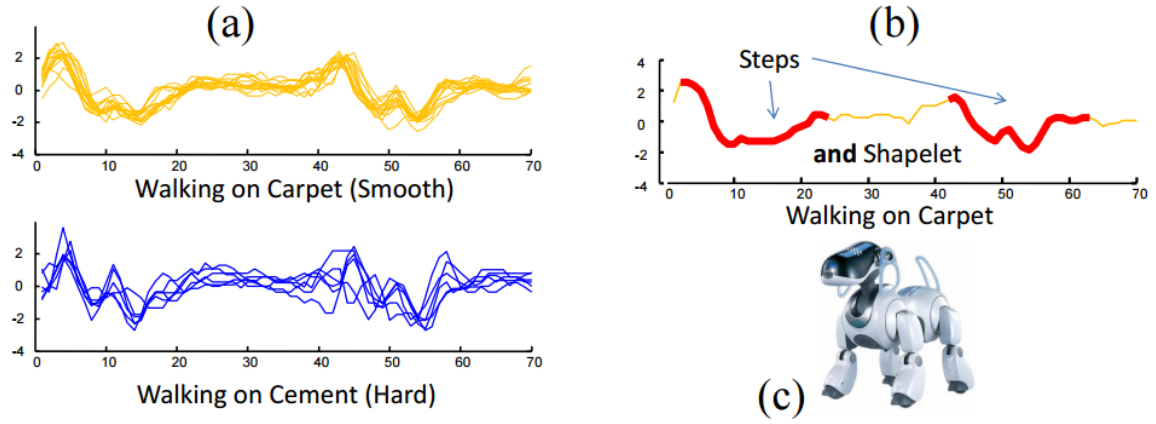


Figure 4.6: (a) Two classes of time series from the Sony AIBO accelerometer. (b) The and-shapelets from the walk cycle on carpet. (c) The Sony AIBO Robot.³

learning a global metric cannot fit well the distance over the data manifold. To overcome this difficulty, they propose to learn a metric on each neighborhood, referred as local metric learning.

Similarly, our M^2TML framework could be extended to learn local combined metric for each neighborhood. The objective is to learn for each n set $Pull_i$ and $Push_i$ (n being the number of samples in the training set) a local metric using the same framework than the one we propose in this work. We obtain n local metrics D_i . Then, to classify a new sample \mathbf{x}_{test} , we compute the n metrics $D_i(\mathbf{x}_{i,test})$ and classify \mathbf{x}_{test} using the k lowest $D_i(\mathbf{x}_{i,test})$.

Re-iteration of the initial metric

Similarly to Large Margin Nearest Neighbors (LMNN) approach proposed by Weinberger & Saul [WS09b], our approach may inherit the same problem of defining the set $Pull_i$ and $Push_i$ according to an initial distance (L2-norm in our work). Other initial distance could have been used. If the initial distance is far away from the optimal solution, the definition of the sets $Pull_i$ and $Push_i$ can impact the convergence to the optimal solution. In same spirit as the multi-pass (LMNN) approach proposed by Weinberger & Saul, we could re-iterate the learning process. At each step, we re-define the sets $Pull_i$ and $Push_i$ using the distance learned at the previous step. Then, we stop the learning when arriving at convergence (*e.g.*, the sets $Pull_i$ and $Push_i$ doesn't evolve anymore between two steps).

Other proposition to define the combined metric

Proposition de Sylvain

Extension to regression problems

For the SVM-based solution, in the pairwise dissimilarity space, each vector \mathbf{x}_{ij} is labeled y_{ij} by following the rule: if \mathbf{x}_i and \mathbf{x}_j are similar, the vector \mathbf{x}_{ij} is labeled -1; and +1 otherwise. For classification problems, the concept of similarity between samples \mathbf{x}_i and \mathbf{x}_j is driven by the class label y_i and y_j in the original space:

$$y_{ij} = \begin{cases} +1 & \text{if } y_i = y_j \\ -1 & \text{if } y_i \neq y_j \end{cases} \quad (4.3)$$

For regression problems, each sample \mathbf{x}_i is assigned to a continuous value y_i . Two approaches are possible to define the similarity concept. The first one discretizes the continuous space of values of the labels y_i to create classes. One possible discretization bins the label y_i into Q intervals as illustrated in Fig. 4.7. Each interval becomes a class which associated value can be set for example as the mean or median value of the interval. Then, the classification framework is used to define the pairwise label y_{ij} .

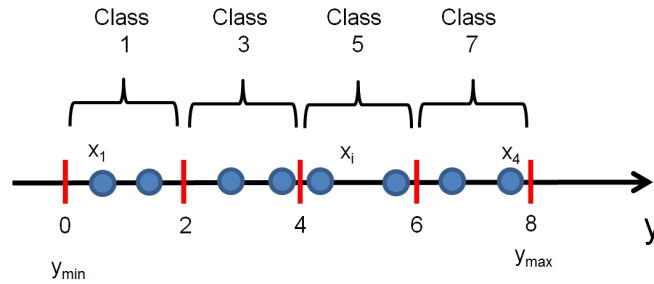


Figure 4.7: Example of discretization by binning a continuous label y into $Q = 4$ equal-length intervals. Each interval is associated to a unique class label. In this example, the class label for each interval is equal to the mean in each interval.

This approach may leads to border effects between the classes. For instance, two samples \mathbf{x}_i and \mathbf{x}_j that are close to a frontier and that are on different sides of the border will be considered as different, as illustrated in Fig 4.8. Moreover, a new sample \mathbf{x}_j will have its labels y_j assigned to a class and not a real continuous value.

The second approach considers the continuous value of y_i , computes a L_1 -norm between the labels $|y_i - y_j|$ and compare this value to a threshold ϵ . Geometrically, a tube of size ϵ around each value of y_i is built. Two samples \mathbf{x}_i and \mathbf{x}_j are considered as similar if the absolute difference between their labels $|y_i - y_j|$ is lower than ϵ (Fig. 4.9):

$$y_{ij} = \begin{cases} -1 & \text{if } |y_i - y_j| \leq \epsilon \\ +1 & \text{otherwise} \end{cases} \quad (4.4)$$

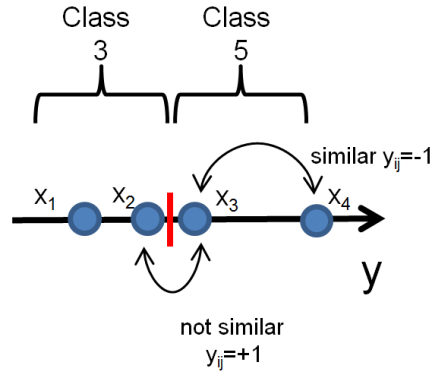


Figure 4.8: Border effect problems. In this example, x_2 and x_3 have closer value labels y_2 and y_3 than x_3 and x_4 . However, with the discretization x_2 and x_3 don't belong to the same class and thus are considered as not similar.

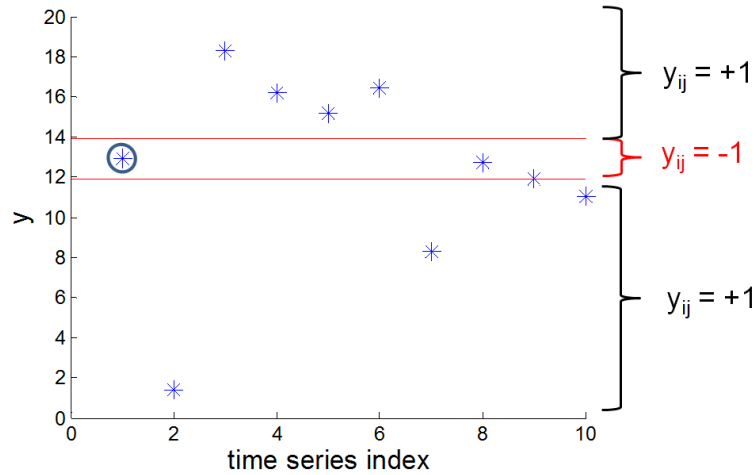


Figure 4.9: Example of pairwise label definition using an ϵ -tube (red lines) around the time series x_i (circled in blue). For, time series x_j that falls into the tube, the pairwise label is $y_{ij} = -1$ (similar) and outside of the tube, $y_{ij} = +1$ (not similar).

Using the learned combined metric in other algorithms

- Perspectives

- **Extension à d'autres modalités:** Le travail effectué s'est intéressé à 3 types de distance (distance euclidienne, corrélation temporelle, distance à base Fourier). Il existe cependant d'autres mesures de distance pour les séries temporelles qu'il serait intéressant d'étudier (wavelets, MFCC, etc.) (cf article).
- **Découpage multi-échelle:** Dans ce travail, nous avons proposé de découper les séries temporelles avec une architecture dichotomique mais d'autres solutions pourraient être proposées: par exemple une fenêtre glissante qui permettrait de localiser plus finement les événements d'intérêt (prendre l'exemple de SonyAIBO comme le

suggère Sylvain).

- **Multi-pass learning:** le framework proposé souffre du même inconvénient que celui de Weinberger. La distance initiale est une norme 2 dans l'espace des paires. Elle permet de définir les voisinages pour la construction des ensembles pull et push qui reste fixe pendant le processus d'optimisation. Une autre solution serait de ré-itérer la distance apprise pour re-définir les ensembles pull et push et ré-apprendre de manière itérative la métrique jusqu'à convergence.
- **Autre proposition pour D :** La proposition de définition de la métrique D pourrait être élargie. Par exemple, on pourrait penser à une solution qui permet d'élargir à des lambdas négatifs. Inconvénient: avec un lambda négatif, la métrique aurait un effet pull et non push, qui risque de binariser la métrique en incluant des distances très proches de 0 pour les plus proches voisins. Or dans une classification kNN, la distance des plus proches voisins nous permet de faire la classification. En ramenant les distances des plus proches voisins, proches de 0, on a un risque de moins bien discriminer.
- **Extension au multivarié:** on pourrait facilement étendre le framework pour des séries temporelles multivariées. Pour chaque dimension de la série multivariée, on calcule les métriques de bases de chaque échelle qui deviennent de nouvelles dimensions basiques pour notre algorithme.
- **Extension à la régression:** Le travail pourrait être étendu à des problèmes de régression. Le seul changement qui est impliqué est la définition du label y_{ij} lorsque le label y_i est une valeur continue et non une classe (cf la proposition).
- **Apprentissage locale de la métrique:** Dans le cadre du metric learning, des auteurs se sont intéressés à l'apprentissage de métrique localisés et non globales. On pourrait penser apprendre une métrique par localité (voisinage) et combinée l'ensemble des métriques locales en une métrique globale.
- **Utilisation de la distance apprise dans d'autres algorithmes:** Dans l'industrie, il est souvent important d'avoir un modèle interprétable. Des algorithmes comme les Arbres de décision apporte une représentation visuelle qui permet de comprendre aisément le modèle appris. Des travaux sur les arbres temporelles ont été effectués dans la littérature (cité papier ahlame) utilisant les métriques classiques (euclidienne, corrélation). On pourrait s'inspirer de ces travaux pour construire un arbre avec la métrique apprise par m2tml.

List of articles

Journal

- C. Do, A. Douzal-Chouakria, S. Marié, M. Rombaut. Multi-modal and Multi-scale Temporal Metric Learning for Time Series Nearest Neighbors Classification, Pattern Recognition Letters 2016 (under submission)

-

papier
Springer?

Conference

- C. Do, A. Douzal-Chouakria, S. Marié, M. Rombaut. Multiple Metric Learning for large margin kNN Classification of time series, EUSIPCO'2015
- C. Do, A. Douzal-Chouakria, S. Marié, M. Rombaut. Temporal and Frequential Metric Learning for Time Series kNN Classification. ECML-PKDD'2015 Workshop on Advanced Analytics and Learning from Temporal Data, 39-45

Detailed presentation of the datasets

Solver library

SVM library

QP resolution

SVM equivalency

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 ¹:

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

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. E.1 for \mathbf{w} and b solves a similar MLD problem in Eq. ?? 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 :

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

$$p_i^- = \frac{L}{2} = \frac{1}{2} \sum_l \frac{1 + y_{il}}{2} \quad (\text{E.3})$$

E.0.0.a Similarities and differences in the constraints

First, we recall the SVM constraints in Eq. E.1:

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

These constraints can be split into two sets of constraints:

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

¹the SVM formulation below divides the loss part into two terms similarly to asymmetric SVM

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:

$$\left\{ \begin{array}{l} \bullet \forall i, j, k, l \text{ such that } y_{ij} = -1, \text{ and } y_{kl} = +1, i \neq j \text{ and } i \neq k : \\ D(\mathbf{x}_k, \mathbf{x}_l) - D(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \frac{\xi_{kl} + \xi_{ij}}{2} \\ \bullet \forall i, j, l \text{ such that } y_{ij} = -1, \text{ and } y_{il} = +1, i \neq j : \\ D(\mathbf{x}_i, \mathbf{x}_l) - D(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \frac{\xi_{il} + \xi_{ij}}{2} \end{array} \right. \quad (\text{E.4})$$

By defining $\xi_{ijl} = \frac{\xi_{ij} + \xi_{il}}{2}$, the second constraint in Eq. E.4 from the MDL formulation is the same as the constraints in the SVM formulation in Eq. ??.

However, an additional set of constraints is present in the SVM formulation (first set of constraints in Eq. E.4) 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 E.1). This is more constraining in the SVM resolution for the resulting metric D especially if the neighborhoods have different spread.

E.0.0.b Similarities and differences in the objective function

Mathematically, from Eq. E.2, we write:

$$\begin{aligned} \sum_{i,l,y_{il}=+1} p_i^+ \xi_{il} &= \sum_{il} p_i^+ \frac{1+y_{il}}{2} \xi_{il} \\ &= \sum_{il} \left(\sum_{j \rightsquigarrow i} \frac{1}{2} \right) \frac{1+y_{il}}{2} \xi_{il} \\ &= \frac{1}{2} \sum_{i,j \rightsquigarrow i,l} \frac{1+y_{il}}{2} \xi_{il} \end{aligned} \quad (\text{E.5})$$

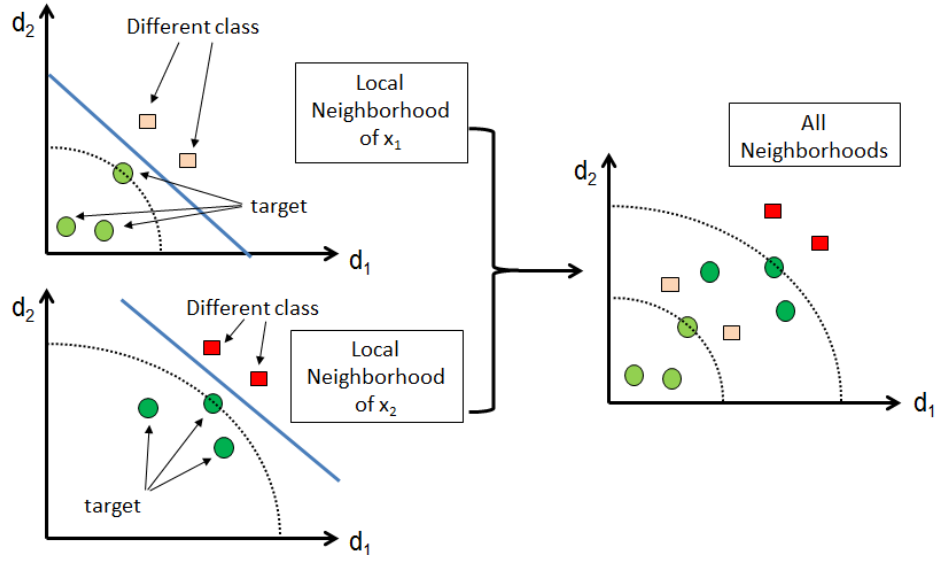


Figure E.1: Geometric representation of the neighborhood of $k = 3$ for two time series \mathbf{x}_1 and \mathbf{x}_2 (left). For each neighborhood, time series of different class are represented by a square and the margin by a blue line. Taking each neighborhood separately, the problem is linearly separable (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 .

And from Eq. E.3, we write:

$$\begin{aligned}
 \sum_{i,j,y_{ij}=-1} p_i^- \xi_{ij} &= \sum_{i,j \rightsquigarrow i} p_i^- \xi_{ij} \\
 &= \sum_{i,j \rightsquigarrow i} \left(\frac{1}{2} \sum_l \frac{1+y_{il}}{2} \right) \xi_{ij} \\
 &= \frac{1}{2} \sum_{i,j \rightsquigarrow i,l} \frac{1+y_{il}}{2} \xi_{ij}
 \end{aligned} \tag{E.6}$$

By replacing Eqs. E.5 and E.6 back into Eq. E.1, the objective function becomes:

$$\begin{aligned}
 \min_{\mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i,j \rightsquigarrow i,l} \frac{1+y_{il}}{2} \frac{\xi_{ij} + \xi_{il}}{2} \\
 \min_{\mathbf{w}, \xi} \quad & \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w}}_{\text{Regularization}} + C \underbrace{\sum_{i,j \rightsquigarrow i,l} \frac{1+y_{il}}{2} \xi_{ijl}}_{\text{Loss}}
 \end{aligned} \tag{E.7}$$

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. E.7), the regularization part tends

to minimize the norm of \mathbf{w} whereas in MLD (Eq. ??), 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.

Bibliography

- [ABR64] M. Aizerman, E. Braverman, and L. Rozonoer. “Theoretical foundations of the potential function method in pattern recognition learning.” In: *Automation and Remote Control* 25 (1964), pp. 821–837 (cit. on p. 19).
- [Alt92] Ns Altman. “An introduction to kernel and nearest-neighbor nonparametric regression.” In: *The American Statistician* 46.3 (1992), pp. 175–185 (cit. on p. 14).
- [AT10] Z. Abraham and P.N. Tan. “An Integrated Framework for Simultaneous Classification and Regression of Time-Series Data.” In: *ACM SIGKDD*. 2010 (cit. on p. 35).
- [BC94] Donald Berndt and James Clifford. “Using dynamic time warping to find patterns in time series.” In: *Workshop on Knowledge Knowledge Discovery in Databases* 398 (1994), pp. 359–370 (cit. on pp. 37, 38).
- [Ben+09] J. Benesty et al. “Pearson correlation coefficient.” In: *Noise Reduction in Speech Processing* (2009) (cit. on p. 35).
- [BGV92] Bernhard E. Boser, Isabelle M. Guyon, and Vladimir N. Vapnik. “A Training Algorithm for Optimal Margin Classifiers.” In: *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*. 1992, pp. 144–152 (cit. on pp. 14, 19).
- [BHS12] Aurélien Bellet, Amaury Habrard, and Marc Sebban. “Good edit similarity learning by loss minimization.” In: *Machine Learning* 89.1-2 (2012), pp. 5–35 (cit. on pp. 50, 51).
- [BHS13] Aurélien Bellet, Amaury Habrard, and Marc Sebban. “A Survey on Metric Learning for Feature Vectors and Structured Data.” In: *arXiv preprint arXiv:1306.6709* (2013), p. 57. arXiv: 1306.6709 (cit. on p. 51).
- [Bis06] Christopher M Bishop. *Pattern Recognition and Machine Learning*. Vol. 4. 4. 2006, p. 738. arXiv: 0-387-31073-8 (cit. on pp. 6, 26).
- [BM67] E. O. Brigham and R. E. Morrow. “The fast Fourier transform.” In: *Spectrum, IEEE* 4.12 (1967), pp. 63 –70 (cit. on p. 34).
- [BMP02] Serge Belongie, Jitendra Malik, and Jan Puzicha. “Shape Matching and Object Recognition Using Shape Contexts.” In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 24 (2002), pp. 509–522 (cit. on p. 14).
- [CH67] T. Cover and P. Hart. “Nearest neighbor pattern classification.” In: *IEEE Transactions on Information Theory* 13.1 (1967), pp. 21–27 (cit. on p. 13).
- [Cha+10] Ratthachat Chatpatanasiri et al. “A new kernelization framework for Mahalanobis distance learning algorithms.” In: *Neurocomputing* 73.10-12 (2010), pp. 1570–1579 (cit. on p. 51).
- [Cha04] Christopher Chatfield. *The analysis of time series : an introduction*. 2004, xiii, 333 p. (Cit. on p. 30).

- [CHL05] Sumit Chopra, Raia Hadsell, and Yann LeCun. “Learning a similarity metric discriminatively, with application to face verification.” In: *CVPR*. Vol. 1. 2005, pp. 539–546 (cit. on p. 43).
- [CHY96] Ming Syan Chen, Jiawei Han, and Philip S. Yu. *Data mining: An Overview from a Database Perspective*. 1996 (cit. on p. 6).
- [Coc77] William C Cochran. “Snedecor G W & Cochran W G. Statistical methods applied to experiments in agriculture and biology. 5th ed. Ames, Iowa: Iowa State University Press, 1956.” In: *Citation Classics* 19 (1977), p. 1 (cit. on p. 11).
- [CS01] Koby Crammer and Yoram Singer. “On the Algorithmic Implementation of Multiclass Kernel-based Vector Machines.” In: *Journal of Machine Learning Research* 2 (2001), pp. 265–292 (cit. on p. 24).
- [CT01] Lijuan Cao and Francis E H Tay. “Financial Forecasting Using Support Vector Machines.” In: *Neural Computing & Applications* (2001), pp. 184–192 (cit. on p. 31).
- [CV95] Corinna Cortes and Vladimir Vapnik. “Support-vector networks.” In: *Machine Learning* 20.3 (1995), pp. 273–297. arXiv: [arXiv:1011.1669v3](#) (cit. on pp. 14, 15).
- [CY11] Colin Campbell and Yiming Ying. *Learning with Support Vector Machines*. Vol. 5. 1. 2011, pp. 1–95 (cit. on pp. 14, 18).
- [DCA11] A. Douzal-Chouakria and C. Amblard. “Classification trees for time series.” In: *Pattern Recognition journal* (2011) (cit. on pp. 35, 40).
- [DCN07] A. Douzal-Chouakria and P. Nagabhushan. “Adaptive dissimilarity index for measuring time series proximity.” In: *Advances in Data Analysis and Classification* (2007) (cit. on p. 37).
- [DD09] Michel Marie Deza and Elena Deza. *Encyclopedia of Distances*. Vol. 2006. 2009, p. 590. arXiv: [0505065](#) [[arXiv:gr-qc](#)] (cit. on pp. 32, 86).
- [Den95] T. Denoeux. “A k-nearest neighbor classification rule based on Dempster-Shafer theory.” In: *IEEE Transactions on Systems, Man, and Cybernetics* 25.5 (1995), pp. 804–813 (cit. on p. 14).
- [DHB95] Thomas G. Dietterich, Hermann Hild, and Ghulum Bakiri. “A comparison of ID3 and backpropagation for English text-to-speech mapping.” In: *Machine Learning* 18.1 (1995), pp. 51–80 (cit. on p. 11).
- [Die97] T. Dietterich. “Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms.” In: (1997) (cit. on p. 11).
- [Din+08] Hui Ding et al. “Querying and Mining of Time Series Data : Experimental Comparison of Representations and Distance Measures.” In: *Proceedings of the VLDB Endowment* 1.2 (2008), pp. 1542–1552. arXiv: [1012.2789v1](#) (cit. on pp. 1, 14, 33).

- [Do+12] Huyen Do et al. “A metric learning perspective of SVM: on the relation of LMNN and SVM.” In: *Proceedings of the 15th International Conference on Artificial Intelligence and Statistics (AISTAS '12)* (2012), pp. 308–317. arXiv: [arXiv:1201.4714v1](#) (cit. on pp. 46, 47).
- [DP12] Robert P W Duin and Elbieta Pçkalska. “The dissimilarity space: Bridging structural and statistical pattern recognition.” In: *Pattern Recognition Letters* 33.7 (2012), pp. 826–832 (cit. on pp. ii, 52).
- [Dud76] Sahibsingh a. Dudani. “DISTANCE-WEIGHTED k-NEAREST-NEIGHBOR RULE.” In: *IEEE Transactions on Systems, Man and Cybernetics* SMC-6.4 (1976), pp. 325–327 (cit. on p. 14).
- [FCH08] RE Fan, KW Chang, and CJ Hsieh. “LIBLINEAR: A library for large linear classification.” In: *The Journal of Machine Learning* (2008) (cit. on pp. 21, 65, 76).
- [FDCG13] C. Frambourg, A. Douzal-Chouakria, and E. Gaussier. “Learning multiple temporal matching for time series classification.” In: *Intelligent Data Analysis*. 2013, pp. 198–209 (cit. on p. 51).
- [FDcG13] Cédric Frambourg, Ahlame Douzal-chouakria, and Eric Gaussier. “Apprentissage de couplages pour la discrimination de séries temporelles.” In: (2013) (cit. on p. 1).
- [Gol+04] Jacob Goldberger et al. “Neighbourhood Components Analysis.” In: *Advances in Neural Information Processing Systems* (2004), pp. 513–520 (cit. on p. 43).
- [HCL08] Chih-Wei Hsu, Chih-Chung Chang, and Chih-Jen Lin. “A Practical Guide to Support Vector Classification.” In: *BJU international* 101.1 (2008), pp. 1396–400. arXiv: [0-387-31073-8](#) (cit. on pp. 12, 24, 65, 76).
- [HHK12] Seok Hwan Hwang, Dae Heon Ham, and Joong Hoon Kim. “Forecasting performance of LS-SVM for nonlinear hydrological time series.” In: *KSCE Journal of Civil Engineering* 16.5 (2012), pp. 870–882 (cit. on p. 31).
- [HHP01] B Heisele, P Ho, and T Poggio. “Face recognition with support vector machines: global versus component-based approach.” In: *IEEE International Conference on Computer Vision, ICCV*. Vol. 2. July. 2001, pp. 688–694 (cit. on p. 14).
- [HWZ13] Jianming Hu, Jianzhou Wang, and Guowei Zeng. “A hybrid forecasting approach applied to wind speed time series.” In: *Renewable Energy* 60 (2013), pp. 185–194 (cit. on p. 31).
- [JJO11] Young-Seon Jeong, Myong K. Jeong, and Olufemi A. Omitaomu. “Weighted dynamic time warping for time series classification.” In: *Pattern Recognition* 44.9 (2011), pp. 2231–2240 (cit. on p. 51).
- [JMF99] a. K. Jain, M. N. Murty, and P. J. Flynn. “Data clustering: a review.” In: *ACM Computing Surveys* 31.3 (1999), pp. 264–323. arXiv: [arXiv:1101.1881v2](#) (cit. on p. 6).

- [Kal60] R E Kalman. “A New Approach to Linear Filtering and Prediction Problems.” In: *Transactions of the ASME Journal of Basic Engineering* 82.Series D (1960), pp. 35–45 (cit. on p. 35).
- [Keo+11] E. Keogh et al. *The UCR Time Series Classification/Clustering Homepage*. 2011 (cit. on p. 73).
- [KGG85] James M. Keller, Michael R. Gray, and James a. Givens. *A fuzzy K-nearest neighbor algorithm*. 1985 (cit. on p. 14).
- [KR04] Eamonn Keogh and Chotirat Ann Ratanamahatana. “Exact indexing of dynamic time warping.” In: *Knowledge and Information Systems* 7.3 (2004), pp. 358–386 (cit. on p. 38).
- [KU02] B Kijssirikul and N Ussivakul. “Multiclass Support Vector Machines using Adaptive Directed Acyclic Graph.” In: *Neural Networks, 2002. IJCNN '02. Proceedings of the 2002 International Joint Conference on* 1 (2002), pp. 980–985 (cit. on p. 24).
- [Lhe+11] S. Lhermitte et al. “A comparison of time series similarity measures for classification and change detection of ecosystem dynamics.” In: *Remote Sensing of Environment* 115.12 (2011), pp. 3129–3152 (cit. on p. 34).
- [Lia+12] Chunquan Liang et al. “Learning very fast decision tree from uncertain data streams with positive and unlabeled samples.” In: *Information Sciences* 213 (2012), pp. 50–67 (cit. on p. 31).
- [MV14] Pablo Montero and José Vilar. “TSclust : An R Package for Time Series Clustering.” In: *Journal of Statistical Software November* 62.1 (2014) (cit. on pp. 1, 33, 86).
- [Naj+12] H. Najmeddine et al. “Mesures de similarité pour l’aide à l’analyse des données énergétiques de bâtiments.” In: *RFIA*. 2012 (cit. on pp. 29, 33, 37).
- [Ngu+12] L. Nguyen et al. “Predicting collective sentiment dynamics from time-series social media.” In: *WISDOM*. 2012 (cit. on p. 30).
- [OE73] Richard O Duda and Peter E Hart. *Pattern Classification and Scene Analysis*. Vol. 7. 1973, p. 482 (cit. on pp. 6, 8, 14).
- [OS06] Jose Oncina and Marc Sebban. “Learning stochastic edit distance: Application in handwritten character recognition.” In: *Pattern Recognition* 39 (2006), pp. 1575–1587 (cit. on p. 51).
- [PAN+08] COSTAS PANAGIOTAKIS et al. “SHAPE-BASED INDIVIDUAL/GROUP DETECTION FOR SPORT VIDEOS CATEGORIZATION.” In: *International Journal of Pattern Recognition and Artificial Intelligence* 22.06 (2008), pp. 1187–1213 (cit. on p. 30).
- [PL12] Zoltán Prekopcsák and Daniel Lemire. “Time series classification by class-specific Mahalanobis distance measures.” In: *Advances in Data Analysis and Classification* 6.3 (2012), pp. 185–200. arXiv: 1010.1526 (cit. on p. 33).
- [PPD02] Elżbieta Pałkalska, Pavel Paclík, and Robert P W Duin. “A Generalized Kernel Approach to Dissimilarity-based Classification.” In: *Journal of Machine Learning Research* 2.2 (2002), pp. 175–211 (cit. on pp. ii, 52).

- [Ram+08] E. Ramasso et al. “Human action recognition in videos based on the transferable belief model : AAAApplication to athletics jumps.” In: *Pattern Analysis and Applications* 11.1 (2008), pp. 1–19 (cit. on p. 30).
- [Rey11] Miguel Reyes. “Feature weighting in dynamic time warping for gesture recognition in depth data.” In: *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on* (2011), pp. 1182–1188 (cit. on p. 51).
- [RJ93] L. Rabiner and B. Juang. *Fundamentals of Speech Recognition*. Vol. 103. 1993 (cit. on p. 38).
- [SC] Stan Salvador and Philip Chan. “FastDTW : Toward Accurate Dynamic Time Warping in Linear Time and Space.” In: () (cit. on p. 38).
- [SC78] H. Sakoe and S. Chiba. “Dynamic Programming Algorithm Optimization for Spoken Word Recognition.” In: *IEEE transactions on acoustics, speech, and signal processing* (1978) (cit. on p. 38).
- [She+02] Noam Shental et al. “Adjustment Learning and Relevant Component Analysis.” In: *European Conference on Computer Vision (ECCV)* 2353 (2002), pp. 776–790 (cit. on p. 43).
- [SJ89] B W Silverman and M C Jones. “E. Fix and J.L. Hodges (1951): An Important Contribution to Nonparametric Discriminant Analysis and Density Estimation: Commentary on Fix and Hodges (1951).” In: *International Statistical Review / Revue Internationale de Statistique* 57.3 (1989), pp. 233–238 (cit. on p. 13).
- [SS12a] M. Sahidullah and G. Saha. “Design, analysis and experimental evaluation of block based transformation in mfcc computation for speaker recognition.” In: *Speech Communication* 54.4 (2012), pp. 543–565 (cit. on p. 1).
- [SS12b] Md Sahidullah and Goutam Saha. “Design, analysis and experimental evaluation of block based transformation in MFCC computation for speaker recognition.” In: *Speech Communication* 54.4 (2012), pp. 543–565 (cit. on p. 34).
- [SS13] Bernhard Schlkopf and Alexander J. Smola. *Learning with Kernels*. Vol. 53. 2013, pp. 1689–1699. arXiv: [arXiv:1011.1669v3](https://arxiv.org/abs/1011.1669v3) (cit. on pp. 14, 20).
- [SSB03] Javad Sadri, Ching Y Suen, and Tien D. Bui. “Application of Support Vector Machines for recognition of handwritten Arabic/Persian digits.” In: *Second Conference on Machine Vision and Image Processing & Applications (MVIP 2003)* 1 (2003), pp. 300–307 (cit. on p. 14).
- [TC98] Christopher Torrence and Gilbert P. Compo. “A Practical Guide to Wavelet Analysis.” In: *Bulletin of the American Meteorological Society* 79.1 (1998), pp. 61–78 (cit. on p. 34).
- [Wan02] Jung-Ying Wang. “Support Vector Machines (SVM) in bioinformatics Bioinformatics applications.” In: *Bioinformatics* (2002), pp. 1–56 (cit. on p. 14).
- [WS09a] K. Weinberger and L. Saul. “Distance Metric Learning for Large Margin Nearest Neighbor Classification.” In: *Journal of Machine Learning Research* 10 (2009), pp. 207–244 (cit. on pp. 43–45, 51, 57).

- [WS09b] Kilian Q Weinberger and Lawrence K Saul. “Distance Metric Learning for Large Margin Nearest Neighbor Classification.” In: *Journal of Machine Learning Research* 10 (2009), pp. 207–244 (cit. on pp. 2, 50, 51, 86, 87).
- [WWK12] Jun Wang, Adam Woznica, and Alexandros Kalousis. “Parametric local metric learning for nearest neighbor classification.” In: *arXiv preprint arXiv:1209.3056* (2012), pp. 1–9 (cit. on p. 86).
- [Xi+06] Xiaopeng Xi et al. “Fast time series classification using numerosity reduction.” In: *Proceedings of the 23rd international conference on Machine learning (ICML)*. 2006, pp. 1033–1040 (cit. on p. 14).
- [XWC12] Zhixiang Xu, Kilian Q. Weinberger, and Olivier Chapelle. “Distance Metric Learning for Kernel Machines.” In: *Arxiv* (2012), pp. 1–17. arXiv: 1208.3422 (cit. on p. 51).
- [YG08] J. Yin and M. Gaber. “Clustering distributed time series in sensor networks.” In: *ICDM*. 2008 (cit. on p. 30).
- [YL99] Yiming Yang and Xin Liu. “A re-examination of text categorization methods.” In: *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval SIGIR 99*. 1999, pp. 42–49 (cit. on p. 14).
- [Zhu07] Xiaojin Zhu. “Semi-Supervised Learning Literature Survey.” In: *Sciences-New York* (2007), pp. 1–59 (cit. on p. 6).
- [ZLL14] X.-L. Zhang, Z.-G. Luo, and M. Li. “Merge-weighted dynamic time warping for speech recognition.” In: *Journal of Computer Science and Technology* 29.6 (2014), pp. 1072–1082 (cit. on p. 51).
- [ZY10] Yu Zhang and Dit-Yan Yeung. “Transfer metric learning by learning task relationships.” In: *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '10* (2010), p. 1199 (cit. on p. 51).
- [G. 06] S. Thiria G. Dreyfus, J.-M. Martinez, M. Samuelides M. B. Gordon, F. Badran. *Apprentissage Apprentissage statistique*. Eyrolles. 2006, p. 471 (cit. on pp. 6, 8).
- [Wie42] Wiener N. *Extrapolation, Interpolation & Smoothing of Stationary Time Series - With Engineering Applications*. Tech. rep. Report of the Services 19, Research Project DIC-6037 MIT, 1942 (cit. on p. 35).

Résumé — Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed non risus. Suspendisse lectus tortor, dignissim sit amet, adipiscing nec, ultricies sed, dolor. Cras elementum ultrices diam. Maecenas ligula massa, varius a, semper congue, euismod non, mi. Proin porttitor, orci nec nonummy molestie, enim est eleifend mi, non fermentum diam nisl sit amet erat. Duis semper. Duis arcu massa, scelerisque vitae, consequat in, pretium a, enim. Pellentesque congue. Ut in risus volutpat libero pharetra tempor. Cras vestibulum bibendum augue. Praesent egestas leo in pede. Praesent blandit odio eu enim. Pellentesque sed dui ut augue blandit sodales. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Aliquam nibh. Mauris ac mauris sed pede pellentesque fermentum. Maecenas adipiscing ante non diam sodales hendrerit. Ut velit mauris, egestas sed, gravida nec, ornare ut, mi. Aenean ut orci vel massa suscipit pulvinar. Nulla sollicitudin. Fusce varius, ligula non tempus aliquam, nunc turpis ullamcorper nibh, in tempus sapien eros vitae ligula. Pellentesque rhoncus nunc et augue. Integer id felis.

Mots clés : Série temporelle, Apprentissage de métrique, k -NN, SVM, classification, régression.

Abstract — The definition of a metric between time series is inherent to several data analysis and mining tasks, including clustering, classification or forecasting. Time series data present naturally several modalities covering their amplitude, behavior or frequential spectrum, that may be expressed with varying delays and at multiple temporal scales—exhibited globally or locally. Combining several modalities at multiple temporal scales to learn a holistic metric is a key challenge for many real temporal data applications. This paper proposes a Multi-modal and Multi-scale Temporal Metric Learning (M²TML) approach for maximum margin time series nearest neighbors classification. The solution refers to embedding time series into a dissimilarity space where a pairwise SVM is used to learn the metric. The M²TML solution is proposed for both linear and non linear contexts. A sparse and interpretable variant of the solution shows the ability of the learned temporal metric to localize accurately discriminative modalities as well as their temporal scales. A wide range of 30 public and challenging datasets, encompassing images, traces and ECG data, that are linearly or non linearly separable, are used to show the efficiency and the potential of M²TML for time series nearest neighbors classification.

Keywords: Time series, Metric Learning, k -NN, SVM, classification.

Schneider Electric
Université Grenoble Alpes, LIG
Université Grenoble Alpes, GIPSA-Lab

