

# Feature extraction over multiple representations for time series classification

Dominique Gay, Romain Guigourès, Marc Boullé, and Fabrice Clérot

Orange Labs

2, avenue Pierre Marzin, F-22307 Lannion Cedex, France

`firstname.name@orange.com`

**Abstract.** We suggest a simple yet effective and parameter-free feature construction process for time series classification. Our process is decomposed in three steps: *(i)* we transform original data into several simple representations; *(ii)* on each representation, we apply a coclustering method; *(iii)* we use coclustering results to build new features for time series. It results in a new transactional (i.e. object-attribute oriented) data set, made of time series identifiers described by features related to the various generated representations. We show that a Selective Naive Bayes classifier on this new data set is highly competitive when compared with state-of-the-art times series classification methods while highlighting interpretable and class relevant patterns.

## 1 Introduction

Time series classification (TSC) has been intensively studied in the past years. The goal is to predict the class of an object (a time series or a curve)  $\tau_i = \langle (t_1, x_1), (t_2, x_2), \dots, (t_{m_i}, x_{m_i}) \rangle$  (where  $x_k, (k = 1..m_i)$  is the value of the series at time  $t_k$ ), given a set of labeled training time series. TSC problems differ from traditional classification problems since there is a time dependence between the variables ; in other terms, the order of the variables is crucial in learning an accurate predictive model. The increasing interest in TSC is certainly due to the wide variety of applications: from e.g., medical diagnosis (like classification of patient electrocardiograms) to the maintenance of industrial machinery. Other domains, where data might be time series, are also concerned: finance, meteorology, signal processing, computer network traffic, . . . The diversity of applications has given rise to numerous approaches (see Section 4 for detailed related work). However, most efforts of the community have been devoted to the following three-step learning process: *(i)* choosing a new data representation, *(ii)* choosing a similarity measure (or a distance) to compare two time series and *(iii)* using the Nearest Neighbor (NN) algorithm as classifier on the chosen representation, using the chosen measure. Wang et al. [17] offer a survey of the various data representations and distances found in the literature and an extensive experimental study using the NN classifier. They conclude that NN classifier coupled with

Euclidean distance (ED) or Dynamic Time Warping (DTW) show the highest predictive performance for TSC problems using the original time domain. Recently, Bagnall et al. [2] experimentally show that the performance of classifiers significantly increases when changing data representation (compared with original temporal domain) ; thus, for a given classifier, there is a high variance of performance depending on the data transformation at use. To alleviate this problem, an ensemble method TSC-ENSEMBLE [2] based on three data representations (plus the original data) and NN algorithm is suggested. The experimental results demonstrate the importance of representations in TSC problems and show that a simple ensemble method based on several data representations provides highly competitive predictive performance. However, with the good performance of NN-based approaches also come the drawbacks of lazy learners: i.e., there is no proper training phase, therefore the training set has to be entirely stored and all the computation time is postponed until deployment phase. Another weakness of the NN approaches is the lack of interpretability; indeed NN only indicates the nearest series w.r.t. the used similarity measure.

The method we suggest takes the pros and leaves the cons of the methods listed above: we come back to the *eager*<sup>1</sup> paradigm, benefit from the combination of multiple representations, build and select valuable features from multiple representations. More precisely, in this paper, we suggest a parameter-free process for constructing valuable features over multiple representations for TSC problems. Our contribution is thus essentially methodological. The next section motivates and describes the three steps of our unsupervised process: *(i)* transformation of original data into several new data representations; *(ii)* coclustering on various data representations ; *(iii)* the exploitation of coclustering results for the construction of new features for the data. The output of the process is then a traditional data set (i.e. labeled objects described by attributes) ready for supervised feature selection and classification. We report the experimental validation of our approach in section 3 and discuss further related work in Section 4 before concluding.

## 2 Feature construction process

**Notations.** In TSC problems, we define a time serie as a pair  $(\tau_i, y_i)$  where  $\tau_i$  is a set of ordered observations  $\tau_i = \langle (t_1, x_1), (t_2, x_2), \dots, (t_{m_i}, x_{m_i}) \rangle$  of length  $m_i$  and  $y_i$  a class value. A time series data set is defined as a set of pairs  $D = \{(\tau_1, y_1), \dots, (\tau_n, y_n)\}$ , where each time series may have a different number of observations (i.e. with different length). Notice that the time series of a data set may also have different values for  $t_k, (k = 1..m_i)$ . The goal is to learn a classifier from  $D$  to predict the class of new incoming time series  $\tau_{n+1}, \tau_{n+2}, \dots$ . To achieve this goal, we suggest the feature construction process summarized as follows:

---

<sup>1</sup> In contrast to lazy learning, eager learning has an explicit training phase and generally deploys faster.

1. We transform original data into multiple data representations.
2. We process a coclustering technique on each representation.
3. We build a set of features from each coclustering result and obtain a new data set gathering the various sets of features.

The new data set is thus object-attribute oriented and ready for supervised classification phase. Since our main contribution is methodological, we will take some time to motivate each step of the process, and when necessary, to make the paper self-contained, we will recall the main principles of the tools used in each step.

## 2.1 Transformations & Representations

Numerous data transformation methods for time series has been suggested in the literature: e.g., polynomial, symbolic, spectral or wavelet transformations, ... (see [17] for a well-structured survey on experienced data representations).

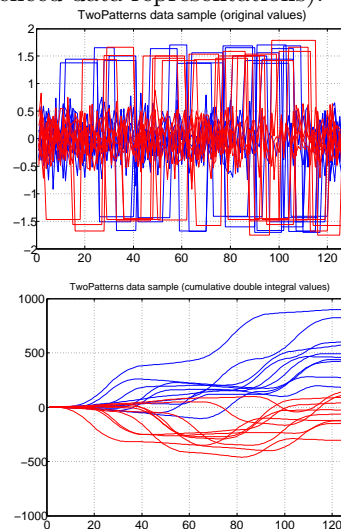
The underlying idea of using data transformation is that transformed data might contain class-characteristic pattern that are easily detectable (i.e. patterns unreachable in the original time domain). The following example illustrates and confirms the relevance of using representations, and highlights simple interpretable features that might arise from data representations.

*Motivating example.* Graphs from figure 1 confirm the relevance of changing data representation: indeed, from original data (a) it is uneasy to separate the two classes (blue/red) whereas simple transformation, like cumulative double integral (b) facilitate class discrimination. For example, after computing the cumulative double integral transformation, we see that curves with some values above 100 are blue and curves with some values below -100 are red. On this tiny example (extracted from the TwoPatterns data set from UCR repository [10]), a simple transformation and two interpretable features are enough to characterize the two classes of curves.

To illustrate and instantiate our process, we use the original representation and we pick six representations among the numerous ones existing in the literature.

**Derivatives : DV et DDV** We use derivatives and double derivatives of original time series (computed between time  $t$  et  $t - 1$ ). These transformations allow us to represent the local evolution (i.e., increasing/decreasing, acceleration/deceleration) of the series.

**Cumulative integrals : IV et IIV** We also use simple and double cumulative



**Fig. 1.** An extract from the TwoPatterns data set (20 series, 2 classes) in its original representation and in double cumulative integral representations

integrals of the series, computed using the trapeze method. These transformations allow us to represent the global (cumulated) evolution of the series.

**Power Spectrum : PS.** A time series can be decomposed in a linear combination of sines and cosines with various amplitudes and frequencies. This decomposition is known as the Fourier transform. And, the Power Spectrum is  $PS(\tau_i) = \langle (f_1, a_1), \dots, (f_{m_i}, a_{m_i}) \rangle$ , where  $f_k$  represent the frequency domain and  $a_k$  the power of the signal (i.e. the sum of the Fourier coefficients squared). This transformation is commonly used in signal processing and plunges the original series into the frequency domain.

**Auto-correlation function : ACF** The transformation by auto-correlation (ACF) is :  $\tau_{i\rho} = \langle (t_1, \rho_1), \dots, (t_{m_i}, \rho_{m_i}) \rangle$  where

$$\rho_k = \frac{\sum_{j=1}^{m_i-k} (x_j - \bar{x}) \cdot (x_{j+k} - \bar{x})}{m \cdot s^2}$$

and where  $\bar{x}$  and  $s^2$  are the mean and variance of the original series. ACF transformation describes the correlation between values of the signal at different times and thus allow us to represent auto-correlation structures like repeating patterns in the time series.

Thus, for a given time series data set  $D_{orig}$ , we build six new data representations:  $D_{DV}$ ,  $D_{DDV}$ ,  $D_{IV}$ ,  $D_{IIV}$ ,  $D_{PS}$  and  $D_{ACF}$  depending on the transformation used. In the following, for the sake of generality, an object from one of these representations will be called “curve” instead of time series since  $D_{PS}$  does not use the time domain.

## 2.2 Coclustering

In classification problems (also in TSC), there might exist intra-class variance, i.e. the variations between objects of the same class might be numerous and of various aspects. Using clustering as a pre-processing step to supervised classification is not new and is a solution to deal with intra-class variance. The idea is to pre-process the data set by grouping together similar objects and to highlight local patterns that might be class-discriminant: e.g., Vilalta et al. [16] suggest a pre-processing step by supervised (per-class) clustering using Expectation Maximization to enhance the predictive performance of Naive Bayes classifier. In order to be able to derive interesting features, we will use an unsupervised coclustering technique as described in the following.

A curve can be seen as a set of points  $(X, Y)$ , described by their abscissa and ordinate values. A set of curves is then also a set of points  $(C_{id}, X, Y)$  where  $C_{id}$  is the curve identifier. This tridimensional representation (one categorical variable and two numerical variables) of a curve data set is needed to apply coclustering methods. Indeed, the goal is to partition the categorical variable and to discretize the numerical variables in order to obtain clusters of curves and intervals for  $X$  and  $Y$ . The result is a tridimensional grid whose cells are defined by a group of curves, an interval for  $X$  and an interval for  $Y$ .

For that purpose, we use the coclustering method KHC [6] (Khiops Coclustering). Originally designed for clustering functional data, it is also suitable for

the particular case of curve data as defined above and it is directly applicable for our pre-processing step. KHC method is based on a piecewise constant non-parametric density estimation and instantiates the generic MODL approach [4] (Minimum Optimized Description Length) – which is similar to a Bayesian Maximum A Posteriori (MAP) approach. The optimal model  $M$ , i.e. the optimal grid, is obtained by optimization of a Bayesian criterion, called *cost*. The *cost* criterion bets on a trade-off between the accuracy and the robustness of the model and is defined as follows:

$$\text{cost}(M) = -\log(\underbrace{p(M | D)}_{\text{posterior}}) = -\log(\underbrace{p(M)}_{\text{prior}}) \times \underbrace{p(D | M)}_{\text{likelihood}}$$

Using a hierarchical prior (on the parameters of a data grid model) that is uniform at each stage of the hierarchy, we obtain an analytic expression for the *cost* criterion:

$$\text{cost}(M) = \log n + 2 \log N + \log B(n, k_C) \quad (1)$$

$$+ \log \binom{N+k-1}{k-1} + \sum_{i_C=1}^{k_C} \log \binom{N_{i_C} + n_{i_C} - 1}{n_{i_C} - 1} \quad (2)$$

$$+ \log N! - \sum_{i_C=1}^{k_C} \sum_{j_X=1}^{k_X} \sum_{j_Y=1}^{k_Y} \log N_{i_C j_X j_Y}! \quad (3)$$

$$+ \sum_{i_C=1}^{k_C} \log N_{i_C}! - \sum_{i=1}^n \log N_i! + \sum_{j_X=1}^{k_X} \log N_{j_X}! + \sum_{j_Y=1}^{k_Y} \log N_{j_Y}! \quad (4)$$

where  $n$  is the number of curves,  $N$  the number of points,  $k_C$  (resp.  $k_X, k_Y$ ) is the number of clusters of curves (resp. the number of intervals for  $X$  and  $Y$ ),  $k$  the number of cells of the data grid,  $n_{i_C}$  the number of curves in cluster  $i_C$ ,  $N_i$  the number of points for curve  $i$  and  $N_{i_C}$  (resp.  $N_{j_X}, N_{j_Y}, N_{i_C j_X j_Y}$ ) is the cumulated number of points for curves of cluster  $i_C$  (resp. for interval  $j_X$  of  $X$ , interval  $j_Y$  of  $Y$ , for cell  $(i_C, j_X, j_Y)$  of the data grid. Notice that  $B(n, k_C)$  is the number of divisions of  $n$  elements into  $k$  subsets. The two first lines stand for the prior and the two last lines relates to the likelihood of the model. Intuitively, low *cost* means high probability ( $p(M | D)$ ) that the model  $M$  arises from the data  $D$ . From an information theory point of view, according to [15], the negative logarithms of probabilities may be interpreted as code length. Thus, the *cost* criterion may also be interpreted as the code length of the grid model plus the code length of data  $D$  given the model  $M$ , according to the Minimum Description Length principle (MDL [9]). Here, low *cost* means high compression of the data using the model  $M$ .

The *cost* criterion is optimized using a greedy bottom-up strategy, (i) starting with the finest grained model, (ii) considering all merges between adjacent clusters or intervals, for the curve and dimension variables, and (iii) performs the best merge if the criterion decreases after the merge. The process loops until no further merge improves the criterion. The obtained grid constitutes a non-parametric estimator of the joint density of the curves and the dimensions of points.

KHC is parameter-free, robust (avoids over-fitting), handles large curve data sets with several millions of data points and its time complexity is  $\Theta(N\sqrt{N} \log N)$  (sub-quadratic) where  $N$  is the number of data points: thus, KHC meets our problem needs (for full details, see [6]).

### 2.3 Feature construction

Feature construction for TSC problems [13] aims at capturing class-relevant properties for describing time series. The generated features goes from simple ones like minimum, maximum, mean, standard deviation of time series to more complex ones like e.g., coefficients of spectral decompositions [12]. The main advantage of feature-based approaches is the final transactional (or vector) representation of the data which is suitable for conventional classifiers like Naive Bayes or decision trees. In our process, we generate features from coclustering results as follows.

For each coclustering result obtained with KHC on a data representation  $(D_{orig}, D_{DV}, D_{DDV}, D_{IV}, D_{IIV}, D_{PS}, D_{ACF})_i$ , we create a set of new features:  $\mathcal{F}_{orig}, \mathcal{F}_{DV}, \mathcal{F}_{DDV}, \mathcal{F}_{IV}, \mathcal{F}_{IIV}, \mathcal{F}_{PS}, \mathcal{F}_{ACF}$ . The new features are the descriptive attributes of the new data set whose objects are curves.

Let  $D_{rep}$  be one of the seven representations described above. Let  $M_{rep} = KHC(D_{rep})$  be the tridimensional optimal grid obtained by coclustering with KHC on  $D_{rep}$ . We denote  $k_C$  the number of clusters of  $M_{rep}$  and  $k_Y$  the number of intervals of  $M_{rep}$  for dimension  $Y$ . We then create similarity-based features and histogram features.

#### Similarity-based features

Considering the good performance of (dis)similarity-based approaches (e.g., ED-NN and DTW-NN), we define a dissimilarity index based on the *cost* criterion.

**Definition 1 (Dissimilarity index).** *The dissimilarity between a curve  $\tau_i$  and a cluster  $c_j$  of the optimal grid  $M_{rep}$  is defined as:*

$$d(\tau_i, c_j) = cost(M_{rep}|\tau_i \cup c_j) - cost(M_{rep})$$

*i.e., the difference of cost between the optimal model  $M_{rep}$  and the model  $M_{rep}|\tau_i \cup c_j$  (the optimal grid in which we add the curve  $\tau_i$  to the cluster of curves  $c_j$ ).*

Intuitively,  $d$  measures the perturbation brought by the integration of a curve into a cluster of curves of the optimal grid (i.e. according to the *cost* criterion used for grid optimization). In terms of code length, if a curve  $\tau_i$  is similar to the curves of cluster  $c_j$ , the total code length of the data is not much different from the total code length of the data plus  $\tau_i$ . Thus, small values of  $d(\tau_i, c_j)$  indicate that  $\tau_i$  is similar to the curves of  $c_j$  whereas high values of  $d$  ( $d(\tau_i, c_j) \gg 0$ ) mean that  $\tau_i$  does not look like the curves of  $c_j$ .

According to the dissimilarity index  $d$ , we generate the following features:

- $k_C$  numerical features (one for each cluster  $c_j$  of curves of  $M_{rep}$ ). The value for a curve  $\tau_i$  is the difference  $d(\tau_i, c_j)$ . Thus, for a given curve  $\tau_i$ , these features tell how  $\tau_i$  is similar to the clusters of curves of the optimal grid (according to  $d$ ).
- One categorical feature indicating the index  $j$  of the cluster of curves that is the closest to a curve  $\tau_i$  according to the dissimilarity  $d$  defined above (i.e.,  $\arg \min_j d(\tau_i, c_j)$ ).

### Histogram features

Taking up the idea of interpretable features (see motivating example and fig.1), we also generate the following features:

- $k_Y$  numerical features (one for each interval  $i_Y$  of  $Y$  from  $M_{rep}$ ) whose value for a curve  $\tau_i$  is the number of points of  $\tau_i$  in interval  $i_Y$ .

These histogram features quantify the presence of a curve in intervals of  $Y$  obtained in the coclustering step.

For a given curve  $\tau_i$ , we now have the following informations provided by the new features (for each representation): (i) the dissimilarity values between  $\tau_i$  and all the clusters of curves, (ii) the index of the closest cluster of curves and (iii) the number of points of  $\tau_i$  in each interval of  $Y$ .

## 2.4 Supervised classification algorithm

We saw that our feature construction process may generate hundreds of new features for each representation. The whole set of features  $\mathcal{F}_{tot}$  for our new data set may contain thousands of attributes. Therefore, the classifier at the end of our process has to be capable of handling a large number of attributes but also selecting the relevant attributes for the classification task. At this stage, we could use conventional classifiers like decision trees or SVM. However, we choose the Selective Naive Bayes classifier (SNB) that meets all the needs, is parameter-free and outperforms classical Naive Bayes [5]. Notice that SNB exploits pre-processing techniques that discretize numerical variables, group values of categorical variables, weight and select features w.r.t. class-relevance by using robust conditional density estimators and following the MODL approach (see [3], [4]). Thus, the generated features benefit from these pre-processing techniques and preserve a potential of interpretability; we lead specific experiments in the next section to support this claim. Moreover, SNB is parameter-free, so is the whole feature construction process. Its time complexity is  $\Theta(KN \log(KN))$ , where  $N$  is the number of objects and  $K$  the number of features.

## 3 Experimental validation

The implementation of the classification process is based on existing tools (KHC for coclustering and SNB for supervised classification<sup>2</sup>). Connections between

<sup>2</sup> KHC and SNB are both available at <http://www.khiops.com>

the tools to handle the whole process, named MODL-TSC, are scripted using MATLAB. The experiments are led to discuss the following questions:

- $\mathcal{Q}_1$  Is MODL-TSC comparable with competitive contenders of the state-of-the-art in terms of accuracy?
- $\mathcal{Q}_2$  MODL-TSC employs and combines several representations. Are they all useful? Do they all bring the same impact?
- $\mathcal{Q}_3$  What kind of data insight do we gain using the coclustering-based features?

### 3.1 Protocol

We experiment our process on 51 time series data sets: 42 data sets are from UCR [10] and 9 new data sets introduced in [2]. The benchmark data sets offer a wide variety in terms of application domains, number of series, length of series and number of class values. We lead experiments in a predefined train-test setting for each data set (see [10]). We compare the predictive performance of our process, called MODL-TSC, with a baseline and three of the most effective alternative approaches:

- ED-NN: the Nearest Neighbor classifier using the Euclidean distance. This approach is considered as a baseline.
- DTW-NN: the Nearest Neighbor classifier using the elastic distance Dynamic Time Warping, considered as hard to beat in the literature (see [18])
- TSC-ENSEMBLE [2] exploits multiple representations via an ensemble method and the NN algorithm. Its performance is comparable to DTW-NN
- FAST-SHAPELETS [14] mines shapelets (i.e., class relevant time series subsequences) that might be embedded in e.g., a decision tree

### 3.2 Results

Due to space limitations, only average accuracy are reported. For full details about predictive performance and running time results, see [1].

**Comparisons with state-of-the-art.** Firstly, global results (mean accuracy, number of wins and mean rank) show that MODL-TSC is very competitive compared to state-of-the-art methods (see table 1). Although ED-NN seems to perform worse than other contenders, the Friedman test (at significance level  $\alpha = 0.05$ ) does not allow us to reject the null hypothesis, i.e., considering the experiments on these 51 benchmark data sets, there is no classifier singled out. We also run Wilcoxon’s sign rank test for pairwise comparisons (also with  $\alpha = 0.05$ ): it appears that MODL-TSC (as well as DTW-NN and TSC-ENS) significantly performs better than ED-NN but other pairwise comparisons conclude that there is no significant difference of performance between MODL-TSC and DTW-NN or TSC-ENS. Table 2 reports the comparison between MODL-TSC and FAST-SHAPELETS over 45 data sets (since accuracy results for FAST-SHAPELETS are available for these data, see [14]). Here Wilcoxon’s sign rank test for pairwise comparisons indicates that MODL-TSC performs significantly better than FAST-SHAPELETS. This



global view of performance results confirms that MODL-TSC is very competitive compared to the most effective contenders of the state-of-the-art.

Data	DTW-NN	ED-NN	TSC-ENS	MODL-TSC
Mean Acc	78.06	73.89	77.45	77.44
#wins	21	2	10	20
MODL-TSC wins vs.	25	34	29	-
Average rank	2.2549	3.1372	2.2745	2.2941

**Table 1.** Comparisons of accuracy results over 51 data sets for MODL-TSC, DTW-NN, ED-NN, TSC-ENSEMBLE.

Data	MODL-TSC	SHAPELETS
Mean Acc	76.68	72.81
#Wins	32	13

**Table 2.** Comparisons of accuracy results over 45 data sets for MODL-TSC and FAST-SHAPELETS.

Secondly, we observe the remarkable performance of MODL-TSC on ARSim, ElectricDevices, FordA and OSULeaf data. On these data, we outperform DTW-NN and TSC-ENSEMBLE: the difference of test accuracy is at least 10. Here, the added-value of the data representations (i.e., the new features) is at work. TSC-ENSEMBLE (exploiting only three representations) and DTW-NN (working in time domain) obtain only poor accuracy results. Conversely the performance of MODL-TSC is very low on Coffee, DiatomSizeReduction, ECGFiveDays and OliveOil data. The difference of test accuracy (about 10 compared with DTW-NN and TSC-ENSEMBLE) is now to our disadvantage. We think that this poor performance might due to one reason: the training set size of these data sets is very small (less than 30 of curves) and could be insufficient for either learning relevant coclusters or learning a predictive model without over-fitting. Indeed, e.g., for OliveOil data, there are only 30 training curves (for 4 classes), no cluster of curves is found by KHC whatever the representation.

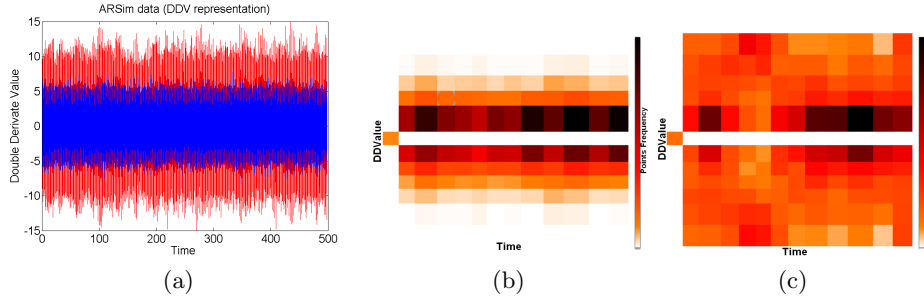
**Added-value of the representations.** In table 3, we also report accuracy results of SINGLE-MODL-TSC using only one representation.

We observe that using one single representation provide poor average accuracy results. Almost always, MODL-TSC using several representations outperforms SINGLE-MODL-TSC on  $\mathcal{F}_{orig}$  (resp.  $\mathcal{F}_{DV}$ ,  $\mathcal{F}_{DDV}$ ,  $\mathcal{F}_{IV}$ ,  $\mathcal{F}_{IIV}$ ,  $\mathcal{F}_{PS}$ ,  $\mathcal{F}_{ACF}$ ). In some cases (e.g., ItalyPowerDemand, MALLAT or MedicalImages), the good performance of MODL-TSC can be attributed to the combination of several representations. Indeed, the gap of test accuracy between any SINGLE-MODL-TSC and MODL-TSC is about 10; thus the combination of features coming from different views of the data improves accuracy results. In other cases (e.g., ARSim or wafer), the good performance seems to be due to only one (or at most two) representation while the other representations are ignored. As an example, for ARSim data, the DDV representation is the most relevant. KHC obtains 43 clusters of curves and 12 intervals for  $Y_{DDV}$ . Most of the clusters are almost pure (only one class of curves per cluster). Moreover, as we can see in figure 2(b) and (c), the number of points in intervals generated by KHC above 6 and below -6 are class-discriminant since curves of class 1 almost never have points in these regions.

	Mean Acc	MODL-TSC wins
MODL-TSC	77.44	-
$\mathcal{F}_{orig}$	65.98	49
$\mathcal{F}_{DV}$	61.20	47
$\mathcal{F}_{DDV}$	55.30	49
$\mathcal{F}_{IV}$	59.26	50
$\mathcal{F}_{IIV}$	53.48	50
$\mathcal{F}_{PS}$	53.27	50
$\mathcal{F}_{ACF}$	57.54	49

**Table 3.** Comparisons of accuracy results for MODL-TSC and SINGLE-MODL-TSC using each single representation.

Thus the combination of features coming from different views of the data improves accuracy results. In other cases (e.g., ARSim or wafer), the good performance seems to be due to only one (or at most two) representation while the other representations are ignored. As an example, for ARSim data, the DDV representation is the most relevant. KHC obtains 43 clusters of curves and 12 intervals for  $Y_{DDV}$ . Most of the clusters are almost pure (only one class of curves per cluster). Moreover, as we can see in figure 2(b) and (c), the number of points in intervals generated by KHC above 6 and below -6 are class-discriminant since curves of class 1 almost never have points in these regions.



**Fig. 2.** ARSim data: (a) double derivate representation, class 1 is blue and class 2 is red. (b) An example of cluster whose curves are mostly of class 1. (c) an example of cluster whose curves are mostly of class 2. The frequency of points are represented for each cell in Time vs. DDValue axis. The stronger the color, the more frequent are the points in a cell.

These experiments recall the very importance of representations in TSC problems and particularly in our feature construction process. Even the simple representations we chose to illustrate our process show good predictive performance. Depending on the application, we may still hope some improvement in performance if we could rely on expert domain knowledge to select relevant representations to use in our generic process.

### 3.3 Interpretation: an example

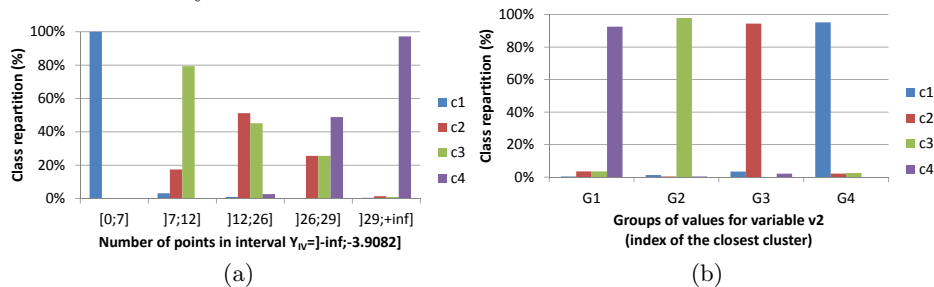
If we consider the cumulative integral (IV) representation of TwoPatterns data, the optimal grid obtained by KHC is made of 224 clusters of curves, 11 intervals for  $X$  and 9 intervals for  $Y_{IV}$ . According to the MODL pre-processing techniques, among all the attributes generated from all representations, the two most relevant attributes are from the IV representation:

1.  $v_1$ , the number of points in interval  $I_{Y_{IV}} = ] - \infty; -3.9082]$
2.  $v_2$ , the index of the closest cluster

In the supervised learning step, the discretization for  $v_1$  and the value grouping for  $v_2$  provide the following contingency tables represented as histograms (see figure 3). We observe (figure 3(a)) that the number of points  $p$  of a curve in interval  $I_{Y_{IV}}$  (i.e. the number of points with value less than  $-3.9082$ ) is class relevant. Indeed, in the learning phase, curves such that  $p \leq 7$  are of class  $c_1$ ; when  $p > 29$  (about 23% of the points of the curve), curves are mostly of class  $c_4$  and when  $7 < p \leq 12$  they are mostly of class  $c_3$ . This type of feature is similar to the ones in the motivating example and figure 1: for a given representation, some regions of  $y$ -axis (delimited by intervals) will be class-discriminant and the number of points of an incoming curve in this interval will also be class-discriminant.

In figure 3(b), we firstly see that, for variable  $v_2$  (“index of the closest cluster”), MODL pre-processing by supervised value grouping provide 4 groups:  $G_1$ , (resp.  $G_2$ ,  $G_3$  et  $G_4$ ) made of 56, (resp. 53, 53 et 62) indexes of clusters that are mostly of class  $c_4$  (resp.  $c_3$ ,  $c_2$ ,  $c_1$ ). The attribute  $v_2$  is then class-relevant. Indeed, for example, if  $j$  is the index of cluster, that is the closest to a curve  $\tau_i$ , and belongs

to  $G_2$  (i.e.  $j \in G_2$ ), then  $\tau_i$  is considered very similar to curves of class  $c_3$ . Moreover, the variable “index of the closest cluster” is an indicator of the relevance of the representation in our process for the current TSC problem. In this example, attribute  $v_2$  alone, is enough to characterize 95% of the data, therefore, IV data representation is very relevant for characterizing the classes of TwoPatterns data. Conversely, for the original representation ( $D_V$ ), the optimal grid obtained with KHC is made of 255 clusters of curves but MODL pre-processing indicates that the variable “index of the closest cluster” is not relevant to characterize the classes of TwoPatterns; as a consequence, SINGLE-MODL-TSC on  $\mathcal{F}_{orig}$  shows bad test accuracy results.



**Fig. 3.** Histogram representation of class repartition for discretization of variable  $v_1$  and value grouping of variable  $v_2$ .

## 4 Related work

In TSC problems, DTW-NN is recognized by the community as a hard-to-beat baseline and it is confirmed by our experiments. However, there exist alternative approaches: besides the numerous similarity measures coupled with Nearest Neighbor algorithm [17], for the sake of interpretability, feature-based approaches have also been intensively studied. Feature-based approaches for TSC aim at extracting class-relevant characteristics of series so that a conventional classifier can be used. A wide variety of features has been studied: e.g., global, trends, symbolic, intervals, distance-based, features coming from spectral transforms [12] or a combination of several types of features [7].

Shapelet-based approaches, a subtopic of feature-based approaches, have drawn much attention in recent years. Shapelets are time series subsequences that are representative of a class. First approaches have embedded extracted shapelets in a decision tree [8, 20], others in a simple rule-based classifier [19], while very recently, Lines et al. [11] have designed a shapelet-based transform.

Our approach generates similarity-based features and histogram features over multiple representations; the former allows us to reach predictive performance comparable to the best similarity-based NN classifiers, with the latter we gain some insight in the data. The closest works are that of Eruhimov et al. [7] who employs a generation of high dimensional feature space and the inspiring work of Bagnall et al. [2] who establish competitive predictive performance by combining multiple representations in an ensemble classifier.

## 5 Conclusion & Perspectives

We have suggested MODL-TSC, a simple yet effective and generic feature construction process for time series classification problems (TSC). Our process is parameter-free, easy to use and the generated features offer a high potential of interpretation. The time complexity of our process is sub-quadratic, thus time-efficient. Experimental results show that the performance of MODL-TSC is highly competitive and comparable with two of the most accurate approaches of the state-of-the-art (namely, DTW-NN and TSC-ENS).

The first results are promising and also confirm the importance of representations in TSC problems. Indeed, depending on the application domain, a particular transformation will facilitate the discovery of class relevant patterns. Moreover, the combination of multiple representations with MODL-TSC leads to highly competitive predictive performance. We have used only a few simple representations in the time, frequency and correlation domains to demonstrate that our feature construction approach is well-founded. The literature offers plenty of relevant data representations (see [17] for a wide view). Notice also that designing new representations is still a hot topic (see e.g., [11]). It gives an *unexplored* potential of improvement for MODL-TSC on data sets and applications where we are less performant.

## References

1. Supporting webpage, <http://sites.google.com/site/tscfeatures/>
2. Bagnall, A., Davis, L.M., Hills, J., Lines, J.: Transformation based ensembles for time series classification. In: SDM'12. pp. 307–318 (2012)
3. Boullé, M.: A bayes optimal approach for partitioning the values of categorical attributes. *Journal of Machine Learning Research* 6, 1431–1452 (2005)
4. Boullé, M.: MODL: A bayes optimal discretization method for continuous attributes. *Machine Learning* 65(1), 131–165 (2006)
5. Boullé, M.: Compression-based averaging of selective naive Bayes classifiers. *Journal of Machine Learning Research* 8, 1659–1685 (2007)
6. Boullé, M.: Functional data clustering via piecewise constant nonparametric density estimation. *Pattern Recognition* 45(12), 4389–4401 (2012)
7. Eruhimov, V., Martyanov, V., Tuv, E.: Constructing high dimensional feature space for time series classification. In: PKDD'07. pp. 414–421 (2007)
8. Geurts, P.: Pattern extraction for time series classification. In: PKDD'01. pp. 115–127 (2001)
9. Grünwald, P.: *The minimum description length principle*. MIT Press (2007)
10. Keogh, E., Zhu, Q., Hu, B., Y., H., Xi, X., Wei, L., Ratanamahatana, C.A.: The UCR time series classification/clustering page (2011), [http://www.cs.ucr.edu/~eamonn/time\\_series\\_data/](http://www.cs.ucr.edu/~eamonn/time_series_data/)
11. Lines, J., Davis, L.M., Hills, J., Bagnall, A.: A shapelet transform for time series classification. In: KDD'12. pp. 289–297 (2012)
12. Mörchen, F.: Time series feature extraction for data mining using DWT and DFT. Tech. rep., Philipps Univeristy Marburg (2003)

13. Nanopoulos, A., Alcock, R., Manolopoulos, Y.: Feature-based classification of time-series data. In: Mastorakis, N., Nikolopoulos, S.D. (eds.) *Information Processing and Technology*, pp. 49–61. Nova Science (2001)
14. Rakthanmanon, T., Keogh, E.: Fast shapelets: a scalable algorithm for discovering time series shapelets. In: *SIAM DM'13* (2013)
15. Shannon, C.E.: A mathematical theory of communication. *Bell System Technical Journal* (1948)
16. Vilalta, R., Rish, I.: A decomposition of classes via clustering to explain and improve naive bayes. In: *ECML'03*. pp. 444–455 (2003)
17. Wang, X., Mueen, A., Ding, H., Trajcevski, G., Scheuermann, P., Keogh, E.: Experimental comparison of representation methods and distance measures for time series data. *Data Mining and Knowledge Discovery* 26(2), 275–309 (2013)
18. Xi, X., Keogh, E.J., Shelton, C.R., Wei, L., Ratanamahatana, C.A.: Fast time series classification using numerosity reduction. In: *ICML'06*. pp. 1033–1040 (2006)
19. Xing, Z., Pei, J., Yu, P.S., Wang, K.: Extracting interpretable features for early classification on time series. In: *SDM'11*. pp. 247–258 (2011)
20. Ye, L., Keogh, E.J.: Time series shapelets: a novel technique that allows accurate, interpretable and fast classification. *Data Mining and Knowledge Discovery* 22(1-2), 149–182 (2011)