

Invariant Factorization of Time Series

Josif Grabocka

Lars Schmidt-Thieme

{josif, schmidt-thieme}@ismll.uni-hildesheim.de
ISMILL, University of Hildesheim, Germany

Abstract

Time-series classification is an important domain of machine learning and a plethora of methods have been developed for the task. In comparison to existing approaches, this study presents a novel method which decomposes a time-series dataset into latent patterns and membership weights of local segments to those patterns. The process is formalized as a constrained objective function and a tailored stochastic coordinate descent optimization is applied. The time-series are projected to a new feature representation consisting of the sums of the membership weights, which captures frequencies of local patterns. Features from various sliding window sizes are concatenated in order to encapsulate the interaction of patterns from different sizes. Finally, a large-scale experimental comparison against 6 state of the art baselines and 43 real life datasets is conducted. The proposed method outperforms all the baselines with statistically significant margins in terms of prediction accuracy.

1 Introduction

Time-series classification is a pillar problem of machine learning and its existence spans over decades of research. Series data emerge in a myriad of application domains, from health-care and astronomy up to economics and botanics. In comparison to other types of data, time series exhibit a high degree of intra-class variability, where patterns occur shifted in time, distorted and scaled. Therefore traditionally strong classifiers, such as Support Vector Machines (SVM), fail to excel in terms of prediction accuracy [14].

A series of attempts have been proposed to address the intra-class variations of time-series patterns. An early pioneer method called Dynamic Time Warping (DTW), (still considered competitive [11, 27]), computes the similarity among series by re-aligning the time indexes. The algorithm explores all the possible relative alignments of time indexes of two series and picks the one yielding the minimum overall distance [18].

The research of time-series classification can be approximately categorized into distance metrics, invariant classifiers, feature extraction and bag-of-patterns streams. Distance metrics focus on defining measurements on the similarity of two series instances [18, 7, 8, 4]. Invariant classifiers, on the other hand, aim at embedding similarities into classifiers. For instance, the invariant kernel functions have

been applied to measure instance similarities in the projected space of a non-linear SVM [32, 14]. Another paper proposes to generate all pattern variations as new instances and inflate the training set [13]. The bag-of-patterns approach splits the time-series into local segments and collects statistics over the segments. Those local segments are converted into symbolic words and a histogram of the words' occurrences is built [22, 23]. Another study constructs a supervised codebook generated from local patterns, which is used to create features for a random forest classifiers [5].

In comparison to existing approaches this study proposes a new perspective. We assume that time series are generated by a set of latent (hidden) patterns which occur at different time stamps and different frequencies across instances. In addition those patterns might be convoluted and/or distorted to produce derived local patterns.

We would like to introduce the concept through the illustration of Figure 1. A synthetic dataset consists of two classes A (green) and B (black), each having two instances. All the time series are composed of three segments of 20 points, while each segment is a convolutional derivative of two latent patterns depicted in red and blue. In other words, each segment is a weighted sum of a single-peaked and double-peaked pattern. The shown coefficients of the convolution are degrees of membership that each local segment has to one of those two latent patterns.

Both Euclidean and DTW based nearest neighbor classifiers have 100% error on a leave-one-out experiment on the dataset of Figure 1. As can be observed, instance A1 is closer to B1 than A2, and the same applies for all other series. In fact the rationale behind this dataset is that A has a higher frequency of the red single-peaked pattern, while B has a higher domination of the blue double-peaked pattern. The method presented in this paper detects the latent patterns, measures the degrees of membership and sums them up into a bag-of-pattern approach. Our approach converts the series of Figure 1 into a new representation F, concretely: $F_{A1} = [1.9, 1.1]$, $F_{A2} = [1.7, 1.3]$, $F_{B1} = [1.3, 1.7]$, $F_{B2} = [1.1, 1.9]$. A nearest neighbor classifier over the new representation F yields 0% error.

In this paper, we will propose a method which detects a set of latent patterns for a time series dataset together with

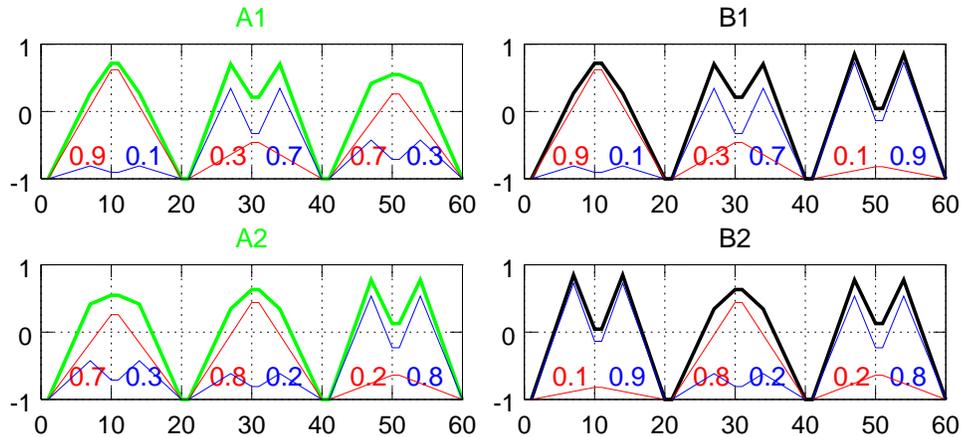


Figure 1: Four series of two classes $A=\{A1,A2\}$ and $B=\{B1,B2\}$, each generated as a convolution of latent patterns

a convolutional degree of membership weights. The product of the membership weights with the patterns approximates the original segments. In contrast to the aforementioned synthetic example, real datasets have segments occurring at arbitrary locations and being of different sizes. Our method employs a sliding window approach to split the series into overlapping local segments and utilizes a factorization model to decompose the segments into latent patterns and weights. We formalize the objective function of the factorization and propose a stochastic coordinate descent technique in order to optimize the objective. The sum of the learned membership degrees is used to project the time series into a new representation. Ultimately, in order to resolve the scale invariance of the patterns, sums of memberships from different sliding window sizes are concatenated.

A throughout experimental comparison is conducted on 43 datasets of the UCR time-series collection against six state of the art baselines. Our method outperforms all the baselines with a statistically significant margin of improvement. Considering the evidence of our experimental results and our survey of related work, we conclude that our prediction accuracy figures are the best published in the realm of time-series classification, regarding the UCR collection of datasets.

2 Related Work

Time-series classification has been elaborated in a vast number of occasions, therefore a complete survey of all the published papers is out of our scope. Instead, we will structure the related work into a set of categories and mention relevant prominent studies.

2.1 Distance Metrics and Invariant Classifiers: A significant portion of time-series research has centralized around the definition of accurate similarity metrics. The most popular of those approaches is the Dynamic Time

Warping (DTW) measure [18], which overcomes deficiencies of the L_2 norm distance by aligning the time indexes of two series instances. The similarity measure is typically plugged into a nearest neighbor classifier. DTW produces competitive prediction accuracies [11, 31] and has been speed up using lower boundary heuristics [27].

Other similarity based distance metrics have extended the edit distance of strings into the time-series domain [7, 8]. Furthermore, the longest common subsequence of time series has also been used as an indication of similarity [28]. Moreover, similarities of sequential data have been measured using sparse spatial kernels [19]. A state of the art method called complexity-invariant distance metric (CID) introduces the total variation regularization for time-series. CID significantly improves the accuracy of DTW [3, 4].

Efforts have been dedicated to incorporating time-series variations into popular classifiers. For instance DTW has been used as a SVM kernel [14], even though the resulting kernel is not positive semi definite. Consecutively, another study has proposed a Gaussian elastic kernel [32]. A method which produces a semi-definite kernel is called global alignment kernels and builds an average statistics from all possible warping paths of time indexes [9]. In addition, another study has inflated the training set by adding new instances that represent variations of original training data [13].

2.2 Feature Extraction and Bag-of-patterns: Other researchers have emphasized the extraction of series features for boosting classification. Dimensionality reduction has been used to project the time series into a low-rank data space [17], while a recent method incorporates class segregation into the projection [12].

However, the most prominent state of the art technique for extracting time-series features is called shapelets mining. Shapelets represent the most discriminative series segment (or set of segments), which yields the maximal prediction

accuracy [26, 24]. A related study detects a set of shapelets and transforms the series data into a new representation, defined by the distance to those shapelets [15].

A recent direction of research has drawn attention on the need to segment the time series into local patterns and measure the frequencies of patterns as classification features. For instance frequencies of time-series motifs have been fed into standard classifiers [6]. Another attempt has focused on building histograms of local patterns represented as symbolic words [23]. Those symbolic words are produced by a piecewise constant approximation technique called SAX [21], while the frequencies of the SAX words are used ultimately for classification [22, 23]. One similar bag-of-words approach has also been applied to long biomedical data [30]. Moreover, a bag-of-patterns study proposes to extract series segments of various lengths and positions and generate a supervised codebook of those patterns [5]. A random forest classifier has been trained over the extracted features. That study demonstrates considerable improvements over baselines in terms of prediction accuracy [5].

2.3 Factorization of Time Series There have been a few attempts in generating invariant time-series features through factorization. A shift-invariant sparse coding of signals has been proposed for reconstructing noisy or missing series segments [20]. In similar domains, sparse coding factorization has been applied for deriving shift and 2D rotation invariant features of hand writing data [2], and also invariant features of audio data [16]. Moreover, a temporal decomposition of multivariate streams has been used to discover patterns in patients’ clinical events [29].

Our method differs from distance metrics principally. Instead of measuring the similarity of series, we project the data into a new representation, where similar instances are positioned close to each other. Furthermore, the proposed method distances away from existing bag-of-patterns methods because we learn a latent decomposition of patterns, instead of counting the occurrence of segments on the original time-series. Finally, our contributions over the existing factorization methods rely on (i) a novel approach in detecting both shift and scale invariant features for time series, and (ii) building a bag-of-patterns representation of the learned invariant features for a classification scenario.

3 Definitions and Notations

1. **Time-series:** A time-series is an ordered sequence of point values. In a dataset of N series instances, where each series has Q points, we denote the series dataset as $T \in \mathbb{R}^{N \times Q}$.
2. **Sliding Window Segment:** A sliding window content of size $L \in \mathbb{N}$, is a series subsequence starting at a position $j \in \{1, \dots, Q - L\}$ of a series i of dataset T , and is denoted as $S_{i,j} \in \mathbb{R}^L$, $S_{i,j} := (T_{i,j}, T_{i,j+1}, \dots, T_{i,j+L-1})$.

3. **All Dataset Segments:** The starting position of each sliding window segment is incremented by an offset $\delta = \{1, \dots, L\}$, therefore the maximum number of segments per series is defined as $M := \frac{Q-L}{\delta}$. All the segments of a time-series datasets are denoted as $S \in \mathbb{R}^{N \times M \times L}$.
4. **Latent Patterns:** Our method mines for K -many latent patterns, each having the same size as one segment, i.e L . So, the latent patterns are denoted as $P \in \mathbb{R}^{K \times L}$.
5. **Degrees of Membership:** Each instance of a dataset will be approximated via the product of latent patterns and the set of membership degrees to those patterns. Each segment of a series will have one membership weight to each of the K latent patterns. Consequently, the degrees of membership of all time-series are defined as $D \in \mathbb{R}^{N \times M \times K}$.

4 Invariant Factorization of Time Series

4.1 Segmenting the Time-Series The series of the dataset are segmented in a sliding window approach having size L and increment δ . The segmentation of each series is described in Algorithm 1. Once derived, the segments are normalized to mean 0 and deviation 1.

Algorithm 1 SegmentSeries

Require: $T \in \mathbb{R}^{N \times Q}$, $L \in \mathbb{N}$, $\delta \in \mathbb{N}$
Ensure: $S \in \mathbb{R}^{N \times M \times L}$

- 1: **for** $i = 1, \dots, N$, $j = 1, \dots, M$ **do**
- 2: **for** $l = 1, \dots, L$ **do**
- 3: $S_{i,j,l} \leftarrow T_{i, \delta(j-1)+l-1}$
- 4: **end for**
- 5: $S_{i,j} \leftarrow \text{normalize}(S_{i,j})$
- 6: **end for**
- 7: **return** S

4.2 Invariant Factorization Objective The objective of the invariant factorization relies on approximating every segment of a series segment as a product of the defined latent patterns P and the membership degrees D . The objective function is described in Equation 4.1.

$$(4.1) \quad \underset{D, P}{\operatorname{argmin}} \sum_{i=1}^N \sum_{j=1}^M \sum_{l=1}^L \left(S_{i,j,l} - \sum_{k=1}^K D_{i,j,k} P_{k,l} \right)^2 + \lambda_P \sum_{k=1}^K \sum_{l=1}^L P_{k,l}^2$$

Subject To:

$$\sum_{k=1}^K D_{i,j,k} = 1, \quad D_{i,j,k} \geq 0, \quad \forall i, j, k$$

The objective function is composed of two loss terms and one constraint. Firstly, the latent patterns P and the

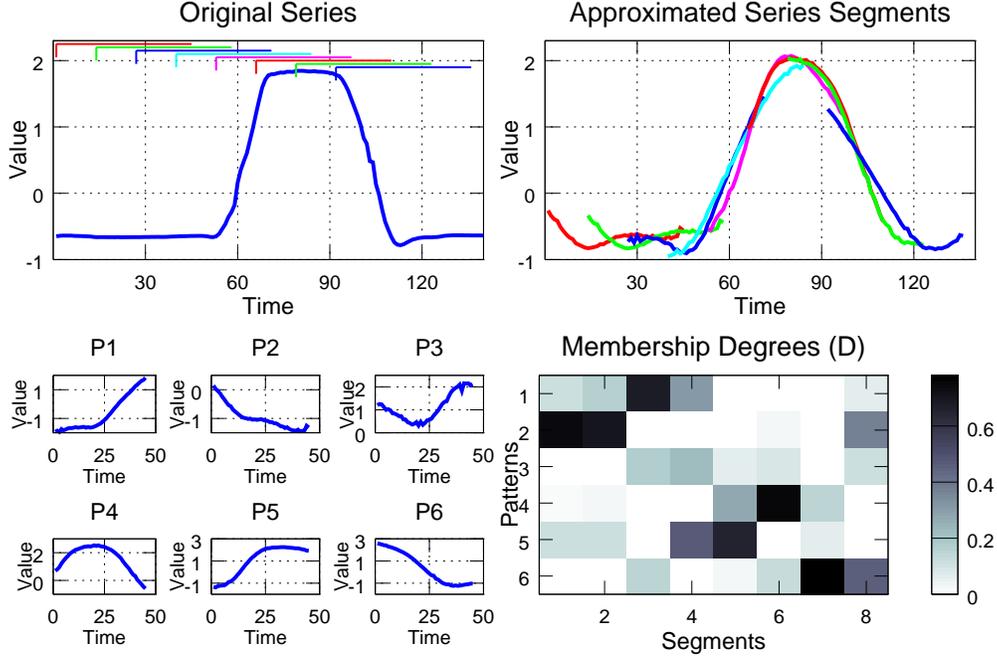


Figure 2: A factorized instance of the Gun Point dataset with parameters $K = 6, L = 45, \delta = 13, \lambda_P = 1$

memberships D should approximate the normalized segments of the series dataset. Therefore, minimizing the L2 norm of the reconstruction error achieves the goal. In addition, a second regularization loss term is added in order to prohibit the patterns P from over-fitting. A hyper-parameter λ_P controls the degree of regularization. Finally, we impose equality and positivity constraints on the membership degrees. The membership degrees of every segment $D_{i,j}$ sum-up to one, because each segment needs to have the same impact factor. Otherwise, in a bag-of-patterns representation of series, different segments would have different scales of memberships. The positivity constraint, on the other hand, prohibits non-interpretable negative memberships.

We would like to illustrate the invariant factorization objective with a concrete illustration, shown in Figure 2. A learned decomposition, as in Equation 4.1, is depicted for the Gun Point dataset. On the left top, a series instance is presented, while the dataset’s latent patterns and the membership degrees of the instance are found below. The product of the patterns and memberships yield the series approximation shown in the right top chart. The series is split into 8 overlapping segments of size 45, each starting at an offset of 13 points. For instance, the 7-th segment starts at 79 and has a high membership value to the 6-th pattern, which matches the descending structure. However, please note that other patterns also contribute with smaller membership degrees (patterns 4 and 5) in order to fit exactly the original segment content.

4.3 Learning the Patterns and Memberships In order to learn the latent patterns and the memberships we are going to optimize the objective function of Equation 4.1 via **stochastic coordinate descent**, which operates by updating each cell of D, P in the direction of the first derivative of the objective.

4.3.1 Update Rules for Latent Patterns In order to compute the update rules for the patterns, we first define the error in approximating a point l of the segment j , in time-series i , as $\xi_{i,j,l}$.

$$\text{Let } \xi_{i,j,l} := S_{i,j,l} - \sum_{k=1}^K D_{i,j,k} P_{k,l}$$

$$P_{k,l}^* := \underset{z}{\operatorname{argmin}} \lambda_P z^2 + \sum_{i,j} (\xi_{i,j,l} + D_{i,j,k} P_{k,l} - D_{i,j,k} z)^2$$

$$2\lambda_P P_{k,l}^* - 2 \sum_{i,j} (\xi_{i,j,l} + D_{i,j,k} (P_{k,l} - P_{k,l}^*)) D_{i,j,k} = 0$$

Subsequently the optimal value of every point l of a latent pattern k is denoted as $P_{k,l}^*$ and is found by solving the first derivative in a coordinate descent way. Therefore the optimal value of $P_{k,l}$ is defined in Equation 4.2. Please note that the error values don’t have to be recomputed for each point, instead we can incrementally update the error terms. Equation 4.3 refreshes the error terms after the change of the pattern value.

$$(4.2) \quad P_{k,l}^* := \frac{\sum_{i,j} (\xi_{i,j,l} + D_{i,j,k} P_{k,l}) D_{i,j,k}}{\lambda_P + \sum_{i,j} D_{i,j,k}^2}$$

$$(4.3) \quad \xi_{i,j,l} \leftarrow \xi_{i,j,l} - (P_{k,l}^* - P_{k,l}) D_{i,j,k}$$

4.3.2 Update Rules for Membership Degrees The update rules for the membership degrees needs to preserve an equality constraint, which enforce the memberships of a segment to sum to one. Therefore any direct update of a membership $D_{i,j,k}$ will violate the constraint. In order to avoid this bottleneck, we propose to update the memberships in pairs, inspired by a similar strategy known as the Sequential Minimal Optimization algorithm [25]. The idea is to draw two random membership weights $D_{i,j,k}, D_{i,j,w}$ and update them such that their sum, denoted $Q = D_{i,j,k} + D_{i,j,w}$, remains equal before and after the updates. In that way, if we increase one membership, then the other would have to decrease and vice versa, while the aim is to find the combination which yield the smallest approximation error. Therefore, the optimal value of $D_{i,j,k}$, will be denoted by $D_{i,j,k}^*$ and is algebraically derived as the solution of the first derivative of our objective function.

$$D_{i,j,k}^* = \underset{z}{\operatorname{argmin}} \sum_l (\xi_{i,j,l} + D_{i,j,k} P_{k,l} + D_{i,j,w} P_{w,l} - Q P_{w,l} + z(P_{w,l} - P_{k,l}))^2$$

$$D_{i,j,k}^* = \frac{-\sum_l (\xi_{i,j,l} - D_{i,j,k} (P_{w,l} - P_{k,l})) (P_{w,l} - P_{k,l})}{\sum_l (P_{w,l} - P_{k,l})^2}$$

Once the optimal value $D_{i,j,k}^*$ is defined then we have to ensure the constraints. Equation 4.4 crops the optimal value to be nonnegative and not exceed the sum of the membership pairs. The error term is refreshed to include the changes of both memberships of the pair in Equation 4.5. As a last step we can commit the optimal values, by preserving their sum before the updates. As Equation 4.6 shows, the best value of $D_{i,j,w}$ can be deduced from the optimal value of $D_{i,j,k}$.

$$(4.4) \quad D_{i,j,k}^* \leftarrow \max(0, \min(Q, D_{i,j,k}^*))$$

$$(4.5) \quad \xi_{i,j,l} \leftarrow \xi_{i,j,l} - (D_{i,j,k}^* - D_{i,j,k}) P_{k,l}, \text{ and}$$

$$\xi_{i,j,l} \leftarrow \xi_{i,j,l} - (Q - D_{i,j,k}^* - D_{i,j,w}) P_{w,l}$$

$$(4.6) \quad D_{i,j,k} \leftarrow D_{i,j,k}^*, \quad D_{i,j,w} \leftarrow Q - D_{i,j,k}^*$$

4.4 Efficient Initialization Since the objective function of Equation 4.1 is nonlinear in terms of P and D together, then a coordinate descent optimization is not guaranteed to avoid local optima. Therefore, good initial values of the patterns and the memberships are crucial for the learning process. The intuition leads into assigning some of the segments as initial patterns, however it is not obvious which of them provide the best initialization.

The answer is addressed via a technique utilized to

find the initial centroids in a clustering setup [1]. The patterns (analogy to centroids) are initialized to segments with a probability proportional to the distance to all the other segments [1]. Therefore, we are assured to pick centroid segments which are evenly distributed across the space of all series segments. The initialization steps are detailed in Algorithm 2. Please note that the first pattern has to be drawn randomly in a uniform distribution, while the other patterns are chosen randomly from the dataset segments based on the probability of their distance to the existing patterns. The function \mathcal{C} measures the distance of a segment to the closest existing pattern.

Algorithm 2 Initialize

Require: $S \in \mathbb{R}^{N \times M \times L}$, $L \in \mathbb{N}$, $K \in \mathbb{N}$

Ensure: $D \in \mathbb{R}^{N \times M \times K}$, $P \in \mathbb{R}^{K \times L}$

1: $P_1 \leftarrow S_{i',j'}$, drawn $i', j' \sim \mathcal{U}(N, M)$

2: **for** $k = 2, \dots, K$ **do**

3: $P_k \leftarrow S_{i',j'}$, with probability weights $\frac{\mathcal{C}(S_{i',j'})^2}{\sum_{i,j} \mathcal{C}(S_{i,j})^2}$

4: **end for**

5: **for** $i = 1, \dots, N$; $j = 1, \dots, M$ **do**

6: $k' = \operatorname{argmin}_{k \in \{1, \dots, K\}} \|S_{i,j} - P_k\|^2$

7: $D_{i,j,k} \leftarrow \begin{cases} 1 & k = k' \\ 0 & k \neq k' \end{cases}, \quad k = 1, \dots, K$

8: **end for**

9: **return** D, P

The initialization of the membership degrees is more trivial than patterns. The degree index k' denotes that pattern $P_{k'}$ is the closest to segment $S_{i,j}$ and its membership $D_{i,j,k'}$ is set to 1, while all the other membership degrees are initialized to zero.

4.5 Learning Algorithm Algorithm 3 finally combines all the steps of the factorization process. In the beginning, the memberships and the patterns are initialized using Algorithm 2. Next the errors are initialized, then the coordinate descent technique updates all the parameters in a number of iterations, denoted as a hyper-parameter \mathcal{I} . Subsequently, the degrees of membership and the patterns are learned by setting the aforementioned optimal values. The membership and pattern indexes are visited in random order to speed up the convergence.

4.6 A New Invariant Representation The final representation will sum the membership degrees in a bag-of-patterns strategy. It enables a quantification of which local patterns appear in a series and how often. The shift invariance is achieved by segmenting the series in a sliding window approach and the scale invariance is addressed using different sliding window sizes. Algorithm 4 describes the algorithmic steps. The algorithm iterates over Φ many different scales of an initial sliding windows size L and solves an invariant fac-

Algorithm 3 InvariantFactorization

Require: $T \in \mathbb{R}^{N \times Q}$, $L \in \mathbb{N}$, $\delta \in \mathbb{N}$, $K \in \mathbb{N}$, $\lambda_P \in \mathbb{R}$, $\mathcal{I} \in \mathbb{N}$

Ensure: $D \in \mathbb{R}^{N \times M \times K}$, $P \in \mathbb{R}^{K \times L}$

```
1:  $S \leftarrow \text{SegmentSeries}(T, L, \sigma)$ 
2:  $(D, P) \leftarrow \text{Initialize}(S, L, K)$ 
3: {Initialize the errors}
4: for  $\forall i \in \mathbb{N}_1^N, \forall j \in \mathbb{N}_1^M, \forall l \in \mathbb{N}_1^L$  do
5:    $\xi_{i,j,l} := S_{i,j,l} - \sum_{k=1}^K D_{i,j,k} P_{k,l}$ 
6: end for
7: {Update the patterns&memberships iteratively}
8: for iteration = 1, ...,  $\mathcal{I}$  do
9:   {Update all degrees of membership}
10:  for  $\forall i \in \mathbb{N}_1^N, \forall j \in \mathbb{N}_1^M$  randomly do
11:    for 1, ...,  $K$ , {Draw  $K$ -many pairs} do
12:       $k, w \sim \mathcal{U}(K, K)$ , s.t.  $D_{i,j,k} + D_{i,j,w} \neq 0$ 
13:       $Q \leftarrow D_{i,j,k} + D_{i,j,w}$ 
14:      {Solve and crop the optimal memberships}
15:       $D_{i,j,k}^* = \frac{-\sum_l (\xi_{i,j,l} - D_{i,j,k} (P_{w,l} - P_{k,l})) (P_{w,l} - P_{k,l})}{\sum_l (P_{w,l} - P_{k,l})^2}$ 
16:       $D_{i,j,k}^* \leftarrow \max(0, \min(Q, D_{i,j,k}^*))$ 
17:      {Update the error terms}
18:      for  $l = 1, \dots, L$  do
19:         $\xi_{i,j,l} \leftarrow \xi_{i,j,l} - (D_{i,j,k}^* - D_{i,j,k}) P_{k,l}$ 
20:         $\xi_{i,j,l} \leftarrow \xi_{i,j,l} - (Q - D_{i,j,k}^* - D_{i,j,w}) P_{w,l}$ 
21:      end for
22:      {Commit the values of the pair}
23:       $D_{i,j,k} \leftarrow D_{i,j,k}^*$ 
24:       $D_{i,j,w} \leftarrow Q - D_{i,j,k}^*$ 
25:    end for
26:  end for
27:  {Update all patterns}
28:  for  $\forall k \in \mathbb{N}_1^K; \forall l \in \mathbb{N}_1^L$ , randomly do
29:     $P_{k,l}^* = \frac{\sum_{i,j} (\xi_{i,j,l} + D_{i,j,k} P_{k,l}) D_{i,j,k}}{\lambda_P + \sum_{i,j} D_{i,j,k}^2}$ 
30:    {Update the error terms}
31:    for  $i = 1, \dots, N; j = 1, \dots, M$  do
32:       $\xi_{i,j,l} \leftarrow \xi_{i,j,l} - (P_{k,l}^* - P_{k,l}) D_{i,j,k}$ 
33:    end for
34:    {Commit the pattern's point value}
35:     $P_{k,l} \leftarrow P_{k,l}^*$ 
36:  end for
37: end for
38: return  $D, P$ 
```

torization from Algorithm 3 per each size. The frequencies of the learned memberships are summed up for all K patterns and the procedure is repeated for every sliding window size. Finally each time series contains $K\Phi$ many features, which denote the frequencies of patterns at different sizes and positions.

The new representation will be used for classification,

instead of the original time series. We deployed a polynomial kernel Support Vector Machines, because we need to capture the interaction among features, i.e. the interaction among patterns of various sizes.

Algorithm 4 InvariantRepresentation

Require: $T \in \mathbb{R}^{N \times Q}$, $L \in \mathbb{N}$, $\delta \in \mathbb{N}$, $K \in \mathbb{N}$, $\lambda_P \in \mathbb{R}$, $\mathcal{I} \in \mathbb{N}$, $\Phi \in \mathbb{N}$

Ensure: $F \in \mathbb{R}^{N \times (K\Phi)}$

```
1: for  $s = 1, \dots, \Phi$  do
2:    $L' \leftarrow L \cdot s$ 
3:    $D \leftarrow \text{InvariantFactorization}(T, L', \delta, K, \lambda_P, \mathcal{I})$ 
4:   for  $i = 1, \dots, N; k = 1, \dots, K$  do
5:      $M \leftarrow \frac{Q - L'}{\delta}$ 
6:      $F_{i,k+(s-1)K} \leftarrow \sum_{j=1}^M D_{i,j,k}$ 
7:   end for
8: end for
9: return  $F$ 
```

4.7 Algorithmic Complexity The run-time complexity of the method is dominated by the updates of memberships and has an order $O(NMKLL)$. Concretely our method needs 48.4 hours to compute on the StarLightCurves (the largest) dataset, while for instance DTW needs 87 hours. The space complexity of our method depends on the storage of the segments S and the memberships D , which is $O(NM \max(K, L))$.

5 Experimental Results

5.1 Baselines We compared the prediction accuracy of our method, denoted Invariant Factorization (**INF**A), against the following six state of the art baselines:

- **TSBF**: The bag-of-features framework for time series (TSBF) uses a supervised codebook to extract features for a random forest classifier [5].
- **SSSK**: Sparse Spatial Similarity Kernel (SSSK) measures sequence similarity through sampling sequence features at different resolutions [19].
- **BOW**: The Bag of Words (BOW) method decomposes the series into local SAX words and uses a histogram representation of words as the new feature representation [22, 23].
- **DTW**: Dynamic Time Warping (DTW) computes the best alignment of time indexes resulting in the minimal distance [18, 27].
- **CID**: The complexity invariant distance (CID) adds a L2-based total variation regularization term into the DTW distance [4].
- **FSH**: Fast shapelet (FSH) extracts the most discriminative segment of the series dataset, such that the distance from the dataset instances to the optimal shapelet can be used as a feature for classification [26].

Table 1: Error Rates - Comparison of Prediction Accuracies on the UCR Collection of Datasets

Dataset	Cls.	Train	Test	Len.	INFA	TSBF	SSSK	BOW	DTW	CID	FSH
50words	50	450	455	270	0.215	0.209	0.488	0.316	0.310	0.226	0.557
Adiac	37	390	391	176	0.315	0.245	0.575	0.325	0.396	0.379	0.514
Beef	5	30	30	470	0.333	0.287	0.633	0.267	0.500	0.467	0.447
CBF	3	30	900	128	0.000	0.009	0.090	0.048	0.003	0.001	0.053
Chlorine	3	467	3840	166	0.451	0.336	0.428	0.405	0.352	0.351	0.417
CinC_ECG	4	40	1380	1639	0.159	0.262	0.438	0.164	0.349	0.054	0.174
Coffee	2	28	28	286	0.000	0.004	0.071	0.036	0.179	0.179	0.068
Cricket_X	12	390	390	300	0.197	0.278	0.585	0.305	0.223	0.249	0.527
Cricket_Y	12	390	390	300	0.208	0.259	0.654	0.313	0.208	0.197	0.505
Cricket_Z	12	390	390	300	0.200	0.263	0.574	0.295	0.208	0.205	0.547
Diatom	4	16	306	345	0.007	0.126	0.173	0.111	0.033	0.065	0.117
ECG200	2	100	100	96	0.140	0.145	0.220	0.110	0.230	0.110	0.227
ECGFiveDays	2	23	861	136	0.000	0.183	0.360	0.164	0.232	0.218	0.004
FaceAll	14	560	1690	131	0.237	0.234	0.369	0.238	0.192	0.144	0.411
FaceFour	4	24	88	350	0.000	0.051	0.102	0.102	0.170	0.125	0.090
FacesUCR	14	200	2050	131	0.083	0.090	0.356	0.137	0.095	0.102	0.328
Fish	7	175	175	463	0.063	0.080	0.177	0.029	0.167	0.154	0.197
Gun_Point	2	50	150	150	0.007	0.011	0.133	0.407	0.093	0.073	0.061
Haptics	5	155	308	1092	0.536	0.488	0.591	0.630	0.623	0.571	0.616
InlineSkate	7	100	550	1882	0.622	0.603	0.729	0.629	0.616	0.586	0.741
ItalyPower	2	67	1029	24	0.049	0.096	0.101	0.044	0.050	0.044	0.095
Lighting2	2	60	61	637	0.148	0.257	0.393	0.328	0.131	0.131	0.295
Lighting7	7	70	73	319	0.233	0.262	0.438	0.370	0.274	0.260	0.403
MALLAT	8	55	2345	1024	0.028	0.037	0.153	0.098	0.066	0.075	0.033
MedicalImages	10	381	760	99	0.275	0.269	0.463	0.401	0.263	0.258	0.433
MoteStrain	2	20	1252	84	0.097	0.135	0.166	0.177	0.165	0.205	0.217
OliveOil	4	30	30	570	0.367	0.090	0.300	0.233	0.133	0.167	0.213
OSULeaf	6	200	242	427	0.190	0.329	0.326	0.153	0.409	0.372	0.359
Sony	2	20	601	70	0.163	0.175	0.376	0.409	0.275	0.185	0.315
SonyII	2	27	953	65	0.075	0.196	0.339	0.154	0.169	0.123	0.215
StarLightCurves	3	1000	8236	1024	0.023	0.022	0.135	0.021	0.093	0.066	0.063
SwedishLeaf	15	500	625	128	0.078	0.075	0.339	0.125	0.210	0.117	0.269
Symbols	6	25	995	398	0.036	0.034	0.184	0.088	0.050	0.059	0.068
synthetic_control	6	300	300	60	0.010	0.008	0.067	0.017	0.007	0.027	0.081
Trace	4	100	100	275	0.000	0.020	0.300	0.000	0.000	0.010	0.002
Two_Patterns	4	1000	4000	128	0.005	0.001	0.087	0.010	0.000	0.004	0.114
TwoLeadECG	2	23	1139	82	0.005	0.046	0.257	0.248	0.096	0.138	0.090
uWaveX	8	896	3582	315	0.177	0.164	0.358	0.242	0.273	0.211	0.293
uWaveY	8	896	3582	315	0.225	0.249	0.493	0.352	0.366	0.278	0.392
uWaveZ	8	896	3582	315	0.229	0.217	0.439	0.325	0.342	0.293	0.364
Wafer	2	1000	6174	152	0.003	0.004	0.029	0.010	0.020	0.006	0.004
WordsSynonyms	25	267	638	270	0.303	0.302	0.553	0.371	0.351	0.243	0.594
yoga	2	300	3000	426	0.112	0.149	0.172	0.145	0.164	0.156	0.269
Absolute Wins					19.33	9.00	0.00	4.33	2.83	7.5	0
INFA One-to-one Wins						26	41	34	34	31	41
INFA One-to-one Draws						0	0	1	1	0	0
INFA One-to-one Losses						17	2	8	8	12	2
Wilcoxon Test (p values) (Stat. significant for $p \leq 0.05$)						0.028	0.000	0.000	0.000	0.005	0.000

5.2 Setup and Reproducibility We conducted a large-scale experimentation in 43 time-series dataset from the UCR collection¹. Our protocol complied to the default train/test split of the data, which is an established benchmark split and is used by the baselines. The metric of comparison is the error rate, i.e. the misclassification rate. Table 1 shows the datasets used for experimentation together with the number of classes, the number of training instances, the number of testing instances and the length of the series.

Our method has a relatively high number of hyperparameters, however most of them can be analytically adjusted. Since all the segments are normalized, then the latent patterns should also have mean 0 and standard deviation 1. Therefore, $\lambda_P = 1$ by the definition of the Tikhonov regularization. We searched for four different sliding windows sizes, i.e. $\Phi = 4$ and $L = 20\%$ of Q , so $L' \in \{20\%, 40\%, 60\%, 80\%\}$ of Q . The number of latent patterns needs to be set large enough to avoid underfitting and was set to $K = 50\%$ of Q . A fine grained sliding window offset was applied as $\delta = 5\%$ of L . However, in order to ensure the scalability for the four largest datasets (Cin_ECG, InlineSkate, MALLAT, StarLightCurves) we set their parameters to $K = 10\%$ of Q and $\delta = 20\%$ of L . The maximum number of iterations was set to $\mathcal{I} = 15$. The applied classifier was a polynomial kernel SVM with a polynomial degree being 3 and the complexity parameter 1, which are competitive SVM settings for the UCR collection [13]. Since the algorithm is based on a probabilistic initialization, it might be possible that it converges to different closeby optima in each execution. However, in our experiments, those optima were very close and the final prediction accuracy results have insignificant differences. **The authors are devoted to promote full reproducibility, therefore the source code, the data and instructions are publicly available².**

5.3 Results The error rate results of the six state of the art baselines and our method INFA are presented in Table 1. The best performing method for each dataset (row) is emphasized in bold. In order to compare multiple classifiers across a large number of datasets we follow the established benchmarks of counting wins and *Wilcoxon's Signed-Rank* test for statistical significance [10]. To be fair with the baselines, we retrieved the results from the baselines' publications [5, 4, 26] over the **same** data splits as INFA. In addition, we verified the published results of the baselines with our own experimental checkups.

Three comparative figures are conducted, the first of which counts the absolute number of wins. Each dataset awards a total value of 1, which is split into equal fractions in case methods have equal error rate scores. The "Absolute

wins" row, in the bottom of the table, counts the datasets where a method has the best prediction accuracy. As can be trivially deduced, our method has a clear superiority in terms of absolute wins, scoring 19.37 wins against 9.00 wins of the second best method. In addition INFA outperforms by large margins all the baselines in an one-to-one comparisons of wins. INFA has more wins, yet the predominant analysis is whether or not those wins represent statistically significant differences. Each cell on the bottom row represents the p value of the Wilcoxon Signed-Rank test on the error rate values of INFA against each baseline. Our method has a statistically significant difference over the error results of all baselines with a two-tailed hypothesis and the standard significance level of 95% confidence ($p \leq 0.05$).

Based on our survey of related work, the results presented in this study are the best published prediction accuracy scores in the realm of time-series classification, with respect to the UCR collection of datasets.

6 Conclusions

In this study we presented Invariant Factorization, a method that initially decomposes the time series into a set of overlapping segments via a sliding window approach. The segments are approximated by learning a set of latent patterns and degrees of memberships of each segment to each pattern. We formalized the factorization as a constraint objective function and proposed a stochastic coordinate descent method to solve it. The new representation of time series are the sums of the membership weights, which represent frequencies of local patterns. Features from various sliding window sizes were concatenated to encapsulate interaction among patterns of various scales. Finally we conducted a thorough experimental comparison against 6 state of the art baselines in 43 real-life time series datasets. Our method outperforms all the baselines with statistically significant margins and marks the best published results in the realm of time-series classification, regarding the UCR collection of datasets.

References

- [1] David Arthur and Sergei Vassilvitskii. k-means++: the advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, SODA '07, pages 1027–1035, Philadelphia, PA, USA, 2007. Society for Industrial and Applied Mathematics.
- [2] Q. Barthelemy, A. Larue, A. Mayoue, D. Mercier, and J.I. Mars. Shift and 2d rotation invariant sparse coding for multivariate signals. *Signal Processing, IEEE Transactions on*, 60(4):1597–1611, 2012.
- [3] Gustavo E. A. P. A. Batista, Xiaoyue Wang, and Eamonn J. Keogh. A complexity-invariant distance measure for time series. In *SDM*, pages 699–710. SIAM / Omnipress, 2011.
- [4] Gustavo E.A.P.A. Batista, Eamonn J. Keogh, Oben Moses Tataw, and Vinicius M.A. Souza. Cid: an efficient

¹www.cs.ucr.edu/~eamonn/time_series_data

²<http://fs.ismll.de/publicspace/>

- complexity-invariant distance for time series. *Data Mining and Knowledge Discovery*, pages 1–36, 2013.
- [5] M. Baydogan, G. Runger, and E. Tuv. A bag-of-features framework to classify time series. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, PP(99):1–1, 2013.
- [6] Krisztian Buza and Lars Schmidt-Thieme. Motif-based classification of time series with bayesian networks and svms. In Andreas Fink et al., editor, *GfKI*, pages 105–114. Springer, 2008.
- [7] Lei Chen and Raymond Ng. On the marriage of lp-norms and edit distance. In *Proceedings of the Thirtieth international conference on Very large data bases - Volume 30, VLDB '04*, pages 792–803. VLDB Endowment, 2004.
- [8] Yueguo Chen, M.A. Nascimento, Beng-Chin Ooi, and A. Tung. Spade: On shape-based pattern detection in streaming time series. In *Data Engineering, 2007. ICDE 2007. IEEE 23rd International Conference on*, pages 786–795, 2007.
- [9] Marco Cuturi. Fast global alignment kernels. In Getoor et al., editor, *Proceedings of the ICML 2011*, ICML 2011, pages 929–936, New York, NY, USA, June 2011. ACM.
- [10] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *J. Mach. Learn. Res.*, 7:1–30, December 2006.
- [11] Hui Ding, Goce Trajcevski, Peter Scheuermann, Xiaoyue Wang, and Eamonn J. Keogh. Querying and mining of time series data: experimental comparison of representations and distance measures. *PVLDB*, 1(2):1542–1552, 2008.
- [12] Josif Grabocka, Alexandros Nanopoulos, and Lars Schmidt-Thieme. Classification of sparse time series via supervised matrix factorization. In Jörg Hoffmann and Bart Selman, editors, *AAAI*. AAAI Press, 2012.
- [13] Josif Grabocka, Alexandros Nanopoulos, and Lars Schmidt-Thieme. Invariant time-series classification. In Peter A. Flach, Tjil De Bie, and Nello Cristianini, editors, *ECML/PKDD (2)*, volume 7524 of *Lecture Notes in Computer Science*, pages 725–740. Springer, 2012.
- [14] Steinn Gudmundsson, Thomas Philip Runarsson, and Sven Sigurdsson. Support vector machines and dynamic time warping for time series. In *IJCNN*, pages 2772–2776. IEEE, 2008.
- [15] J. Hills, J. Lines, E. Baranauskas, J. Mapp, and A. Bagnall. Classification of time series by shapelet transformation. *Data Mining and Knowledge Discovery, accepted subject to minor corrections*, 2013.
- [16] Po-Sen Huang, Jianchao Yang, Mark Hasegawa-Johnson, Feng Liang, and Thomas S. Huang. Pooling robust shift-invariant sparse representations of acoustic signals. In *INTERSPEECH*. ISCA, 2012.
- [17] Eamonn Keogh, Kaushik Chakrabarti, Michael Pazzani, and Sharad Mehrotra. Dimensionality reduction for fast similarity search in large time series databases. *Knowledge and Information Systems*, 3(3):263–286, 2001.
- [18] Eamonn J. Keogh and Michael J. Pazzani. Scaling up dynamic time warping for datamining applications. In *KDD*, pages 285–289, 2000.
- [19] P.P. Kuksa and V. Pavlovic. Spatial representation for efficient sequence classification. In *Pattern Recognition (ICPR), 2010 20th International Conference on*, pages 3320–3323, 2010.
- [20] Michael S. Lewicki and Terrence J. Sejnowski. Coding time-varying signals using sparse, shift-invariant representations. In *Proceedings of NIPS*, pages 730–736, Cambridge, MA, USA, 1999. MIT Press.
- [21] Jessica Lin, Eamonn Keogh, Li Wei, and Stefano Lonardi. Experiencing sax: a novel symbolic representation of time series. *Data Min. Knowl. Discov.*, 15(2):107–144, October 2007.
- [22] Jessica Lin, Rohan Khade, and Yuan Li. Rotation-invariant similarity in time series using bag-of-patterns representation. *J. Intell. Inf. Syst.*, 39(2):287–315, 2012.
- [23] Jessica Lin and Yuan Li. Finding structural similarity in time series data using bag-of-patterns representation. In *Proceedings of the 21st International Conference on Scientific and Statistical Database Management, SSDBM 2009*, pages 461–477, Berlin, Heidelberg, 2009. Springer-Verlag.
- [24] Abdullah Mueen, Eamonn J. Keogh, and Neal Young. Logical-shapelets: an expressive primitive for time series classification. In Chid Apté, Joydeep Ghosh, and Padhraic Smyth, editors, *KDD*, pages 1154–1162. ACM, 2011.
- [25] John C. Platt. Advances in kernel methods. chapter Fast training of support vector machines using sequential minimal optimization, pages 185–208. MIT Press, Cambridge, MA, USA, 1999.
- [26] T. Rakthanmanon and E. Keogh. Fast shapelets: A scalable algorithm for discovering time series shapelets. *Proceedings of the 13th SIAM International Conference on Data Mining*, 2013.
- [27] Thanawin Rakthanmanon, Bilson Campana, Abdullah Mueen, Gustavo Batista, Brandon Westover, Qiang Zhu, Jesin Zakaria, and Eamonn Keogh. Searching and mining trillions of time series subsequences under dynamic time warping. In *Proceedings of the 18th ACM SIGKDD*, KDD 2012, pages 262–270, New York, NY, USA, 2012. ACM.
- [28] M. Vlachos, G. Kollios, and D. Gunopulos. Discovering similar multidimensional trajectories. In *Data Engineering, 2002. Proceedings. 18th International Conference on*, pages 673–684, 2002.
- [29] Fei Wang, Noah Lee, Jianying Hu, Jimeng Sun, and Shahram Ebadollahi. Towards heterogeneous temporal clinical event pattern discovery: a convolutional approach. In *Proceedings of ACM SIGKDD*, KDD '12, pages 453–461, New York, NY, USA, 2012. ACM.
- [30] Jin Wang, Ping Liu, Mary F.H. She, Saeid Nahavandi, and Abbas Kouzani. Bag-of-words representation for biomedical time series classification. *Biomedical Signal Processing and Control*, 8(6):634 – 644, 2013.
- [31] Xiaoyue Wang, Abdullah Mueen, Hui Ding, Goce Trajcevski, Peter Scheuermann, and Eamonn Keogh. Experimental comparison of representation methods and distance measures for time series data. *Data Mining and Knowledge Discovery*, 26(2):275–309, 2013.
- [32] Dongyu Zhang, Wangmeng Zuo, David Zhang, and Hongzhi Zhang. Time series classification using support vector machine with gaussian elastic metric kernel. In *ICPR*, pages 29–32. IEEE, 2010.