Similarity Preserving Representation Learning for Time Series Analysis

Qi Lei[‡] Jinfeng Yi[†] Roman Vaculin[†] Lingfei Wu[†] Inderjit S. Dhillon[‡]

The University of Texas at Austin, Austin, TX, USA

†IBM Thomas J. Watson Research Center, Yorktown Heights, NY, USA
leiqi@ices.utexas.edu, {jinfengy, vaculin, wuli}@us.ibm.com, inderjit@cs.utexas.edu

Abstract

A considerable amount of machine learning algorithms take instance-feature matrices as their inputs. As such, they cannot directly analyze time series data due to its temporal nature, usually unequal lengths, and complex properties. This is a great pity since many of these algorithms are effective, robust, efficient, and easy to use. In this paper, we bridge this gap by proposing an efficient representation learning framework that is able to convert a set of time series with equal or unequal lengths to a matrix format. In particular, we guarantee that the pairwise similarities between time series are well preserved after the transformation. The learned feature representation is particularly suitable to the class of learning problems that are sensitive to data similarities. Given a set of n time series, we first construct an $n \times n$ partially observed similarity matrix by randomly sampling $O(n \log n)$ pairs of time series and computing their pairwise similarities. We then propose an extremely efficient algorithm that solves a highly non-convex and NP-hard problem to learn new features based on the partially observed similarity matrix. We use the learned features to conduct experiments on both data classification and clustering tasks. Our extensive experimental results demonstrate that the proposed framework is both effective and efficient.

1 Introduction

Modeling time series data is an important but challenging task. It is considered by [Yang and Wu, 2006] as one of the 10 most challenging problems in data mining. Although time series analysis has attracted increasing attention in recent years, the models that analyze time series data are still much fewer than the models developed for static data. The latter category of models, which usually take instance-feature matrices as their inputs, cannot directly analyze time series data due to its temporal nature, usually unequal lengths, and complex properties [Längkvist et al., 2014]. This is a great pity since many static models are effective, robust, efficient, and easy to use. Introducing them to time series analysis can greatly enhance the development of this domain.

In this work, we bridge this gap by proposing an efficient unsupervised representation learning framework that is able to convert a set of time series data with equal or unequal lengths to a matrix format. In particular, the pairwise similarities between the raw time series data are well preserved after the transformation. Therefore, the learned feature representation is particularly suitable to the similarity-based models in a variety of learning problems such as data clustering, classification, and learning to rank. Notably, the proposed framework is flexible to any time series distance or similarity measures such as Mikowski distance, cross-correlation, Kullback-Leibler divergence, dynamic time warping (DTW) similarity, and short time series (STS) distance. In this work, we use DTW similarity by default since it is known as the best measure for time series problems in a wide variety of domains [Rakthanmanon et al., 2013].

Given a total of n time series, our first step is to generate an $n \times n$ similarity matrix A with A_{ij} equaling to the DTW similarity between the time series i and j. However, computing all the pairwise similarities requires to call the DTW algorithm $O(n^2)$ times, which can be very timeconsuming when n is large. As a concrete example, generating a full similarity matrix when n = 150,000 takes more than 28 hours on an Intel Xeon 2.40 GHz processor with 256 GB of main memory. In order to significantly reduce the running time, we follow the setting of matrix completion [Sun and Luo, 2015] by assuming that the similarity matrix A is of low-rank. This is a very natural assumption since DTW algorithm captures the co-movements of time series, which has shown to be driven by only a small number of latent factors [Stock and Watson, 2005: Basu and Michailidis, 2015]. According to the theory of matrix completion, only $O(n \log n)$ randomly sampled entries are needed to *perfectly* recover an $n \times n$ low-rank matrix. This allows us to only sample $O(n \log n)$ pairs of time series to generate a partially observed similarity matrix A. In this way, the time spent on generating similarity matrix is significantly reduced by a factor of $O(n/\log n)$, a number that scales almost linearly with respect to n. When n = 150,000, it only takes about 3 minutes to construct a partially observed similarity matrix with $[20n \log n]$ observed entries, more than

¹Since this similarity matrix is symmetric, we only need to call the DTW algorithm about $[10n \log n]$ times to generate this matrix.

500 times faster than generating a full similarity matrix.

Given the generated partially observed similarity matrix $\hat{\bf A}$, our second step learns a new feature representation for n time series such that their pairwise DTW similarities can be well approximated by the inner products of new features. To this end, we solve a symmetric matrix factorization problem to factorize $\hat{\mathbf{A}}$, i.e., learning a feature representation $\mathbf{X} \in \mathcal{R}^{n \times d}$ such that $P_{\Omega}(\tilde{\mathbf{A}}) \approx P_{\Omega}(\mathbf{X}\mathbf{X}^{\top})$, where P_{Ω} is a matrix projection operator defined on the observed entry set Ω . Despite its relatively simple formulation, this optimization problem is hard to solve since it is highly non-convex and NP-hard. To address this challenge, we propose a very efficient exact cyclic coordinate descent algorithm. By wisely updating variables and taking the advantage of sparse observed entries in Ã, the proposed algorithm incurs a very low computational cost, and thus can learn new feature representations in an extremely efficient way.

To evaluate the performance of the learned feature representation, we use more than 10 real-world data sets to conduct experiments on data classification and clustering tasks. Our results show that classical static models that are fed with our learned features outperform the state-of-the-art time series classification and clustering algorithms in both accuracy and computational efficiency. In summary, our main contributions of this work are two-fold:

- We bridge the gap between time series data and a great amount of static models by learning a new feature representation of time series. The learned feature representation preserves the pairwise similarities of the raw time series data, and is general enough to be applied to a variety of learning problems.
- We propose the first, to the best of our knowledge, optimization parameter free algorithm to solve symmetric matrix factorization problem on a partially observed matrix. The proposed algorithm is highly efficient and converges at a very fast rate.

2 Related Work

In this section, we review the existing work on learning feature representations for time series data. Among them, a family of methods use a set of derived features to represent time series. For instance, [Nanopoulos et al., 2001] proposed to use the mean, standard deviation, kurtosis, and skewness of time series to represent control chart patterns. The authors in [Wang et al., 2006] introduced a set of features such as trend, seasonality, serial correlation, chaos, nonlinearity, and self-similarity to partition different types of time series. [Deng et al., 2013] used some easy to compute features such as mean, standard deviation and slope temporal importance curves to guide time series classification. In order to automate the selection of features for time series classification, the authors in [Fulcher and Jones, 2014] proposed a greedy forward method that can automatically select features from thousands of choices. Besides, several techniques have been proposed to represent time series by a certain types of transformation, such as discrete Fourier transformation [Faloutsos et al., 1994], discrete wavelet transformation [Chan and Fu, 1999], piecewise aggregate approximation [Keogh *et al.*, 2001], and symbolic aggregate approximation [Lin *et al.*, 2007]. In addition, deep learning models such as Elman recurrent neural network [Elman, 1990] and long short-term memory [Schmidhuber, 2015] are capable of modeling complex structures of time series data and learn a layer of feature representations. Due to their outstanding performance on a number of applications, they have become increasingly popular in recent years.

Despite the remarkable progress, the feature representations learned by these algorithms are usually problem-specific, and are not general enough for applications in multiple domains. Besides, the learned features cannot preserve similarities of the raw time series data, thus they are not suitable to the problems that are sensitive to the data similarity. These limitations inspire us to propose a problem-independent and similarity preserving representation learning framework for time series data.

3 Similarity Preserving Representation Learning for Time Series Analysis

In this section, we first present a general approach of our similarity preserving time series representation learning framework. We then propose an extremely efficient algorithm that is significantly faster than a naive implementation.

3.1 Problem Definition and General Framework

Given a set of n time series $\mathcal{T} = \{T_1, \dots, T_n\}$ with equal or unequal lengths, our goal is to convert them to a matrix $\mathbf{X} \in \mathcal{R}^{n \times d}$ such that the time series similarities are well preserved after the transformation. Specifically, we aim to learn a mapping function $f: T \to \mathcal{R}^d$ that satisfies

$$S(T_i, T_j) \approx \langle f(T_i), f(T_j) \rangle \ \forall i, j \in [n],$$
 (1)

where $\langle\cdot,\cdot\rangle$ stands for the inner product, one of the most commonly used similarity measure in analyzing static data. $S(\cdot,\cdot)$ denotes the pairwise time series similarity that can be computed by a number of functions. In this work, we use dynamic time warping (DTW) algorithm to measure the similarity by default. By warping sequences non-linearly in the time dimension, the DTW algorithm can calculate an optimal match between two given temporal sequences with equal or unequal lengths. Due to its superior performance, DTW has been successfully applied to a variety of applications, including computer animation [Müller, 2007], surveillance [Sempena et al., 2011], gesture recognition [Celebi et al., 2013], signature matching [Efrat et al., 2007], protein sequence alignment [Vial et al., 2009], and speech recognition [Muda et al., 2010].

Normally, DTW algorithm outputs a pairwise distance between two temporal sequences, thus we need to convert it to a similarity score. Since the inner product space can be induced from the normed space using $\langle x,y\rangle=(\|x\|^2+\|y\|^2-\|x-y\|^2)/2$ [Adams, 2004], we generate the DTW similarity by

$$S(T_i, T_j) = \frac{DTW(T_i, \theta)^2 + DTW(T_j, \theta)^2 - DTW(T_i, T_j)^2}{2},$$

where θ denotes the length one time series with entry θ . We note that the similarity computed via the above equation is a

more robust choice than some other similarity measures such as the reciprocal of distance. This is because when two time series are almost identical, their DTW distance is close to 0 and thus its reciprocal tends to infinity.

In order to learn the matrix \mathbf{X} , an intuitive idea is to factorize the similarity matrix $\mathbf{A} \in \mathcal{R}^{n \times n}$ where $\mathbf{A}_{ij} = \mathrm{S}(T_i, T_j)$. In more detail, this idea consists of two steps, i.e., a similarity matrix construction step and a symmetric matrix factorization step. In the first step, it constructs a similarity matrix \mathbf{A} as follows

$$\mathbf{A}_{ij} = \begin{cases} [b_i^2 + b_j^2 - \text{DTW}^2(T_i, T_j)]/2, & \text{if } i < j \\ (b_i^2 + b_j^2)/2, & \text{if } i = j \\ \mathbf{A}_{ji} & \text{if } i > j, \end{cases}$$
 (2)

where $b_i = DTW(T_i, 0)$, $i = 1, \dots, n$. In the second step, it learns an optimal data-feature matrix **X** by solving the following optimization problem

$$\min_{\mathbf{X}=(\mathbf{x}_1,\cdots,\mathbf{x}_n)} \sum_{i=1}^n \sum_{j=i+1}^n \|\mathbf{A}_{ij} - \langle \mathbf{x}_i, \mathbf{x}_j \rangle\|^2,$$
 (3)

which can be further expressed in a matrix form

$$\min_{\mathbf{X} \in \mathcal{R}^{n \times d}} \|\mathbf{A} - \mathbf{X} \mathbf{X}^{\top}\|_F^2. \tag{4}$$

The problems (3) and (4) can be solved by performing eigendecomposition on A, i.e.,

$$\mathbf{X} = \mathbf{Q}_{1:n,1:d} \times \sqrt{\mathbf{\Lambda}_{1:d,1:d}} \;,$$

where $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{\top}$ and the notation $\mathbf{Q}_{1:k,1:r}$ represents the sub-matrix of \mathbf{Q} that includes its first k rows and the first r columns.

Although the inner products of the learned data points $\mathbf{x}_1, \cdots, \mathbf{x}_n$ can well approximate the DTW similarities of the raw time series, the idea described above is not practical since both two steps are painfully slow when n is large. In order to generate a $n \times n$ similarity matrix, we need to call the DTW algorithm $O(n^2)$ times. In addition, a naive implementation of eigen-decomposition takes $O(n^3)$ time. Although we can reduce this cost by only computing the d largest eigenvalues and the corresponding eigenvectors, it is still computational intensive with a large n. As a concrete example, when n=150,000 and the length of time series is 30, it takes more than 28 hours to generate the similarity matrix and more than 3 days to compute its 30 largest eigenvalues on an Intel Xeon 2.40 GHz processor with 256 GB of main memory.

3.2 A Parameter-free Scalable Algorithm

In this subsection, we propose an extremely efficient approach that significantly reduces the computational costs of both steps while learns the new feature representation with a high precision. In addition to efficiency, another nice property is that the proposed approach is *optimization parameter free*, that is, the user does not need to decide any optimization parameter such as step length or learning rate.

To significantly improve the efficiency of the first step, we make the key observation that the similarity matrix **A** should be of low-rank. This is due to the fact that the DTW algorithm measures the level of co-movement between time series, which has shown to be dictated by only a small number of latent factors [Stock and Watson, 2005; Basu and Michailidis, 2015]. Indeed, we can verify the low-rankness in another way. Since the matrix **A** is a special case of Wigner random matrix, the gaps between its consecutive eigenvalues should not be small [Marčenko and Pastur, 1967]. This implies the low-rank property since most of its energy is concentrated in its top eigenvalues [Erdős *et al.*, 2009].

Based on the theory of matrix completion [Sun and Luo, 2015], only $O(n \log n)$ randomly sampled entries are needed to perfectly recover an $n \times n$ low-rank matrix. Thus, we don't need to compute all the pairwise DTW similarities. Instead, we randomly sample only $O(n \log n)$ pairs of time series, and then compute the DTW similarities only within the selected pairs. In other words, we generate a partially observed similarity matrix $\tilde{\bf A}$ with $O(n \log n)$ observed entries as

$$\tilde{\mathbf{A}}_{ij} = \begin{cases} \mathbf{S}(T_i, T_j) & \text{if } \Omega_{ij} = 1\\ \text{unobserved} & \text{if } \Omega_{ij} = 0, \end{cases}$$
 (5)

where $\Omega \in \{0,1\}^{n \times n}$ is a binary matrix indicating the indices of sampled pairs. In this way, the running time of the first step is significantly reduced by a factor of $O(n/\log n)$. Since this factor scales almost linearly with n, we can greatly save the running time when n is large. For instance, when n=150,000, it only takes 194 seconds to construct a partially observed similarity matrix with $[20n\log n]$ observed entries, more than 500 times faster than generating a full similarity matrix.

Given the partially observed similarity matrix $\tilde{\mathbf{A}}$, our second step aims to learn a new feature representation matrix \mathbf{X} . Instead of first completing the full similarity matrix \mathbf{A} and then factorize it, we propose an efficient symmetric factorization algorithm that is able to directly factorize the partially observed similarity matrix $\tilde{\mathbf{A}}$, i.e., find a $\mathbf{X} \in \mathcal{R}^{n \times d}$ that minimizes the following optimization problem

$$\min_{\mathbf{X} \in \mathcal{R}^{n \times d}} \|P_{\Omega} \left(\tilde{\mathbf{A}} - \mathbf{X} \mathbf{X}^T \right)\|_F^2, \tag{6}$$

where $P_{\Omega}: \mathcal{R}^{n \times n} \to \mathcal{R}^{n \times n}$ is a projection operator defined as

$$[P_{\Omega}(\mathbf{B})]_{ij} = \begin{cases} \mathbf{B}_{ij} & \text{if } \Omega_{ij} = 1\\ 0 & \text{if } \Omega_{ij} = 0. \end{cases}$$
 (7)

The objective function (6) does not have a regularization term since it already bounds the Frobenius norm of **X**. Despite its relatively simple formulation, solving problem (6) is non-trivial since its objective function is highly non-convex and the problem is NP-hard. To address this issue, we propose a very efficient optimization algorithm that solves problem (6) based on exact cyclic coordinate descent (CD). Although [Ge *et al.*, 2016] shows that the symmetric matrix factorization problem can be solved via (stochastic) gradient descent algorithm as well, our proposed coordinate descent algorithm has the following two advantages: (i) for each iteration, our CD algorithm directly updates each coordinate to

²Indeed, we need to call the DTW algorithm at least n(n+1)/2 times to generate full similarity matrix **A**. This number is achieved when we pre-compute all the b_i , $i=1,\cdots,n$.

the optimum. Thus, we do not need to decide any optimization parameter; and (ii) by directly updating coordinates to the optimums using the most up-to-date information, our CD algorithm is highly efficient and converges at a very fast rate.

At each iteration of the exact cyclic CD method, all variables but one are fixed, and that variable is updated to its optimal value. One of the main strengths of our algorithm is its capacity to update variables in an extremely efficient way. Besides, the proposed algorithm takes the advantage of sparse observed entries in $\tilde{\bf A}$ to further reduce the computational cost. To be precise, our algorithm consists of two loops that iterate over all the entries of $\bf X$ to update their values. The outer loop of the algorithm traverses through each column of $\bf X$ by assuming all the other columns known and fixed. At the i-th iteration, it optimizes the i-th column $\bf X_{1:n,i}$ by minimizing the following subproblem

$$\|\mathbf{R} - P_{\Omega}(\mathbf{X}_{1:n,i}\mathbf{X}_{1:n,i}^T)\|_F^2, \tag{8}$$

where **R** is the residual matrix defined as $\mathbf{R} = P_{\Omega}(\tilde{\mathbf{A}} - \sum_{j \neq i} \mathbf{X}_{1:n,j} \mathbf{X}_{1:n,j}^T)$.

In the inner loop, the proposed algorithm iterates over each coordinate of the selected column and updates its value. Specifically, when updating the j-th entry \mathbf{X}_{ji} , we solve the following optimization problem

$$\min_{\mathbf{X}_{ji}} \|\mathbf{R} - P_{\Omega}(\mathbf{X}_{1:n,i}\mathbf{X}_{1:n,i}^{T})\|_{F}^{2}$$

$$\iff \min_{\mathbf{X}_{ji}} \|\mathbf{R}\|_{F}^{2} - 2\langle \mathbf{R}, P_{\Omega}(\mathbf{X}_{1:n,i}\mathbf{X}_{1:n,i}^{T})\rangle$$

$$+ \|P_{\Omega}(\mathbf{X}_{1:n,i}\mathbf{X}_{1:n,i}^{T})\|_{F}^{2}$$

$$\iff \min_{\mathbf{X}_{ji}} \mathbf{X}_{ji}^{4} + 2(\sum_{k \in \Omega_{j}, k \neq j} \mathbf{X}_{ki}^{2} - \mathbf{R}_{jj})\mathbf{X}_{ji}^{2}$$

$$- 4(\sum_{k \in \Omega_{j}, k \neq j} \mathbf{X}_{ki}\mathbf{R}_{jk})\mathbf{X}_{ji} + C$$

$$\iff \min_{\mathbf{X}_{ji}} \psi(\mathbf{X}_{ji}) = \mathbf{X}_{ji}^{4} + 2p\mathbf{X}_{ji}^{2} + 4q\mathbf{X}_{ji} + C,$$

where Ω_i , $i=1,\cdots,n$ contains the indices of the observed entries in the *i*-th row of matrix $\tilde{\mathbf{A}}$. C is a constant that is independent of \mathbf{X}_{ji} . Since the $\psi(\mathbf{X}_{ji})$ is a fourth-degree polynomial function, \mathbf{X}_{ji} , as will be shown later, can be updated in a very efficiently way. Algorithm 1 describes the detailed steps of the proposed exact cyclic CD algorithm.

The proposed algorithm incurs a very low computational cost in each iteration. Lines 7-11 of the algorithm can be computed in $O(n\log n)$ operations. This is because the costs of computing p and q are only proportional to the cardinality of Ω_j . Besides, the derivative $\nabla \psi(\mathbf{X}_{ji})$ is a third-degree polynomial, thus its roots can be computed in closed form. By using Cardano's method [Cardano and Witmer, 1993], the optimal solution of \mathbf{X}_{ji} can be calculated in a constant time given the computed p and q. Likewise, lines 6 and 12 of the algorithm also take $O(n\log n)$ time since matrix \mathbf{R} can be updated by only considering the observed entries. To sum up, the proposed algorithm has a per-iteration cost of $O(dn\log n)$, significantly faster than the recent matrix factorization algorithms that take at

Algorithm 1 Efficient Exact Cyclic Coordinate Descent Algorithm for Solving the Optimization Problem (6)

1: Inputs:

- $\tilde{\mathbf{A}} \in \mathbb{R}^{n \times n}$: partially observed similarity matrix
- Ω_i , $i = 1, \dots, n$: indices of the observed entries in the *i*-th row of matrix $\tilde{\mathbf{A}}$
- *I*: number of iterations
- d: dimension of features

2: Initializations:

```
\bullet \ \mathbf{X}^{(0)} \leftarrow \mathbf{0_{n \times d}}
                    • \mathbf{R} \leftarrow P_{\Omega}(\tilde{\mathbf{A}} - \mathbf{X}^{(0)}\mathbf{X}^{(0)\top}) = P_{\Omega}(\tilde{\mathbf{A}})
   3: for t = 1, \dots, I do
4: \mathbf{X}^{(t)} \leftarrow \mathbf{X}^{(t-1)}
   5:
                   for i=1,\cdots,d do
                          \begin{aligned} \mathbf{R} \leftarrow \mathbf{R} + P_{\Omega}(\mathbf{X}_{1:n,i}^{(t)}\mathbf{X}_{1:n,i}^{(t)\top}) \\ \mathbf{for} \ j = 1, \cdots n \ \mathbf{do} \end{aligned}
   6:
   7:
                                   p \leftarrow \sum_{k \in \Omega_j}^{-1} \mathbf{X}_{ki}^{(t)2} - \mathbf{X}_{ji}^{(t)2} - \mathbf{R}_{jj}
                                    q \leftarrow -\sum_{k \in \Omega_j} \mathbf{X}_{ki}^{(t)} \mathbf{R}_{jk} + \mathbf{X}_{ji}^{(t)} \mathbf{R}_{jj}
   9:
                           \mathbf{X}_{ji}^{(t)} \leftarrow \operatorname{argmin}\{\mathbf{X}_{ji}^{t} + 2p\mathbf{X}_{ji}^{2} + 4q\mathbf{X}_{ji}\}
end for
\mathbf{R} \leftarrow \mathbf{R} - P_{\Omega}(\mathbf{X}_{1:n,i}^{(t)}\mathbf{X}_{1:n,i}^{(t)\top})
10:
11:
12:
13:
                    end for
14: end for
15: Output: \mathbf{X}^{(I)}
```

least $O(dn^2)$ time in each iteration [Vandaele *et al.*, 2015; Yu *et al.*, 2012].

In addition to a low per-iteration cost, we also expect that the proposed algorithm yields a fast convergence. This is because our algorithm always uses the newest information to update variables and each variable is updated to the optimum in a single step. This hypothesis is verified by a convergence test conducted on the UCR Non-Invasive Fetal ECG Thorax1 testbed [Chen et al., 2015]. This testbed contains a total of 3,765 time series with a length of 750. In this test, we generate a full similarity matrix A by computing the DTW similarities between all the time series pairs, and then randomly sample $[20n \log n]$ of its entries to generate a partially observed matrix A. We call the proposed algorithm to factorize matrix $\tilde{\mathbf{A}}$ by setting d=30. To measure the performance of the proposed method, we compute two error rates, i.e., the observed error $\|P_{\Omega}(\tilde{\mathbf{A}} - \mathbf{X}\mathbf{X}^T)\|_F / \|P_{\Omega}(\tilde{\mathbf{A}})\|_F$ and the underlying true error $\|\mathbf{A} - \mathbf{X}\mathbf{X}^T\|_F / \|\mathbf{A}\|_F$, at each iteration. Figure 1 shows how they converge as a function of time. This figure clearly demonstrates that the proposed exact cyclic CD algorithm converges very fast – it only takes 1 second and 8 iterations to converge. Besides, the construction accuracy is also very encouraging. The observed error and the underlying true error rates are close to each other and both of them are only about 0.1%. This result not only indicates that the inner products of the learned features well approximate the pairwise DTW similarities of the raw time series, but also

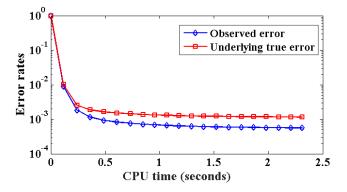


Figure 1: Two error rates as a function of CPU time on UCR Non-Invasive Fetal ECG Thorax1 data set

verifies that we can learn accurate enough features by only computing a small portion of pairwise similarities. In addition, this test validates the low-rank assumption. It shows that a $3,765\times3,765$ DTW similarity matrix can be accurately approximated by a rank 30 matrix.

4 Experiments

In this section, we evaluate the proposed framework, i.e., Similarity PreservIng RepresentAtion Learning (SPIRAL for short), on both classification and clustering tasks. For both tasks, we learn new feature representations by setting the number of features d=30, the number of iterations I=20, and the sample size $|\Omega| = [20n \log n]$. We then fed the learned features into some static models and compare them with the state-of-the-art time series classification and clustering algorithms. In our experiments, we set a same DTW window size min(40, [average time series length/10]) for all the DTW-related algorithms evaluated here. Although we can improve the DTW performance by tuning a problemdependent warping window size [Mueen and Keogh, 2016], we skip this step in our experiments to ensure that the superior performance is achieved by the proposed representation learning framework instead of tuning the DTW algorithm. All the results were obtained on a Linux server with an Intel Xeon 2.40 GHz CPU and 256 GB of main memory.

4.1 Classification Task

Six real-world time series data sets are used in our classification analysis. As a research institute, we have partnered with one of the world's largest online brokers on a challenge problem of predicting its clients' propensity of trading options. We are provided with a historical records data of 64, 523 sampled clients with the lengths of the time series range from 1 month to 67 months. The data contains 76 dynamic attributes with none of them related to option trading. Given this data, our task is to predict whether the clients will trade options in the next 3 months. Besides, we also conduct experiments on 5 benchmark data sets, i.e., FordA, ECG5000, PhalangesOutlinesCorrect (POC), ItalyPowerDemand (IPD), and HandOutlines (HO), from the UCR Time Series Repository [Chen et al., 2015]. The data sets used in our experiments have widely

Table 1: Statistics of classification data sets

Data sets	# training	# testing	length of	
	time series	time series	time series	
Online Broker	51,618	12,905	1 – 67	
ECG5000	1,320	3,601	500	
FordA	1,320	3,601	500	
POC	1,800	858	80	
IPD	67	1,029	24	
НО	370	1,000	2,709	

Table 2: Average AUC Scores of the classification algorithms

Data sets	SPIRAL	SPIRAL	NN	LSTM
	-XGB	-LR	-DTW	LSTW
Online Broker	0.84	0.83	0.76	0.79
ECG5000	0.91	0.85	0.76	0.76
FordA	0.73	0.66	0.56	0.59
POC	0.77	0.66	0.68	0.66
IPD	0.94	0.97	0.95	0.91
НО	0.87	0.81	0.77	0.61

varying sizes and encompass multiple domains such as finance, healthcare, and manufacture. Table 1 summarizes the statistics of these data sets.

Given the feature representations learned by the proposed SPIRAL framework, we fed them into the algorithms of XGBoost [Chen and Guestrin, 2016] and ℓ_2 -regularized logistic regression that is implemented in LibLinear [Fan et al., 2008]. These approaches, denoted as SPIRAL-XGB and SPIRAL-LR, respectively, are compared with the state-of-theart time series classification algorithms NN-DTW [Wang et al., 2013] and Long Short-Term Memory (LSTM) [Schmidhuber, 2015]. For the XGBoost algorithm, we use the default parameters specified in its source code. The parameter of logistic regression is determined by a 5-fold cross validation.

We use AUC (area under the ROC Curve) to measure the classification performance. All the experiments in this study are repeated five times, and the AUCs averaged over the five trials are reported in Table 2. From this table, we first observe that XGBoost and logistic regression algorithms that are fed with our learned features outperform the two baseline algorithms on all the data sets. In particular, the method SPIRAL-XGB yields the best performance on five out of six data sets, and the method SPIRAL-LR performs the best on the other data set. More encouragingly, SPIRAL-LR can achieve better classification performance than NN-DTW and LSTM on almost all the data sets. This verifies that the feature representation learned by the proposed method is powerful – even classical models such as logistic regression can achieve satisfactory performance on it.

On the other hand, although LSTM has been successfully applied to a number of sequence prediction tasks, it fails to deliver strong performance in our empirical study. We conjecture that this is because the data sets are not large enough for LSTM to model complex structures. Besides, NN-DTW is usually a very strong baseline that is considered as hard-

to-beat in the literature [Xi et al., 2006; Wang et al., 2013; Mueen and Keogh, 2016]. However, it also yields an overall worse performance than SPIRAL-LR and SPIRAL-XGB. One possible reason is that NN-DTW uses 1-nearest neighbor classifier, thus is sensitive to noises and outliers. This observation shows another advantage of the proposed method: by learning a feature representation instead of directly developing a time series model, our method is more flexible and can exploit the strengths of different learning algorithms.

In addition to the superior performance, SPIRAL-LR and SPIRAL-XGB are very efficient as well. They have a significantly lower running time than that of the baseline algorithms. For example, it takes NN-DTW more than 2 days to classify the FordA data set while the running time for SPIRAL-LR and SPIRAL-XGB is only about 5 minutes.

4.2 Clustering Task

Similar to time series classification, time series clustering is also an important task that has found numerous applications. To further test our learned features on the clustering task, we conduct experiments on another 7 UCR time series data sets. Since data clustering is an unsupervised learning task, we merge their training and testing sets together. Statistics of these data sets are summarized in Table 3.

We fed the features learned by the proposed framework into the kMeans algorithm as our clustering method, and compare it to the state-of-the-art time series clustering algorithm k-Shape [Paparrizos and Gravano, 2015], which has been shown to outperform many state-of-the-art partitional, hierarchical, and spectral clustering approaches. Besides, we also compare our method with clustering algorithms kMeans-DTW and CLDS [Li and Prakash, 2011] since our ideas are similar in some respects. kMeans-DTW is a popular time series clustering algorithm that uses DTW algorithm to measure pairwise distances between data points. Although it looks similar to the idea of our SPIRAL-kMeans that also utilizes the DTW and kMeans algorithms, it is less desirable than SPIRAL-kMeans mainly because: (i) kMeans-DTW suffers from a very high computational cost since it needs to compute the pairwise DTW distances between all the time series and all the cluster centers at each iteration; and (ii) the DTW distance does not satisfy the triangle inequality, thus can make the cluster centers computed by averaging multiple time series drift out of the cluster [Niennattrakul and Ratanamahatana, 2007]. By designing an efficient algorithm that only needs to call the DTW function $O(n \log n)$ times and by embedding time series data to the Euclidean space while preserving their original similarities, the proposed method SPIRAL successfully addresses both these issues. Similar to the idea of representation learning presented in this paper, CLDS also learns a new feature representation of the time series data, and then partition the learned representation via an EM-like algorithm.

In this study, we use the normalized mutual information (NMI for short) to measure the coherence between the inferred clustering and the ground truth categorization. NMI scales from 0 to 1, and a higher NMI value implies a better partition. Each experiment is repeated five times, and the performance averaged over the five trials is reported in Ta-

Table 3: *Statistics of clustering data sets*

Data sets	number of	length of	number of
	time series	time series	clusters
Swedish Leaf	1,125	128	15
Cricket_X	780	300	12
uWGLX	4,478	315	8
50words	905	270	50
SLC	9,236	1,024	3
SC	600	60	6
ED	16,637	96	7

Table 4: Average NMI Scores of all the clustering algorithms. N/A indicates that the clustering task cannot be completed due to memory limitations..

Data sets	SPIRAL -kMeans	k-Shape	CLDS	kMeans -DTW
Swedish Leaf	0.61	0.56	0.60	0.50
Cricket_X	0.32	0.31	0.26	0.22
uWGLX	0.47	0.38	0.28	0.32
50words	0.68	0.50	0.56	0.53
SLC	0.60	0.56	0.41	N/A
SC	0.80	0.78	0.51	0.70
ED	0.35	0.23	0.2	0.24

ble 4. Compared to all the baseline algorithms, our clustering method SPIRAL-kMeans yields the best performance on all the seven data sets, indicating that it delivers the state-of-theart performance. In addition, SPIRAL-kMeans has a significantly lower running time than all the baseline clustering algorithms evaluated here. For instance, clustering the ED data set takes k-Shape, CLDS, and kMeans-DTW 20, 57, and 114 minutes, respectively. As a comparison, our clustering algorithm SPIRAL-kMeans only spends one minute and a half to partition the ED data set.

We finally note that kMeans is just one choice of clustering algorithms that can take our learned features as the input. By replacing it with some more advanced clustering algorithms, it is expected to achieve an even better clustering performance.

5 Conclusions

In this paper, we propose a similarity preserving representation learning framework for time series analysis. Given a set of n time series, the key idea is to first generate a partially observed similarity matrix with $O(n\log n)$ observed DTW similarities, and then factorize this matrix to learn a new feature representation. To this end, we propose the first optimization parameter free algorithm for factorizing partially-observed symmetric matrix. The proposed algorithm updates variables via exact cyclic coordinate descent, and incurs a very low computational cost in each iteration. The feature representation learned by the proposed framework preserves the DTW similarities of the raw time series data, and is general enough to be applied to a variety of learning problems. Our empirical studies on both classification and clustering tasks verify the

effectiveness and efficiency of the proposed method.

References

- [Adams, 2004] Colin Conrad Adams. *The knot book: an elementary introduction to the mathematical theory of knots.* American Mathematical Soc., 2004.
- [Basu and Michailidis, 2015] S. Basu and G. Michailidis. Low-rank and sparse modeling of high-dimensional vector autoregressions. 2015.
- [Cardano and Witmer, 1993] G. Cardano and T R. Witmer. *Ars magna or the rules of algebra*. 1993.
- [Celebi *et al.*, 2013] S. Celebi, A. Aydin, T. Temiz, and T. Arici. Gesture recognition using skeleton data with weighted dynamic time warping. In *VISAPP* (1), pages 620–625, 2013.
- [Chan and Fu, 1999] Kin-Pong Chan and Ada Wai-Chee Fu. Efficient time series matching by wavelets. In *ICDE*, pages 126–133. IEEE, 1999.
- [Chen and Guestrin, 2016] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. *arXiv* preprint *arXiv*:1603.02754, 2016.
- [Chen et al., 2015] Y. Chen, E. Keogh, B. Hu, N. Begum, A. Bagnall, A. Mueen, and G. Batista. The ucr time series classification archive, July 2015. www.cs.ucr.edu/~eamonn/time_series_data/.
- [Deng et al., 2013] H. Deng, G. Runger, E. Tuv, and M. Vladimir. A time series forest for classification and feature extraction. *Information Sciences*, 239:142–153, 2013.
- [Efrat *et al.*, 2007] A. Efrat, Q. Fan, and S. Venkatasubramanian. Curve matching, time warping, and light fields: New algorithms for computing similarity between curves. *Journal of Mathematical Imaging and Vision*, 27(3):203–216, 2007.
- [Elman, 1990] Jeffrey L Elman. Finding structure in time. *Cognitive science*, 14(2):179–211, 1990.
- [Erdős *et al.*, 2009] L. Erdős, B. Schlein, and H. Yau. Local semicircle law and complete delocalization for wigner random matrices. *Communications in Mathematical Physics*, 287(2):641–655, 2009.
- [Faloutsos et al., 1994] C. Faloutsos, M. Ranganathan, and Y. Manolopoulos. Fast subsequence matching in timeseries databases, volume 23. ACM, 1994.
- [Fan et al., 2008] R. Fan, K. Chang, C. Hsieh, X. Wang, and C. Lin. Liblinear: A library for large linear classification. JMLR, 9(Aug):1871–1874, 2008.
- [Fulcher and Jones, 2014] Ben D Fulcher and Nick S Jones. Highly comparative feature-based time-series classification. *IEEE Transactions on Knowledge and Data Engineering*, 26(12):3026–3037, 2014.
- [Ge *et al.*, 2016] R. Ge, J. Lee, and T. Ma. Matrix completion has no spurious local minimum. In *NIPS*, pages 2973–2981, 2016.

- [Keogh *et al.*, 2001] E. Keogh, K. Chakrabarti, M. Pazzani, and S. Mehrotra. Dimensionality reduction for fast similarity search in large time series databases. *Knowledge and information Systems*, 3(3):263–286, 2001.
- [Längkvist et al., 2014] M. Längkvist, L. Karlsson, and A. Loutfi. A review of unsupervised feature learning and deep learning for time-series modeling. Pattern Recognition Letters, 42:11–24, 2014.
- [Li and Prakash, 2011] Lei Li and B Aditya Prakash. Time series clustering: Complex is simpler! In *ICML*, pages 185–192, 2011.
- [Lin et al., 2007] J. Lin, E. Keogh, L. Wei, and S. Lonardi. Experiencing sax: a novel symbolic representation of time series. *Data Mining and Knowledge Discovery*, 15(2):107–144, 2007.
- [Marčenko and Pastur, 1967] V. Marčenko and L. Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of the USSR-Sbornik*, 1(4):457, 1967.
- [Muda et al., 2010] L. Muda, M. Begam, and I Elamvazuthi. Voice recognition algorithms using mel frequency cepstral coefficient (mfcc) and dynamic time warping (dtw) techniques. arXiv preprint arXiv:1003.4083, 2010.
- [Mueen and Keogh, 2016] A. Mueen and E. Keogh. Extracting optimal performance from dynamic time warping. In *SIGKDD*, pages 2129–2130. ACM, 2016.
- [Müller, 2007] Meinard Müller. Dtw-based motion comparison and retrieval. *Information Retrieval for Music and Motion*, pages 211–226, 2007.
- [Nanopoulos *et al.*, 2001] A. Nanopoulos, R. Alcock, and Y. Manolopoulos. Feature-based classification of timeseries data. *International Journal of Computer Research*, 10(3):49–61, 2001.
- [Niennattrakul and Ratanamahatana, 2007] V. Niennattrakul and C. Ratanamahatana. Inaccuracies of shape averaging method using dynamic time warping for time series data. In *International conference on computational science*, pages 513–520. Springer, 2007.
- [Paparrizos and Gravano, 2015] J. Paparrizos and L. Gravano. k-shape: Efficient and accurate clustering of time series. In *SIGMOD*, pages 1855–1870. ACM, 2015.
- [Rakthanmanon *et al.*, 2013] T. Rakthanmanon, B. Campana, A. Mueen, G. Batista, B. Westover, Q. Zhu, J. Zakaria, and E. Keogh. Addressing big data time series: mining trillions of time series subsequences under dynamic time warping. *TKDD*, 7(3):10, 2013.
- [Schmidhuber, 2015] Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117, 2015.
- [Sempena *et al.*, 2011] S. Sempena, N. Maulidevi, and P. Aryan. Human action recognition using dynamic time warping. In *International Conference on Electrical Engineering and Informatics*, pages 1–5. IEEE, 2011.
- [Stock and Watson, 2005] J. Stock and M. Watson. Implications of dynamic factor models for var analysis. Technical report, National Bureau of Economic Research, 2005.

- [Sun and Luo, 2015] Ruoyu Sun and Zhi-Quan Luo. Guaranteed matrix completion via nonconvex factorization. In *FOCS*, pages 270–289. IEEE, 2015.
- [Vandaele *et al.*, 2015] Arnaud Vandaele, Nicolas Gillis, Qi Lei, Kai Zhong, and Inderjit S. Dhillon. Coordinate descent methods for symmetric nonnegative matrix factorization. *CoRR*, abs/1509.01404, 2015.
- [Vial et al., 2009] Jérôme Vial, Hicham Noçairi, Patrick Sassiat, Sreedhar Mallipatu, Guillaume Cognon, Didier Thiébaut, Béatrice Teillet, and Douglas N Rutledge. Combination of dynamic time warping and multivariate analysis for the comparison of comprehensive two-dimensional gas chromatograms: application to plant extracts. *Journal of Chromatography A*, 1216(14):2866–2872, 2009.
- [Wang et al., 2006] X. Wang, K. Smith, and R. Hyndman. Characteristic-based clustering for time series data. *Data mining and knowledge Discovery*, 13(3):335–364, 2006.
- [Wang et al., 2013] X. Wang, A. Mueen, H. Ding, G. Trajcevski, P. Scheuermann, and E. Keogh. Experimental comparison of representation methods and distance measures for time series data. *Data Mining and Knowledge Discovery*, 26(2):275–309, 2013.
- [Xi et al., 2006] X. Xi, E. Keogh, C. Shelton, L. Wei, and C. A. Ratanamahatana. Fast time series classification using numerosity reduction. In *ICML*, pages 1033–1040, 2006.
- [Yang and Wu, 2006] Qiang Yang and Xindong Wu. 10 challenging problems in data mining research. *International Journal of Information Technology & Decision Making*, 5(04):597–604, 2006.
- [Yu *et al.*, 2012] H. F. Yu, C.J. Hsieh, S. Si, and I. Dhillon. Scalable coordinate descent approaches to parallel matrix factorization for recommender systems. In *ICDM*, pages 765–774, 2012.