# A pattern representation of stock time series based on DTW

Tian Han, Qinke Peng *, Zhibo Zhu, Yiqing Shen, Huijun Huang,
Nahiyoon Nabeel Abid

*Systems Engineering Institute, Xi'an Jiaotong University, Xi'an, China*

## ARTICLE INFO

## ABSTRACT

Time series analysis based on pattern discovery has received a lot of interests in the fields of economic physics and machine learning due to its simplicity and ability to reveal complex nonlinear behavior in stock market. Dynamic Time Warping (DTW) is a useful tool to extract morphological characteristics of time series for its capacity to cope with time shifts and warpings. In this paper, we propose a new time series representation method for stock time series based on dynamic time warping (DTW) called PR-DTW. A combinatorial optimization model with strict constraints is built to get the pattern representation of stock time series. To simplify the calculation, we construct another unconstrained global optimization problem whose optimal solution includes the optimal solution of the original combinatorial optimization problem based on a theorem proved in this paper. Particle Swarm Optimization algorithm is used to solve the global optimization problem, then the results can be converted into the optimal solution of the combinatorial optimization problem through a few simple formulas given in the theorem. The results of three classifiers (1NN, Decision Tree, Multi-layer Perceptron) implemented on 15 sectors in Chinese A-share market unanimously demonstrate that PR-DTW has the capability of extracting time series short-term patterns which is widely regarded as difficulty. And we conclude that PR-DTW has the capability of prevention of End Effect, anti-noise and segmentation. Moreover, by extracting the top ten patterns predicting stock's rise and fall in short term (10 days) according to the ranking of stock's rising probability in the next three days, we find out the short-term patterns obtained by PR-DTW have prospective directive to the stock trend analysis in short term.

© 2020 Elsevier B.V. All rights reserved.

## 1. Introduction

In the analysis of financial time series, technicians believe that recognizable (and predictable) price developing on a pattern due to investor behavior repeats itself so often as a proverb called "History tends to repeat itself". The study of the patterns, which assists stock investors to make reasonable decision [1,2], has been proposed as a technique for evaluating the trend of a given stock time series [3,4]. As one of the seminal works, 'Technical Analysis of Stock Trends' published by Robert D. Edwards and John Magee focuses on trend analysis and remains in use to present [5]. The wave theory named after Ralph Elliot, who contended that stock market tends to move in discernible and predictable patterns, is another significant theory for traders to analyze stock trend [6]. In details, 53 well-known stock patterns concluded by Bulkowski [7] have a great effect on the research of stock pattern, and they are compiled with comprehensive formal specification by Yuqing Wan in 2018 [8]. However, these patterns are confined to the traders' priori knowledge of stock

---

* Corresponding author.
  *E-mail address:* qkpeng@mail.xjtu.edu.cn (Q. Peng).

price in history. Meanwhile, for the noise in stock market, it is very hard to extract the short-term patterns to guide investment in stock market. Therefore, in order to support short-term investment, there is a need to present a suitable representation method to find meaningful patterns in short-term stock time series.

There are various existing representation methods for time series [9]. Generally speaking, the representations of time series can be divided into three categories: generative model based representations, transformation based representations and time-domain based representations [10]. Generative model based representations use the parameters in different models to represent the time series, including hidden Markov model based representation [11–13], Bayesian network based representation [14–16] and deep learning based representation [17,18], etc. Transformation based representations, aiming at transferring original data into the feature space to represent time series, includes discrete Fourier Transform (FT) [19–21], Discrete Wavelet Transform (DWT) [22–24], Discrete Cosine Transform (DCT) [25–27] and Symbolic Aggregate Approximation (SAX) [28], etc. Time-domain based representations, whose core is to obtain representative points from the original time series, contains Piecewise Linear Approximation (PLA) [29], Adaptive Piecewise Constant Approximation(APCA) [30], Piecewise Aggregate Approximation (PAA) [31], Piecewise Vector Quantized Approximation (PVQA) [32] and Indexable Piecewise Linear Approximation (IPLA) [33], etc. Compared with the other two categories of representation methods, time-domain based representations, as a kind of simplest and most effective representations, can maintain various properties such as fluctuations, shapes and trends in the patterns to provide more information [34].

Reducing the dimension by preserving the salient points is a promising idea for time-domain based representation methods and has been widely used in financial applications [35]. However, considering the high volatility and risk in stock market, the representative points chosen from stock time series should be considered from a global perspective rather than local perspective. In other words, the importance of a point not only depends on its value, but also its representativeness throughout the whole time series [36,37]. As a method measuring similarity of time series with character of pattern expression, dynamic time warping (DTW) can be used to evaluate the quality of different representation methods. The shorter distance between the representation sequence and the original sequence, the better the representation methods are. Inspired by DTW, we present a pattern representation method based on DTW called PR-DTW aiming to find a subsequence from original time series which is closest to original sequence in DTW measurement. Thus a combinatorial optimization model is built to get the pattern representation. Considering the strict constraints in solving combinatorial optimization problem, we provide a fast calculation method solved by particle swarm optimization algorithm to adapt to the real-time analysis of stock market. In order to evaluate the effectivity of our representation, three classifiers (1NN, Decision Tree, Multi-layer Perceptron) are used to experiment on the stocks in 15 sectors in the A-share market. The results demonstrate that PR-DTW is an effective method for extracting the pattern of time series from global perspective. Furthermore, it is able to overcome the impact of noise in time series and maintain patterns on the time scale. Finally, we sum up three advantages of PR-DTW and make statistical analysis on 20 typical patterns from the whole A-share market to conduct statistical analysis.

The rest of paper is organized as follows: The related works are briefly reviewed in the next section. Our proposed method and its fast approximation are described in Section 3. Section 4 presents the experimental results to evaluate the effectiveness of PR-DTW. Section 5 analyzes 20 typical patterns obtained by PR-DTW in A-share Market. Finally we conclude this paper in Section 6.

## 2. Related work

In this section, we overview previous work on the time-domain based representation methods and present the calculation of DTW.

### 2.1. Time-domain based representation methods

Representing a time series by the significant points is one of the basic ideas for time-domain based representation methods. The selected salient points may contribute critically on the overall shape of the time series [38]. In the stock market, the common technical patterns are typically characterized by some important points such as Head-and-Shoulders, Double Tops and Pennants [39]. The identification of the important points was first introduced by Chung et al. and had raised strong interests due to its effectiveness and interpretability [40]. Man and Wong proposed a lattice structure to represent the identified peaks and troughs (called control points) in the time series [41]. Pratt and Fink compressed the time series by selecting only certain important extrema [42]. A high-level representation based on a sequence of critical points were proposed by Bao and Yang for financial data analysis [43]. Cun Ji presented a piecewise linear representation based on importance data points for time series data called PLR_IDP, which can hold the main characteristic with small fitting error of segments and single points [44]. Yu Wei proposed Important Turning Points Set (ITPS) to represent motion datas and had received efficient results [45].
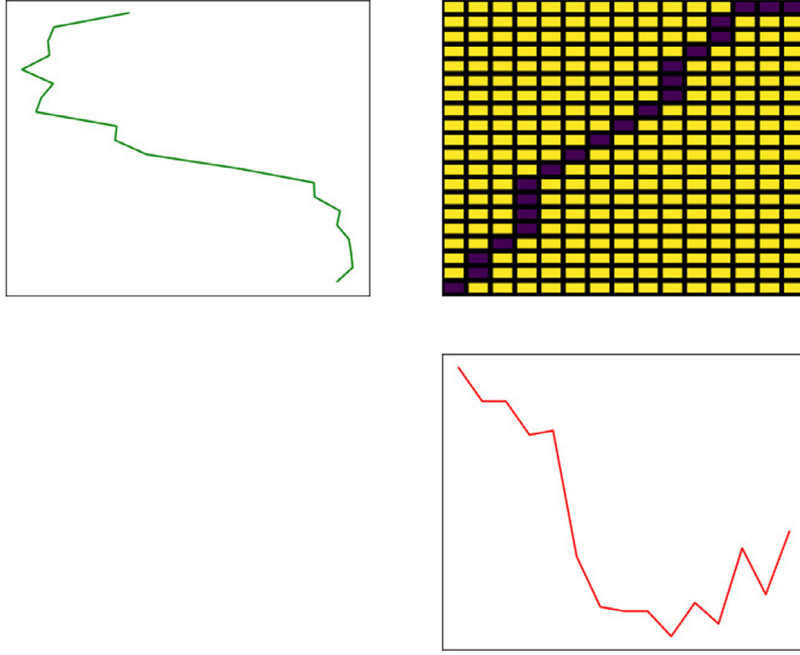
**Fig. 1.** The optimal warping path in distance matrix of two stock time series.

## 2.2. DTW

DTW is a kind of elastic measurement for calculating time series similarity [46,47]. Suppose there are two time series $S = \{s_1, s_2, \ldots, s_m\}$ and $Q = \{q_1, q_2, \ldots, q_n\}$ with the length of $m$ and $n$, DTW minimizes the distance between two time series by constructing an optimal warping path $P$ which shows the corresponding relationship between each other. A warping path denoted as $P = \{p_1, p_2, \ldots, p_n\}$ means the mapping information derived from the two time series. The parameter $K \in [\max(m, n), m + n - 1]$ is an integer representing length of the path, and $d(p_k)$ denotes the measurement between two points $s_i$ and $q_j$. Here we choose $d(p_k) = d(i, j) = |s_i - q_j|, i \in [1, m], j \in [1, n]$. At the same time, the warping path must satisfy at least three constraints named boundary conditions, continuity and monotonicity.

- Boundary condition: $p_1 = (1, 1)$ and $p_K = (m, n)$;
- Continuity: Given $p_k = (a, b)$ and $p_{k-1} = (a', b')$ where $a - a' \leq 1$ and $b - b' \leq 1$;
- Monotonicity: Given $p_k = (a, b)$ and $p_{k-1} = (a', b')$ where $a - a' \geq 0$ and $b - b' \geq 0$;

The path we want is an optimal solution with minimal warping cost, i.e. The calculation formula of DTW is as follows.

$$DTW(S, Q) = \min \sum_{k=1}^{K} d(p_k) \tag{1}$$

The best warping path can be found by using dynamic programming, which calculates the cumulative distance $CD(i, j)$ with distance $d(i, j)$ by adding the minimum of the cumulative distance restricted by three adjacent elements, i.e.

$$CD(i, j) = d(i, j) + \min(CD(i, j - 1), CD(i - 1, j), CD(i - 1, j - 1)) \tag{2}$$

Where $CD(0, 0) = 0, CD(i, 0) = CD(0, j) = \infty$. Fig. 1 shows two different time series' corresponding relationship when calculating DTW.

## 3. Pattern representation of stock time series based on DTW

In this section, we introduce the main contents of pattern representation based on DTW

### 3.1. PR-DTW

Suppose we have a time series $Q = \{q_1, q_2, \ldots, q_n\}$ and an integer $m \in [1, n]$ representing the length of extracted pattern. In PR-DTW, we aim to extract $m$ points from the normalized original time series as the pattern, which has the

shortest DTW distance to the original normalized time series. The pattern is denoted as $Q' = \left\{ q_{t_1}, q_{t_2}, \ldots, q_{t_m} \right\}$ and the optimization model of PR-DTW can be written as

$$\min_{Q'} f\left(Q'\right) = DTW\left(Q, Q'\right) \tag{3}$$

$$s.t. \begin{cases} Q' \subseteq Q \\ t_i < t_j, \forall 1 \leq i \leq j \leq m \end{cases}$$

Since the optimization model is a kind of NP-hard combination optimization problem, we want to relax it into an unconstrained global optimization model to simplify the calculation. Therefore, we propose a fast PR-DTW to find the optimal solution in a reasonable computing time. The details of calculation is introduced in 3.2.

### 3.2. Fast PR-DTW calculation

In this section, we propose a method to obtain a fast PR-DTW solution to meet the practical requirements.

Firstly, we construct another global optimization model which solution has a close connection to that of (3). Suppose a time series $Q'' = \{x_1, x_2, \ldots, x_m\}$ $(m < n)$, the constructed optimization model is as follows:

$$\min_{Q''} f\left(Q''\right) = DTW\left(Q, Q''\right) \tag{4}$$

$$s.t. \quad Q'' \subseteq R$$

Different from (3), here we want to find a pattern $Q'' = \{x_1, x_2, \ldots, x_m\}$ $(m < n)$ in a global perspective which is closest to the original time series $Q$ in DTW measurement.

Secondly, we give a theorem on the relationship of the optimal solution of (3) and that of (4).

**Theorem.** *The optimal solution of* (3) *is a subset of the optimal solution of* (4)

The proof process is as follows:

- Step1: For $m = 1$, the result of the optimization problem (4) is determined by the middle points of the sequence generated by the ascending order of $Q$.

Supposing $Q_{rank} = \left\{ q_{rank_1}, q_{rank_2}, \ldots, q_{rank_n} \right\}$ is a sequence whose points are sorted in ascending order from original series $Q$.

For $m = 1$, in this case $Q'' = \{x_1\}$, Eq. (4) can be written as:

$$\min_{x_1} f(x_1) = \sum_{i=1}^{n} |x_1 - q_i| = \sum_{i=1}^{n} \left| x_1 - q_{rank_i} \right| \tag{5}$$

While, $x_1 < q_{rank_1}$, then (5) can be rewritten as:

$$\min_{x_1} f(x_1) = (-n) x_1 + \sum_{i=1}^{n} q_{rank_i} \tag{6}$$

It is obvious that $f(x_1)$ is monotone decreasing.

Then considering $x_1 \geq q_{rank_1}$:

If there are $h$ points in $Q$ should satisfy the following formula.

$$x_1 - q_{rank_i} > 0, \quad h \in [1, n] \tag{7}$$

And (5) can be rewritten as:

$$\min_{x_1} f(x_1) = (2h - n) x_1 + \left( \sum_{i=h+1}^{n} q_{rank_i} - \sum_{i=1}^{h} q_{rank_i} \right) \tag{8}$$

Then we discuss (8) in two situations.

(1) When $n$ is odd:

- (a) When $x_1 < q_{rank_{(n+1)/2}}$: $0 \leq h < (n+1)/2$. We can see that $2h - n < 0$, so $f(x_1)$ is monotone decreasing.
- (b) When $x_1 > q_{rank_{(n+1)/2}}$: $(n+1)/2 < h \leq n$ We can see that $2h - n > 0$ so $f(x_1)$ is monotone increasing.

In this situation, $f(x_1)$ gets minimum value while $x_1 = q_{rank_{(n+1)/2}}$.

(2) When $n$ is even:

- (a) When $x_1 \leq q_{rank_{n/2}}$: $0 \leq h < n/2$. We can see that $2h - n < 0$, so $f(x_1)$ is monotone decreasing.

(b) When $q_{rank_{n/2}} < x_1 < q_{rank_{(n/2)+1}}$: $h = n/2$. We can see that $2h - n = 0$, so $f(x_1)$ is a constant.

(c) When $x_1 \geq q_{rank_{(n/2)+1}}$: $n/2 < h \leq n$. We can see that $2h - n > 0$ so $f(x_1)$ is monotone increasing.

In this situation, $f(x_1)$ gets minimum value while $x_1 \in \left[ q_{rank_{n/2}}, q_{rank_{(n/2)+1}} \right]$

So for $m = 1$, the sequence $Q'' = \{x_1\}$ which is closest to the original time series $Q$ in DTW measurement, is determined by the following formula:

$$x_1 = \begin{cases} q_{rank_{(n+1)}} \\ \\ t \quad t \in \left[ q_{rank_{n/2}}, q_{rank_{(n/2)+1}} \right] \end{cases} \tag{9}$$

From formula (9) we can see that the points we get are determined by the middle points of the sequences in ascending order of original time series $Q$. It can be explained that DTW has a good effect on retaining the pattern of time series.

• Step2: The optimal solution of (3) belong to the optimal solution of (4)

It is known that $Q''$ is the global optimum solution of (4). When we calculating the DTW distance between $Q''$ and $Q$, the distance matrix can be written as:

$$P = \begin{bmatrix} p_{11} & \cdots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{m1} & \cdots & p_{mn} \end{bmatrix} \tag{10}$$

The set of nodes contained in the optimal path can be written as:

$$optimized\_path = \{p_{11}, p_{12}, \ldots, p_{mn}\} \tag{11}$$

For each $i \in [1, m]$, we can get a corresponding points subscript set and a corresponding points set.

$$node\_index(i) = \left\{ j \,|\, p_{ij} \in optimized\_path \right\} \tag{12}$$

$$corresponding\_points\_set = \left\{ q_j \,|\, j \in node\_index(i) \right\} \tag{13}$$

Each point in $Q'' = \{x_1, x_2, \ldots, x_m\}$ corresponds to a set of continuous partitioning points in the original time series $Q$. Then the DTW distance between $Q''$ and $Q$ can be written as:

$$DTW\left(Q'', Q\right) = \sum_{i=1}^{m} DTW\left(x_i, corresponding\_points\_set(i)\right) \tag{14}$$

The $corresponding\_points\_set(i)$ is a points set corresponding to each point $x_i$ in the optimal solution. For each $corresponding\_points\_set(i)$, $x_i$ must satisfy the formula (9). Here we can prove the above conclusion by using reduction to absurdity. If there exists $x_i$ does not satisfy the formula (9), there must exist $x_i'$ that satisfies (9) to make $DTW\left(Q'', Q\right)$ smaller, against the hypothesis that $Q''$ is the optimal solution.

Since the solution of combinatorial optimization model in PR-DTW is a subset of that of constructed global optimization model, and we have found out a constructed solution of the combinatorial optimization model belonging to the optimal solution of constructed global optimization model. The solution we find must be the optimal solution of the combinatorial optimization. In general, the optimal solution of combinatorial optimization problem must be a subset of the optimal solution of global optimization problem.

Finally, the original combinatorial optimization problem has been relaxed into the constructed unconstrained global optimization problem. Particle Swarm Optimization (PSO), as a widely used heuristic algorithm with the advantage of fast searching speed and high efficacy, is used to solve the constructed global optimization model. In the experiments, we choose the population number as 200, the maximum iteration as 150 and the inertia weight as 0.6. After we get the global optimal solution $Q'' = \{x_1, x_2, \ldots, x_m\}$, the optimal solution of combinatorial optimization model $Q' = \left\{ q_{t_1}, q_{t_2}, \ldots, q_{t_m} \right\}$ can be constructed in this way.

$$q_{t_i} = \begin{cases} x_i, & length(corresponding\_points\_set(i)) \text{ is odd} \\ t, t \in \left\{ q_{rank_{n/2}}, q_{rank_{(n/2)+1}} \right\} & length(corresponding\_points\_set(i)) \text{ is even} \end{cases} \tag{15}$$

## 4. Evaluation experiments

In this section, we will verify the effectiveness of PR-DTW. We believe that after representation, the sequence can filter out noise and extract key information for better prediction performance. Therefore, we can evaluate time series representations according to prediction ability of different classification method. By using three classifiers (1NN, Decision Tree, Multi-layer Perceptron), we conduct experiments on stocks' closing price of nearly 500 trading days in 15 datasets from different stock sectors.

**Table 1**
The specific information of the selected stocks.

| Stock sector | Number of stocks | Number of samples |
|---|---|---|
| Glass | 17 | 485 |
| Oil exploitation | 13 | 347 |
| Tourism service | 14 | 384 |
| Dyes and coatings | 12 | 349 |
| Publish industry | 12 | 434 |
| Culture & education | 20 | 563 |
| Aviation | 18 | 515 |
| Commerce | 22 | 594 |
| Road and bridge | 18 | 356 |
| Supermarket | 10 | 276 |
| Machine tool (M_T) manufacturing | 9 | 257 |
| Housewear & furnishings (H & F) | 18 | 490 |
| Rubber | 7 | 214 |
| Fishery | 8 | 198 |
| Diversified Financial | 19 | 513 |

### 4.1. Data description

The daily close prices of 223 stocks from 21 March 2016 to 9 April 2018 are obtained from TDX software by different sectors of Chinese A-share market. The specific information of the selected stocks are shown in Table 1.

In the experiments, all stock price sequence are segmented into sub-sequences by the sliding time window with time span of 10 and time step of 10 ($n = 10$).

$$S_i(k) = \{s_i(k), s_i(k+1), \ldots, s_i(k+n-1)\} \tag{16}$$

$S_i(k)$ denotes the $k$th sub-sequence of the $i$th stock close price. For highlighting the advantages of PR-DTW's pattern extraction in short term, the return ($r_i(k)$) of $S_i(k)$ is chosen as the label of sub-sequence and the calculation is as follows.

$$r_i(k) = s_i(k+n)/s_i(k+n-1) - 1 \tag{17}$$

Then Z-score normalization is chosen as the normalization method in experiments.

Considering 0.01 is the low level of interest rate [48], the label of each data is chosen as 1 if $r_i(k) \geq 0.01$ and chosen as 0 if $r_i(k) \leq -0.01$. Therefore, stock data are transformed into sequence samples with same length and their labels. And then different represented sequence samples are obtained by different representation methods. Finally, we will evaluate time series representation methods according to prediction results of their represented sequence samples using different binary classification methods. The better representation method is, the better prediction result of represented sequence samples is. The methods are compared by their classification accuracy (ACC) rate in 15 datasets. The ACC rate is defined as:

$$ACC = Number\ of\ samples\ correctly\ classified/Total\ number\ of\ samples \tag{18}$$

### 4.2. Experimental results

Tables 2–4 show the results of three representation methods using three different classifiers based on 5-fold cross validation respectively. The patterns obtained by PR-DTW has higher classification accuracy with an appropriate length of extracted sequence (parameter $m$) in most stock sectors using three different classifiers. Thus we can conclude that PR-DTW can preserve the morphological information of the original sequence effectively and filter out the noise properly. If $m$ is appropriate, PR-DTW can improve the accuracy of computation. Meanwhile, the representation may cause damage to the pattern of the original time series when $m$ is inappropriate. That is to say, the choice of parameter $m$ has great influence on the final representation.

Compare the patterns obtained by PR-DTW with the original sequences, we find that the computing time of classifier decreases obviously for the patterns obtained by PR-DTW due to their reduction in data size. Therefore, PR-DTW can meet the requirements of real-time trading in stock market.

### 4.3. The three advantages of PR-DTW

In this section, we conclude three advantages of PR-DTW.

(1) Prevention of End Effect: As a key problem in the field of EMD [49–51], End Effect can also make damage on time series representation. For the traditional time series representation methods such as PLA and PIP, they have to add the end points of time series in the representation sequence as the initial condition, which makes great restriction and impact

**Table 2**
The classification ACC rate of the patterns obtained by PR-DTW and Original Sequence (OR) using three classifiers on different sectors.
The ACC rate of PR-DTW and Original Sequence (OR) on different sectors using three classifiers

| Sector | 1NN classifier | | Decision Tree classifier | | MLP classifier | |
|---|---|---|---|---|---|---|
| | PR-DTW (m) | OR(10) | PR-DTW (m) | OR(10) | PR-DTW (m) | OR(10) |
| Glass | **0.5567 (8)** | 0.5443 | **0.5278 (3)** | 0.5216 | **0.5484 (3)** | 0.4948 |
| Oil exploitation | 0.5591 (7) | **0.5675** | **0.5733 (6)** | 0.5447 | **0.5618 (3)** | 0.5447 |
| Tourism service | **0.5730 (6)** | 0.5312 | **0.5678 (4)** | 0.4896 | **0.5912 (7)** | 0.5471 |
| Dyes & coatings | **0.5557 (4)** | 0.5014 | **0.5645 (7)** | 0.4840 | **0.5330 (3)** | 0.5301 |
| Publish industry | **0.5737 (8)** | 0.5712 | **0.5623 (6)** | 0.5391 | **0.5530 (8)** | 0.5368 |
| Culture & education | **0.5611 (8)** | 0.5345 | 0.5487 (6) | **0.5488** | **0.5506 (7)** | 0.5041 |
| Aviation | 0.6427 (8) | **0.6834** | **0.6135 (8)** | 0.6116 | **0.5378 (6)** | 0.5048 |
| Commerce | **0.5134 (3)** | 0.5032 | **0.5741 (4)** | 0.5369 | **0.5791 (8)** | 0.5319 |
| Road and bridge | 0.5281 (3) | **0.5474** | 0.5730 (8) | **0.5760** | **0.5590 (4)** | 0.5364 |
| Supermarket | **0.5727 (8)** | 0.5361 | **0.5687 (3)** | 0.5324 | **0.5793 (4)** | 0.5181 |
| M_T manufacturing | **0.5720 (4)** | 0.4944 | **0.5917 (8)** | 0.5379 | **0.5760 (7)** | 0.5034 |
| H&F | **0.5244 (3)** | 0.5183 | **0.5999 (7)** | 0.5693 | **0.5775 (6)** | 0.5530 |
| Rubber | **0.6630 (7)** | 0.6215 | **0.5840 (5)** | 0.5281 | **0.5653 (4)** | 0.5373 |
| Fishery | **0.5705 (7)** | 0.4742 | **0.5702 (6)** | 0.4647 | **0.5703 (6)** | 0.5451 |
| Diversified Finance | **0.5341 (4)** | 0.5303 | **0.5458 (5)** | 0.5204 | **0.5498 (5)** | 0.5300 |

**Table 3**
The classification ACC rate of the patterns obtained by PR-DTW and PIP on different sectors using three classifiers on different sectors.
The ACC rate of PR-DTW and PIP on different sectors using three classifiers

| Sector | 1NN classifier | | Decision Tree classifier | | MLP classifier | |
|---|---|---|---|---|---|---|
| | PR-DTW (m) | PIP (m) | PR-DTW (m) | PIP (m) | PR-DTW (m) | PIP (m) |
| Glass | **0.5567 (8)** | 0.5216 (7) | **0.5278 (3)** | **0.5278 (7)** | **0.5484 (3)** | 0.5216 (6) |
| Oil exploitation | 0.5591 (7) | **0.5822 (4)** | **0.5733 (6)** | 0.5506 (4) | **0.5618 (3)** | 0.5503 (6) |
| Tourism service | **0.5730 (6)** | 0.5157 (5) | **0.5678 (4)** | 0.5495 (6) | **0.5912 (7)** | 0.5861 (6) |
| Dyes & coatings | **0.5557 (4)** | 0.5102 (5) | **0.5645 (7)** | 0.4959 (4) | **0.5330 (3)** | 0.4984 (6) |
| Publish industry | **0.5737 (8)** | 0.5462 (7) | **0.5623 (6)** | 0.5391 (4) | **0.5530 (8)** | 0.5229 (4) |
| Culture education | **0.5611 (8)** | 0.5348 (5) | **0.5487 (6)** | 0.5363 (6) | **0.5506 (7)** | **0.5506 (6)** |
| Aviation | **0.6427 (8)** | 0.6214 (6) | **0.6135 (8)** | 0.6038 (5) | **0.5378 (6)** | 0.5242 (4) |
| Commerce | 0.5134 (3) | **0.5638 (4)** | **0.5741 (4)** | 0.5470 (3) | **0.5791 (8)** | 0.5706 (4) |
| Road and bridge | 0.5281 (3) | **0.5584 (5)** | 0.5730 (8) | **0.6068 (7)** | **0.5590 (4)** | 0.5506 (5) |
| Supermarket | **0.5727 (8)** | 0.5578 (7) | 0.5687 (3) | **0.5723 (7)** | 0.5793 (4) | **0.5867 (8)** |
| M_T manufacturing | 0.5720 (4) | **0.5915 (5)** | 0.5917 (8) | **0.5955 (3)** | **0.5760 (7)** | 0.5565 (8) |
| H&F | 0.5244 (3) | **0.5510 (4)** | **0.5999 (7)** | 0.5734 (3) | **0.5775 (6)** | 0.5734 (6) |
| Rubber | **0.6630 (7)** | 0.6172 (7) | **0.5840 (5)** | 0.5704 (8) | **0.5653 (4)** | 0.5607 (6) |
| Fishery | **0.5705 (7)** | 0.5553 (4) | 0.5702 (6) | **0.5755 (3)** | **0.5703 (6)** | 0.5701 (8) |
| Diversified Finance | **0.5341 (4)** | 0.5263 (8) | **0.5458 (5)** | 0.5380 (4) | 0.5498 (5) | **0.5751 (3)** |

**Table 4**
The classification ACC rate of the patterns obtained by PR-DTW and PLA on different sectors using three classifiers on different sectors.
Comparison of PR-DTW and PLA using three classifiers

| Sector | 1NN classifier | | Decision Tree classifier | | MLP classifier | |
|---|---|---|---|---|---|---|
| | PR-DTW (m) | PLA (m) | PR-DTW (m) | PLA (m) | PR-DTW (m) | PLA (m) |
| Glass | **0.5567 (8)** | 0.5195 (3) | **0.5278 (3)** | 0.5216 (4) | **0.5484 (3)** | 0.5154 (6) |
| Oil exploitation | 0.5591 (7) | **0.5736 (3)** | **0.5733 (6)** | 0.5648 (3) | 0.5618 (3) | **0.5846 (4)** |
| Tourism service | **0.5730 (6)** | 0.5181 (6) | **0.5678 (4)** | 0.5547 (5) | **0.5912 (7)** | 0.5652 (4) |
| Dyes & coatings | **0.5557 (4)** | 0.5187 (7) | **0.5645 (7)** | 0.5443 (3) | **0.5330 (3)** | 0.5299 (6) |
| Publish industry | **0.5737 (8)** | 0.5621 (3) | **0.5623 (6)** | 0.5438 (6) | **0.5530 (8)** | 0.5391 (6) |
| Culture education | **0.5611 (8)** | 0.5345 (7) | 0.5487 (6) | **0.5541 (6)** | **0.5506 (7)** | 0.5454 (8) |
| Aviation | **0.6427 (8)** | 0.6368 (6) | **0.6135 (8)** | **0.6135 (5)** | **0.5378 (6)** | 0.5242 (5) |
| Commerce | 0.5134 (3) | **0.5335 (7)** | **0.5741 (4)** | 0.5487 (7) | **0.5791 (8)** | 0.5740 (6) |
| Road and bridge | 0.5281 (3) | **0.5450 (3)** | **0.5730 (8)** | 0.5591 (5) | 0.5590 (4) | **0.5644 (6)** |
| Supermarket | 0.5727 (8) | **0.5979 (3)** | 0.5687 (3) | **0.5688 (8)** | **0.5793 (4)** | 0.5722 (5) |
| M_T manufacturing | **0.5720 (4)** | 0.5638 (6) | **0.5917 (8)** | 0.5881 (3) | 0.5760 (7) | **0.6187 (7)** |
| H&F | **0.5244 (3)** | **0.5244 (5)** | **0.5999 (7)** | 0.5897 (7) | 0.5775 (6) | **0.5857 (4)** |
| Rubber | **0.6630 (7)** | 0.6075 (6) | **0.5840 (5)** | 0.5652 (8) | **0.5653 (4)** | 0.5560 (7) |
| Fishery | **0.5705 (7)** | 0.5557 (3) | 0.5702 (6) | **0.5721 (4)** | **0.5703 (6)** | 0.5549 (3) |
| Diversified Finance | **0.5341 (4)** | 0.5087 (7) | **0.5458 (5)** | 0.5321 (4) | 0.5498 (5) | **0.5750 (3)** |

on the pattern expression of time series. Owing to the representative points are selected from the whole series, PR-DTW can avoid the limitation of End Effect in time series representation.
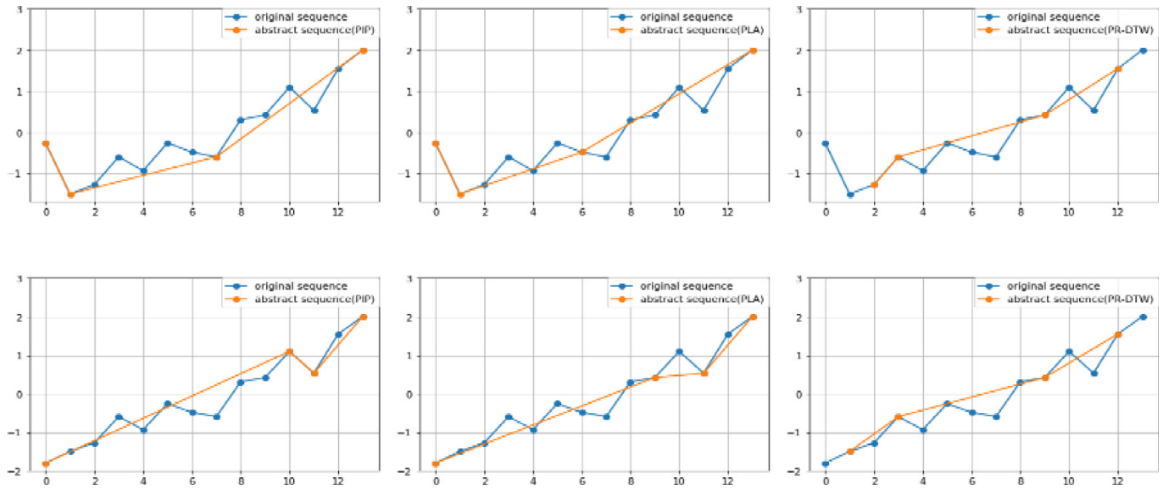
**Fig. 2.** The comparison of three representation methods before and after value of the first data point changing.

**Table 5**
The comparison of three methods and the original sequence using Decision Tree Classifier.

| Comparison of three methods and the original sequence using Decision Tree classifier | | | |
|---|---|---|---|
| Sector | Or_DTW (10) | PR-DTW (m) | PIP (m) | PLA (m) |
| Glass | 0.5216 | **0.5278 (3)** | **0.5278 (7)** | 0.5216 (4) |
| Oil exploitation | 0.5447 | **0.5733 (6)** | 0.5506 (4) | 0.5648 (3) |
| Tourism service | 0.4896 | **0.5678 (4)** | 0.5495 (6) | 0.5547 (5) |
| Dyes & coatings | 0.4840 | **0.5645 (7)** | 0.4959 (4) | 0.5443 (3) |
| Publish industry | 0.5391 | **0.5623 (6)** | 0.5391 (4) | 0.5438 (6) |
| Culture education | 0.5488 | 0.5487 (6) | 0.5363 (6) | **0.5541 (6)** |
| Aviation | 0.6116 | **0.6135 (8)** | 0.6038 (5) | **0.6135 (5)** |
| Commerce | 0.5369 | **0.5741(4)** | 0.5470(3) | 0.5487(7) |
| Road and bridge | 0.5760 | 0.5730 (8) | **0.6068 (7)** | 0.5591 (5) |
| Supermarket | 0.5324 | 0.5687(3) | **0.5723(7)** | 0.5688(8) |
| M_T manufacturing | 0.5379 | 0.5917(8) | **0.5955(3)** | 0.5881(5) |
| H&F | 0.5693 | **0.5999 (7)** | 0.5734 (3) | 0.5897(7) |
| Rubber | 0.5281 | **0.5840 (5)** | 0.5704 (8) | 0.5652 (8) |
| Fishery | 0.4647 | 0.5702 (6) | **0.5755 (3)** | 0.5721 (4) |
| Diversified Finance | 0.5204 | **0.5458 (5)** | 0.5380 (4) | 0.5321 (4) |

Fig. 2 shows the representation result using PLA, PIP and PR-DTW before and after the value of the first data point changes. It is obvious that the pattern obtained by PLA and PIP has changed a lot before and after we change the value of the first data point. By contrast, PR-DTW can keep the basic shape of time series and avoid End Effect in time series representation.

(2) Anti-noise: In stock market, the noise has nothing to do with the stocks' basic value and even may lead investors to make wrong judgments on the pattern of stocks.

Table 5 shows that the patterns obtained by PR-DTW has higher accuracy than the general methods using Decision Tree Classifier, and it means PR-DTW can obtain high quality patterns by filtering the noise of the short-term stock time series in most of the sectors. Meanwhile, Fig. 2 shows that PR-DTW has the properties of anti-noise especially for the end points. We can conclude that PR-DTW can avoid such short-term noise and maintain the pattern of time series, which can help us to shield noise and better analyze the pattern of stock time series.

(3) Segmentation: In stock time series pattern discovery, as a pre-processing step to reduce the dimension of data, segmentation extracts important data points to represent the time series [52]. In [53], the stock time series are segmented based on PIP and the segmentation is typically characterized by a few critical PIPs [30]. The more important the data point is, the earlier it is chosen as a cutting point. Moreover, different definitions like critical points [40] and key points [54], which were predefined by one's prior knowledge instead of the time series' nature structure, were presented. From the theorem in 3.2, we can know that PR-DTW can segment the time series without prior knowledge.

The purpose of PR_DTW is to find a sub-sequence that has the shortest distance to the original sequence in DTW measurement, thus PR_DTW has natural advantages in pattern representation. Fig. 3 shows the representation of PR-DTW has deep connection to the original sequence. Each point in the extracted pattern corresponds to a segment of the original sequence and the segmentation is automated. For stock time series, segmentation can also be considered as a division of the original sequence based on the trend consistency. PR-DTW can segment the stock time series by their trend in the
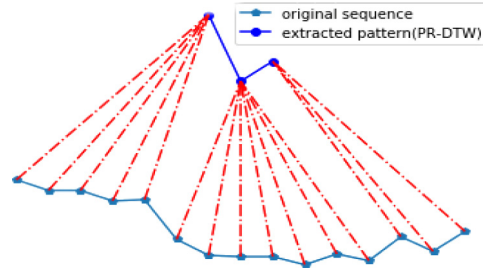
**Fig. 3.** The corresponding relationship between the original sequence and the extracted pattern by PR-DTW.
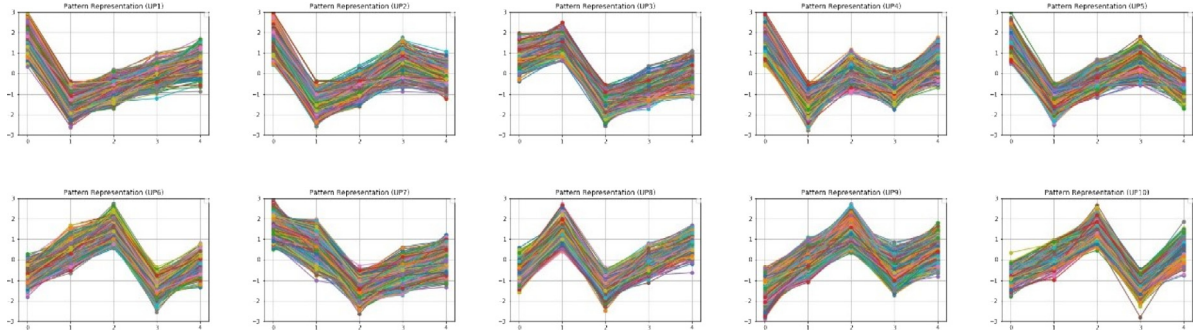


**Fig. 4.** Stock time series top 10 patterns obtained by PR-DTW for positive rate of return in short term named UP1, UP2, UP3, UP4, UP5, UP6, UP7, UP8, UP9 and UP10.
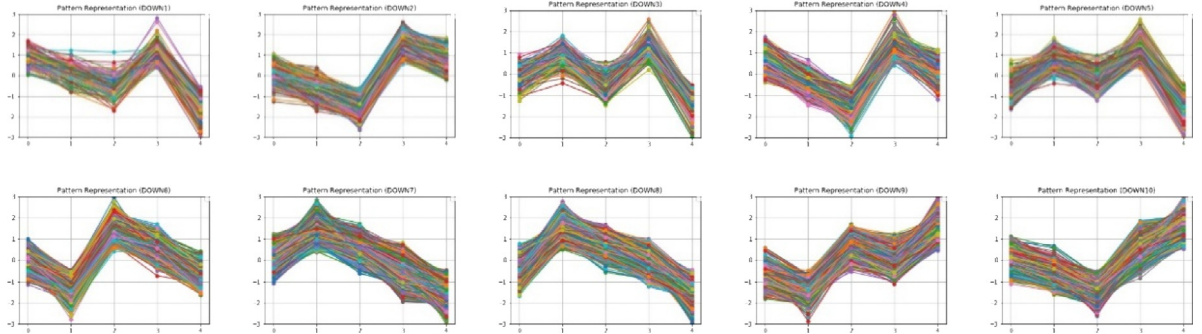


**Fig. 5.** Stock time series top 10 patterns obtained by PR-DTW for negative rate of return in short term named DOWN1, DOWN2, DOWN3, DOWN4, DOWN5, DOWN6, DOWN7, DOWN8, DOWN9 and DOWN10.

whole time scale accurately. Therefore, the segmentation by PR_DTW has a significant power on studying the volatility of stocks.

## 5. Empirical analysis in the Chinese A-share market

The 53 chart patterns, which have been summarized by Wiley John from over 38,500 pattern samples of 500 stocks with the evaluation indexes like statistics of failure rates and breakout, has contribution on the pattern analysis [55]. The works are helpful for us to make an analysis on stock time series with long period. However, they may lose its guiding significance for short-term stock time series with stochasticity and heavy noise.

In order to find out the patterns of stock market in short term and verify the effectiveness of our method, we validate PR-DTW in Chinese A-share market. For the purpose of studying the pattern of stock price in a short term to support real-time transactions, we select all closing price data of 3588 stocks listed in A-share market for nearly 500 trading days from 21 March 2016 to 9 April 2018. For the short term in the stock market is generally considered as 10 trading days or more [2,56,57], all stock price sequence are normalized and segmented into sub-sequences by the sliding time window with time span of 10 and time step of 10. Z-score normalization is chosen as the normalization method in experiments.

**Table 6**
The performance analysis of top 10 patterns for positive rate of return obtained by PR-DTW.

| Pattern name | Rise probability (%) | Average rise (%) |
| --- | --- | --- |
| UP1 | 0.8336 (516/619) | 2.1559 |
| UP2 | 0.7962 (461/579) | 2.2207 |
| UP3 | 0.7804 (231/296) | 1.9916 |
| UP4 | 0.7690 (393/511) | 2.2685 |
| UP5 | 0.7557 (359/475) | 1.7885 |
| UP6 | 0.7465 (268/359) | 2.1826 |
| UP7 | 0.7420 (1844/2485) | 1.6944 |
| UP8 | 0.7416 (178/240) | 1.4556 |
| UP9 | 0.7274 (1145/1574) | 1.8951 |
| UP10 | 0.7231 (175/242) | 1.8539 |

**Table 7**
The performance analysis of top 10 patterns for negative rate of return obtained by PR-DTW.

| Pattern name | Fall probability (%) | Average rise (%) |
| --- | --- | --- |
| DOWN1 | 0.8188 (104/127) | −1.9969 |
| DOWN2 | 0.8081 (299/370) | −2.0621 |
| DOWN3 | 0.8009 (177/221) | −2.3007 |
| DOWN4 | 0.7988 (409/512) | −2.0637 |
| DOWN5 | 0.7823 (302/386) | −2.6300 |
| DOWN6 | 0.7786 (475/610) | −1.9532 |
| DOWN7 | 0.7782 (3039/3905) | −2.8719 |
| DOWN8 | 0.7717 (1065/1380) | −2.3469 |
| DOWN9 | 0.7713 (577/748) | −1.8978 |
| DOWN10 | 0.7697 (846/1099) | −2.0338 |

Zhao-Rong Lai concluded that the most important and reliable information affecting stock prices can be obtained in 5 trading time points in stock market [58], so we choose $m = 5$. Finally, we get 141708 sub-sequence samples. PR-DTW, with the parameter $m$ chosen as 5, is used to extract the stock patterns in short term. Considering the stock market volatility, we choose different kinds of labels for the patterns predicting the rise of stocks and the patterns predicting the fall.

$$LABEL_{up} = \max \left( next\ three\ days'\ close\ prices \right) / recent\ close\ price - 1 \tag{19}$$

$$LABEL_{fall} = \min \left( next\ three\ days'\ close\ prices \right) / recent\ close\ price - 1 \tag{20}$$

Then the samples are clustered by the ranking of pattern points' values, and it means the patterns with the same order of value ranking belong to a cluster. Then each class is sorted in descending order by the probability of stock rising or fall in the next three days. Accordingly, the top ten patterns predicting rise of stocks and the top ten patterns predicting fall of stocks in the whole Chinese A-share market are found out and shown in Figs. 4 and 5. The rise/fall probability and average rise are used as evaluation indexes of pattern performance. The rise/fall probability index represents the proportion of rise/fall of samples to all samples in the same class, and average rise index means the average of the patterns' return in the same class.

From Tables 6 and 7, with the representation of PR-DTW, the patterns we discovered have the capability of prediction in short term and the guiding significance for investment in short term.

## 6. Conclusions and discussions

In this paper, by constructing optimization models, we propose a pattern representation of stock time series based on DTW called PR-DTW with the advantages of avoiding the End Effect, anti-noise and segmentation based on the trend. In order to simplify the calculation, we propose a theorem to illustrate the inclusion relationship between the optimal solution set of (3) and the optimal solution set of (4) and the optimal solution of global optimization problem can be mapped into that of combinatorial optimization problem according to the function (15). The results of three classifiers (1NN, Decision Tree, Multi-layer Perceptron) implemented on 10 sectors in Chinese A-share market unanimously shows that PR-DTW has the capability of extracting time series short-term patterns which is widely regarded as difficulty.

For mining the typical patterns of short-term stock time series, we make statistical analysis on the short term patterns in two situations with two recognized indicators called rise probability (fall probability) and average rise. The rise probability (fall probability) of the patterns are more than 70%, which can make effective predictions for future earnings. Moreover, it is obvious that the 10 patterns predicting stock fall are more significant than the 10 patterns predicting stock

rise, which is quite suitable to the situation of Chinese A-share market. Due to the influence of retail investors, the prices of stocks are more sensitive in the bearish market. On the whole, PR-DTW, as an effective pattern representation method, is helpful for short-term investors to predict the trend in stock market.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgments

## References

[1] John J. Murphy, Technical Analysis of the Financial Markets, New York Institute of Finance, ISBN: 0-7352-0066-1, 1999, pp. 1–5, 24–31.
[2] Anne-Marie Baiynd, The Trading Book: A Complete Solution to Mastering Technical Systems and Trading Psychology, McGraw-Hill, 2011, p. 272.
[3] P. Blakey, Pattern recognition techniques in stock prices and volumes, Microw. Mag. IEEE 3 (1) (2002) 28–33.
[4] S. Kim, H. Lee, H. Ko, et al., Pattern matching trading system based on the dynamic time warping algorithm, Sustainability 10 (12) (2018) 4641.
[5] R.D. Edwards, J. Magee, W.H.C. Bassetti, Technical analysis of stock trends, in: Technical Analysis of Stock Trends, tenth ed., AMACOM, 2007.
[6] G.S. Atsalakis, E.M. Dimitrakakis, C.D. Zopounidis, Elliott wave theory and neuro-fuzzy systems, in: Stock Market Prediction: The WASP System, 2011.
[7] J.S. Chou, T.K. Nguyen, Forward forecast of stock price using sliding-window metaheuristic-optimized machine learning regression, IEEE Trans. Ind. Inf. PP (99) (2018) 1.
[8] Y. Wan, Y.W. Si, A formal approach to patterns classification in financial time series, Inform. Sci. 411 (2017) 151–175.
[9] X. Wang, A. Mueen, H. Ding, et al., Experimental comparison of representation methods and distance measures for time series data, Data Min. Knowl. Discov. 26 (2) (2013) 275–309.
[10] Yang Yun, Unsupervised Ensemble Learning and Its Application to Temporal Data Clustering, University of Manchester, 2011.
[11] M. Azzouzi, I.T. Nabney, Analysing time series structure with hidden Markov models, in: Neural Networks for Signal Processing VIII Proceedings of the IEEE Signal Processing Soc, 1998, pp. 402–408.
[12] J. Meng, L.X. Wu, X.K. Wang, et al., Granulation-based symbolic representation of time series and semi-supervised classification, Comput. Math. Appl. 62 (9) (2011) 3581–3590.
[13] F. Angeletti, E. Bertin, P. Abry, Random vector and time series definition and synthesis from matrix product representations: From statistical physics to hidden Markov models, IEEE Trans. Signal Process. 61 (21) (2013) 5389–5400.
[14] K.P. Murphy, Dynamic Bayesian Networks: Representation, Inference and Learning, University of California, Berkeley, 2002.
[15] A. Tucker, X. Liu, A Bayesian Network Approach to Explaining Time Series with Changing Structure, IOS Press, 2004.
[16] Z.Y. Zhao, Bayesian Multiregression Dynamic Models with Applications in Finance and Business (Dissertations & theses - Gradworks), 2015.
[17] Y. Yuan, G. Xun, Q. Suo, et al., Wave2Vec: Deep representation learning for clinical temporal data, Neurocomputing (2018).
[18] E. Choi, M.T. Bahadori, E. Searles, et al., Multi-Layer Representation Learning for Medical Concepts, 2016.
[19] R. Agrawal, C. Faloutsos, A. Swami, Efficient similarity search in sequence databases, in: International Conference on Foundations of Data Organization and Algorithms, Springer, Berlin, Heidelberg, 1993.
[20] K. Samiee, P. Kovacs, M. Gabbouj, Epileptic seizure classification of EEG time-series using rational discrete short-time Fourier transform, IEEE Trans. Biomed. Eng. 62 (2) (2015) 541–552.
[21] M.W. Reed, Fourier series representation of industrial load profiles, in: Symposium on Simulation, IEEE Computer Society Press, 1987.
[22] Z.R. Struzik, A. Siebes, The Haar wavelet transform in the time series similarity paradigm, in: Principles of Data Mining & Knowledge Discovery, Third European Conference, Pkdd 99, September, DBLP, Prague, Czech Republic, 1999.
[23] M. Fuad, M. Marwan, Aggressive pruning strategy for time series retrieval using a multi-resolution representation based on vector quantization coupled with discrete wavelet transform, Expert Syst. 34 (1) (2017) e12171.
[24] A.E. Milne, R.M. Lark, Wavelet transforms applied to irregularly sampled soil data, Math. Geosci. 41 (6) (2009) 661–678.
[25] F. Korn, H.V. Jagadish, C. Faloutsos, Efficiently supporting Ad Hoc queries in large datasets of time sequences, SIGMOD Rec. 26 (2) (1997) 289–300.
[26] Y. Liu, Z. Yang, L. Yang, Online signature verification based on DCT and sparse representation, IEEE Trans. Cybern. 45 (11) (2014).
[27] A. Gonzalez-Vidal, P. Barnaghi, A.F. Skarmeta, BEATS: Blocks of eigenvalues algorithm for time series segmentation, IEEE Trans. Knowl. Data Eng. (2018) 1.
[28] J. Lin, E.J. Keogh, S. Lonardi, et al., A symbolic representation of time series, with implications for streaming algorithms, in: Proceedings of the 8th ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery, DMKD 2003, June 13, 2003, ACM, San Diego, California, USA, 2003.
[29] E. Keogh, S. Chu, D. Hart, M. Pazzani, An online algorithm for segmenting time series, in: Data Mining, 2001 ICDM 2001, Proceedings IEEE International Conference on, IEEE, 2001, pp. 289–296.
[30] D. Bao, A generalized model for financial time series representation and prediction, Appl. Intell. 29 (1) (2008) 1–11.
[31] N.Q.V. Hung, D.T. Anh, An Improvement of PAA for Dimensionality Reduction in Large Time Series Databases, 2008.
[32] V. Megalooikonomou, G. Li, Q. Wang, A dimensionality reduction technique for efficient similarity analysis of time series databases, in: Proceedings of the 2004 ACM CIKM International Conference on Information and Knowledge Management, November, ACM, Washington, DC, USA, 2004, pp. 8–13.
[33] Q. Chen, L. Chen, X. Lian, et al., Indexable PLA for efficient similarity search, in: Proceedings of the 33rd International Conference on Very Large Data Bases, September 23–27, 2007, University of Vienna, Austria, 2007, VLDB Endowment.
[34] C.A. Ratanamahatana, J. Lin, D. Gunopulos, E. Keogh, M. Vlachos, G. Das, Mining time series data, in: O. Maimon, L. Rokach (Eds.), Data Mining and Knowledge Discovery Handbook, Springer, Berlin, 2010, pp. 1049–1077, http://dx.doi.org/10.1007/0-387-25465-x_36.
[35] T.C. Fu, Y.K. Hung, F.L. Chung, Improvement algorithms of perceptually important point identification for time series data mining, in: 2017 IEEE 4th International Conference on Soft Computing & Machine Intelligence, ISCMI, IEEE, 2017.

[36] E. Lefevre, J.D. Markman, P.T. Jones, et al., Reminiscences of a Stock Operator: With New Commentary and Insights on the Life and Times of Jesse Livermore, John Wiley & Sons, 2010.
[37] H. Li, L. Yang, Time series visualization based on shape features, Knowl.-Based Syst. 41 (2013) 43–53.
[38] S.J. Wilson, Data representation for time series data mining: time domain approaches, Wiley Interdiscip. Rev. Comput. Stat. 9 (1) (2017).
[39] S. Liu, L. Lu, G. Liao, et al., Pattern discovery from time series using growing hierarchical self-organizing map, in: Neural Information Processing, Springer, Berlin Heidelberg, 2006.
[40] F.L. Chung, T.C. Fu, R. Luk, V. Ng, Flexible time series pattern matching based on perceptually important points. in: International Joint Conference on Artificial Intelligence Workshop on Learning from Temporal and Spatial Data, 2001, pp. 1–7.
[41] P.W.P. Man, M.H. Wong, ACM Press the Tenth International Conference - Atlanta, Georgia, USA (2001.10.05-2001.10.10), Proceedings of the Tenth International Conference on Information and Knowledge Management, - Cikm\01 - Efficient and Robust Feature Extraction and Pattern Matching of Time Series by a Lattice Structure, 2001, p. 271.
[42] K.B. Pratt, E. Fink, Search for patterns in compressed time series, Int. J. Imag. Graph. 2 (1) (2002) 89–106.
[43] D. Bao, Z. Yang, Intelligent stock trading system by turning point confirming and probabilistic reasoning, Expert Syst. Appl. 34 (1) (2008) 620–627.
[44] C. Ji, S. Liu, C. Yang, et al., A piecewise linear representation method based on importance data points for time series data, in: 2016 IEEE 20th International Conference on Computer Supported Cooperative Work in Design, CSCWD, IEEE, 2016.
[45] Y. Wei, Motion representation based on important turning points set and its application in dance training, in: 2015 Second International Conference on Soft Computing and Machine Intelligence, ISCMI, IEEE, 2015.
[46] A. Mueen, N. Chavoshi, N. Abu-El-Rub, et al., Speeding up dynamic time warping distance for sparse time series data, Knowl. Inf. Syst. (2017).
[47] A. Sharabiani, H. Darabi, S. Harford, et al., Asymptotic dynamic time warping calculation with utilizing value repetition, Knowl. Inf. Syst. (12) (2018) 1–30.
[48] T. Berrada, Incomplete information, heterogeneity, and asset pricing, J. Financ. Econ. 4 (1) (2006) 136–160.
[49] T.C. Fu, A review on time series data mining, Eng. Appl. Artif. Intell. 24 (1) (2011) 164–181.
[50] J. Han, J. Qian, X. Dong, Suppression of end-effect in empirical mode decomposition by mirror extension and radial basis function neural network prediction, J. Vib. Meas. Diagn. 30 (4) (2010) 414–417.
[51] C. Zhong, Z. Shixiong, Analysis on end effects of EMD method, J. Data Acquis. Process. (2003).
[52] Y. Wan, X. Gong, Y.W. Si, Effect of segmentation on financial time series pattern matching, Appl. Soft Comput. 38 (2016) 346–359.
[53] X.D. Zhang, A. Li, R. Pan, Stock trend prediction based on a new status box method and AdaBoost probabilistic support vector machine, Appl. Soft Comput. 49 (2016).
[54] F.L. Chung, T.C. Fu, R. Luk, V. Ng, Flexible time series pattern matchingbased on perceptually important points. in: International Joint Conference onArtificial Intelligence Workshop on Learning from Temporal and Spatial Data, 2001, pp. 1–7.
[55] M. Leng, X. Lai, G. Tan, et al., Time series representation for anomaly detection, in: IEEE International Conference on Computer Science & Information Technology, IEEE, 2009.
[56] T.N. Bulkowski, Encyclopedia of Chart Patterns, second ed., Wiley John Sons, 2011.
[57] J.L. Wang, S.H. Chan, Stock market trading rule discovery using pattern recognition and technical analysis, Expert Syst. Appl. 33 (2) (2007) 304–315.
[58] L. Zhao-Rong, Y. Pei-Yi, W. Xiaotian, et al., A kernel-based trend pattern tracking system for portfolio optimization, Data Min. Knowl. Discov. (2018).