

Development of Patternrecognition and Representation Algorithm for Time Series Datasets

Emmanuel N. Nwajiobi, Sylvanus .O. Anigbogu, Kenechukwu Anigbogu

Dept of Computer Science Nwafor Orizu College of Education, Nsugbe, Anambra State, Nigeria

Submitted: 15-07-2021

Revised: 29-07-2021

Accepted: 31-07-2021

ABSTRACT

In diverse areas of human endeavor such as business, industry, sciences and so on, massive amount of time series data are generated daily and due to the fact that time series data are typically very large, discovering information from such massive datasets therefore becomes a major challenge. A number of algorithms have been introduced to represent, classify, cluster, segment, index, detect motifs and anomalies in a time series data. In view of the above, this paper proposes a robust algorithm for pattern recognition and representation of a time series. The algorithm first normalises a time series dataset into the range [0,1]. The normalized version is now used for pattern identification and representation. In the proposed algorithm, we pre-defined patterns as up, down and flat patterns, and having equal length (three, five or ten data points). Each pattern represents a segment (subsequence) of the time series. The algorithm was tested with historical time series datasets obtained online from (a) Dow Jones Industrial Average (b) Nasdaq, and (c) S&P 500 via yahoo finance. Each dataset consisted of 5158 data points, covering the period 2000-2020. The algorithm captured all the pre-defined patterns in the datasets and was able to represent the patterns in the entire historical datasets with symbols. The algorithm is a veritable tool for time series data mining operations. Object-Oriented Analysis and Design Methodology (OOADM) and prototyping methodology were used to design the system; while PHP, MYSQL, HTML and CSS were used to develop the system. The system was well tested and the outputs were excellent.

Key terms: Pattern recognition, time series, representation, algorithm, datamining

I. INTRODUCTION

At present, majority of activities in companies and organisations generate large amounts of data which are typically saved in databases. However, the question of what to do

with such huge amounts of data is not always easily obvious or answered in most situations by owners of such large databases. Though large computer storage disks make the storage of huge data possible, computational algorithms are also needed to analyse the data. Massive data sets are rarely profitable; their real worth lies in the possibility to extract useful information for making decisions or for understanding the phenomena that generated such data.

To this extent, information retrieval is no longer enough anymore for decision-making. Thus, the availability of these huge collections of data now created new needs that will help us make better and informed decisions, including making predictions about the future. These new needs include automatic summarization of data, extraction of information buried in stored data, and the discovery of patterns in raw data. With the availability of these enormous amounts of data stored in files, databases, and other repositories, it is therefore very important and necessary to develop improved means of analysing and interpreting the data, as well as extracting interesting knowledge and patterns that could help in decision-making and prediction.

To make these large data sets more useful, we need techniques to analyse them with a view to finding out something surprising and interesting from the gathered data. In this regard, we are faced with the problem of how to find patterns from the datasets and show that the patterns are useful, informative and important. Data mining techniques can be used to discover patterns from large datasets, including time series datasets.

Time series is a collection of observations made sequentially in time (Abdullah, 2016). It is an ordered sequence of values (real numbers) of a variable or variables measured, observed or calculated at regular time intervals over a period of time. According to Pohl and Bouchachia (2012),

the following activities can be performed on a time series data: detecting motifs, recognizing and extracting patterns, finding correlation between time series or finding similar time series. Similarly, analysis of a time series can be said to comprise

three processing steps, namely: (a) Abstraction (or representation), (b) Mining and Discovery of trends and patterns, and (c) Prediction (Pohl and Bouchachia, 2012).



Figure. 1: Processing stages in time series analysis (Source: Pohl and Bouchachia (2012))

Any information of the sequential nature can be processed by pattern recognition algorithms to make the sequences comprehensible for its practical use. The term pattern recognition connotes automatic discovery of regularities in data through the use of computer algorithms and with the use of these regularities to take such actions as classifying the data into different categories (Bishop, 2006). These regularities in data can be referred to as patterns.

According to Raj et al (2015), pattern recognition is a multi-disciplinary subject covering the following fields: statistics, engineering, artificial intelligence, computer science, psychology and physiology, etc. They posited that computer-based automated pattern recognition systems are required when: (a) the human senses fail to recognize patterns, (b) there is need to automate and speed up the recognition process. And considering the volume of data generated by businesses and companies these days, it is obvious that pattern recognition is inevitable in exploring the data for information buried in the data. For instance, people measure things like blood pressure, annual rainfall, value of stock, etc., and as such time series occur in virtually every medical, scientific and business domain. Time series reveals the temporal behaviour of the underlying mechanism that produced the data.

However, as the amount of data generated by business houses increases, there is therefore the need to explore new ideas and algorithms to analyse it to gather information necessary for decision making and predictions. The type of time series data considered in this paper are mostly those that can generate forecasts, such as stock closing price. Based on the foregoing, this research paper, propose a new and novel pattern recognition algorithm/model to (a) efficiently represent time series dataset, and(b) detect and extract patterns of interest buried in time series datasets. Time series datasets collected via yahoo finance website from different sources were used to test and validate the model.

II. LITERATURE REVIEW

Time series according to Nguyen and Duong (2007) is a sequence of real numbers, each number representing a value at a time point. It is a collection of observations made sequentially in time (Abdullah, 2016). A time series is an ordered sequence of values of a variable (univariate) or many variables (multivariate) measured, observed or calculated at equally spaced time intervals over a period of time. It consists of a sequence of values and their corresponding timestamps (i.e. the time at which the values were observed or measured). Figure 2 illustrates a typical plot of a time series dataset.

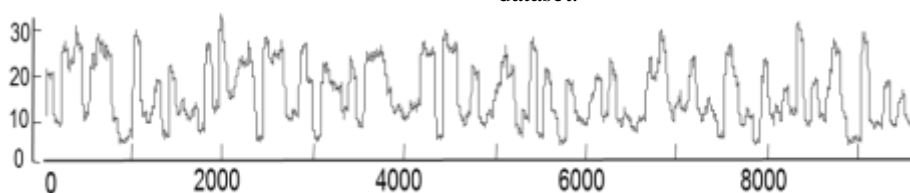


Figure 2: A typical illustration of a time series plot. On the Y-axis are the values, and along the X-axis are time stamps; Source Lin et al, (2002)

Time series abstraction is concerned with finding a suitable way to represent time series for further computational analysis. The process creates

a more compact representation of the time series, while at the same time preserving the information content of the original time series. Thus, one of the

cardinal objectives of time series representation is to find a lower dimensionality that preserves the fundamental characteristics of the original data. Again, good representation will foster further time series analysis towards discovering patterns and making informed decisions.

A good number of representation and dimension reduction techniques have been proposed in the literature for representing and summarizing time series data. Among them is the symbolic representation approach, which constitutes a simple way of reducing the dimensionality of the time series data by turning it into sequences of symbols. Once a time series is available in a string symbolic form, string analysis techniques can be used to analyse the data faster and efficiently. Time series representation techniques already developed include, but not limited to the work of (Ding, 2008; Han & Kamber, 2001; Jingpei, et al 2013). These are discrete Fourier Transformation (DFT), Piecewise Aggregate Approximation (PAA), Adaptive Piecewise Constant Approximation (APCA), Symbolic Aggregate approximation (SAX), Indexable Piecewise Linear Approximation, Independent Component Analysis (ICA), Principal Component Analysis (PCA), Piecewise Linear Approximation (PLA), Discrete Wavelet Transformation (DWT), Single Value Decomposition (SVD), Discrete Cosine Transformation.

Singh (2000) addressed the problem of time series representation by creating an algorithm called binary representation, in which “1” was used to represent increase and “0” was used to represent decrease. It partially solved the problem of time series representation by transforming it into strings of ones and zeros for further processing. It did not however address the issue of patterns in time series. What about a situation where there exists consecutive increases or decreases? The model was silent on that.

Álvarez (2010) proposed a clustering approach to find patterns in electricity time series. He applied K-means, Expectation Maximisation (EM) and Fuzzy C-Means (FCM) clustering techniques to find patterns in stock market data and electricity pricing data. The model proposed can be used to forecast a stock market and electricity pricing time series as recorded in the study. The approach did not delve into the use of large historical data to find patterns necessary for pattern extraction and prediction. No definite means of extracting patterns from a historical database.

Jiangling et al, (2011) developed a novel time series segmentation method that was based on turning points to extract trends from the maximum or minimum points of the time series. It was a very solid and useful idea for detecting patterns in a time series. It segmented time series into up and down structure that minimized destruction of the original underlying trends in the dataset. It did not address the issue of how to symbolically (or otherwise) represent the time series or the discovered trends.

Prasanna, S. and Ezhilmaran, D. (2013) performed analysis of past and present financial data to generate patterns and decision making algorithms using artificial intelligence and data mining techniques. The study was able to establish that data mining can be applied in evaluating past stock prices and acquire valuable information.

The weakness of the study was inability to define the type of patterns that can be generated and how to represent them.

Badhiye, et al (2015) addressed time series representation to facilitate data mining of large time series databases. The method used symbolic piecewise trend approximation to represent the original dataset. It achieved dimensional reduction, and was able to symbolically represent time series dataset. The shortcoming of the approach was classification of trend into two: up and down only. It ignored the existence of flat trend, and lacked the ability to predict future trend.

Nguyen and Duong, 2007 proposed the use of Piecewise Linear Approximation (PLA) to segment time series, as a preprocessing approach necessary for further analysis. The approach represents a time series with straight lines. PLA refers to the approximation of a time series T , of length n with k straight lines (where $k < n$) (Nguyen and Duong, 2007). The PLA is composed of a series of segments representing the trend (up and down) of the raw data. Thus, PLA approximates a time series into a representation of linear segments that is efficient to manipulate and faster to process than the raw data. The linear segments can be visualized in an identical way to the original data in a time series plot, while the number of data points is significantly reduced without losing the intrinsic nature of the underlying activity (Berlin, 2009).

The algorithms for implementing piecewise linear approximation are sliding window,

bottom-up and top-down. Applications of PLA include pattern matching and prediction of trading points in the stock market (Zhang, et al, 2010), Wu at al, 2004). One of the problems of the three basic PLA approaches (sliding window, bottom up and top down) is the design of a stopping condition (usually denoted as a user-defined threshold value) as each of them heavily depend on the threshold to stop.

Keogh Eamonn and Jessica Lin in 2002 invented SAX. SAX stands for Symbolic Aggregate Approximation. It was the first, and a novel symbolic representation for a time series. SAX is a symbolization method that involves placing a symbol for each segment obtained by using PAA, since it is based on the Piecewise Aggregate Approximation (PAA) representation.

The PAA representation is merely an intermediate step required to obtain SAX. SAX is a process that maps the PAA representation of a time series into a sequence of discrete symbols (Nguyen and Duong, 2007). In other words, SAX uses alphabet symbols (a – z) to represent segments obtained through PAA. In order to place the symbols, it is essential to specify the number of symbols to be used and the intervals (or breakpoints) of the values for each symbol. To this end, Burcu, et al (2011) stated that the number of symbols to be used is generally determined by an expert having knowledge about the application domain under study. However, to help solve the problem of specifying the intervals (breakpoints) for each symbol, Burcu, et al (2011) suggested the use of histograms of the data values as shown in figure 3.

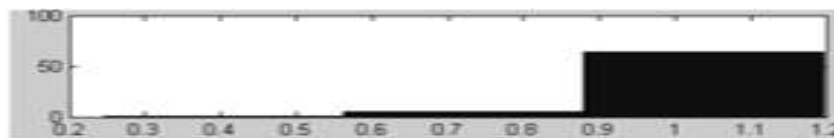


Figure 3: Histogram of segment values to help determine breakpoints

Another way to surmount the problem of determining breakpoints is to make use of the predefined statistical table. Table 2.1 shows a typical predefined lookup statistical table for 3 to 10 alphabets.

α	3	4	5	6	7	8	9	10
β_1	-0.43	-0.67	-0.84	-0.97	-1.07	-1.15	-1.22	-1.28
β_2	0.43	0	-0.25	-0.43	-0.57	-0.67	-0.76	-0.84
β_3		0.67	0.25	0	-0.18	-0.32	-0.43	-0.52
β_4			0.84	0.43	0.18	0	-0.14	-0.25
β_5				0.97	0.57	0.32	0.14	0
β_6					1.07	0.67	0.43	0.25
β_7						1.15	0.76	0.52
β_8							1.22	0.84
β_9								1.28

Table 1: Lookup table from a pre-defined statistical table that contains the breakpoints ($\beta_1 \dots \beta_9$) for alphabet size $a = 3$ to 10 that divides a Gaussian distribution into an arbitrary number. (Source: Lin et al, 2003).

Another attribute of SAX, in addition to the use of PAA technique, is normalization in order to transform the series to a Gaussian distribution so that the breakpoints can be determined from the curve in accordance with the required alphabet size. SAX has also the potential for dimensionality reduction. Thus, it can reduce a time series of arbitrary length n to a symbolic string of arbitrary length w ($w < n$), with the string composed of z different symbols ($z > 2$), Sant'Anna and Wickström (2011).

III. SYSTEM METHODOLOGY

This work approached the issue of pattern recognition and representation of time series from the data mining perspective, rather than from the statistical point of view. This is because, statistical tools will not suffice for large time series datasets analysis as it concerns pattern identification. As a result, the pattern recognition approach applied was an unsupervised learning since there was no prior labelling or classes of patterns unto which new patterns can be mapped to. However, to facilitate the task of pattern recognition, patterns were defined as either Up, Down or Flat, with fixed lengths of either 3, 5 or 10 data points (days with their respective timestamps).

The algorithm begins with a historical time series dataset which it receives as input. Prior to that, the dataset should have been preprocessed by removing blank cells of data and transforming (normalization process) the dataset into the range [0,1], such that the highest value is 1 and the least value in the series is 0. After this normalisation process, the pattern recognition algorithm can be applied to the resulting dataset to fish out patterns of interest and thus represent them with symbols. All patterns identified were symbolized and stored in a database for future uses and manipulation. Object-Oriented Analysis and Design Methodology (OOADM) and prototyping methodology were used to design the system; while PHP, MYSQL, HTML and CSS were used to

develop the system. The system was well tested and the outputs were excellent.

IV. SYSTEM DESIGN AND IMPLEMENTATION

4.1 The Proposed Pattern Recognition and Representation Algorithm

In this work, we propose an algorithm for pattern detection, extraction and representation (using symbols). For ease of identification of patterns, extraction and representation, we pre-defined patterns as either Up (U), Down (D) or Flat (F). Each pattern represents a segment, and can be drawn as shown in figure 4.

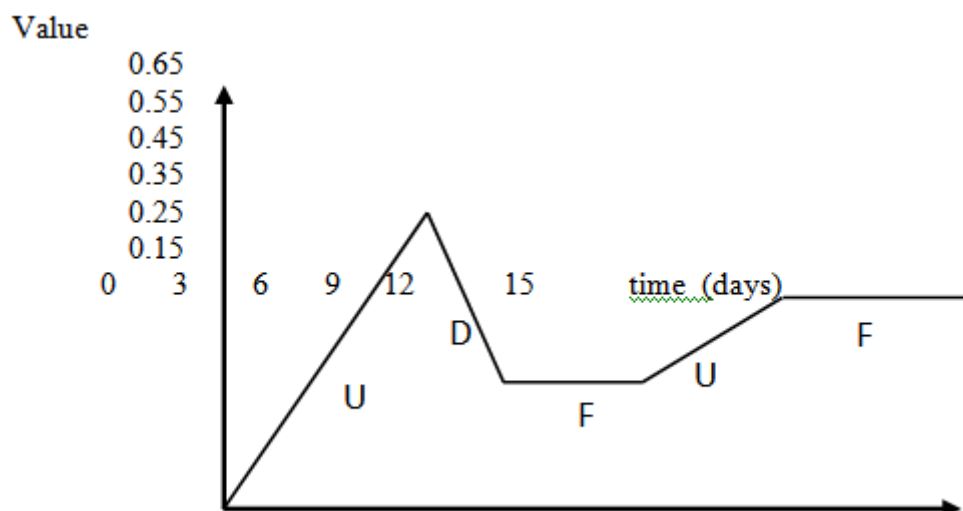


Figure 4: Visualisation of patterns of a time series, showing the up, down and flat patterns. (value = average values of data points for each segment: Up, Down or Flat pattern)

The symbolic representation of the time series as shown in figure 4 is UDFUF. Therefore, a time series of length 25 (data points) has been reduced to a string of UDFUF (which is five characters).

The algorithm for pattern identification, extraction and symbolic representation is hereby presented.

Input: S, Segment_size

Output: Pattern string symbols (for Up, Down and Flat patterns)

Repeat

Initialize tup = tdn = 0;

For (i = 0; i ≤ Segment_size - 1, i++) {
df = (i + 1) - i;

if df is positive, tup++ //augment increasing pattern variable

if df is negative, tdn++ //augment decreasing pattern variable

}
If tup = Segment_size - 1 or tdn = Segment_size - 1 then, pattern is Up or Down respectively otherwise pattern is Flat.

Calculate segment average; //Call SegmentAvg () function

Store segment MinDate, MinValue MaxDate, MaxValue, Segment_Symbol (U,D,F), SegmentAvg

Until end of S is reached.

4.2 System Implementation

The model cum algorithm was tested with real-valued discrete univariate time series data, mostly stock market data, obtained online from Yahoo website. The algorithm achieved 100% success in detecting and symbolizing the three

types of pre-defined patterns and equally symbolized them. Table 2 presents some of the outputs from the system.

Record No	Date	Value	Normalised Value	Year
1	2000-01-03	11357.5097656250	0.8104355931	2000
2	2000-01-04	10997.9296875000	0.6239265800	2000
3	2000-01-05	11122.6503906250	0.6886174083	2000
4	2000-01-06	11253.2597656250	0.7563626170	2000
5	2000-01-07	11522.5595703125	0.8960445523	2000
6	2000-01-10	11572.2001953125	0.9217924476	2000
7	2000-01-11	11511.0800781250	0.8900903463	2000
8	2000-01-12	11551.0996093750	0.9108479023	2000
9	2000-01-13	11582.4296875000	0.9270983338	2000
10	2000-01-14	11722.9804687500	1.0000000000	2000
11	2000-01-18	11560.7197265625	0.9158377051	2000
12	2000-01-19	11489.3603515625	0.8788245916	2000

Table 2: Sample raw and normalized time series data.

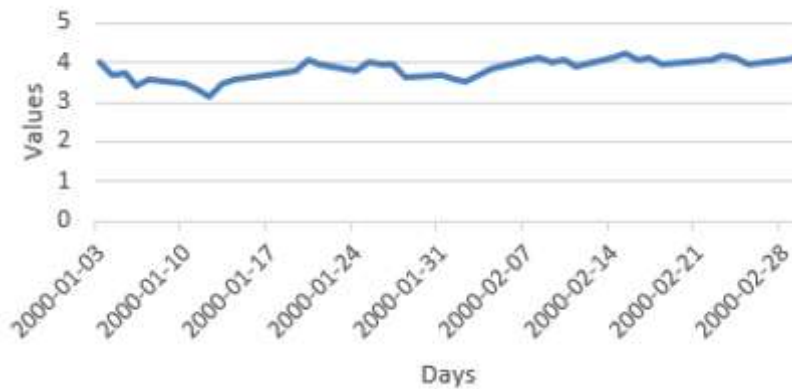


Figure 5a: Raw data plotted without patterns extracted for two (2) months

Again, figure 5b shows the normalized plot of the extracted patterns for the same two (2) months.

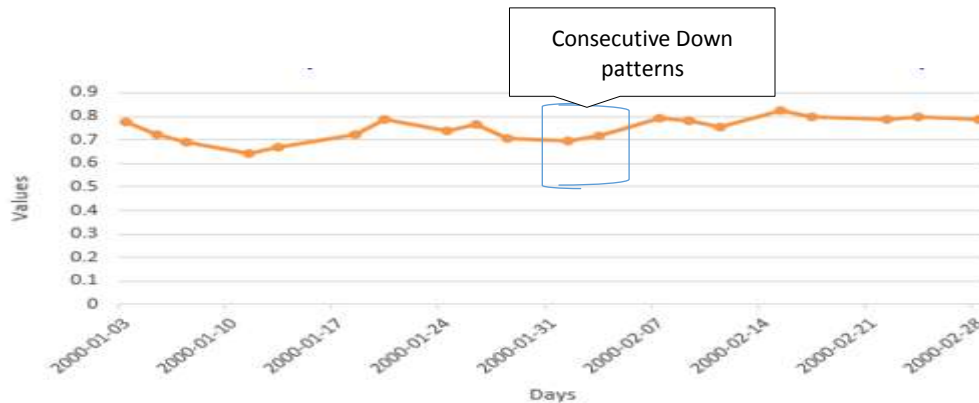


Figure 5b: The normalized data plot of the extracted patterns for the same 2 months.
 Figure 6 shows the symbolic representation of the extracted patterns.

```

FUFUDDUDUFDUDFDUUDUFFUFDFFFDFUDFDFFFUFD
FFFUFDUFFFFFUDUUUDFDFUUFUFFFFUUFFDFDDFFDFUDDF
FFFUUDFDDFFDFUDUDUUFFUUFFDFDFDFDFDFDFDUU
UFFFFDFUFFFFUUDFUDFDFUUDFFUUDFDFUDDFFDFUFFFFDU
UDFFFFDFDFUDDDDUFFUFDUFFFFUFDUUFFUFFFFFUDDFUU
UFUUFUDDDFUFDUDDFFUFDUFFFFUFDUDDUDFFFFDFDUFFF
FFUDUDDFFDFDFDFDDDFDFDFDFDUUDUFFFFDFUFDFFF
FDDFUUFUFFFFUUDFUFFFFDFDFDFDFUUDDFDUFFDFDF
FFFFFDUUUFDUFFFFFUDUDDFFDFUUFUUFUDFFFFFUDU
UDFFFFDFUUFUFDUFFFFFDFUUFFDFDFUUFFDFDFUUU
UFFFFUUFFFFFFUFDFFFUFFFFDFDFDFDFUUFFDFUFDUDU
UFDUDDUFDFFUFDUFDUFDUFDUFDUFDUFDUFDUFDUFDU
UFUFDUFFFFDFDDUUFUUDUUFUFDUFDUFDUFDUFDUFDU

```

Figure 6: The symbolic representation of the extracted patterns.

DD means consecutive Down patterns, likewise FFF or UU. With the patterns symbolically represented, further programming work can be done to compare time series of similar months of several years to find out areas of similarity.

V. RESULTS AND DISCUSSION

Table 2 presented a portion of the historical dataset, which has 5158 records (data points). The table showed both the raw and normalized values and their timestamps (date and year).

As already noted, figure 5a presented the raw data plotted without pattern extracted for the two (2) months; while figure 5b presented the normalized data of the extracted patterns for the same two months. In comparing figure 5a and 5b, the patterns in figure 5b looked straight lines in Up, Down and Flat positions. Furthermore, figure 5b conspicuously showed patterns in the series, unlike in figure 5a which showed everything as a curve with no indication of where there is an occurrence of patterns.

Again in figure 6, DD means consecutive Down patterns, likewise FFF or UU. And with the patterns symbolically represented, further programming work can be done to compare time series of similar months of several years to find out areas of similarity.

VI. CONCLUSION

Time series analysis cuts across the following activities: time series representation, pattern recognition, similarity search and prediction. This work explored development of time series representation and pattern recognition algorithms from the data mining perspective, and successfully developed an algorithm to mine time series datasets for patterns and also represented the

patterns using symbols (U, D, F) for Up, Down and Flat respectively. Our algorithm is very efficient, easy to understand and implement towards finding patterns in a time series. Researchers in time series analysis from different perspectives like statistics, economics and data mining will benefit immensely from the contributions of this work. Further research can be carried out to incorporate prediction ability into the algorithm.

REFERENCES

- [1]. Abdullah, M; Suman, N and Jie, L (2016). Similarity Search on TS Data: Past, Present and Future. CIKM2016 Tutorial, obtained from <http://www.cs.unm.edu/~mueen/Tutorial/CIKM2016Tutorial.pdf> on 17/4/2017.
- [2]. Álvarez, F.M. (2010) Pattern sequence analysis to forecast time series. Unpublished thesis work. Universidad Pablo De Olavide De Sevilla
- [3]. Badhiye, S.S; Hatwar, K, S. and Chatur, P.N (2015). Trend based Approach for Time Series Representation. International Journal of Computer Applications, Volume 113, No. 16, (0975 – 8887).
- [4]. Bishop, C.M (2006). Pattern Recognition and Machine Learning. Singapore: Springer Science+Business Media, LLC.
- [5]. Burcu, K; Serhan, O and Bora K (2011). Application of Symbolic Piecewise Aggregate Approximation (PAA) Analysis To ECG Signals. Artificial Intelligence & Design Lab, Mechanical Engineering Department, Computer Engineering Department, Izmir Institute of Technology, 35430 Izmir, Turkey. A pdf file obtained 15/4/2017
- [6]. Ding, H; Hui D.; Goce T.; Scheuermann P.; Xiaoyue W., and Keogh E. (2008) "

- Querying and Mining of Time Series Data: experimental comparison of representations and distance measures". Proceedings of the VLDB Endowment VLDB Endowment, Volume 1 Issue 2, pp 1542-1551.
- [7]. Han, J. and Kamber, M. (2001). Data mining concepts and techniques. San Francisco: Morgan Kaufman
- [8]. Jiangling, Y; Yain-Whar, Si and Zhiguo, G (2011). Financial Time Series Segmentation Based On Turning Points. Proceedings of 2011 International Conference on System Science and Engineering, Macau, China - June 2011.
- [9]. Jingpei, D; Weiren, S; Fangyan, D and Kaoru, H (2013). Piecewise Trend Approximation: A Ratio-Based Time Series Representation. Journal of Abstract and Applied Analysis, Volume 2013. <http://dx.doi.org/10.1155/2013/603629>
- [10]. Keogh, E. (2010) Data Mining Time Series Data, in Lovrić (Ed.), International Encyclopaedia of Statistical Science. New York, USA: Springer
- [11]. Keogh, E. (2007). Mining Shape and Time Series Databases with Symbolic Representations. SIGKDD 2007 Tutorial.
- [12]. Keogh, E., Chakrabarti, K., Pazzani, M. & Mehrotra, S. (2001). Locally adaptive dimensionality reduction for indexing large time series databases. In proceedings of ACM SIGMOD Conference on Management of Data. Santa Barbara, CA, May 21-24. pp 151-162.
- [13]. Keogh, E.J; Chu, S; Hart, D and Pazzani, M.J (2001). An Online Algorithm for Segmenting Time Series. In ICDM Proceedings of the IEEE International Conference on Data Mining. Washington DC, USA: IEEE Computer Society, pp. 289 - 296.
- [14]. Keogh E., S. Lonardi, Ratanamabatana C.A (2001). Towards parameter-free data mining. In proceedings of Tenth ACM SIGKDD international Conference on Knowledge Discovery and Data Mining
- [15]. Lin, J; Keogh, E; Lonardi, S and Patel, P. (2002) Finding Motifs in Time Series. ACM SIGKDD, July 23-26,Edmonton, Alberta, Canada.
- [16]. Nguyen, Q.V.H and Duong, T.A (2007). Combining SAX and Piecewise Linear Approximation to Improve Similarity Search on Financial Time Series. International Symposium on Information Technology Convergence, IEEE Computer Society.
- [17]. Pohl, D and Bouchachia, A (2012). Financial Time Series Processing: A Roadmap of Online and Offline Methods. A pdf file downloaded on April 5, 2017.
- [18]. Prasanna, S. and Ezhilmaran, D. (2013). An analysis on Stock Market Prediction using Data Mining Techniques; International Journal of Computer Science & Engineering Technology (IJCSET), Vol. 4 No. 02; p. 49-51.
- [19]. Raj, M.P; Swaminarayan, P.R; Saini, J.R and Parmar, D.k (2015). Applications of Pattern Recognition Algorithms in Agriculture: A Review. Int. J. of Advanced Networking and Applications Volume: 6 Issue: 5 Pp: 2495-2502
- [20]. Singh, S (2000). Pattern Modelling in Time-Series Forecasting. Cybernetics and Systems - An International Journal, vol. 31, issue 1
- [21]. Wu, H; Salzberg, B and Zhang, D (2004). Online event-driven subsequence matching over financial data streams, in Proceedings of the 2004 ACM SIGMOD International Conference on Management of Data, Paris pp. 23-34.
- [22]. Zhang, Z; Jiang, J; Liu, X; Lau, W.C; Wang, H; Wang, S; Song, X and Xu, D (2010). Pattern recognition in stock data based on a new segmentation algorithm, in Lecture Notes in Computer Science, pp. 520-525, Springer Berlin / Heidelberg.