# Predominant Pattern Distribution Model for Noise Distributed Time Series Database

[1]B. Sujatha and [2]S. Chenthur Pandian
[1]Department of CSE, Sengunthar Engineering College Tiruchengode, TamilNadu, India
[2]Dr. Mahalingam College of Engineering and Technology, Pollachi, TamilNadu, India

**Abstract:** A time series is collection of well-defined data sets obtained through repeated measurements of time. Extraction of periodic pattern in a time series database is significant one in data mining problem that predicts and forecasts the future behavior of the data at regular time interval. Periodic pattern mining involves several applications such as prediction, forecasting, detection of unusual activities. The difficulty is not trivial because the data to be examined are regularly noisy and diverse periodicity types (that is symbol, sequence and segment) are to be examined. The whole time series or in a subsection of it to effectively handle various types of noise (to a definite degree) and at the same time to detect different types of periodic patterns. The existing suffix tree based periodic pattern mining algorithm can detect symbol, sequence and segment periodicity in time series data with noise filters for diverse noise kinds. But the running time desired to identify the patterns without redundancy is high. So, to overcome this issue, in this study, Predominant Pattern Distribution Model is introduced with which redundant and unwanted noisy patterns are identified and discarded from the time series data. Predominant patterns are extracted with automatic or user defined threshold of pattern of interest, generated from the dynamic online time series data. Experiments conducted on both synthetic and real data sets of research repositories including protein sequences. Performance of proposed framework is measured and evaluated in terms of periodic pattern mining accuracy, noise distribution rate and predominant pattern occurrence.

**Key words:** Time series, periodic pattern mining, periodicity types, suffix tree, Predominant Pattern Distribution Model

## INTRODUCTION

A time series is a collection of data sets accumulated at regular interval time to indicate certain behavior of an entity. Real time examples of time series data, e.g., meteorological data, network delays, power consumption, stock growth, gene expression data analysis, etc. Data mining is the process of extracting patterns and trends from large amounts of data stored in data repositories that uses statistical and mathematical techniques. Study in time series data mining has concentrated on extracting different types of patterns.

Periodicity detection is a process of discovering temporal regularities of data within the time series and the objective of analyzing a time series is to find whether and how a periodic pattern is repeated within the series. A time series is described by a set of repeating cycles.

A time series is discretized (Rasheed *et al.*, 2011; Elfeky *et al.*, 2005; Chen, 2010) before it is analyzed. Let $T = e_0, e_1, e_2, ..., e_{n-1}$ be a time series with n events where $e_i$ denotes the event recorded at time i, time series T may be discretized by considering m different ranges such that all values in the same range are represented by $\sum$. For example, consider the time series including the hourly number of transactions in a superstore the discretization procedure may define the following mapping by taking the distinguished possible range of transactions {0} transactions: a, {1-100} transactions: b, {101-200} transactions: c, {201-300} transactions: d, {>300} transactions: e, Based on this mapping, the time series T = 147, 162, 185, 296, 76, 0, 0 and 95 can be discretized into T' = cccdbaab. In general, three types of periodic patterns can be detected in a time series: symbol periodicity, sequence periodicity or partial periodic patterns and segment or full-cycle periodicity. A time series is known as symbol periodicity if at least one symbol is repeated periodically. For e.g., in time series T = abd acb cab abc, symbol a is periodic with periodicity p = 3, starting at position zero (stPos = 0). Next, a pattern consisting of more than one symbol may be periodic in a time series and this is known as partial periodic patterns. If the whole time series can be

**Corresponding Author:** B. Sujatha, Department of CSE, Sengunthar Engineering College Tiruchengode, TamilNadu, India

represented as a repetition of pattern or segment and then this kind of periodicity is called segment or full-cycle periodicity. For example, the time series T = abcab abcab abcab has segment periodicity of 5 (p = 5) starting at the first positions (stops = 0). That is T consists of only three occurrences of the segment abcab.

To conclude, time series exists often in the daily life and their analysis could lead to valuable discoveries. Those are identifying three types of periodic patterns, handling asynchronous periodicity by locating periodic patterns that may drift from their expected positions and finding periodic pattern in the whole time series as well as in a subsection of the time series using suffix tree. The algorithm proposed in this study can identify and discard the unwanted and redundant data from the time series data and predominant pattern are extracted with automatic or user defined threshold. It is applicable to biological datasets and it has been analyzed for periodic pattern mining accuracy, noise distribution rate and predominant pattern occurrence.

**Literature review:** Existing work on time series examination approximately faces two kinds of algorithms. The primary group comprises algorithms that need the user to identify the period and then appear only for patterns happening with that period. The second class, conversely is algorithms which appear for all probable periods in the time series. The algorithm group it performs more than the other algorithms by appearing for all probable periods initiating from all potential locations inside a pre-specified range whether the complete time series or a section of the time series. Rasheed *et al.* (2011) presented an algorithm which can identify sign, series (partial) and section (full cycle) periodicity in time series. The algorithm employs suffix tree as the fundamental data structure this permits us to plan the algorithm such that it's most terrible case difficulty of the time series. The algorithm is noise pliant; it has been productively established to effort with substitute, addition, crossing out or a combination of these kinds of noise.

Classification has been employed for representing numerous types of data sets, counting deposits of items, text credentials, graphs and systems. Representing such data is helpful with the budding GPS and RFID knowledge and is significant for efficient haulage and traffic setting up. Lee *et al.* (2011) considered techniques for categorizing routes on road networks. Periodicity mining is employed for forecasting development in time series data. Determining the time at which the time series is cyclic has always been an obstruction for completely

mechanized periodicity mining. In the crisis of noticing, the periodicity time of a database is addressed (Elfeky *et al.*, 2005). Two types of periodicities are distinct and a scalable, computationally proficient algorithm is planned for every type. Conventional outline growth-based strategies for chronological pattern mining obtain patterns supported on the estimated databases recursively. At every level of recursion, they unidirectional produce the span of noticed patterns by one all along the suffix of noticed patterns which wants k stages of recursion to discover a pattern (Chen, 2010).

To decrease the number of models and develop the efficiency of the salgorithm, Lo *et al.* (2007) have also commenced mining blocked iterative patterns, i.e., maximal patterns devoid of any super-pattern containing the similar support. Lo *et al.* (2011) to officially intensify study on iterative pattern mining, the researchers initiate mining iterative generators, i.e., negligible patterns devoid of any sub-pattern containing the similar support. Periodic pattern mining or periodicity detection has numerous applications such as prediction, forecasting, detection of unusual events, etc. The periodic patterns are detected in a time-series database depending on the time intervals (Obulesu *et al.*, 2012).

Rasheed *et al.* (2007) described a dynamic periodicity discovery algorithm to determine periodicity in DNA series. The algorithm supported suffix tree as the fundamental data structure. The proposed strategy by Rasheed and Alhajj (2010) believes the periodicity of substitute substrings, besides allowing for active window to notice the periodicity of convinced illustrations of substrings. Nevertheless, even devoid of nested patterns, the lingo is influential sufficient to detain different forward-temporal stipulations from numerous release source requests (Lo *et al.*, 2007; Lo and Maoz, 2008). In the Move Mine System a position of normally employed touching object mining purposes are constructed and a user-friendly boundary is presented to make possible interactive examination of touching object data mining and bendable alteration of the mining restraints and parameters. Analyzing such data consumes more time and noises on the set of data proceed in it.

Research in time series data mining has concerted on determining diverse kinds of patterns: chronological patterns sequential patterns, episodic connection rules, incomplete cyclic patterns and astonishing patterns to forename a few. This periodicity pulling out methods need the user to identify a time that decides the time at which the time series is cyclic. They imagine that users either recognize the value of the period earlier or are enthusiastic to attempt different period values in anticipation of

acceptable cyclic patterns emerge. Since, the mining process must be performed frequently to attain good results, this trial-and-error method is obviously not competent. Even in the study of time series information with a priori recognized periods, there might be incomprehensible periods and as a result, motivating periodic patterns that will not be exposed. The clarification to these troubles is to invent methods for determining possible periods in time series data. In this research, a highly efficient Periodic Pattern Mining Model is presented for determining the unwanted data in the distributed time series data and processed the database in terms of user specified threshold value.

## MATERIALS AND METHODS

### Proposed predominant pattern distribution model for noise distributed time series database
**Time series database:** A time-series database is a compilation of data values grouped commonly at consistent period of time to reproduce convinced actions of an entity. In genuine globe, there are numerous examples of time-series for instance weather circumstances of a distinct position, expensing patterns, stock development, communication in a superstore, network holdup, power utilization, computer network burden examination and safety contravene discovery, earthquake calculation. The periodicity recognition is a procedure of identifying the temporal regularities inside the time-series and the objective of examining a time-series database is to discover how common a periodic prototype (full or partial) is repetitive inside time intervals. In common, there are three kinds of periodic patterns can be noticed in a time series database. They are declared:

- Symbol periodicity
- Sequence periodicity or partial periodic patterns
- Segment or full-cycle periodicity

**Symbol periodicity:** A time-series is supposed to be a sign periodicity, if no less than one symbol is repetitive occasionally. For instance, in a time-series, let T = abdacbabdabc, symbol a is cyclic inside periodicity p = 3, opening at position zero (StPos = 0).

**Sequence periodicity:** A time-series is supposed to be a sequence periodicity, if more than one symbol might be cyclic and it is also termed as limited periodic patterns. For example, in a time-series database let T = bbaaabbdabcaabbcabcd then the series ab is cyclic inside periodicity p = 4 initiating at location 4 (StPos = 4).

**Segment priodicity:** A time-series is supposed to be a segment periodicity if the entire time-series can be typically symbolized as a replication of a model or segment and it is also recognized as full-cycle periodicity. For example, in a time-series database let T = abcababcababcab contain Segment periodicity of 5 (p = 5) initiating at the first location (stPos = 0), i.e., T comprises of only three incidences of the segment abcab.

It is not essential to forever have ideal periodicity in a time series as in the exceeding three examples. Generally, the degree of excellence is symbolized by confidence which is 100% in the three examples. Alternatively, real-life instances are typically distinguished by the occurrence of noise in the data and this depressingly concerns the confidence level. The confidence of a prototype is termed as the part of its genuine occurrence in the series over its projected ideal occurrence in the series. The genuine and estimated perfect frequency are both the similar in the specified three instances nevertheless, in the time series $T_x$ = abefd abcde acbcd abefa, the pattern ab initiating at location 0 with p = 5 contain four and five as its genuine and estimated ideal frequencies, correspondingly; so the self-reliance of ab is 4/5.

The first part of the research concentrate on building highly efficient noisy removal to the unwanted data distributed across the time series. The architecture diagram of the proposed Predominant Pattern Distribution Model for Noise Distributed time series database [PPDMND] is shown in Fig. 1.

From the Fig. 1, it is being observed that predominant patterns are extracted with automatic or user defined threshold of pattern of interest, generated from the dynamic online time series data. The predominant pattern distribution model is introduced with which redundant and unwanted noisy patterns are identified and discarded from the time series data.
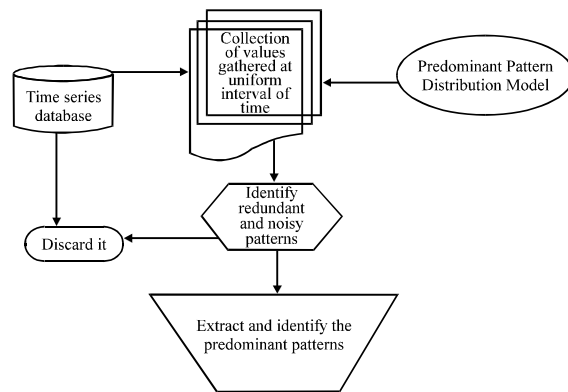


Fig. 1: Architecture diagram of the proposed PPDMND

**Periodicity detection in presence of noise:** Three types of noise usually measured in time series data are substitution, addition and deletion noise. In substitution noise, various symbols in the dishonored time series are restored at arbitrary with further symbols. In case of addition and deletion noise, several symbols are interleaved or deleted, correspondingly, arbitrarily at diverse locations (or time values). Noise is also a combination of these three kinds; for example, substitution type noise resources the consistent combination of Replacement (R) and addition (I) noise. When the time series is moreover completely cyclic or includes only substitution noise and achieves inadequately in the occurrence of addition or deletion noise. This is since insertion and deletion noise develops or deals the time axis important to move of the imaginative time series values.

For instance, time series T = abcabcabc after adding symbol b at locations 2 and 6 would be T' = abbcabcabbc. The event for sign a in T is (0; 3; 6), as it is (0; 4; 7) in T'. It is very obvious that when the time series is indistinct by addition and or removal noise, STNR do not execute well. To enhance the noise detection process, this system presents a Predominant Distribution Model to remove the redundant and unwanted noisy distribution of data in the time series dataset.

With the set of patterns obtained from the time series database, identify the predominant patterns. Before that set a user defined threshold value T for all the patterns derived. The threshold values are used to identify the interesting patterns without any redundancy in data. After that Perdominant Distribution Model is presented to identify the unwanted data distributed on the time series dataset by using the threshold values. The procedure below describes the Predominant Pattern Distribution Model (Fig. 2).

With the above process, a highly efficient Predominant Pattern Distribution Model is presented for removal of noise and redundancy for the data being observed in the time series data.

**Experimental evaluation:** In this study, researchers provide the outcome of numerous experiments that have been processed using both synthetic and real data. Researchers also account the outcome of trying different individuality of the proposed PPDMND algorithm against other existing algorithms like STNR (Suffix-Tree-based Noise-Resilient algorithm). There are two kinds of algorithms explained in the literature: algorithms in the initial group discover periodic patterns for a precise period value and those in the other group ensure the time series for each and every time. The first set of test is
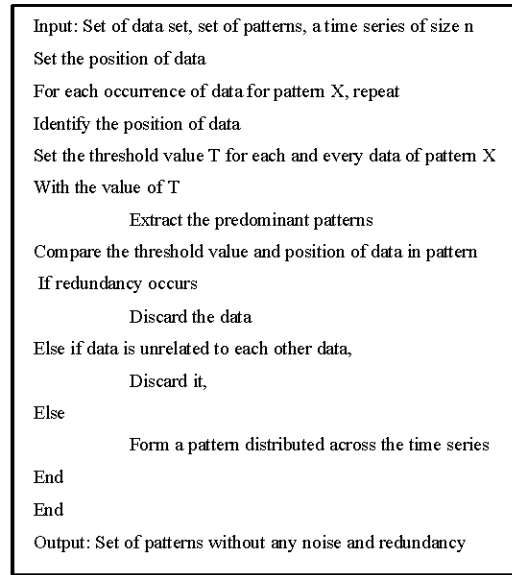
```
Input: Set of data set, set of patterns, a time series of size n
Set the position of data
For each occurrence of data for pattern X, repeat
Identify the position of data
Set the threshold value T for each and every data of pattern X
With the value of T
        Extract the predominant patterns
Compare the threshold value and position of data in pattern
If redundancy occurs
        Discard the data
Else if data is unrelated to each other data,
        Discard it,
Else
        Form a pattern distributed across the time series
End
End
Output: Set of patterns without any noise and redundancy
```

Fig. 2: Process of PPDMND

Table 1: Periodic pattern in packet data

| Pattern | Period | Confidence |
| --- | --- | --- |
| aa | 2 | 0.42 |
| aa | 2 | 0.42 |
| aa | 2 | 0.47 |
| aaa | 3 | 0.44 |
| aaa | 3 | 0.43 |
| aaa | 3 | 0.47 |
| aaa*a | 5 | 0.44 |
| aaaa* | 5 | 0.36 |
| aaa** | 5 | 0.44 |

committed to reveal the comprehensiveness of PPDMND in the sense that it should be capable to discover a time once it subsists in the time series. Researchers check how PPDMND suits this on both synthetic and real data.

The synthetic data was produced in the similar way as done. The parameters prescribed through data creation are data allocation (uniform or normal), alphabet extent (number of exclusive signs in the data), dimension of the data (quantity of symbols in the data), time period size, and the style and quantity of noise in the data. A datum might hold substitution, addition and removal of noise or any combination of these kinds of noise. For real data experiments, the packet data have been employed where the researchers established the intense periodic patterns, i.e., the area where the time series is typically cyclic the three regions which were originated cyclic for symbol a and its patterns are shown in Table 1.

The PPDMND algorithm does discover all the periods in the intense periodic section, generally at superior poise level. Since, the confidence level of the periodicity patterns increases, the noise level in the periodic data is low compared to an existing STNR technique. The

performance of the proposed Predominant Pattern Distribution Model for Noise Distributed time series database [PPDMND] is measured in terms of:

- Periodic pattern mining accuracy
- Noise distribution rate
- Execution time

**RESULTS AND DISCUSSION**

In this research, researchers have seen how the noise has been removed from the set of periodic patterns in the time series dataset. At first, periodic pattern distribution model applied to the time series database to remove the redundant and unwanted noisy patterns present in it. Experiments have been conducted with several time series dataset and the performance has also been evaluated. Table 1 and Fig. 3 describe the performance of the proposed Predominant Pattern Distribution Model for noise distributed time series database and compared the results with an existing STNR technique (Rasheed *et al.*, 2011).

The accuracy of periodic patterns in the time series dataset is analyzed based on the number of patterns formed is illustrated in the Table 2 for the proposed PPDMND and existing STNR.

Figure 3 describes the accuracy of the periodic patterns in the time series dataset is analyzed based on the number of patterns formed. The accuracy of the patterns is measured based on the confidence level of the patterns. The confidence of a pattern is termed as the ratio of its genuine occurrence in the series over its estimated ideal frequency in the series. Since, the proposed PPDMND set a user defined threshold value for patterns generated from the dynamic time series data, the predominant patterns are automatically extracted. With that, redundant and unwanted data in the patterns are easily identified by building a Predominant Pattern Distribution Model. So, the accuracy of the periodic patterns in the dataset is high in the proposed PPDMND. Compared to an existing STNR (Suffix-Tree-based Noise-Resilient algorithm), the proposed PPDMND provides high level of accuracy on periodic patterns for the given dataset and the variance is 35-45% high in it. The noise distribution rate for the periodic patterns present in the time series dataset is illustrated in the Table 3.

Figure 4 describes the noise distribution rate for the periods obtained for assigning the patterns present in the time series dataset. The noise distribution rate is measured in terms of presence of noise in the specified data distributed in the given dataset. Since, the proposed PPDMND presented a predominant pattern distribution model, the noisy and redundant data are removed from the distributed dataset by setting a user defined threshold. With the threshold value, the mechanism makes insertion,deletion and modification of patterns in the given dataset. While changing the patterns, the user defined threshold value might change. So, the process of noise distribution is less in the proposed PPDMND.

Table 2: No. of patterns vs. periodic pattern mining accuracy

| | Periodic pattern mining accuracy (%) | |
|---|---|---|
| No. of patterns | Proposed PPDMND | Existing STNR |
| 25 | 54 | 30 |
| 50 | 60 | 33 |
| 75 | 68 | 38 |
| 100 | 73 | 40 |
| 125 | 78 | 43 |
| 150 | 83 | 46 |
| 175 | 85 | 50 |

Table 3: Periods vs. noise distribution rate

| | Noise distribution rate | |
|---|---|---|
| Periods | Proposed PPDMND | Existing STNR |
| 20 | 10 | 20 |
| 40 | 15 | 25 |
| 60 | 13 | 30 |
| 80 | 18 | 28 |
| 100 | 20 | 35 |
| 120 | 25 | 33 |
| 140 | 23 | 40 |

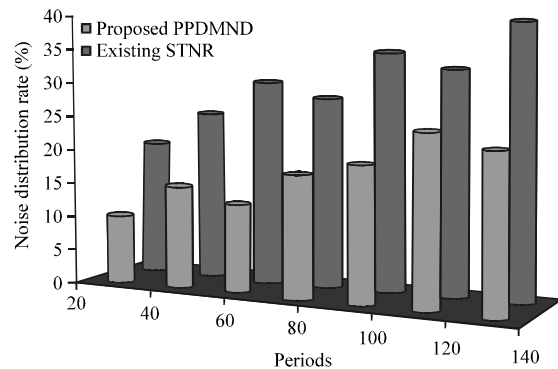

Fig. 3: No of patterns vs. periodic pattern mining accuracy



Fig. 4: Periods vs. noise distribution rate

Fig. 5: Noise ratio vs. execution time

Table 4: Noise ratio vs. execution time

| | Execution time (sec) | |
| Noise ratio | Proposed PPDMND | Existing STNR |
| --- | --- | --- |
| 0.1 | 0.5 | 2.0 |
| 0.2 | 1.0 | 2.3 |
| 0.3 | 0.9 | 2.5 |
| 0.4 | 0.5 | 2.0 |
| 0.5 | 0.7 | 1.9 |
| 0.6 | 0.9 | 2.2 |
| 0.7 | 1.0 | 2.6 |

Compared to an existing STNR (Suffix-Tree-based Noise-Resilient algorithm), the proposed PPDMND provides less noise distribution on the given dataset and the variance is 45-55% low in it (Table 4).

The next set of researches determines the impact of noise ratio on the time presentation of the proposed PPDMND. For this trial, researchers set the time series span, size, alphabet size and data allocation and considered the force of unstable noise ratio on time recital of the algorithm. The results are plotted in Fig. 5. These experiments have been processed employing the two algorithms: STNR and PPDMND. It is true that PPDMND consumes more time when the noise ratio lies between 10-15% but when the noise ratio is level is high, PPDMND tends to obtain the analogous time. As the results show, the noise ratio does not concern the time presentation of existing STNR.

Finally, it is being observed that the proposed PPDMND efficiently removed the noise and discard the unwanted data distributed in the time series dataset by adapting the predominant pattern distribution model. At first, predominant patterns are extracted from the dataset by setting a user threshold pattern of interest generated from the dynamic online time series data. Then, the periodic patterns are extracted devoid of any noise in the generated patterns.

## CONCLUSION

In this study, PPDMND is presented as a Predominant Pattern Distribution Model for removal of noise in time series data. The algorithm is noise-resilient and processes the time series dataset efficiently even in the worst case. The single algorithm can discover symbol, series (partial periodic) and section (full cycle) periodicity in the time series. It can also discover the periodicity inside a part of the time series. A highly efficient noisy removal strategy is introduced to eradicate the unwanted data distributed across the time series. A Predominant Pattern Distribution Model identified and discarded the redundant and unwanted noisy patterns from the time series data. Predominant patterns are extracted with automatic or user defined threshold of pattern of interest, generated from the dynamic online time series data. Several experiments are performed to show the time behavior, accuracy and noise resilience characteristics of the data. The algorithm is processed on both real and synthetic data. The reported results demonstrated the efficiency and accuracy of predominant patterns in the time series dataset.

## REFERENCES

Chen, J., 2010. An updown directed acyclic graph approach for sequential pattern mining. IEEE Trans. Knowledge Data Eng., 22: 913-928.

Elfeky, M.G., W.G. Aref and A.K. Elmagarmid, 2005. Periodicity detection in time series databases IEEE Trans. Knowl. Data Eng., 17: 875-887.

Lee, J.G., J. Han, X. Li and H. Cheng, 2011. Mining discriminative patterns for classifying trajectories on road networks. IEEE Trans. Knowledge Data Eng., 23: 713-726.

Lo, D. and S. Maoz, 2008. Mining scenario-based triggers and effects. Proceedings of the 23rd IEEE/ACM International Conference on Automated Software Engineering, September 15-19, 2008, L'Aquila, pp: 109-118.

Lo, D., J. Li, I. Wong and S.C. Khoo, 2011. Mining iterative generators and representative rules for software specification discovery. IEEE Trans. Knowl. Data Eng., 23: 282-296.

Lo, D., S. Maoz and S.C. Khoo, 2007. Mining modal scenario-based specifications from execution traces of reactive systems. Proceedings of the 22nd IEEE/ACM International Conference on Automated Software Engineering, March 11-15, 2007, Seoul, Korea.

Obulesu, O., A.R.M. Reddy and K. Suresh, 2012. Finding maximal periodic patterns and pruning strategy in spatiotemporal databases. Int. J. Adv. Res. Comput. Sci. Software Eng., 2: 423-426.

Rasheed, F. and R. Alhajj, 2010. STNR: A suffix tree based noise resilient algorithm for periodicity detection in time series databases. Applied Intell., 32: 267-278.

Rasheed, F., M. Alshalalfa and R. Alhajj, 2007. Adapting machine learning technique for periodicity detection in nucleosomal locations in sequences. Proceedings of the 8th International Conference Intelligent Data Engineering and Automated Learning, December 16-19, 2007, Birmingham, UK., pp: 870-879.

Rasheed, F., M. Alshalalfa and R. Alhajj, 2011. Efficient periodicity mining in time series databases using suffix trees. IEEE Trans. Knowledge Data Eng., 23: 79-94.