

K-Shape Efficient and Accurate Clustering Algorithm Using ARIMA Model

Mr. B. Venkateswarlu¹, M. Anitha²

¹Assoc. Prof, Dept. of Information Technology, RVR&JCCE, AP,
Guntur, 522017, India

²Student, Dept. of Information Technology, RVR&JCCE, AP, Guntur,
522017, India

Abstract – *k*-Shape depends on a ascendable unvaried refinement procedure, that creates uniform and well-separated clusters. As its distance live, *k*-Shape uses a normalized version of the cross-correlation live so as to think about the shapes of your time series whereas examination them. supported the properties of that distance live, we tend to develop a way to work out cluster centroids, that ar employed in each iteration to update the assignment of your time series to clusters. To demonstrate the lustiness of *k*-Shape, we tend to perform an in depth experimental analysis of our approach against partitional, ranked, and spectral agglomeration ways, with combos of the foremost competitive distance measures. *k*-Shape outperforms all ascendable approaches in terms of accuracy. *k*-Shape emerges as a domain-independent, extremely correct, and extremely economical agglomeration approach for statistic with broad applications.

INTRODUCTION

Most time-series analysis techniques, together with bunch, critically rely upon the selection of distance live. A key issue once examination 2 time-series sequences are a way to handle the variability of distortions.

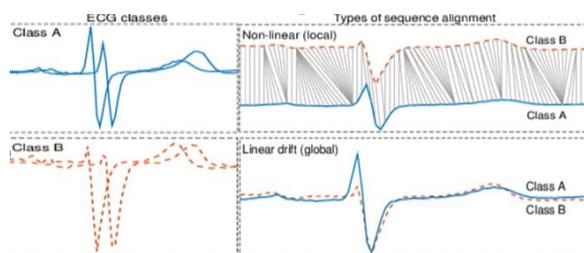


Figure 1: ECG sequence examples and types of alignments for the two classes of the ECGFiveDays

To illustrate now, contemplate the well-known ECG FiveDays dataset [1], with electrocardiogram sequences recorded for identical patient on 2 totally different days. Where as the sequences appear similar overall, they exhibit patterns that belong in one in every of the 2 distinct categories (see Figure 1): category A is

characterized by a pointy rise, a drop, and another gradual increase whereas category B is characterized by a gradual increase, a drop, and another gradual increase. Time-Series Invariance's: Based on the domain, sequences area unit typically distorted in how, and distance measures have to be compelled to satisfy variety of invariance so as to match sequences meaningfully. During this section, we have a tendency to review common time-series distortions and their invariances. Scaling and translation invariance's: In several cases, it's helpful to acknowledge the similarity of sequences despite variations in amplitude (scaling) and offset (translation). In alternative words, reworking a sequence $x \rightarrow$ as $x \rightarrow = ax \rightarrow + b$, wherever a and b area unit constants, shouldn't similarity to alternative sequences. as an example, these invariances can be helpful to investigate differences due to the season in currency values on interchange markets while not being biased by inflation.

PROBLEM STATEMENT

We address the problem of domain-independent, accurate, and scalable clustering of time series into k

clusters, for a given value of the target number of clusters k . Even though different domains might require different invariances to data distortions, we focus on distance measures that offer invariances to scaling and shifting, which are generally sufficient. Furthermore, to easily adopt such distance measures, we focus our analysis on raw based clustering approaches. Next, we introduce k-Shape, our novel clustering algorithm.

EXISTING SYSTEM

Euclidean Distance:

Euclidean distance is one amongst the foremost used distance metric. it's calculated exploitation scientist Distance formula by setting p 's worth to 2. This could update the house'd' formula as below:

$$ED(\vec{x}, \vec{y}) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}$$

Let's stop for a while! can this formula look familiar? Well affirmative, we tend to tend to easily saw this formula on high of throughout this text whereas discussing "Pythagorean Theorem". Euclidean distance formula is also accustomed calculate the house between a pair of information points in associate degree extremely plane.

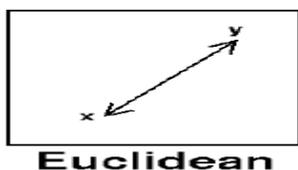


Figure 2: Euclidean Distance

Dynamic Time Warping:

Dynamic Time warp (DTW) is one amongst the algorithms for measure similarity between a pair of temporal sequences, which may vary in speed. as associate degree example, similarities in walking is also detected exploitation DTW, tho' one person was walking faster than the alternative, or if there are accelerations Associate in Nursing decelerations throughout the course of associate degree observation. DTW has been applied to temporal sequences of video, audio, and graphics information therefore, any information which can be become a linear sequence is also analyzed with DTW.

DTW is also a method that calculates Associate in Nursing best match between a pair of given sequences (e.g. time series) Every index from the first sequence

ought to be matched with one or plenty of indices from the alternative sequence, and also the different means around. The initial index from the first sequence ought to be matched with the first index from the alternative sequence (but it does not have to be compelled to be its alone match). The last index from the first sequence ought to be matched with the last index from the alternative sequence (but it does not have to be compelled to be its alone match). The mapping of the indices from the first sequence to indices from the alternative sequence ought to be monotonically increasing, and also the different means around.

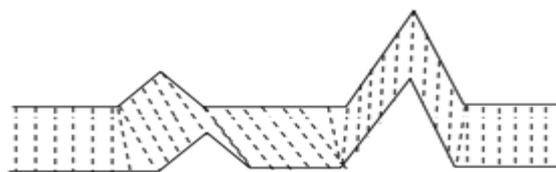


Figure 3: Dynamic Time Warping

Applications of Dynamic Time Warping:

1. Automatic speech recognition- to address totally different speaking speeds.
2. Speaker recognition and on-line signature recognition- It can even be employed in partial form matching application.

PROPOSED SYSTEM

Autoregressive integrated moving average (ARIMA) model is to boot a generalization of Associate in Nursing autoregressive moving average (ARMA) model. Every of these models unit of activity fitted to info information either to higher understand the information or to predict future points among the series (forecasting). ARIMA models unit of activity applied in some cases where information show proof of non-stationary, where Associate in Nursing initial differencing step (corresponding to the "integrated").The AR a locality of ARIMA indicates that the evolving variable of interest is regressed on its own lagged (i.e., prior) values. The MA [*fr1] indicates that the regression error is totally a linear combination of error terms whose values occurred contemporaneously and at varied times among the past. The I (for "integrated") indicates that the information values unit replaced with the excellence between their values and additionally the previous values (and this differencing technique may unit performed quite once). the aim of each of these picks is to create the model

match the information likewise as realizable. Non-seasonal ARIMA models unit of activity typically denoted ARIMA(p,d,q) where parameters p, d, and letter unit of activity non-negative integers, p is that the order (number of a jiffy lags) of the autoregressive model, d is that the degree of differencing (the sort of times the information have had past values subtracted), and letter is that the order of the moving-average model. seasonal ARIMA models unit of activity typically denoted ARIMA(p,d,q)(P,D,Q)m, where m refers to the quantity of periods in each season, and additionally the majuscule P,D,Q see the autoregressive, differencing, and moving average terms for the seasonal a locality of the ARIMA model.

4.1 Auto-Regressive Model:

Models future values as a operate of recent past consecutive values

$$Y_t = f(Y_{t-1}, Y_{t-2}, \dots, Y_{t-p}, Y_t)$$

Representation: Associate in Nursing AR model with past p values is denoted as AR(p). $Y_t = \emptyset_0 + \emptyset_1 Y_{t-1} + \emptyset_2 Y_{t-2} + \dots + \emptyset_p Y_{t-p} + Y_t$

4.2 Moving Average Model:

Models future values as a operate of recent past consecutive error terms

$$Y_t = f(\epsilon_t, \epsilon_{t-1}, \epsilon_{t-2}, \dots, \epsilon_{t-q})$$

Representation: Associate in Nursing MA model with past q values is denoted as MA(q)

4.3 ARIMA MODEL:

Auto regressive Moving Average (ARMA) model:

Models future values as a operate of recent past consecutive values and error terms

$$Y_t = f(\epsilon_t, \epsilon_{t-1}, \epsilon_{t-2}, \dots, \epsilon_{t-q}, Y_{t-1}, Y_{t-2}, \dots, Y_{t-p})$$

Representation: ARMA(p, q) model

$$Y_t = \emptyset_0 + \emptyset_1 Y_{t-1} + \emptyset_2 Y_{t-2} + \dots + \emptyset_p Y_{t-p} + \epsilon_t + \$1\epsilon_{t-1} + \$2\epsilon_{t-2} + \dots + \$q\epsilon_{t-q}$$

An ARIMA model is to boot a class of arithmetic models for analyzing and statement info information. It expressly caters to a gaggle of traditional structures in info information, and in and of itself provides a

straightforward but powerful technique for making skillful info forecasts. ARIMA is Associate in nursing kind that stands for Autoregressive Integrated Moving Average. it is a generalization of the less complicated Autoregressive Moving Average and adds the notion of integration.

AR: Auto regression. A model that uses the dependent relationship between Associate in Nursing observation and kind of favor of lagged observations.

I: Integrated. The use of differencing of raw observations (e.g. subtracting Associate in Nursing observation from Associate in Nursing observation at the previous time step) thus on manufacture the information stationary.

MA: Moving Average. A model that uses the dependency between Associate in Nursing observation and a residual error from a moving average model applied to lagged observations.

Each of these elements unit of activity expressly per the model as a parameter. Associate in Nursing everyday notation is employed of ARIMA(p,d,q) where the parameters unit of activity substituted with vary values to quickly indicate the actual ARIMA model being utilized.

The parameters of the ARIMA model unit of activity written as follows:

p: the quantity of lag observations enclosed among the model, to boot remarked as a results of the lag order.

d: the quantity of times that the raw observations unit of activity differenced, to boot remarked as a results of the degree of differencing.

q: the size of the moving average window, to boot remarked as a results of the order of moving average.

A regression toward the mean model is formed beside the specified varied and kind of terms, and additionally the information is prepared by a degree of differencing thus on manufacture it stationary, i.e. to induce eliminate trend and seasonal structures that negatively have an impact on the regression model. A value of zero unit typically used for a parameter, that indicates to not use that a district of the model. This way, the ARIMA model unit typically organized to perform the operate of Associate in Nursing ARMA model, and even a straightforward AR, I, or MA model. Adopting Associate in Nursing ARIMA model for Associate in

nursing info assumes that the underlying technique that generated the observations is Associate in Nursing ARIMA technique. This may seem obvious, but helps to encourage the requirement to substantiate the assumptions of the model among the raw observations and among the residual errors of forecasts from the model.

METHODOLOGY

Our objective is to develop a domain-independent, accurate, and climbable algorithmic rule for time-series clump, with a distance live that is invariant to scaling and shifting. we tend to propose k-type, a novel centroid-based clump algorithmic rule which can preserve the shapes of time-series sequences. Specifically, we tend to initial discuss our distance live, that depends on the cross-correlation live. supported this distance live, we tend to propose a method to reason centroids of time-series clusters. Finally, we tend to explain our k-Shape clump algorithmic rule, that depends on degree unvarying refinement procedure that scales linearly at intervals the vary of sequences and generates homogeneous and well separated clusters.

Algorithm 3: $[IDX, C] = k\text{-Shape}(X, k)$

Input: X is an n -by- m matrix containing n time series of length m that are initially z -normalized.
 k is the number of clusters to produce.
Output: IDX is an n -by-1 vector containing the assignment of n time series to k clusters (initialized randomly).
 C is a k -by- m matrix containing k centroids of length m (initialized as vectors with all zeros).

```

1  iter ← 0
2  IDX' ← []
3  while IDX' ≠ IDX' and iter < 100 do
4    IDX' ← IDX
5    // Refinement step
6    for j ← 1 to k do
7      X' ← []
8      for i ← 1 to n do
9        if IDX(i) = j then
10       X' ← [X'; X(i)]
11     C(j) ← ShapeExtraction(X', C(j)) // Algorithm 2
12 // Assignment step
13 for i ← 1 to n do
14   mindist ← ∞
15   for j ← 1 to k do
16     [dist, x] ← SBD(C(j), X(i)) // Algorithm 1
17     if dist < mindist then
18       mindist ← dist
19       IDX(i) ← j
20   iter ← iter + 1

```

CONCLUSION

k-Shape compares statistic with efficiency and computes centroids effectively below the scaling and shift invariances. Our in depth analysis shows that k-Shape outperforms all progressive partitional, stratified, and spectral agglomeration approaches, with only 1 methodology achieving similar performance. Curiously, this methodology is 2 orders of magnitude slower than k-Shape and its distance live needs calibration, not like that for k-Shape. Overall, k-Shape may be a domain-independent, accurate, and scalable approach for time-series agglomeration.

ACKNOWLEDGEMENT

The authors would like to thank and acknowledge the DST-FIST (Govt. of India) for funding to set up the research computing facilities at RVR & JC College of Engineering.

REFERENCES

- [1] A. J. Bagnall and G. J. Janacek. agglomeration time series from ARMA models with clipped knowledge. In *KDD*, pages 49–58, 2004. The UCR statistic Classification/Clustering Homepage. http://www.cs.ucr.edu/~eamonn/time_series_data.
- [2] D. Aloise, A. Deshpande, P. Hansen, and P. Popat. NP-hardness of geometer sum-of-squares agglomeration. *Machine Learning*.
- [3] A. J. Bagnall and G. J. Janacek. agglomeration statistic from ARMA models with clipped knowledge. In *KDD*, pages 49–58, 2004. The UCR statistic Classification/ClusteringHomepage. http://www.cs.ucr.edu/~eamonn/time_series_data. Accessed: could 2014.
- [4] G. E. Batista, E. J. Keogh, O. M. Tataw, and V. M. de Souza. CID: associate degree economical complexity-invariant distance for statistic. data processing and data Discovery, pages 1–36, 2013.
- [5] D. J. Berndt and J. Clifford. victimization dynamic time distortion to seek out patterns in statistic. In *AAAI Workshop on KDD*, pages 359–370, 1994.