

Time-series classification using matrix-based methods: Application to blackhole state identification of RXTE satellite data

First A. Author, *Fellow, IEEE*, Second B. Author, and Third C. Author, Jr., *Member, IEEE*

Abstract—Across diverse domains such as medicine, weather, finance, agriculture, astronomy, etc., it is required to deal with timeseries of measurements. Classification of timeseries as stochastic (noise-like) or non-stochastic (which has a well-defined structure), helps understand the underlying phenomenon. The methods used to accomplish this classification are either : (i) Correlation Integral (CI)-based or (ii) Entropy-based approaches, both of which are computationally expensive. In this work, we propose two matrix-based methods to achieve stochastic vs non-stochastic classification, without requiring the computation-intensive phase space. The proposed matrix-based methods are: (a) SVD-decomposition followed by topological analysis (using Betti number descriptors) (b) PCA-based technique. The proposed methods have been applied to synthetic data, as proof of concept. The utility of the methods is illustrated on astronomy data which are 12 categories of timeseries pertaining to blackhole *GRS 1915 + 105*, obtained from RXTE satellite. Comparisons of obtained results with those in literature are also presented. The order of computational complexity using the proposed approaches is $XXXXXX$ of N , where N is the length of the timeseries. In contrast, CI based approaches require $XXXX$, while Entropy-based approaches need $XXXXX$. It is found that among the proposed matrix based methods, SVD analysis concurs with CI based analysis on all 12 categories of time series utilized. However, the inference using PCA based approach illustrates that one class among the 12 turns out to be inconsistent with the other approaches. Investigation into these (in)consistencies is expected to have long standing implications in astrophysics and otherwise.

Index Terms—Timeseries classification, stochastic, non-stochastic, SVD analysis, PCA analysis

I. INTRODUCTION

Several real-world phenomena are studied by collecting associated measurements over time, popularly called as time-series. Timeseries classification as stochastic (noise-like) or non-stochastic (which has a well-defined structure), is the first step in understanding the underlying physical phenomenon. Standard stochastic signals such as white noise, pink noise, etc. exhibit characteristics such as nearly zero auto-correlation coefficients for all possible values of lags and a power spectral

density that decays with frequency. The rate of decay determines the kind of noise. On the other hand, standard non-stochastic signals such as Logistic map (at growth rate = 4), Lorenz system result in timeseries that exhibit a well-defined structure, such as having a certain number of fixed points. For such phenomena, computing parameters such as Correlation Dimension helps in revealing the underlying dynamics. However, for stochastic timeseries the Correlation Dimension never saturates. Hence if the goal of the study is to check if the timeseries is stochastic, then such computations must be avoided.

In literature, methods that accomplish this classification can be broadly categorised as : (i) Correlation Integral based (ii) Entropy-based. Correlation Integral based approach was proposed in [?]. It is a computation-intensive process, since the Correlation Integral (CI) needs to be computed for different choices of Embedding dimension, which can only be approximated using the autocorrelation plot. Besides, it is well-known that this value of correlation dimension does not saturate for a stochastic time series. Hence to establish if the considered timeseries is stochastic, this computation needs to be repeated for a large range of values of Embedding Dimension, making the order of computations needed greater by that factor. Entropy-based approaches [?], [?], [?] utilize concepts of phase-space. This is also a computation-intensive process, with several assumptions about the phenomenon that resulted in the timeseries. TODO write something more on Entropy based methods.

The problem of stochastic vs non-stochastic classification is important for one of the challenging problems in astrophysics, which could lead to the understanding of black holes. As a black hole cannot be seen directly, to identify it, one has to look for its environment forming a disc-like structure by the infalling matter called accretion disc. In this work, we focus on the black hole source *GRS 1915+105*, which presents several intriguing facets. It has been divided into 12 different temporal categories: α , β , γ , δ , λ , κ , μ , ν , ρ , ϕ , χ and θ [?], with their respective distinct timeseries. One fundamental aspect of the understanding is to determine if the black hole source is stochastic or non-stochastic (implying turbulent nature of the system). There are studies reported that utilize the Correlation Integral (CI) approach to determine the characterization of this specific black hole data [?], [?]. However, in this work, we propose to utilize matrix-based methods, Principal Component Analysis (PCA) and Singular Value Decomposition (SVD), to understand the same data. It is useful to compare the inferences

This paragraph of the first footnote will contain the date on which you submitted your paper for review. It will also contain support information, including sponsor and financial support acknowledgment. For example, “This work was supported in part by the U.S. Department of Commerce under Grant BS123456.”

The next few paragraphs should contain the authors’ current affiliations, including current address and e-mail. For example, F. A. Author is with the National Institute of Standards and Technology, Boulder, CO 80305 USA (e-mail: author@boulder.nist.gov).

S. B. Author, Jr., was with Rice University, Houston, TX 77005 USA. He is now with the Department of Physics, Colorado State University, Fort Collins, CO 80523 USA (e-mail: author@lamar.colostate.edu).

obtained using these distinct approaches; the implications of the (dis)similarities in inferences, if any, could lead to better understanding of the temporal dynamics of the system.

It is widely known that the true nature of the source is understood by studying both temporal and spectral features. If the source radiation is temperature dependent, it is called multicolour blackbody or “diskbb” [?]. On the other hand, the temperature independent radiation consists of power-law tail (“PL”) [?], [?]. The difference in their implications lies in the fact that the former leads the underlying accretion disc around the black hole to be geometrically thin, while the latter leads to a geometrically thick disc. Studies in literature combine these observations into four possible black hole states [?]:

- 1) Non-stochastic and diskbb: Keplerian disc [?].
- 2) Non-stochastic and PL: Advection Dominated Accretion Flow (ADAF) [?].
- 3) Stochastic and diskbb: Slim disc [?].
- 4) Stochastic and PL: General Advective Accretion Flow (GAFF) [?], [?].

The contributions of this paper:

- SVD decomposition of the data matrix is used for identifying temporal dynamics of the time series as in [?]. The novelty in the work is that this decomposition is followed by topological analysis of the plot involving the top two right singular vectors for classification of a time series as stochastic or non-stochastic.
- PCA, which is widely used approach for decorrelating features and dimensionality reduction, is utilized to classify a time series as stochastic vs non-stochastic. We propose a novel approach that hierarchically splits the timeseries, computing eigenvalue ratios of covariance matrix of data sub-intervals, exhaustively. Multiple features are devised from the eigenvalue ratios, which are utilized as input features to a linear SVM for classification as stochastic or non-stochastic.

II. RELATED WORK

Several groups have worked on distinguishing between stochastic and non-stochastic time series. The idea of utilizing Permutation Entropy (PE) to determine the complexity measure of a time series was explored in [?]. In the work reported in [?], PE was used to parameterize a given time series followed by classification using Neural Network. The paper explored the idea of utilizing PE of a time series to determine if it is strongly correlated with known stochastic signals (noise). The claim was that for non-stochastic signals the deviation of the parameter is relatively large as compared to that of the parameter of a stochastic signal. Another set of reported studies are based on graph theory. In the work reported in [?], the authors have utilized the horizontal visibility algorithm in order to distinguish between stochastic and non-stochastic processes. A recent work, reported in [?], mapped time series into graphs and computed various topological properties, which they called *NetF*, capturing measures such as centrality, distance, connectivity etc. PCA was applied on the *NetF* feature matrix and clustering was performed on the principal components.

In the approach outlined in [?], the authors combined the idea of sparsity and machine learning with non-linear dynamical systems, in order to determine the governing dynamics. Sparse regression was used to determine the fewest terms in the equations that govern the dynamics of the phenomenon. The user-defined dictionary of basis functions consists of well-known functions such as polynomials, trigonometric functions and exponentials. However, the optimal choice of dictionary for a specific choice of problem remains a challenge.

In this work, we propose to utilize classical matrix based methods which do not require any assumptions about the underlying phenomenon.

III. PROPOSED METHOD

In this work, we propose two different matrix based approaches to characterize time series as stochastic vs non-stochastic. They are 1) SVD decomposition followed by topological analysis (using Betti number descriptors) and 2) PCA derived features followed by SVM classification. Proof of Concept on synthetic signals is also presented.

A. SVD based approach

In this approach, we form uncorrelated observation vectors from the raw time series data by choosing an approximate value of embedding dimension [?] using autocorrelation plot. Data matrix, D , is formed with each row as the time shifted version of the original time series. The time shift is chosen to be large enough so that each column can be viewed as a different observation vector of the same time-evolving phenomenon. Temporal dynamics is understood by utilizing the right singular vectors of the SVD decomposition of D as given in equation (1) below. We consider the top two right singular vectors, $E1$ and $E2$, for our analysis.

$$D = U\Sigma V^T. \quad (1)$$

We observe the topology of the plot $E1$ vs $E2$. For non-stochastic time series this plot is expected to show a specific pattern (attractor behavior, where the plot follows a structured trajectory leaving well-defined voids). On the other hand, $E1$ Vs $E2$ plot for a stochastic time series, appears as a single blob without any voids. The topology of the $E1$ vs $E2$ plot is captured using Betti numbers [?]. Betti number descriptor for a d -dimensional manifold is a vector of d integers which is represented as $\beta = (\beta_0, \beta_1 \dots \beta_{d-1})$. Here β_0 is the number of blobs (connected components) and β_k represents number of k -dimensional holes for $k > 0$. The $E1$ vs $E2$ plots are 2-D manifolds, which are described by $\beta = (\beta_0, \beta_1)$. For a stochastic time series the values of β_0 and β_1 are expected to be 1 and 0 respectively, as the $E1$ vs $E2$ plot consists of one single blob. Hence the $L1$ -norm of a stochastic time series will be 1. However, for a non-stochastic time series, we observe that the value of β_0 can be greater than 1 and the value of β_1 is always greater than 0 due to the attractor behavior. Hence the $L1$ -norm of non-stochastic time series will always be greater than 1. In this work, we utilize the $L1$ -norm of the $E1$ Vs $E2$ plot of a given time series to classify it as stochastic or non-stochastic.

B. PCA Based approach

PCA decomposition is carried out to infer if the given time series possesses a dominant orientation or not. This is computed by hierarchically splitting the time series into two halves, and computing the covariance matrix of this split observations. The eigenvalues of this 2×2 covariance matrix will show one of the signatures: If the data indeed show any dominant direction (as in non-stochastic time series), then the larger eigenvalue will be significantly greater than the other. This will lead to a large eigen value ratio. On the other hand, if the data do not show any dominant direction (as in stochastic time series), then the two eigenvalues of the covariance matrix will be comparable. This will lead to small eigen value ratio. This observation is utilized in devising features for stochastic Vs non-stochastic classification. The steps are outlined as below.

For a time series consisting of n values $z_1, z_2 \dots z_n$.

- Split the series into two halves $(z_1, z_2 \dots z_{\lfloor \frac{n}{2} \rfloor})$ and $(z_{\lfloor \frac{n}{2} \rfloor + 1}, \dots z_n)$.
- Compute covariance matrix, C , by treating the samples in two halves as $\lfloor \frac{n}{2} \rfloor$ observations of 2-D vectors.
- Find eigenvalues of C , λ_1 and λ_2 ; the eigenvalue ratio is computed as λ_1/λ_2 where $\lambda_1 > \lambda_2$ (eigenvalues of a covariance matrix are real).

If the eigenvalue ratio for an interval is greater than a value of threshold, τ (computing optimal value of τ is described in subsection III-B1 below), the interval is further split into two sub-intervals of equal size. Subsequently, the eigenvalue ratio for each sub-interval is computed. The process is repeated as long as the length of the sub-interval is greater than a predefined number of samples (here taken as 100). For a fixed value of τ , the following features are derived

- **Variance of Eigenvalue Ratio (VER)**: This is the variance of the eigenvalue ratios of covariance matrices across sub-intervals in the entire time series.
- **Area Under the Eigenvalue Ratio curve (AUER)**: This measure captures the area under the curve of the eigenvalue ratio for the entire time series.

TODO should we show an eigen vale ratio curve to show the features.

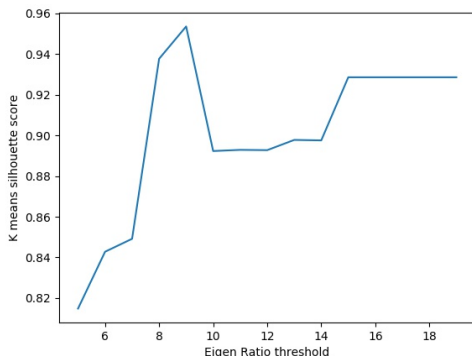


Fig. 1. Plot of Silhouette score vs eigenratio threshold

1) *Computing optimal value of threshold τ* : In order to arrive at the optimal value of τ , we observe the plot of the Silhouette score of K-Means clustering, with $K = 2$ (stochastic and non-stochastic), performed using the above feature set, as a function of the threshold value. The value of the threshold that results in the best Silhouette clustering score is taken as τ . This process is illustrated in the Silhouette score plot shown in Fig1. For the time series considered here for illustration, it is evident from the plot that the best clustering is obtained at threshold value 9, resulting in the maximum value of Silhouette score. Hence we use the corresponding $\tau = 9$ to arrive at the optimal hierarchical splitting and subsequent computing of the devised PCA-based features, VER and AUER.

C. Proof of Concept on Synthetic Data

The proposed approaches have been applied to standard synthetic signals. For stochastic class of signals, white noise and pink noise are considered; for non-stochastic class of signals, Lorenz system and Logistic map (for growth rate = 4) are considered.

SVD-Decomposition based technique : The SVD decomposition of the data is computed, followed by the plot of the top two right singular vectors. This plot is utilized to determine the Betti descriptors. The plot in Fig. 2 corresponds to the E1 vs E2 plot for a realization of white noise, which is known to be stochastic. The plot shows one single blob implying Betti descriptor of (1 0), which has L1-norm of 1. As discussed before, L1-norm of 1 implies that the timeseries is labelled as stochastic. On the other hand, the plot in Fig. 3 corresponds to the E1 vs E2 plot of a realization of Lorentz system, which is known to be non-stochastic. The plot shows two distinct voids, implying Betti descriptor to be (1 2), which has L1-norm of 3, implying that the timeseries is labelled as non-stochastic. This inference mechanism has been utilized on real data described in section IV. TODO Table with series name, betti descriptor, L1 norm, Label PCA-Decomposition based technique :PCA-

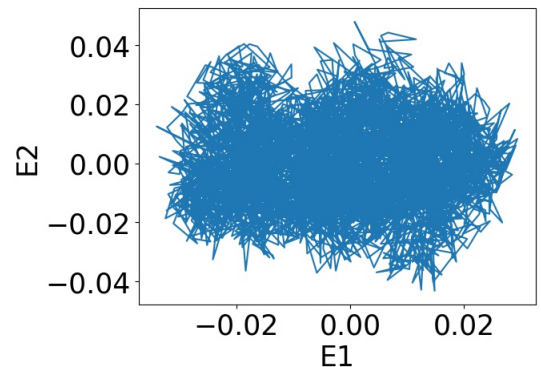


Fig. 2. Plot of Top-2 Right singular vectors for a stochastic timeseries

based features (i) VER and (ii) AUER are computed for the considered synthetic signals which are 24 different realizations of white noise and logistic map (growth rate = 4). The scatter plot of these features for these realizations is shown Fig 4 below. The plot makes it evident that in this feature space the two classes of timeseries, stochastic and non-stochastic,

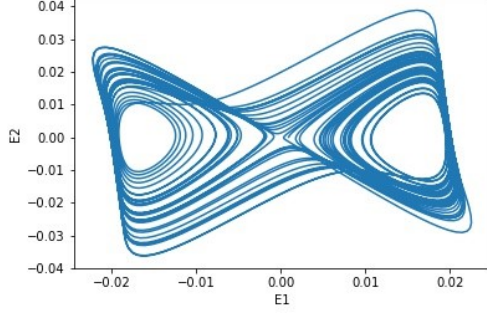


Fig. 3. Plot of Top-2 Right singular vectors for a non-stochastic timeseries

become linearly separable. Hence a linear SVM classifier is utilized. For training the SVM, computed features from white noise and Logistic map are utilized. For validating the trained SVM, 12 realizations of pink noise and Lorentz system are used. The classification on all 12 realizations of pink noise, yields the label stochastic, and the classification label for Lorentz system is obtained as non-stochastic, leading to perfect validation accuracy. This trained SVM is used to classify real data as described in section IV

TODO Table with series name, PCA features, Label

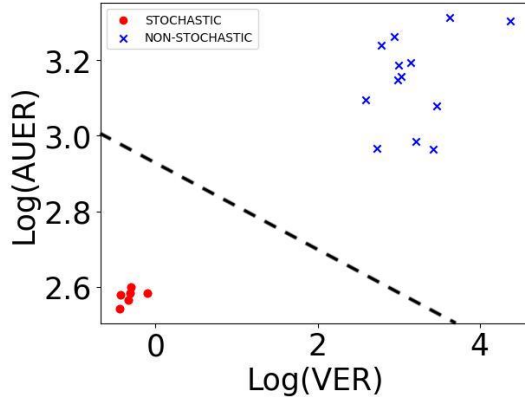


Fig. 4. Scatter plot of PCA-based features for synthetic data

IV. RESULTS AND DISCUSSION

In this section, we present the real data used, results obtained using proposed approaches and comparison with results in literature.

A. Real Data

The proposed approaches are illustrated on the publicly available data of *GRS 1915 + 105* taken from website [?]. 12 distinct categories of time series are utilized from the available data. All these time series are re-sampled with a sampling interval of 0.1 second. These datasets were also used in the work reported in [?], where the authors use CI based approaches, leading us to be able to compare our obtained results with theirs. TODO..consistently ..categories of time series in place of classes of time series TODO..Topological

analysis ...Betti descriptors are used interchangeably..we need to check

B. Results of SVD based analysis

From SVD decomposition of the data matrix, we plot the top 2 right singular vectors (E1 vs E2) to understand the temporal dynamics for each time series. Figure ?? shows representative E1 vs E2 plots for time series that are classified as non-stochastic and Figure ?? shows the corresponding plots for time series that are classified as stochastic. The Betti number descriptors for each of the E1 vs E2 plots are tabulated in Table II under the column *Betti descriptor*. In order to infer the label of the time series from the Betti descriptors, we use the L1-norm of β , $\|\beta\|_1$. If $\|\beta\|_1 > 1$, the time series is classified as non-stochastic, else the time series is stochastic.

TODO :Table of Betti number - Time series, Betti Descriptor, l1 norm , SVD label

C. Results of PCA based analysis

Figures ?? and ?? show the eigenvalue ratio plots for stochastic time series shown in Figure ?? and non-stochastic time series shown in Figure ?? respectively. We compute PCA-derived features, VER and AUER. These features are input to the SVM classifier to obtain the class labels.

- VER: For a stochastic signal since the variation in the eigenvalue ratios is typically small, the computed variance across the values is small. On the other hand, for a non-stochastic signal, since the eigenvalue ratios occupy diverse values, VER is typically high.
- AUER: For a stochastic time series, since the eigenvalue ratios are small across the entire span, the value of AUER is also small. However, for a non-stochastic signal, the eigenvalue ratios remain high for longer time intervals. Hence the value of AUER is significantly higher.

D. Consolidated Results

1) *Comparison of Results*: Table II tabulates the computed features and the respective inferences using the proposed approaches. Comparison of our results with CI based approach [?] is also presented. The columns of the table are described below.

- 1) Column 1 (*Class*) gives the class of the time series [?].
- 2) Column 2 (*diskbb*) and column 3 (*PL*) give quantities *diskbb* and *PL*, respectively, which indicate the spectral states of the black hole [?].
- 3) Column 4 (*CI Inference*) gives the inference about the state of the time series using CI approach [?].
- 4) Column 5 SVD based inference.
- 5) Our inference using these PCA features is given in column XX.
- 6) Finally the last column gives if there is a match between all three inferences.

TABLE I

TIMESERIES: COMPARISON BETWEEN CI BASED LABEL AND INFERENCE USING PROPOSED APPROACHES. THE MISMATCHED TIME SERIES CLASS, δ , IS SHOWN IN BOLD. (HERE NS STANDS FOR NON-STOCHASTIC AND S STANDS FOR STOCHASTIC)

| Class | CI Label | Betti Norm | SVD Label | VER | AUER | PCA Label | Match |
|-----------|----------|------------|-----------|------|------|-----------|-------|
| β | NS | 4 | NS | 483 | 43 | NS | Yes |
| θ | NS | 5 | NS | 778 | 58 | NS | Yes |
| λ | NS | 4 | NS | 6782 | 314 | NS | Yes |
| κ | NS | 4 | NS | 5199 | 144 | NS | Yes |
| μ | NS | 2 | NS | 51 | 12 | NS | Yes |
| ν | NS | 7 | NS | 32 | 16 | NS | Yes |
| α | NS | 6 | NS | 1.9 | 27.7 | NS | Yes |
| ρ | NS | 2 | NS | 147 | 35 | NS | Yes |
| δ | S | 1 | S | 9.7 | 26.2 | NS | NO |
| ϕ | S | 1 | S | 0.5 | 15 | S | YES |
| γ | S | 1 | S | 1 | 16 | S | YES |
| χ | S | 1 | S | 0.25 | 6.05 | S | YES |

TABLE II

BLACKHOLE STATE INFERENCE COMPARISON ACROSS CI-BASED AND PROPOSED APPROACHES:

| Name | diskbb | PL | State by CI | State by SVD | State by PCA | Match |
|-----------|--------|----|-------------|--------------|--------------|-------|
| β | 46 | 52 | ADAF | ADAF | ADAF | Yes |
| θ | 11 | 88 | ADAF | ADAF | ADAF | Yes |
| λ | 54 | 46 | Keplerian | Keplerian | Keplerian | Yes |
| κ | 59 | 51 | Keplerian | Keplerian | Keplerian | Yes |
| μ | 56 | 41 | Keplerian | Keplerian | Keplerian | Yes |
| ν | 28 | 72 | ADAF | ADAF | ADAF | Yes |
| α | 23 | 77 | ADAF | ADAF | ADAF | Yes |
| ρ | 28 | 72 | ADAF | ADAF | ADAF | Yes |
| δ | 48 | 50 | ADAF | ADAF | GAAF | NO |
| ϕ | 50 | 34 | slimdisc | slimdisc | slimdisc | YES |
| γ | 60 | 31 | slimdisc | slimdisc | slimdisc | YES |
| χ | 09 | 89 | GAAF | GAAF | GAAF | YES |

E. Identification of black hole states

We observe that SVD based analysis results in classification are consistent with CI based results for all the 12 categories of time series. However, with the PCA based approach the inference for δ time series is not consistent with the other two approaches. We observe that the PCA based features, VER and AUER, make the space of timeseries, linearly separable, as shown in Figure ???. This feature space is utilized to train a linear SVM. According to the CI based analysis δ turns out to be in between states slim disc and GAAF [?]. However, the present analysis shows that δ falls in between ADAF and Keplerian disc.

approach with those obtained using the proposed matrix based methods. Of the 12 categories of time series analysed, a mismatch is observed in the PCA based inference of only one class, while all other classes concur.

V. CONCLUSION

Exploring different techniques in order to have a conclusive inference for black hole systems turns out to be indispensable. We explore two different classical matrix based techniques to identify states of *GRS 1915+105* black hole using the time series obtained from *RXTE* satellite data. Based on our analysis, we are able to identify two extreme temporal dynamical classes of accretion around black holes. In the first approach we extend SVD decomposition to understand temporal dynamics, by adding topological descriptors, to classify time series as stochastic vs non-stochastic. In yet another approach, a novel application of PCA to characterize the time series is proposed. We compare inferences of the CI based