

Symbolic Analysis-based Reduced Order Markov Modeling of Time Series Data

Devesh K. Jha^{a,1}, Nurali Virani^{a,2}, Jan Reimann^b, Abhishek Srivastav^c, Asok Ray^{a,b}

Keywords: Symbolic Analysis, Markov Modeling, Order reduction, Combustion Instability

Abstract

This paper presents a technique for reduced-order Markov modeling for compact representation of time-series data. In this work, symbolic dynamics-based tools have been used to infer an approximate generative Markov model. The time series data are first symbolized by partitioning the continuous measurement space of the signal and then, the discrete sequential data are modeled using symbolic dynamics. In the proposed approach, the size of temporal memory of the symbol sequence is estimated from spectral properties of the resulting stochastic matrix corresponding to a first-order Markov model of the symbol sequence. Then, hierarchical clustering is used to represent the states of the corresponding full-state Markov model to construct a reduced-order (or size) Markov model with a non-deterministic algebraic structure. Subsequently, the parameters of the reduced-order Markov model are identified from the original model by making use of a Bayesian inference rule. The final model is selected using information-theoretic criteria. The proposed concept is elucidated and validated on two different data sets as examples. The first example analyzes a set of pressure data from a swirl-stabilized combustor, where controlled protocols are used to induce flame instabilities. Variations in the complexity of the derived Markov model represent how the system operating condition changes from a stable to an unstable combustion regime. In the second example, the data set is taken from NASA's data repository for prognostics of bearings on rotating shafts. We show that, even with a very small state-space, the reduced-order models are able to achieve comparable performance and that the proposed approach provides flexibility in the selection of a final model for representation and learning.

1. MOTIVATION AND INTRODUCTION

Hidden Markov model (HMM) is a widely used statistical learning tool for modeling uncertain dynamical systems [5], where the associated temporal data are used to infer a Markov chain with unobserved states. In this setting, the learning task is to infer the states and the corresponding parameters of the Markov chain. In addition to HMM, several other nonlinear techniques have been proposed for Markov modeling of time-series data. Symbolic time-series analysis-based Markov modeling is a recently proposed technique [24] where the states of a Markov chain are represented as a collection of words (i.e., symbol blocks, also referred to as memory words) of different lengths, which can be identified from the time-series data on a discrete space with finite cardinality [7], [18], [24], [28]. The symbols are created from the continuously varying time-series data by projecting the data to a set with finite cardinality. A common ground among all these tools of Markov modeling as discrete sequences, is that the Markov chain is induced by probabilistic representation of a deterministic finite state automaton (DFSA), often called probabilistic finite state automata (PFSA) [31]. While the PFSA-based inference provides a consistent, deterministic graph structure for learning, the deterministic algebraic structure is generally not a very compact representation and may often lead to large number of states in the induced Markov model. To circumvent this problem attempts have been made to reduce the state-space by merging statistically similar states of the model [18]. The problem is, however, that as these models are constructed by partitioning of phase space of the dynamical system, merging states that are statistically similar leads to algebraic inconsistency. On the other hand, if the states are merged to preserve the algebraic consistency, it leads to statistical impurity in the final models (i.e., states which have different statistics could be merged together). Other approaches for state aggregation in Markov chains could be found in [9], [32], [35]. However, these papers do not consider inference of the Markov model from the data which may not be suitable for analysis of data-driven systems [8].

The state space for Markov models, created by using symbolic analysis, increases exponentially with increase in memory or order of the symbolic sequence. Estimating the right memory is critical for temporal modeling of patterns observed in the sequential data. However, some of the states may be statistically similar and thus merging them can reduce the size of state-space. This paper presents reduced-order Markov modeling of time-series data to capture temporal patterns, where we estimate the size of temporal memory of the symbolic data using the spectral properties of a PFSA whose states are words of length one [13], [29]. The constraint of deterministic algebraic structure is not imposed by the end objective, but due to

^a Devesh K. Jha, N. Virani and Asok Ray are with Mechanical & Nuclear Engineering Department, Pennsylvania State University, University Park, PA 16802, USA, and are partially supported by the U.S. Air Force Office of Scientific Research under Grant No. FA9550-15-1-0400; {dkj5042, nnv105, axr2}@psu.edu

^b Jan Reimann and Asok Ray are with Department of Mathematics, Pennsylvania State University, University Park, PA, 16802, USA; {jan.reimann, axr2}@psu.edu. Jan Reimann was partially supported by NSF Grant DMS-1201263

^c Abhishek Srivastav is with AI and Machine Learning Lab, GE Global Research Center, San Ramon, CA, USA; srivastav@ge.com

¹ Currently with Mitsubishi Electric Research Laboratories, Cambridge, MA 02139

² Currently with AI and Machine Learning Lab, GE Global Research Center, Niskayuna, NY

the choice of the data representation model. Thus we propose to merge the states and remove the constraint of deterministic algebraic properties associated with PFSA, where the states of the Markov chain are now collection of words from its alphabet of length estimated in the last step. This state aggregation induces a non-determinism in the finite state model. The parameters of the reduced-order Markov model are estimated by a Bayesian inference technique from the parameters associated with the higher-order Markov model. The final model for data representation is selected using information-theoretic criteria, and thus, we get a unique stopping point to terminate the state-merging procedure. We also present a bound on the distortion of the predictive capability of the models up on reduction in the size of the state-space. The final model obtained is a generative model for the data; however, some predictive capability is lost as we remove the deterministic algebraic structure of a DFSA.

The proposed technique of state merging is inspired by time-critical applications where it is imperative to arrive at a reliable decision quickly as the dynamics of the process being monitored is really fast. In such applications, there are strict constraints on accuracy as well as the time needed to come to a decision. In this paper, we illustrate the concepts using two different datasets. We discuss in detail the example of combustion instability which is a highly nonlinear and complex phenomena and results in severe structural degradation in jet turbine engines. Some good surveys on the current understanding of the mechanisms for the combustion instability phenomena could be found in [6], [11], [17], [21], [27]. Active combustion instability control (ACIC) with fuel modulation has proven to be an effective approach for reducing pressure oscillations in combustors [3], [4]. Based on the work available in literature, one can conclude that the performance of ACIC is primarily limited by the large delay in the feedback loop and the limited actuator bandwidth [3], [4]. Early detection of combustion instability can potentially alleviate the problems with delay in the ACIC feedback loop and thus possibly improve the performance. Some recent work for detection and prediction of combustion instabilities could be found in [12], [19], [20], [25], [33]. While the results in these papers are encouraging, there is no interpretation of the expected changes in the data-driven model that could be observed during changes in the operating regime of the underlying process. In contrast to the work reported in literature, we have presented an overall idea of changes in the underlying stochastic model structure and parameters during the complex instability phenomenon.

Contributions. This paper presents a technique for Markov modeling of time series data using a PFSA with nondeterministic algebraic structure. Nondeterminism is induced by merging states of a PFSA with deterministic algebraic structure inferred from discrete sequential data, which in turn allows very compact representation of temporal data. In contrast to the approach in [18], we present a method to use information-theoretic criteria to arrive at a consistent stopping criterion for model selection. The resulting reduced-order model has fewer parameters to estimate; this in turn leads to faster convergence rates and thus faster decisions during test (or operation). We also present a bound on the distortion in the predictive capability of the models due to state-space reduction using Hamming distance between the sequences generated by the original and final model. The algorithms presented in the paper are validated on two different datasets— pressure data obtained from a swirl-stabilized combustor to monitor thermo-acoustic instability and a public data set for bearing prognostics. We show changes in the complexity of the pressure data as the process moves from stable to unstable through the transient phase which is then used to arrive at a criterion that provides perfect class separability. Apart from the results on Markov modeling, the results on combustion instability could be of independent interest in combustion community.

2. BACKGROUND AND MATHEMATICAL PRELIMINARIES

Symbolic analysis of time-series data is a recent approach where continuous sensor data are converted to symbol sequences via partitioning of the continuous domain [14], [24]. The dynamics of the symbols sequences are then modeled as a probabilistic finite state automaton (PFSA), which is defined as follows:

Definition 2.1 (PFSA): A probabilistic finite state automaton (PFSA) is a tuple $G = (\mathcal{Q}, \mathcal{A}, \delta, \mathbf{M})$ where

- \mathcal{Q} is a finite set of states of the automata;
- \mathcal{A} is a finite alphabet set of symbols $a \in \mathcal{A}$;
- $\delta : \mathcal{Q} \times \mathcal{A} \rightarrow \mathcal{Q}$ is the state transition function;
- $\mathbf{M} : \mathcal{Q} \times \mathcal{A} \rightarrow [0, 1]$ is the $|\mathcal{Q}| \times |\mathcal{A}|$ emission matrix. The matrix $\mathbf{M} = [m_{ij}]$ is row stochastic such that m_{ij} is the probability of generating symbol a_j from state q_i .

Remark 2.1: The PFSA defined above has a deterministic algebraic structure which is governed by the transition function δ ; thus a symbol emission from a particular state will lead to a fixed state. However, the symbol emissions are probabilistic (represented by the emission matrix). On the other hand, the transition function for a non-deterministic finite state automaton is given by a map, $\delta : \mathcal{Q} \times \mathcal{A} \rightarrow 2^{\mathcal{Q}}$ where, $2^{\mathcal{Q}}$ denotes the power set of \mathcal{Q} and includes all subsets of \mathcal{Q} . The idea is also presented in Figure 1 where we show that the same symbol can lead to multiple states, however in a probabilistic fashion. This allows more flexibility in modeling at the expense of some predictive accuracy.

For symbolic analysis of time-series data, a class of PFSA's called the D -Markov machine have been proposed [24] as a sub-optimal but computationally efficient approach to encode the dynamics of symbol sequences as a finite state machine.

Definition 2.2: (D -Markov Machine [18], [24]) A D -Markov machine is a statistically stationary stochastic process $S = \dots a_{-1}a_0a_1 \dots$ (modeled by a PFSA in which each state is represented by a finite history of D symbols), where the probability of occurrence of a new symbol depends only on the last D symbols, i.e.,

$$\Pr(s_n \mid \dots s_{n-D} \dots s_{n-1}) = \Pr(s_n \mid s_{n-D} \dots s_{n-1})$$

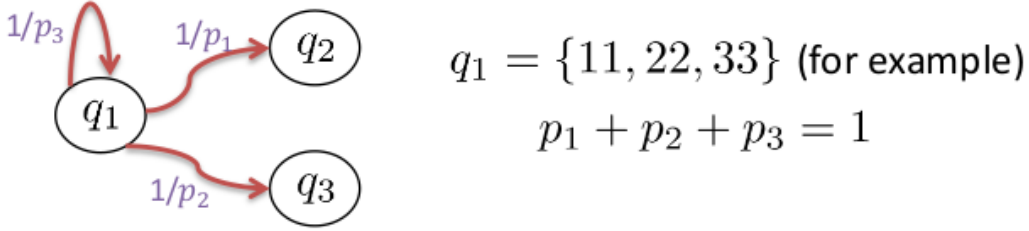


Fig. 1: Graphical model showing non-determinism in a PFSA. The symbol 1 emitted from state q_1 leads to different states with fixed probabilities indicating non-deterministic behavior.

where D is called the depth of the Markov machine.

A D -Markov machine is thus a D^{th} -order Markov approximation of the discrete symbolic process. For most stable and controlled engineering systems that tend to forget their initial conditions, a finite length memory assumption is reasonable. The D -Markov machine is represented as a PFSA and states of this PFSA are words over alphabet \mathcal{A} of length D (or less); the state transitions are described by a sliding block code of memory D and anticipation length of one [15].

For systems with fading memory it is expected that the predictive influence of a symbol progressively diminishes. In this context, depth is defined as follows.

Definition 2.3 (Depth): Let $\vec{s} = s_1 \dots s_k s_{k+1} s_{k+2} \dots$ be the observed symbol sequence where each $s_j \in \mathcal{A} \forall j \in \mathbb{N}$. Then, the depth of the process generating \vec{s} is defined as the length D such that:

$$\Pr(s_k | s_{k-1}, \dots, s_1) = \Pr(s_k | s_{k-1}, \dots, s_{k-D}) \quad (1)$$

An accurate estimation of depth for the symbolic dynamical process is required for the precise modeling of the underlying dynamics of the discrete sequence. Next we introduce an information-theoretic metric which is used for merging the states of the Markov model later in next section.

Definition 2.4 (Kullback-Leibler Divergence): [10] The Kullback-Leibler (K-L) divergence of a discrete probability distribution P from another distribution \tilde{P} is defined as follows.

$$D_{\text{KL}}(P \| \tilde{P}) = \sum_{x \in X} p_X(x) \log \left(\frac{p_X(x)}{\tilde{p}_X(x)} \right)$$

It is noted that K-L divergence is not a proper distance as it is not symmetric. However, to treat it as a distance it is generally converted into symmetric divergence as follows, $d(P, \tilde{P}) = D_{\text{KL}}(P \| \tilde{P}) + D_{\text{KL}}(\tilde{P} \| P)$. This is defined as the K-L distance between the distributions P and \tilde{P} .

This distance is used to find out the structure in the set of the states of the PFSA-based Markov model whose states are words, over the alphabet of the PFSA, of length equal to the depth estimated for the discretized sequence.

3. TECHNICAL APPROACH

In this section, we present the details of the proposed approach for inferring a Markov model from the time series data. As discussed earlier, the first step is the discretization of the time-series data to generate a discrete symbol sequence. While it is possible to optimize the symbolization of time-series using some optimization criterion, we do not discuss such a technique here. The data is discretized using the unbiased principle of entropy maximization of the discrete sequence using maximum entropy partitioning (MEP) [23]. The proposed approach for Markov modeling then consists of the following four critical steps

- Estimate the approximate size of temporal memory (or order) of the symbol sequence.
- Cluster the states of the high-order Markov model.
- Estimate the parameters of the reduced-order Markov model (i.e., the transition matrix).
- Select the final model using information theoretic scores (described below, Section 3-C).

Memory of the discrete sequence is estimated using a recently introduced method based on the spectral analysis of the Markov model with depth 1, induced by a PFSA [13], [29]. It is noted that these steps are followed during training to estimate the approximate model for data and during test, the parameters are estimated for the reduced-order model. The key ideas behind these steps are explained in the next section.

A. Estimation of Reduced-Order Markov Model

Depth D of a symbol sequence has been redefined in [29] as the number of time steps after which probability of current symbol is independent of any past symbol i.e.:

$$\Pr(s_k | s_{k-n}) = \Pr(s_k) \quad \forall n > D \quad (2)$$

Note that dependence in the proposed definition (eq. 2) is evaluated on individual past symbols using $\Pr(s_k|s_{k-n})$ as opposed to the assessing dependence on words of length D using $\Pr(s_k|s_{k-1}, \dots, s_{k-D})$. It is shown that if the observed process is *forward causal* then observing any additional intermediate symbols $s_{k-1}, \dots, s_{k-n+1}$ cannot induce a dependence between s_k and s_{k-n} if it did not exist on individual level [29].

Let $\mathbf{\Pi} = [\pi_{ij}^{(1)}]$ be the one-step transition probability matrix of the PFSA G constructed from this symbol sequence i.e.

$$\mathbf{\Pi} = \Pr(s_k|s_{k-1}) \quad (3)$$

Then using the distance of the transition matrix after steps from the stationary point, depth can be defined as a length D such that

$$|\text{trace}(\mathbf{\Pi}^n) - \text{trace}(\mathbf{\Pi}^\infty)| \leq \sum_{j=2}^J |\lambda_j|^n < \epsilon \quad \forall n > D \quad (4)$$

where J is number of non-zero eigenvalues of $\mathbf{\Pi}$. Thus, the depth D of the symbol sequence is estimated for a choice of ϵ by estimating the stochastic matrix for the one-step PFSA. Next, another pass of data is done to estimate the PFSA parameters whose states are words over \mathcal{A} of length D , i.e., $\mathbf{\Pi} = \Pr(s_k|s_{k-1}, \dots, s_{k-D})$. It is noted that this step is critical for modeling accuracy.

The states of the reduced-order Markov model are then estimated by partitioning the set of words over \mathcal{A} of length D estimated in the last step. This is done by using an agglomerative hierarchical clustering approach. The advantage of using the hierarchical clustering approach is that it helps visualize the structure of the set of the original states using an appropriate metric. Agglomerative hierarchical clustering is a bottom-up clustering approach [34] that generates a sparse network (e.g., a binary tree) of the state set \mathcal{Q} (where $|\mathcal{Q}| = |\mathcal{A}|^D$) by successive addition of edges between the elements of \mathcal{Q} . Initially, each of the states q_1, q_2, \dots, q_n is in its own cluster C_1, C_2, \dots, C_n where $C_i \in \mathcal{C}$, which is the set of all clusters for the hierarchical cluster tree. The distance between any two states in \mathcal{Q} is measured using the K-L distance between the symbol emission probabilities conditioned on them, i.e.,

$$d(q_i, q_j) = D_{\text{KL}}(\Pr(\mathcal{A}|q_i) \parallel \Pr(\mathcal{A}|q_j)) + D_{\text{KL}}(\Pr(\mathcal{A}|q_j) \parallel \Pr(\mathcal{A}|q_i)) \quad (5)$$

where the terms on the right have the following meaning.

$$\begin{aligned} & D_{\text{KL}}(\Pr(\mathcal{A}|q_i) \parallel \Pr(\mathcal{A}|q_j)) \\ &= \sum_{s \in \mathcal{A}} \Pr(s|q_i) \log \left(\frac{\Pr(s|q_i)}{\Pr(s|q_j)} \right) \end{aligned}$$

In terms of the distance measured by eq. (5), the pair of clusters that are nearest to each other are merged and this step is repeated till only one cluster is left. The tree structure displays the order of splits in the state set of the higher-order Markov model and is used to aggregate the states close to each other. For clarification of presentation, we show an example of a Markov chain with 27 states and 3 symbols on a simplex plane in Figure 2, where each **red pentagon** on the simplex represents one row of the symbol emission matrix. The hierarchical clustering is used to find the structure of the state set on the simplex plane using the K-L distance. The set of states clustered together could be obtained based on the number of final states required in the final Markov model.

The overall algorithm is presented as a pseudo-code in Algorithm 1. This algorithm is used to find the parameters of the models during training. The parameters during test are estimated using the clustering map $f_{N_{\text{max}}}$ and is further discussed in next section. In the later sections we show how an information theoretic criterion could be used to select the appropriate model to terminate the state merging algorithm or select a final model from the set of reduced-order models. Through numerical experiments using two different data-sets we also illustrate the main motivation of this work that although the right memory is required for accurate modeling of the symbolic process, the state-space not necessarily consist of all words corresponding to the estimated memory and we can achieve sufficiently-high predictive accuracy even with a smaller state-space. We are able to achieve this trade-off between the model complexity and predictive modeling accuracy using the information-theoretic criteria.

B. Parameter Estimation of the Reduced-Order Markov Model

The parameters of the Markov model obtained after clustering the states of the original PFSA with $|\mathcal{A}|^D$ states is obtained using a Bayesian inference technique using the parameters estimated for the PFSA. In this proposed approach, the state transition matrix $\mathbf{\Pi}$, the emission matrix \mathbf{M} , and the state probability vector \mathbf{p} of the original PFSA model G are available, along with the deterministic assignment map $f: \mathcal{Q} \rightarrow \tilde{\mathcal{Q}}$ of the state in \mathcal{Q} (i.e., state set of original model) to one of the state in $\tilde{\mathcal{Q}}$ (i.e., state set of the reduced order model). Since the reduced order model can be represented by the tuple $\tilde{G} = (\tilde{\mathcal{Q}}, \tilde{\mathbf{\Pi}})$, where $\tilde{\mathbf{\Pi}} = [\tilde{\pi}_{ij}]$ is the state transition matrix, we employ a Bayesian inference technique to infer the individual values of transition probabilities $\tilde{\pi}_{ij} = \Pr(\tilde{q}_{k+1} = j \mid \tilde{q}_k = i)$ for all $i, j \in \tilde{\mathcal{Q}}$.

Let Q_k be the random variable denoting the state of PFSA model at some time step $k \in \mathbb{N}$ and S_k denotes the symbol emitted from that state, this probabilistic emission process is governed by the emission matrix \mathbf{M} . The state of the reduced

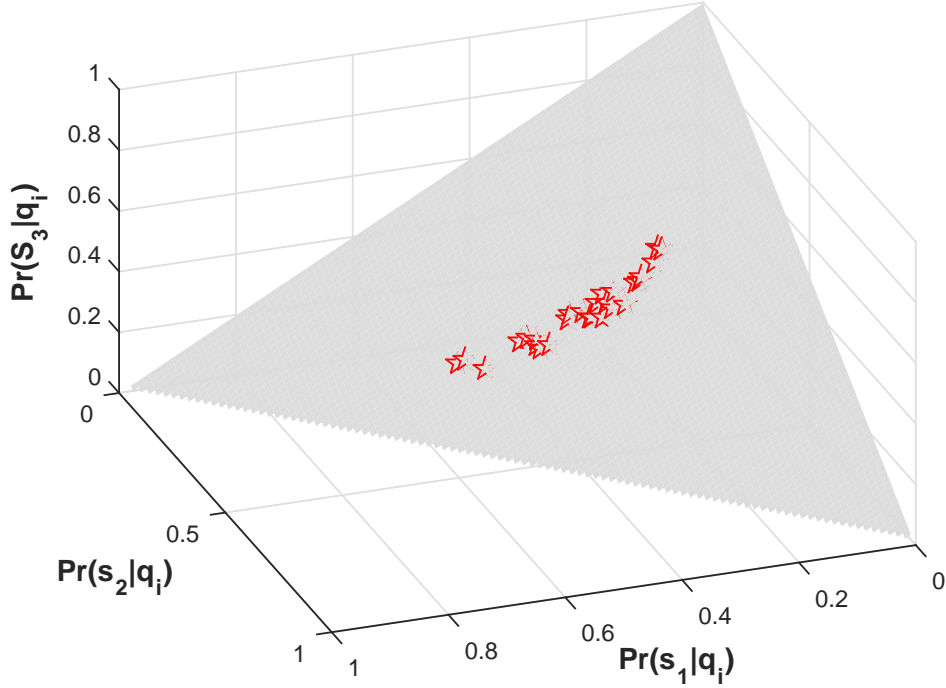


Fig. 2: The symbol emission probabilities for a Markov chain with 3 symbols are shown on a simplex. Symmetric K-L distance is used to find the structure in the state-set in the information space and the states are clustered based on the revealed structure.

Algorithm 1: Reduced Order Markov Modeling

- Input:** The observed symbol sequence $\vec{s} = \{ \dots s_1 s_2 s_3 \dots | s_i \in \mathcal{A} \}$
Output: The final Markov model, $\mathcal{M} = (\tilde{\mathcal{Q}}, \tilde{M}, \tilde{\Pi})$
- 1 Estimate the Π matrix for 1-step Markov model using frequency counting with an uniform prior;
 - 2 Estimate the size of temporal memory, $D(\epsilon)$ for \vec{s} using equation (4);
 - 3 Estimate M and Π for the $D(\epsilon)$ -Markov model using frequency counting with an uniform prior;
 - 4 $\mathcal{C}_{|\mathcal{Q}|} = \{ q_i | q_i \in \mathcal{Q} \}$;
 - 5 **for** $i = |\mathcal{Q}| - 1, \dots, 1$ **do**
 - 6 find distinct clusters $A, B \in \mathcal{C}_{i+1}$ minimizing $d(A \cup B)$;
 - 7 $\mathcal{C}_i := (\mathcal{C}_{i+1} \setminus \{A, B\}) \cup \{A \cup B\}$
 - 8 **return** $\mathcal{C}_1, \dots, \mathcal{C}_{|\mathcal{Q}|}$ and $f_i : \mathcal{Q} \rightarrow \mathcal{C}_i \forall i \in \{1, \dots, |\mathcal{Q}|\}$
 - 9 Calculate the parameters of reduced model using $\tilde{\mathcal{Q}} = \mathcal{C}_{N_{\max}}, f_{N_{\max}}$ and equations (7) through (13);
 - 10 Calculate the Log-likelihood for models with Equation (16);
 - 11 The final model is selected using the AIC or BIC criteria explained in Section 3-C;

order model is obtained from a deterministic mapping of the state of the PFSA model, thus the state of this model is also a random variable, which is denoted by $\tilde{Q}_k = f(Q_k)$. The Bayesian network representing the dependencies between these variables is shown in the recursive as well as unrolled form in the Figure 3. The conditional density $\Pr(\tilde{Q}_k = \tilde{q} | Q_k = q)$ can be evaluated by checking if state q belongs to the state cluster \tilde{q} and assigning the value of 1 if true, else assign it the value of 0. Since we know that $\tilde{\mathcal{Q}}$ partitions the set \mathcal{Q} , the conditional density is well-defined. Thus, it can be written as

$$\Pr(\tilde{Q}_k = \tilde{q} | Q_k = q) = I_{\tilde{q}}(q), \tag{6}$$

where I is the indicator function with $I_{\tilde{q}}(q) = 1$, if element q belongs to the set \tilde{q} , else it is 0. The derivation of the Markov

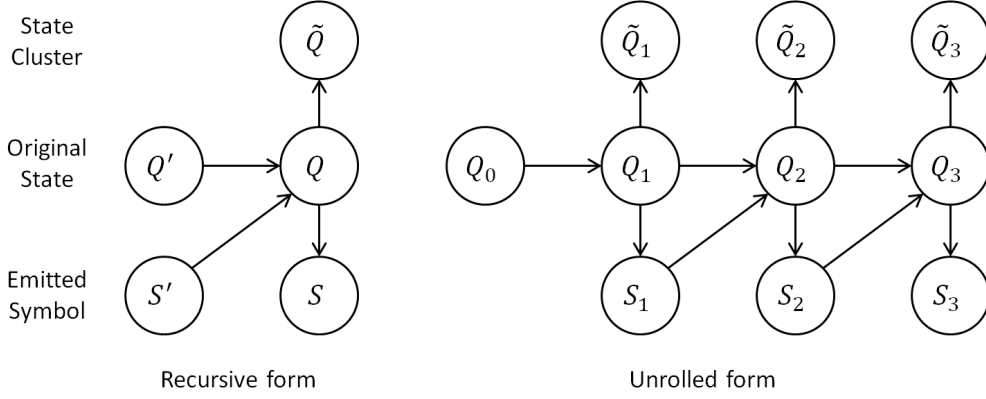


Fig. 3: Graphical models representing the dependencies between the random variables

model $\Pr(\tilde{Q}_{k+1} | \tilde{Q}_k)$ using $\Pr(Q_{k+1} | Q_k)$, stationary probability vector \mathbf{p} , and assignment map f is shown ahead.

$$\Pr(\tilde{Q}_{k+1} | \tilde{Q}_k) = \sum_{q \in \mathcal{Q}} \Pr(\tilde{Q}_{k+1}, Q_{k+1} = q | \tilde{Q}_k) \quad (7)$$

(Marginalization)

$$= \sum_{q \in \mathcal{Q}} \Pr(Q_{k+1} = q | \tilde{Q}_k) \Pr(\tilde{Q}_{k+1} | Q_{k+1} = q) \quad (8)$$

(Factorization using Figure 3)

$$= \sum_{q \in \mathcal{Q}} \Pr(Q_{k+1} = q | \tilde{Q}_k) \mathbb{I}_{\tilde{Q}_{k+1}}(q) \quad (9)$$

(using (6))

$$= \sum_{q \in \tilde{\mathcal{Q}}_{k+1}} \Pr(Q_{k+1} = q | \tilde{Q}_k). \quad (10)$$

We can obtain $\Pr(Q_{k+1} | \tilde{Q}_k)$ from Bayes' rule as

$$\Pr(Q_{k+1} | \tilde{Q}_k) = \frac{\Pr(\tilde{Q}_k | Q_{k+1}) \Pr(Q_{k+1})}{\sum_{q \in \mathcal{Q}} \Pr(\tilde{Q}_k | Q_{k+1} = q) \Pr(Q_{k+1} = q)}. \quad (11)$$

Following the steps to obtain (10), we also derive

$$\Pr(\tilde{Q}_k | Q_{k+1}) = \sum_{q \in \tilde{\mathcal{Q}}_k} \Pr(Q_k = q | Q_{k+1}). \quad (12)$$

We can obtain $\Pr(Q_k | Q_{k+1})$ from Bayes' rule as

$$\Pr(Q_k | Q_{k+1}) = \frac{\Pr(Q_{k+1} | Q_k) \Pr(Q_k)}{\sum_{q \in \mathcal{Q}} \Pr(Q_{k+1} | Q_k = q) \Pr(Q_k = q)}. \quad (13)$$

Note that, for the distribution $\Pr(Q_k)$ and $\Pr(Q_{k+1})$, we use the stationary probability \mathbf{p} . Using the equations (10), (11), (12), and (13) together, one can easily obtain the desired state transition matrix $\tilde{\mathbf{\Pi}}$ of the reduced order model. Once the state cluster set $\tilde{\mathcal{Q}}$ and state transition matrix $\tilde{\mathbf{\Pi}}$ are available, the reduced order model is completely defined.

C. Model Selection using information theoretic criteria

In this section, we describe the model selection process during the underlying state merging process for model inference. We compute “penalized” likelihood estimates for different models. Then, the model with the lowest score is selection as the optimal model.

The (unpenalized) log-likelihood of a symbol sequence \vec{s} given a Markov model G is computed as follows:

$$\mathcal{L}(\vec{s}|G) \cong \sum_{k=1}^N \log \Pr(s_k | q_k) \quad (14)$$

where the effects of the initial state are ignored because they become negligible for long statistically stationary symbol sequences. It is noted that with a finite symbol sequence, the log-likelihood is always finite. Furthermore, with the Markov

models considered in this paper, the sum is simplified to the following form.

$$\mathcal{L}(\vec{s}|G) \cong \sum_{k=D+1}^N \log \Pr(s_k | s_{k-1}, \dots, s_{k-D}) \quad (15)$$

As discussed earlier, the states are merged using hierarchical clustering and thus, for every desired number of final states we get the deterministic map $f_{N_{\max}}$ which determines how the original states are partitioned using the hierarchical clustering. This map is known for every terminal number of states and thus, we can find the log-likelihood of the symbol sequence using the following relationship.

$$\mathcal{L}(\vec{s}|\tilde{G}) \cong \sum_{k=D+1}^N \log \Pr(s_k | \tilde{q}_k = f_{N_{\max}}(q_k)) \quad (16)$$

where, \tilde{q}_k is the state of the reduced model and q_k is the state of the original full-order model.

In the next step of the model selection process, a ‘‘complexity penalty’’ is added to the log-likelihood estimates, thereby balancing goodness of fit against the complexity of the model (and hence trying to prevent overfitting). We apply two widely-used such model selection functions, namely the Akaike information criterion (AIC) [2] and the Bayesian information criterion (BIC) [26]:

- 1) $\mathcal{M}_{\text{BIC}} = -2\mathcal{L}(\vec{s}|\tilde{G}) + K \log(N)$, where K is the number of free parameters and N is the number of observations.
- 2) $\mathcal{M}_{\text{AIC}} = -2\mathcal{L}(\vec{s}|\tilde{G}) + 2K$, where K is the number of free parameters.

The number of free parameters to be estimated from the data is the parameters of the symbol emission parameters, i.e., $K = |\mathcal{A} \times \tilde{\mathcal{Q}}|$. It is noted that this allows model selection for individual symbol sequences. The criterion here allows a terminal condition for state merging; however, different symbol sequences can have different models. The model with the minimum score is selected as the best model. Through the results presented in next sections we illustrate the fact that most of the temporal and predictive capabilities can be preserved for the models with a very small number of states when compared to the original model.

Remark 3.1: The final Markov model is a finite depth approximation of the original time-series data. However, compared to the PFSA-based D-Markov machines in [18], [24], the current aggregated model has a non-deterministic algebraic structure, i.e., the same symbol emissions from a state can lead to different states. While this leads to some loss in predictive capability as compared to the models in [18], [24], this allows us to compress the size of the model as per the requirement at hand. This allows faster convergence rates for the symbol emission probabilities as we only require fewer parameters to estimate from data, which might lead to faster decisions during testing.

In the rest of the paper, we will present a Hamming distance-based bound for distortion in the predictive capabilities of reduced models and demonstrate the utility of these models in practical problems of fault/anomaly detection from time-series data.

4. ANALYSIS OF THE PROPOSED ALGORITHM

In this section, we will present a bound on the distortion of the model due to the reduction of state-space of the Markov model using Hamming distance between two symbol sequences. We first present the Pinsker’s inequality [10] which relates the information divergence with the variational distance between probability measures defined on arbitrary spaces. This is followed by another theorem which can be used to derive Hamming distance bounds using the informational divergence.

Theorem 4.1 (Pinsker’s inequality): [10] Let P and Q be two probability distributions on a measurable space (\mathbb{X}, Σ) . Then, the following is true

$$d_{TV}(P, Q) \leq \sqrt{\frac{1}{2} D_{\text{KL}}(P||Q)} \quad (17)$$

where $d_{TV}(P, Q) = \sup_{A \in \Sigma} \{|P(A) - Q(A)|\}$ is the total variation distance.

Theorem 4.2: [16] Let \mathbb{X} be a countable set and let us denote by x^n the sequence $(x_1, x_2, \dots, x_n) \in \mathbb{X}^n$. Let q^n be a Markov measure on \mathbb{X}^n , that is, $q(x^n) = q(x_1) \prod_{i=2}^n q_i(x_i | x_{i-1})$. Then for any probability measure p^n on \mathbb{X}^n , the following is true

$$\bar{d}(p^n, q^n) \leq \left[\frac{1}{2n} D_{\text{KL}}(p^n || q^n) \right]^{1/2} \quad (18)$$

where, \bar{d} denotes the normed Hamming distance on $\mathbb{X}^n \times \mathbb{X}^n$:

$$\bar{d}(x^n, y^n) = n^{-1} \sum_{i=1}^n d(x_i, y_i), \quad (19)$$

where $d(x_i, y_i) = 1$ if $x_i \neq y_i$ and 0 otherwise. The \bar{d} -distance between p^n and q^n is

$$\bar{d}(p^n, q^n) = \min E \bar{d}(\hat{X}^n, X^n), \quad (20)$$

where min is taken over all joint distributions with marginals $p^n = \text{dist} \hat{X}^n$ and $q^n = \text{dist} X^n$ and E denotes the expectation operator.

The above theorem provides us a way to bound Hamming distance between sequences generated by two different distributions. Thus, using the above theorem, we find a bound on the Hamming distance between the symbol sequences generated by the reduced-order Markov model and the original model by estimating the K-L distance between the measure on symbol sequences induced by these models. An approximate estimate of the K-L distance between the original and a reduced model could be expressed and estimated as shown in the following.

Let the original D-Markov model be denoted by \mathcal{M} and the reduced-order model by $\hat{\mathcal{M}}$. The Markov measure on the probability space (S^n, \mathcal{E}, P) where the set S^n consists of sequences of length n from an alphabet \mathcal{A} could be estimated using the symbol emission probabilities. More explicitly, the Markov measure of a sequence S_n on S^n induced by \mathcal{M} is given by $P_{\mathcal{M}}(S_n) = \Pr(q_1) \prod_{i=D+1}^n \Pr(s_i | q_i)$ (where D is the depth of the model). Then, the K-L divergence between \mathcal{M} and $\hat{\mathcal{M}}$ is given by the following expression.

$$D_{\text{KL}}(P_{\mathcal{M}}^n \| P_{\hat{\mathcal{M}}}^n) = \sum_{S_n \in S^n} P_{\mathcal{M}}(S_n) \log \left(\frac{P_{\mathcal{M}}(S_n)}{P_{\hat{\mathcal{M}}}(S_n)} \right) \quad (21)$$

Then, the above expression can be simplified as follows.

$$\log \left(\frac{P_{\mathcal{M}}(S_n)}{P_{\hat{\mathcal{M}}}(S_n)} \right) = \sum_{i=D+1}^n \log(\Pr(s_i | q_i)) - \log(\Pr(s_i | \hat{q}_i)),$$

where, \hat{q} is the merged state and q is the original state. Then the expression on the right could be further bounded using the Lipschitz constant for the logarithm function and under the assumption that $\log(\Pr(s_j | q_i)) \neq 0 \forall q_i \in \mathcal{Q}$ and all $s_j \in \mathcal{A}$.

$$\sum_{i=D+1}^n \log(\Pr(s_i | q_i)) - \log(\Pr(s_i | \hat{q}_i)) \quad (22)$$

$$\leq \sum_{i=D+1}^n \left(\frac{\Pr(s_i | q_i) - \Pr(s_i | \hat{q}_i)}{\Pr(s_i | q_i)} \right) \quad (23)$$

$$\leq (n - D - 1)\kappa \quad (24)$$

where, $\kappa = \max_{q \in \mathcal{Q}, s \in \mathcal{A}} \frac{\Pr(s|q) - \Pr(s|\hat{q})}{\Pr(s|q)}$. In the above inequalities, equation (23) is obtained from equation (22) by using the observation that $\Pr(s_i | \hat{q}_i) = \Pr(s_i | q_i) + \eta$, where η is the perturbation in the symbol emission probability from q_i when it is clustered into a new state \hat{q}_i . Hence, the K-L distance in equation (21) could be bounded by the following term.

$$\begin{aligned} D_{\text{KL}}(P_{\mathcal{M}}^n \| P_{\hat{\mathcal{M}}}^n) &\leq \sum_{S_n \in S^n} P_{\mathcal{M}}(S_n) (n - D - 1)\kappa \\ &= (n - D - 1)\kappa \sum_{S_n \in S^n} P_{\mathcal{M}}(S_n) \\ &= (n - D - 1)\kappa \end{aligned} \quad (25)$$

Thus, a uniform bound on the Hamming distance between the original and the final model could then be obtained as follows,

$$\bar{d}(P_{\mathcal{M}}(S_n), P_{\hat{\mathcal{M}}}(S_n)) \leq \sqrt{\frac{(n - D - 1)\kappa}{2n}} \quad (26)$$

The above inequality thus, allows us to compare models with different state-space based on the predictive accuracy of a reduced model when compared to the original model. As compared to the earlier information theoretic criteria, which were based on the efficiency of data compression by different models, the inequality in (26) allows to compare them based on their symbol emission statistics and thus, is computationally efficient. It is possible to find a rather tighter bound in an expected sense by using the stationary distribution of the two Markov chains to find an expected bound on Hamming distance. However, finding the same is left as an exercise for future work. Using the above bound for selection of models could be more efficient than the information theoretic metrics (as it can be estimated by using the symbol emission probabilities instead of the penalized likelihoods); however, finding a penalized version of the bound for model selection is also left as a future exercise.

5. DESCRIPTION OF EXPERIMENTATION AND DATA SETS

In this section, we briefly describe the two different data-sets which have been used in this paper to illustrate and validate the proposed concepts. Specifically, we will describe the experiments done at Penn State to investigate instability in lean-premixed combustion and another benchmark data-set for anomaly detection in bearings. An important point to be noted here is that the numerical experiments we present in the following sections is to justify the fact that the reduced-order models obtained by the proposed algorithms are able to achieve the trade-off between predictive accuracy and model complexity. Further results

for classification and anomaly detection are to illustrate that this proposed approach of model learning can still achieve good performance for machine learning objectives of class separability and anomaly detection.

A. Combustion

A swirl-stabilized, lean-premixed, laboratory-scale combustor was used to perform the experimental study. Tests were conducted at a nominal combustor pressure of 1 atm over a range of operating conditions, as listed in Table I.

TABLE I: Operating conditions

Parameters	Value
Equivalence Ratio	0.525, 0.55, 0.60, 0.65
Inlet Velocity	25-50 m/s i m/s increments
Combustor Length	25-59 inch in 1 inch increments

In each test, the combustion chamber dynamic pressure and the global OH and CH chemiluminescence intensity were measured to study the mechanisms of combustion instability. The measurements were made simultaneously at a sampling rate of 8192 Hz (per channel), and data were collected for 8 seconds, for a total of 65536 measurements (per channel). A total of 780 samples of data were collected from all the tests where in every test the combustion process was driven from stable to unstable by changing either the equivalence ratio, ϕ . However, as the accurate model of the process is not available, an accurate label of transition of the process to unstable phase is not available. It is noted that the data consists the behavior of the process over a large number of operating condition and thus provides a rich set of data to test the efficacy of the algorithm in detecting classes irrespective of the underlying operating conditions.

B. Bearing Prognostic Data

This test data has been picked from NASA’s prognostics data repository [1], [30]. A detailed description of the experiments could be found in [22]. The bearing test rig hosts four test bearings on one shaft which is driven by an AC motor at a constant speed. A constant force is applied on each of the bearings and accelerometer data is collected at every bearing at a sampling rate of 20 kHz for about 1 s. The tests are carried for 35 days until a significant amount of debris was found in the magnetic plug of the test bearing. A defect in at least one of the bearings is found at the end of every test. In this paper, we will use the data from a bearing which shows anomalous behavior in the later parts of test. In particular, out of the three data sets, we use set one where an inner race fault occurred on Bearing 3. In the analysis, we use data from Bearing 3.

6. MARKOV MODELING

In this section, we present results for modeling and analysis of the time-series data which are presented in this paper.

A. Combustion

Time-series data is first normalized by subtracting the mean and dividing by the standard deviation of its elements; this step corresponds to bias removal and variance normalization. Data from engineering systems is typically oversampled to ensure that the underlying dynamics can be captured (in the current experiments, it was 8192 Hz). Due to coarse-graining from the symbolization process, an over-sampled time-series may mask the true nature of the system dynamics in the symbolic domain (e.g., occurrence of self loops and irrelevant spurious transitions in the Markov chain). Time-series is first down-sampled to find the next crucial observation. The first minimum of auto-correlation function generated from the observed time-series is obtained to find the uncorrelated samples in time. The data sets are then down-sampled by this lag. The autocorrelation function for the time-series data during unstable case is shown in Figure 4 where the data are downsampled by the lag marked in **red** rectangle in Figure 4. To avoid discarding significant amount of data due to downsampling, down-sampled data using different initial conditions is concatenated. Further details of this preprocessing can be found in [29].

The continuous time-series data set is then partitioned using maximum entropy partitioning (MEP), where the information rich regions of the data set are partitioned finer and those with sparse information are partitioned coarser. In essence, each cell in the partitioned data set contains (approximately) equal number of data points under MEP. A ternary alphabet with $\mathcal{A} = \{0, 1, 2\}$ has been used to symbolize the continuous combustion instability data. As discussed in section 5, we analyze data sets from different phases, as the process goes from stable through the transient to the unstable region (the ground truth is decided using the RMS-values of pressure).

In Figure 5a, we show the observed changes in the behavior of the data as the combustion operating condition changes from stable to unstable. A change in the empirical distribution of data from unimodal to bi-modal is observed as the system moves from stable to unstable. We selected 150 samples of pressure data from the stable and unstable phases each to analyze and compare. First, we compare the expected size of temporal memory during the two stages of operation. There are changes in the eigenvalue decomposition rate for the 1-step stochastic matrix calculated from the data during the stable and unstable behavior, irrespective of the combustor length and inlet velocity. During stable conditions, the eigenvalues very quickly go to zero as

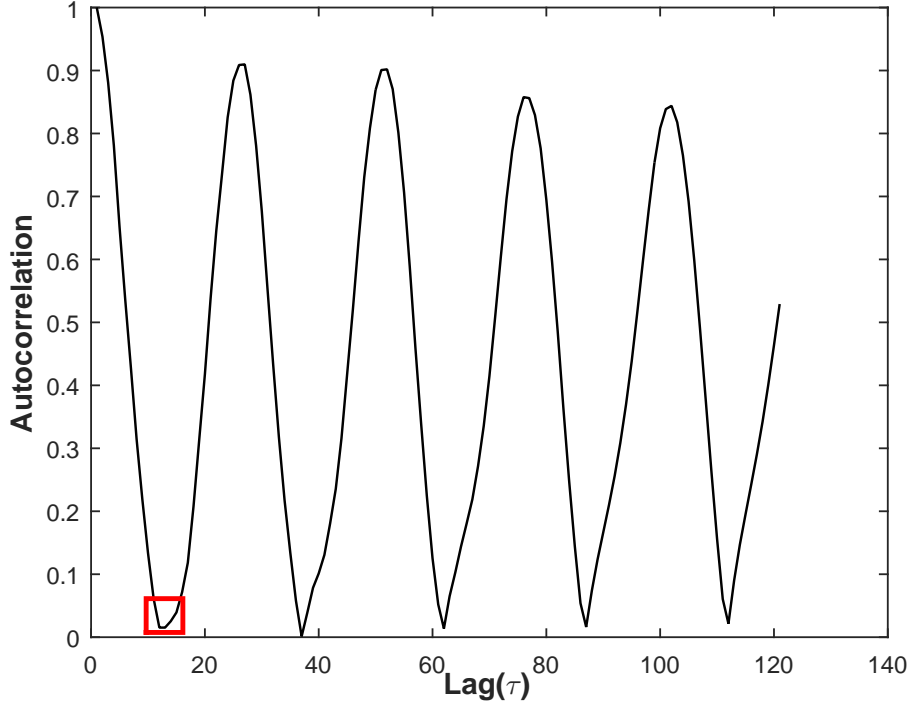
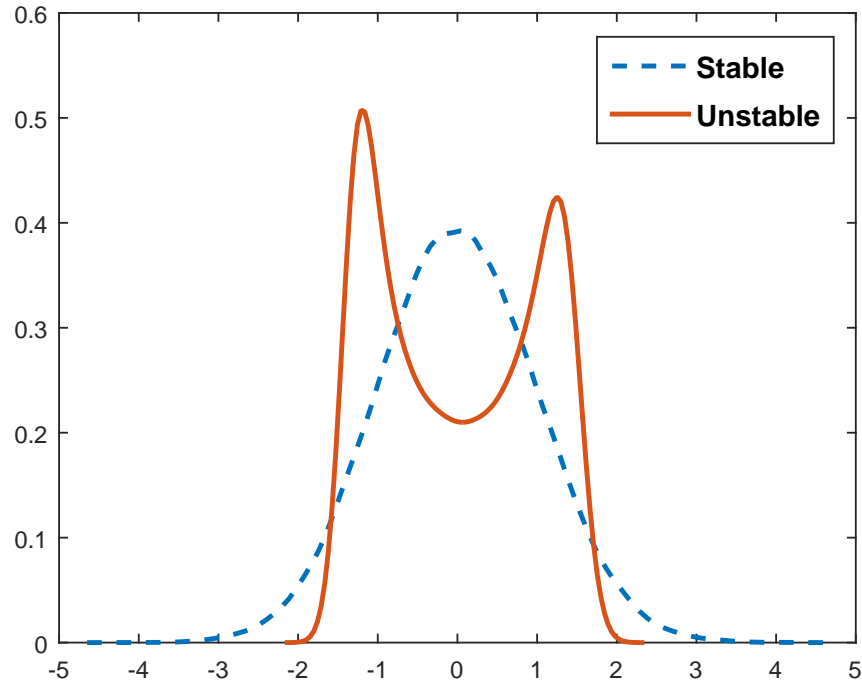


Fig. 4: Autocorrelation function of time-series data during the unstable phase of combustion. The time-series data is down-sampled by the lag marked in red square. It is noted that the individual time-series have their own down-sampling lags.

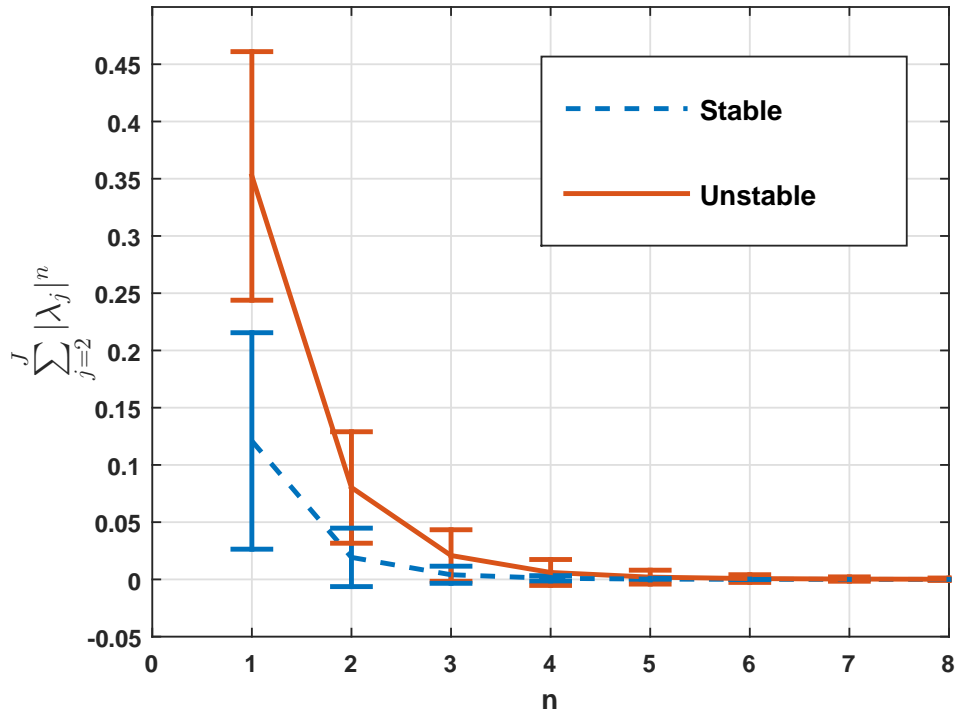
compared to the unstable operating condition (see Figure 5b). This suggests that the size of temporal memory of the discretized data increases as we move to the unstable operating condition. This indicates that under the stable operating condition, the discretized data behaves as symbolic noise as the predictive power of Markov models remain unaffected even if we increase the order of the Markov model. On the other hand, the predictive power of the Markov models can be increased by increasing the order of the Markov model during unstable operating condition, indicating more deterministic behavior. An $\epsilon = 0.05$ is chosen to estimate the depth of the Markov models for both the stable and unstable phases. Correspondingly, the depth was calculated as 2 and 3 for the stable and unstable conditions (see Figure 5). The corresponding $D(\epsilon)$ is used to construct the Markov models next. First a PFSA whose states are words over \mathcal{A} of length $D(\epsilon)$ is created and the corresponding maximum-likely parameters (\mathbf{M} and $\mathbf{\Pi}$) are estimated. Then, the hierarchical clustering algorithm using K-L distance is used to cluster and aggregate the states. It is noted that we create individual models for every sample, i.e., every sample is partitioned individually so that the symbols will have different meaning (i.e., they represent different regions in the measurement space of the signals) for every sample. Consequently, each sample will have a different state-space when viewed in the continuous domain. Thus, we do not show the mean behavior of the samples during any operating regime as the state-space would be inconsistent (even though the cardinality could be the same).

In Figure 6, we show the hierarchical cluster tree which details the structure of the state-space for the PFSA with depth $D(\epsilon)$ for a typical sample during stable and unstable behavior. The cluster tree also suggests the symbolic noise behavior of the data during the stable regime (the states are very close to each other based on the K-L distance). However, clearly a coarse clustering of states in the model during the unstable behavior would lead to significant information loss (as the states are statistically different). However, to compare the two Markov models, we keep the cardinality of the final models the same. For example, the algorithm is terminated with 3 states in the final Markov model during the stable as well as the unstable regime, and the final aggregated states are the three clusters depicted in the Figure 6. Once the final aggregated states are obtained, we estimate the parameters of the model using the Bayesian inference discussed in section 3-B.

Next, we present some results for model selection using the information-theoretic criteria discussed earlier in section 3-C. BIC and AIC are used to select the model which achieves the minimum score. As seen in the Figures 7a through 7b, the model with 5 states is selected for stable as well as for the unstable case (note that the original model for the stable class had 9 states for depth 2 and the unstable model had 27 states for a depth of 3). In contrast to cross-validation, the two criteria provide an unsupervised way for model selection. Thus we see that we need much smaller state-space to preserve the temporal statistics of the data and AIC and BIC provide us with a technique to select the compact model. In Figure 8, we show the Hamming distance between the sequences generated by the original model and the reduced models for a typical sample each from stable and unstable combustion. The box-plots are generated by simulating the original model and the reduced-order model to generate symbol sequences of length 1000 from 100 different initial states (i.e., a total of 100 strings are generated)

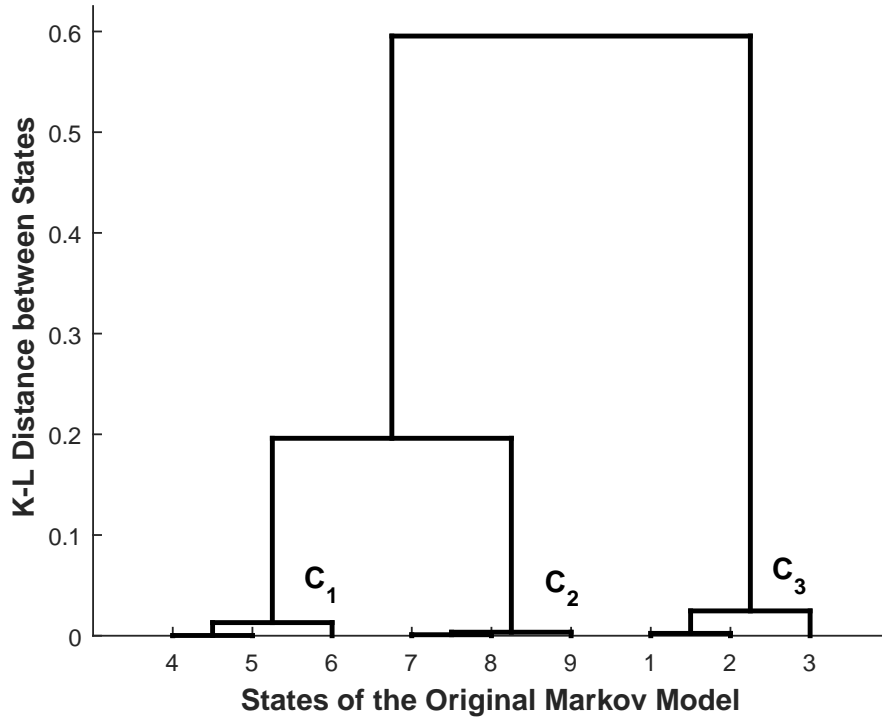


(a) Probability density function for the pressure time series data

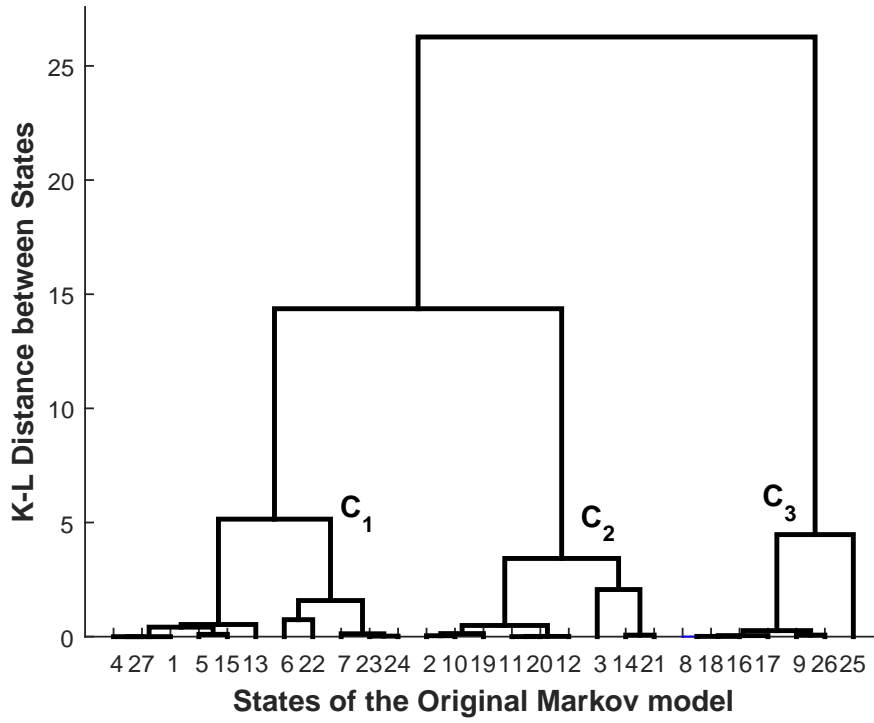


(b) Spectral decomposition of the stochastic matrix for 1-step Markov model

Fig. 5: The first plate in the above Figure shows the change in the empirical density calculated for the pressure time-series data as the process deviates from the stable operating condition to unstable operating condition. The second plate shows the spectral decomposition of the 1-step stochastic matrix for the data under stable and unstable operating conditions.

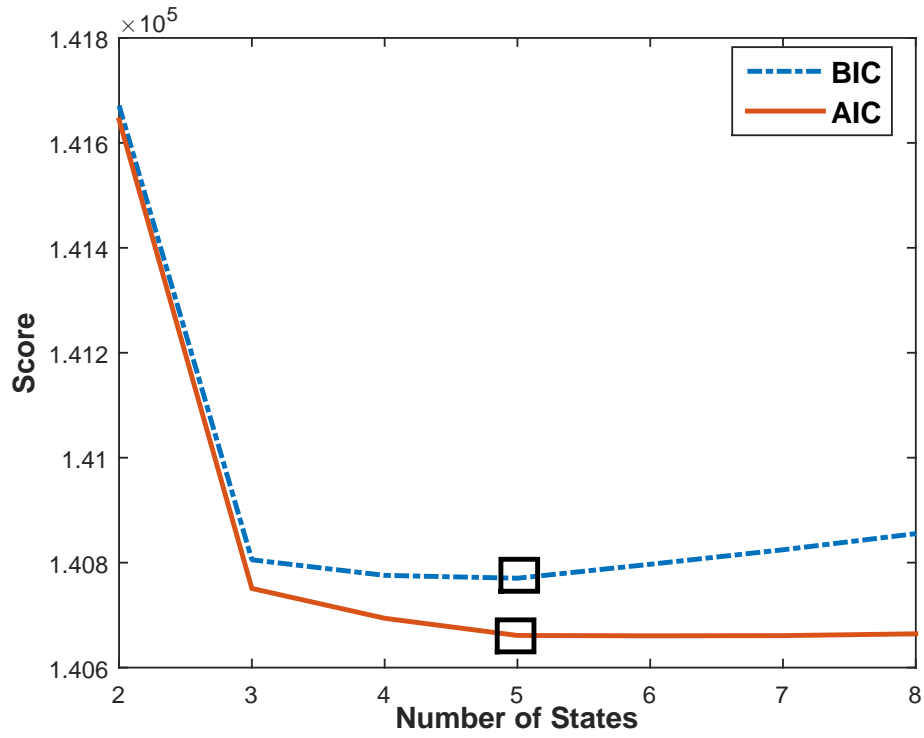


(a) Hierarchical cluster tree of stable states

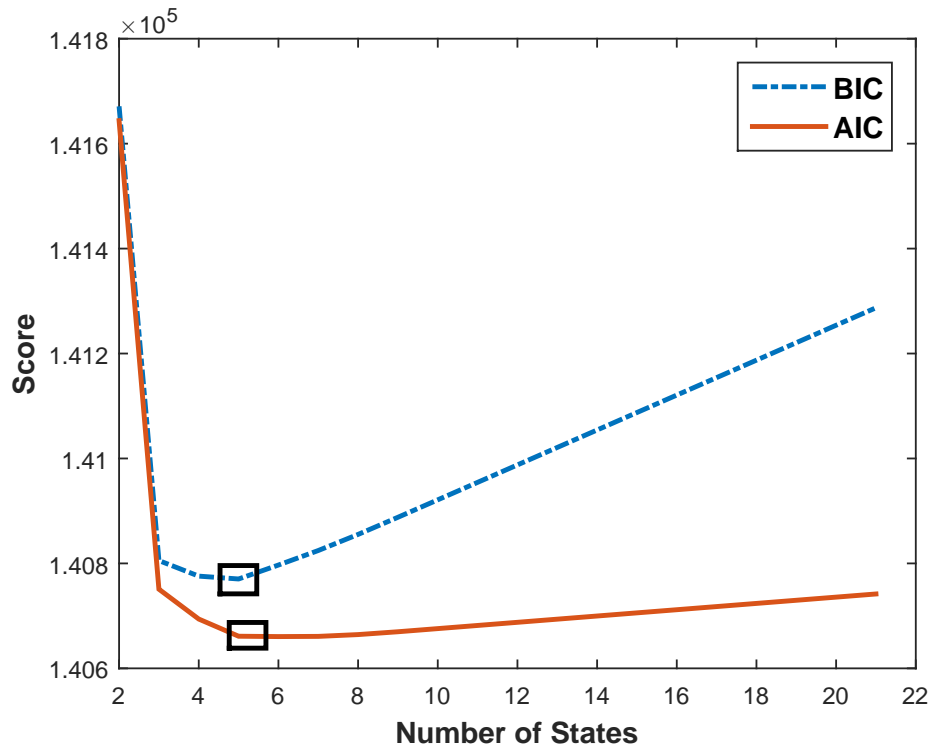


(b) Hierarchical cluster tree of unstable states

Fig. 6: State clustering under stable and unstable conditions.



(a) Model scores using the BIC and AIC criterion during a typical stable condition



(b) Model scores using the BIC and AIC criterion during a typical unstable condition

Fig. 7: Unsupervised model selection under stable and unstable conditions.

and the Hamming distance between them is calculated. A bound on the Hamming distance between the sequences generated by the original model and final model is also calculated using the inequality (26). The results are shown in Figure 8. It is possible to use the proposed Hamming distance metric to select a final model; however, this measures the distance between the distributions induced by the Markov models, and model selection using it is left as a future work. It is noted that the bounds on Hamming distance can provide a computationally convenient way to select model scores as it can be found from the symbol emission probabilities of the model instead of explicitly looking at the predictive capability by looking at the likelihoods of the symbol sequences.

B. Bearing

The same procedure of downsampling and depth estimation is followed for analysis of bearing data as was described in the previous section for combustion. A ternary alphabet is again chosen to discretize the continuous data after downsampling and the maximum entropy partitioning is used to find the partitions. Using the spectral method, a depth of 2 (i.e., a total of 9 states) is estimated for an $\epsilon = 0.02$ (we skip the plot of spectral decomposition plot for brevity). The BIC and AIC score for the different models is shown in Figure 9 and the model with five states is selected using the obtained scores (marked in black rectangle). In Figure 10, we show the Hamming distance between the sequences generated by the original model (with 9 states) and the reduced models and the corresponding bounds obtained by inequality (26).

7. CLASSIFICATION AND ANOMALY DETECTION RESULTS

In this section, we present some results for anomaly detection and classification using the pressure time-series data to infer the underlying reduced-order Markov model. As we discussed earlier in section 5-A, the exact transition point of the system from stable to unstable is unknown, we first present results on anomaly detection and clustering of the data into different clusters which can be then associated with the stable and unstable class. We will present two different metrics for anomaly detection that allows models of different state-space and structure to be compared. It is noted that the word metric is used here in a loose sense; it is meant to be a distance that could be used to compare two different Markov models.

A. Anomaly Detection

As individual time-series have different state-space, we define some metrics to compare them. These metrics reflect changes in the information complexity of Markov models and reveal different behavior of combustion process based on the changes in the inferred data model. In particular, the following two metrics are defined.

- 1) Cluster Divergence: This measure is defined for individual Markov models based on the cluster structure of the state-space of the model. Physically, it represents the maximum statistical difference between the states of the Markov model measures using K-L distance. It is calculated for a particular model \mathcal{M} as follows

$$\Delta_{\mathcal{M}} = \max_{q_i, q_j \in \mathcal{Q}} d(q_i, q_j) \quad (27)$$

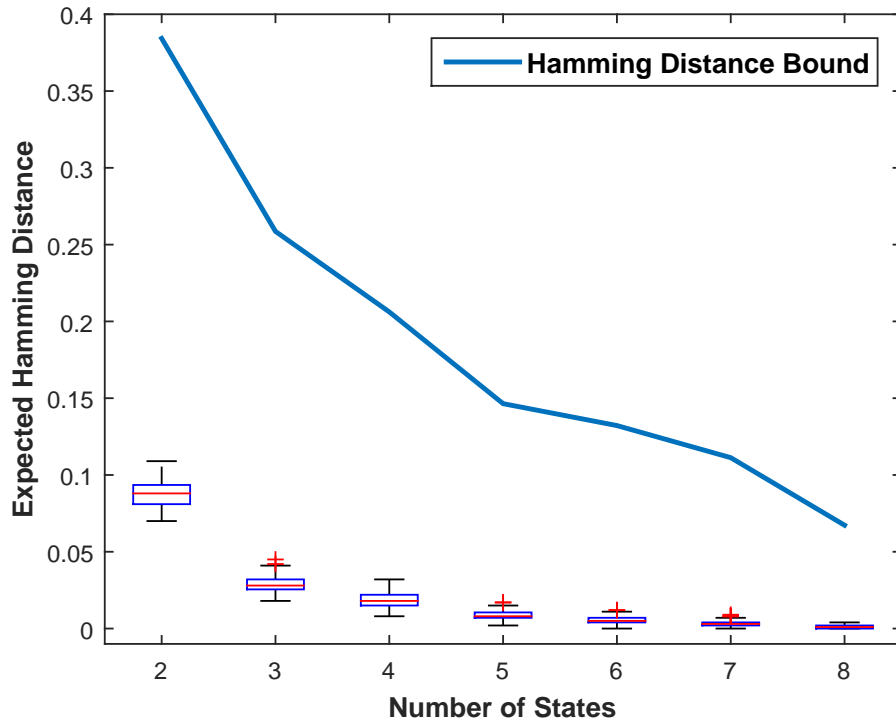
where d is defined by equation (5).

- 2) Discrepancy Statistics: We measure the discrepancy between the i.i.d. statistics and the Markov statistics for the discretized data. This could be also interpreted as the information gain for Markov models. This measure also represents the information complexity of the data. If the i.i.d. statistics and the Markov statistics are very close, then the data has no temporal statistics; however, an increase in this measure would indicate the information gain by creating a temporal Markov model for the data. This is measured by the following equation.

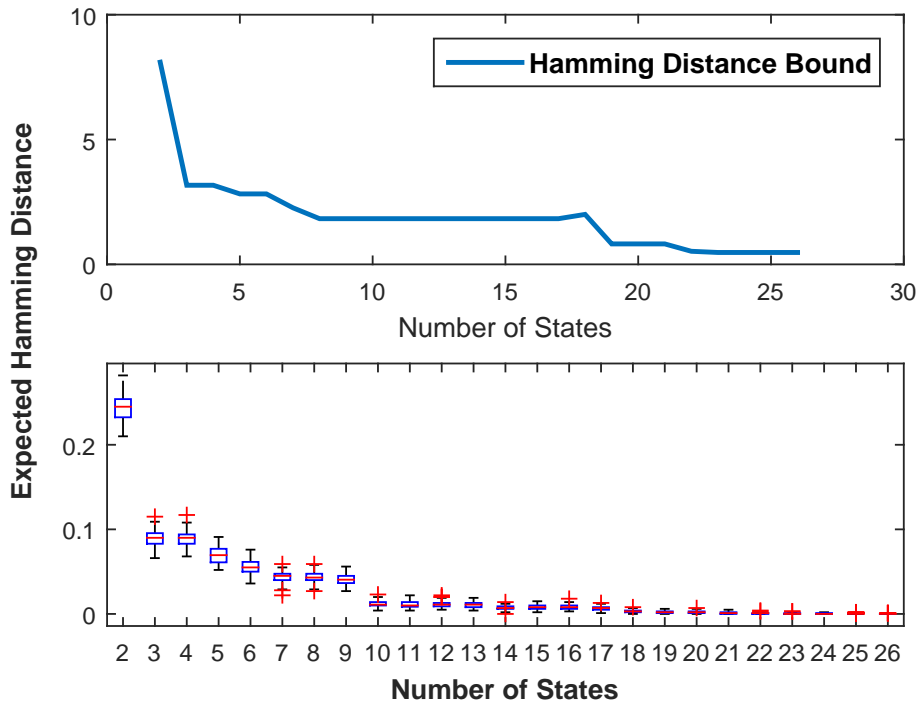
$$H_{\mathcal{M}} = \sum_{q \in \mathcal{Q}} \Pr(q) D_{KL}(\Pr(\mathcal{A} | q) || \Pr(\mathcal{A})) \quad (28)$$

where $\Pr(\mathcal{A} | q)$ represents the symbol emission probability conditioned on a state q of the Markov model and $\Pr(\mathcal{A})$ represents the marginal symbol emission probability. The term D_{KL} represents the symmetric K-L distance between the two distributions.

In Figure 10, we present some results to show the behavior of $\Delta_{\mathcal{M}}$ with increasing pressure fluctuations. It is noted that every model has been created in an unsupervised fashion by first discretizing and then, estimating the memory of the discrete sequence. As seen in Figure 11a, there are three distinct behaviors that can be associated with $\Delta_{\mathcal{M}}$. With low pressure fluctuations, the metric is very close to 0, indicating that the states of the model are very similar statistically. This is seen until data number 200 with corresponding $P_{rms} \sim 0.065$ psig, which leads to a gradual change to a point where the measure saturates with $P_{rms} \sim 0.12$ psig (when the process becomes unstable). Thus, with this gradual trend with increasing pressure fluctuations, we associate different behaviors with the process. However, as is seen in the Figure 11a, the transition from stable to unstable behavior is not clearly defined and is very difficult to label during the experiments as the process is very fast. We show the pressure signals from the three different clusters in Figure 11b where it could be seen that the sample number 250 could be seen to approach an approximate limit cyclic behavior (and thus, could be loosely classified as transient stage). An important point to note at this point is that this measure is independent of any operating conditions and only depends on stability (or



(a) Hamming distance between the original and final models for a typical stable combustion sample



(b) Hamming distance between the original and final models for a typical unstable combustion sample

Fig. 8: Box plot for Hamming distance between the original and reduced-order models obtained after merging based on the results in Section 4

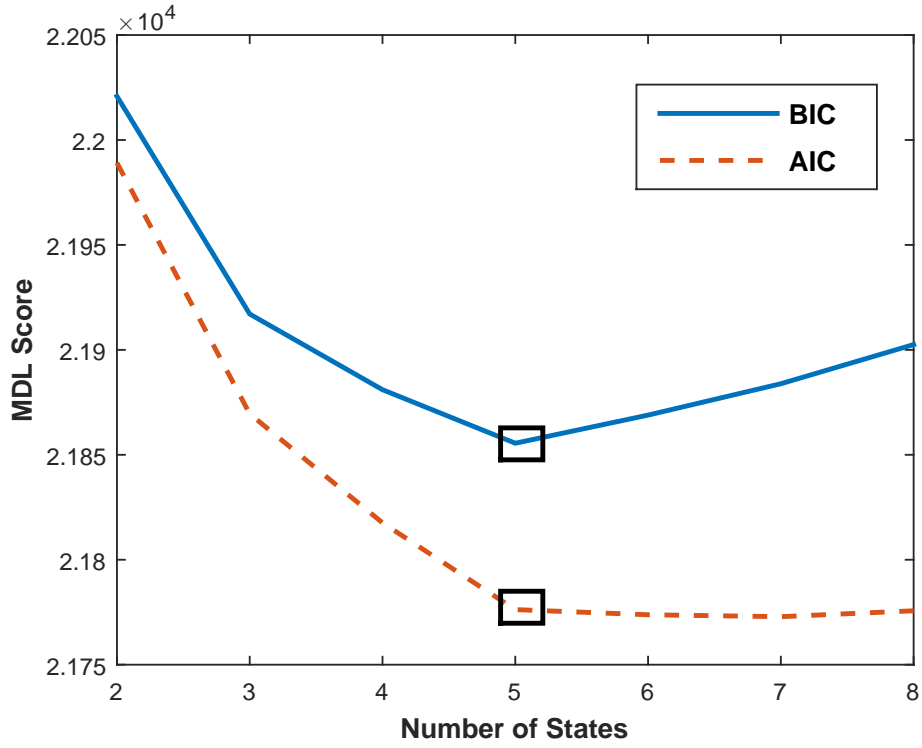


Fig. 9: Model scores using the BIC and AIC criteria; selected models are depicted by black rectangles.

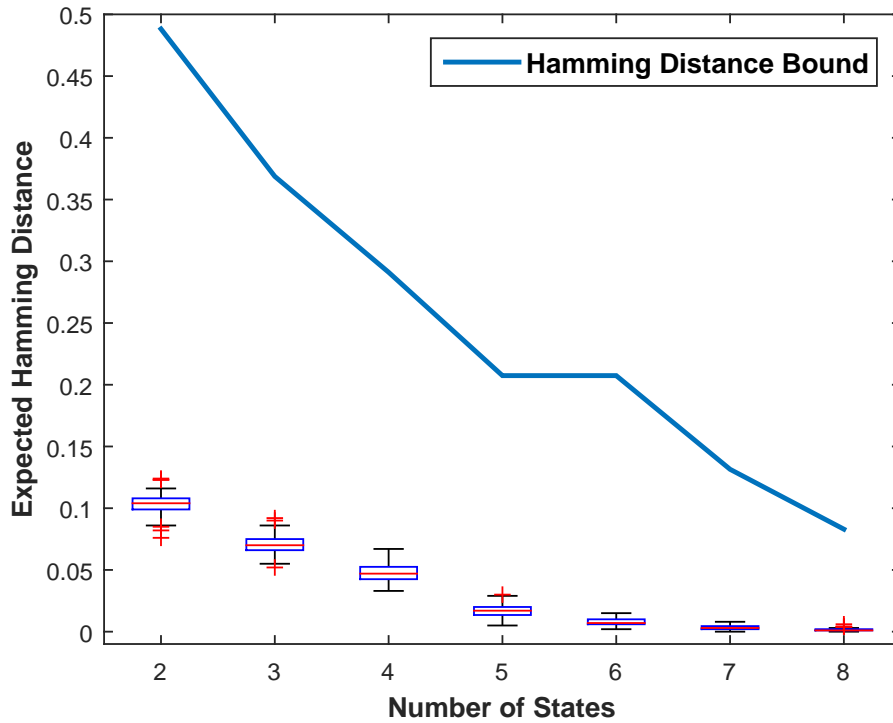


Fig. 10: Box plot of the Hamming distance between the original and reduced-order models along with the analytical bound presented in Section 4.

TABLE II: Performance of classifiers with different number of states. Mean Error= Lower is better.

Number of States	Classifier	Classification Error (%)
9	SVM	3.48 ± 0.74
	DT	9.83 ± 3.24
8	SVM	3.62 ± 0.71
	DT	9.38 ± 3.11
7	SVM	2.87 ± 0.68
	DT	7.70 ± 2.61
6	SVM	2.48 ± 0.61
	DT	7.00 ± 2.55
5	SVM	2.05 ± 0.54
	DT	6.10 ± 2.17
4	SVM	1.86 ± 0.43
	DT	4.72 ± 2.29
3	SVM	1.69 ± 0.45
	DT	5.56 ± 1.90
2	SVM	1.67 ± 0.43
	DT	4.83 ± 1.80

instability) of the process. This metric is thus used for anomaly detection. In Figure 10c, we show the statistics of $\Delta_{\mathcal{M}}$ with four states. We see that there is some loss of information up on merging states in the unstable class; the stable cluster remains unchanged implying that the states are statistically similar and the model distortion up on merging of states is insignificant. Thus, $\Delta_{\mathcal{M}}$ can be reliably used to detect departure from stable behavior.

The statistics for the discrepancy measure for the full state models is shown in Figure 11. The plot in Figure 11 also agrees qualitatively with the earlier results on $\Delta_{\mathcal{M}}$. From these plots, we can infer that the Markov statistics for the stable cluster is very similar to the i.i.d. statistics and thus the data is very much independently distributed and conditioning on the inferred states of the Markov models doesn't improve predictability (or information complexity) of the temporal model. Thus, these two measures help infer the changes in the behavior of the data during the combustion process and are useful for anomaly detection.

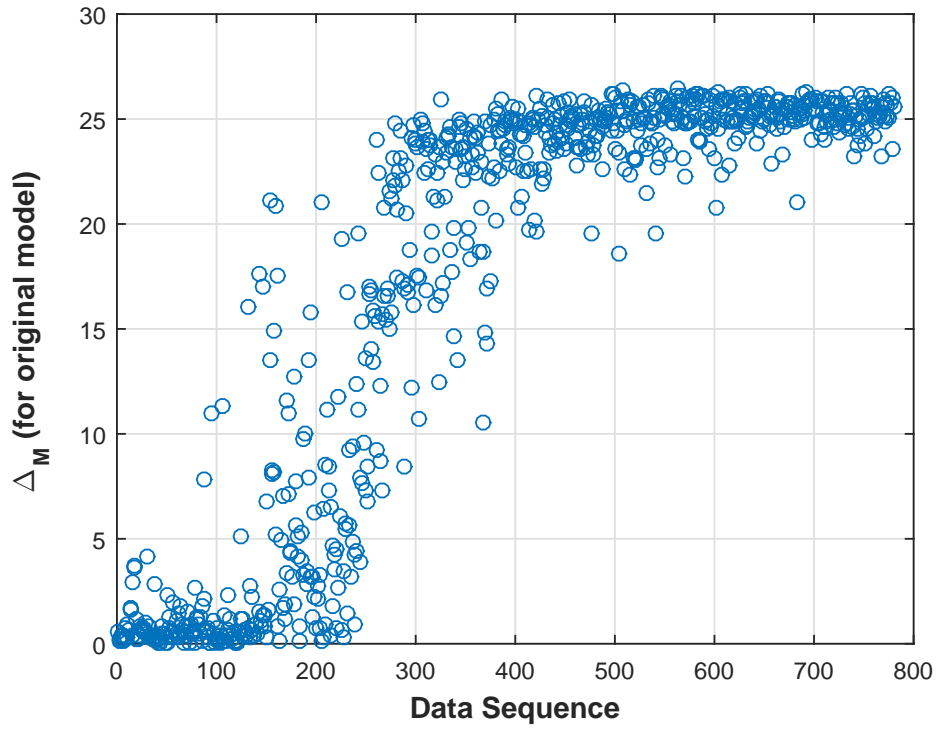
To see more explicitly the changes in the underlying models, the models during stable and unstable phases are visualized in the information space. To do this, we reduce the state space of the models to just 2 states and estimate the corresponding emission parameters. As the models have three symbols, the emission matrix has 2 rows and each row corresponds to the symbol emission probabilities conditioned on the two states. Each of these rows for 100 cases from stable and 100 cases from unstable are plotted on a single simplex plane which is shown in Figure 12. The Figure shows the clusters of stable and unstable cases in the information space and that the model with even 2 states are clustered separately. This shows that there is a structured change in the temporal dynamics of the data at the two phases and that the inferred Markov models are able to capture this change. Furthermore, the distinctive features of the models are sufficiently retained even after significant reduction in the state-space of the models.

B. Classification

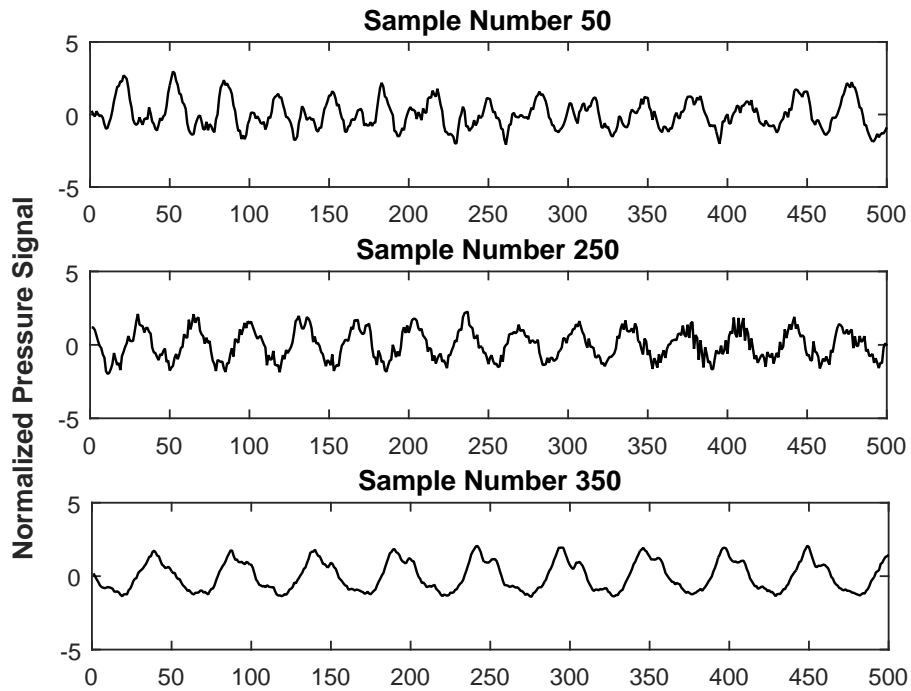
These models are then used to train classifiers using support vector machines (SVM) and decision trees (DT) [5]. The rationale behind using multiple classifier is to show that the performance of the Markov models is independent of the classification technique (i.e., it works equally well with maximum margin classifiers or decision tree classifiers). The SVM classifier is trained using a radial basis function kernel while the decision tree is trained using the standard euclidean distance. The classifiers are trained with 100 data points from each class and are tested on the remaining data (around 80 and 380 for stable and unstable respectively). The tests are repeated for 100 different train and test data sets from the total data. The results of classification accuracy are listed in Table II. The SVM classifier is able to achieve around 1.67% error using models with 2 states while the decision tree classifier is able to achieve around 4.70% error using models with 4 states. This provides another way of selecting the final model for state merging in a supervised learning setting. It is noted that the original models contain 9 states for stable and 27 for unstable class.

8. SUMMARY, CONCLUSIONS AND FUTURE WORK

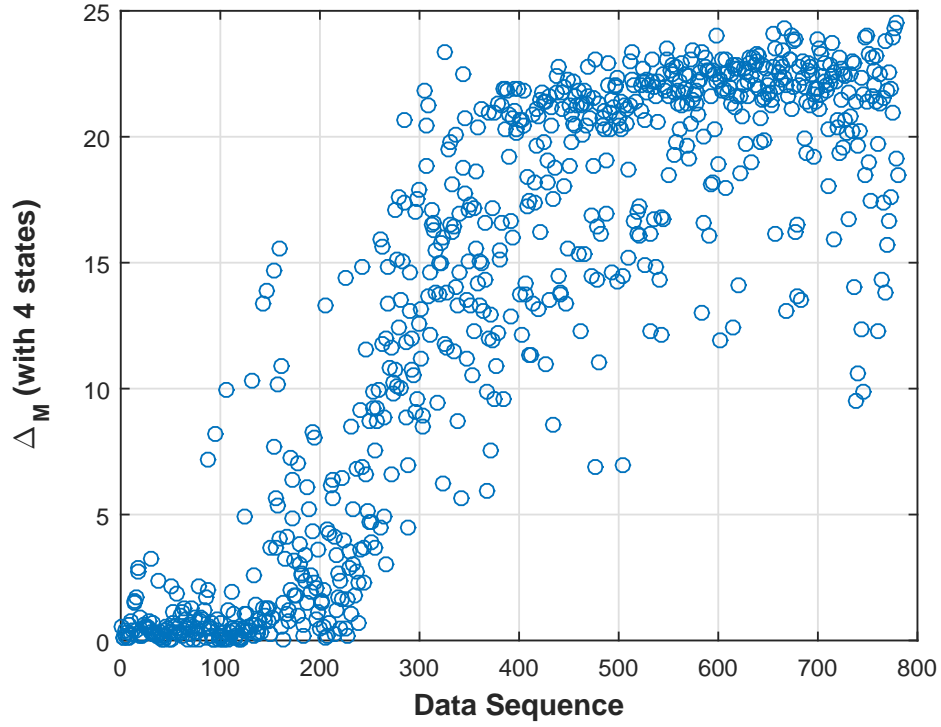
In recent times the idea of representation learning has become very popular in the machine learning literature as it allows decoupling of data for model learning from the end-objectives like classification or clustering. In this paper, we presented a technique for Markov modeling of time-series data using concepts of symbolic dynamics which allows inference of model structure as well as parameters for compact data representation. In the proposed technique we first estimate the memory size of the discretized time-series data. The size of memory is estimated using spectral decomposition properties of the one-step Markov model created from the symbol sequence. Then, a second pass of data is made to infer the model with the right memory and the corresponding symbol emission matrix is estimated. Then the equivalence class of states based on K-L distance between the states are estimated using hierarchical clustering of the corresponding states of the Markov model. The proposed concepts were validated using two different datasets– combustion instability and bearing. Modeling of combustion instability



(a) Δ_M for the full state model for the time-series data with increasing pressure root mean square



(b) Typical pressure signals from the three clusters seen in Figure 11a



(c) $\Delta_{\mathcal{M}}$ for models with 4 states for the time-series data with increasing pressure root mean square

Fig. 10: Anomalous behavior of data in the combustion process

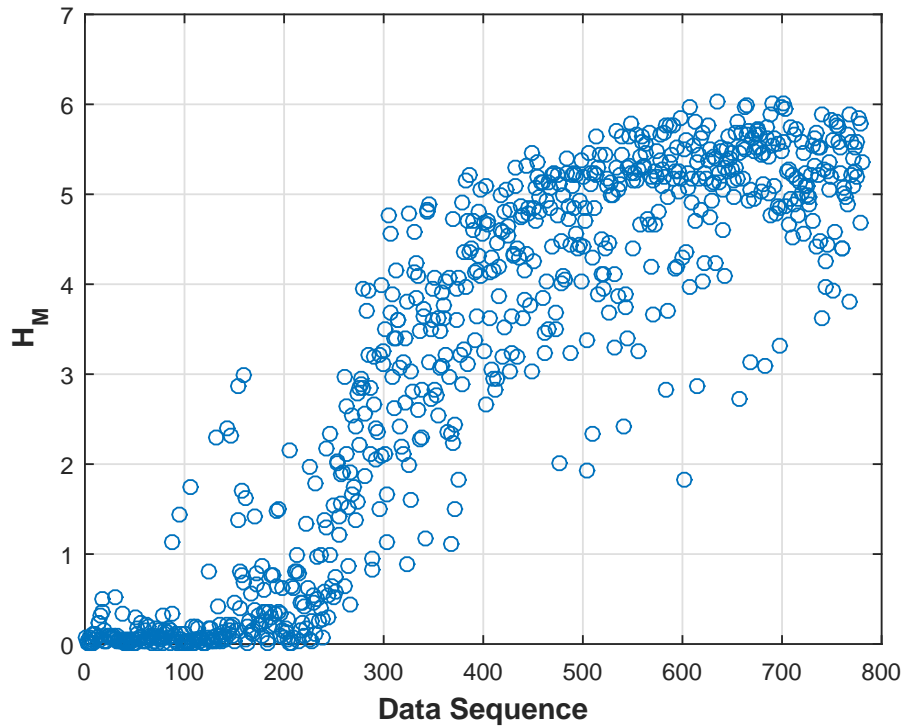


Fig. 11: Variation of discrepancy statistics $H_{\mathcal{M}}$ with increasing pressure fluctuations. This also shows an anomaly around the point 200 and qualitatively agrees to the behavior of $\Delta_{\mathcal{M}}$.

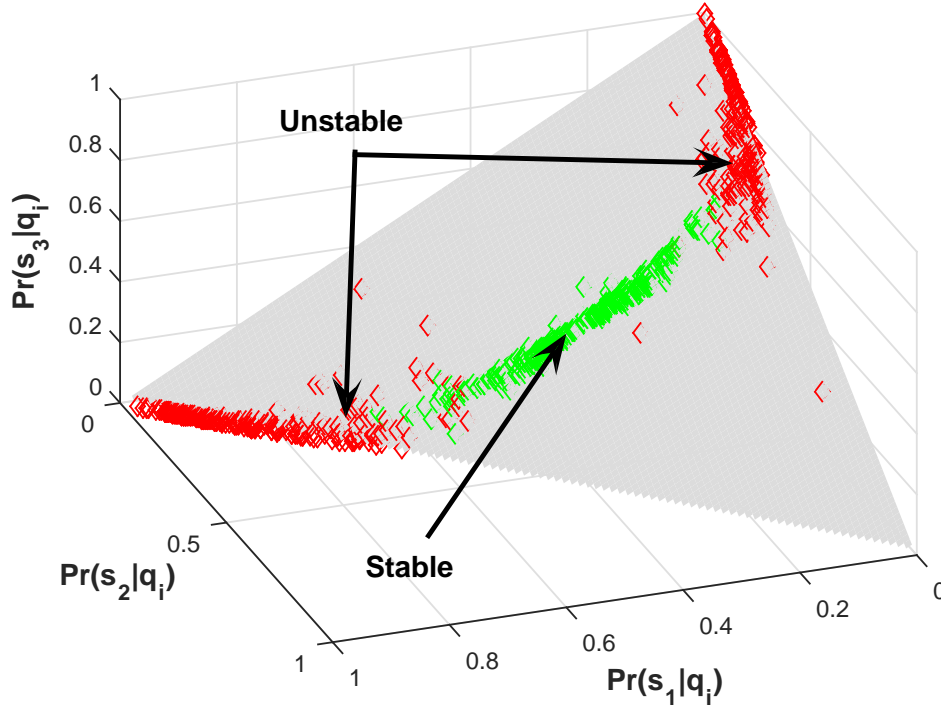


Fig. 12: Cluster of stable and unstable phase in information space. Each point is a row of the emission matrix for the reduced Markov model with 2 states. The plot shows the change in the Markov model as the process moves from stable and unstable. **Red diamonds** represent the unstable phase while **green diamonds** represent the stable phase.

still remains a puzzle in the combustion community. The Markov modeling technique was used to analyze the problem of combustion instability. The proposed ideas were tested on experimental data from a swirl-stabilized combustor used to study unstable thermo-acoustic phenomenon during combustion process. The proposed approach allows us to infer the complexity of the time-series data based on the inferred Markov model. Two different metrics were proposed for anomaly detection and classification of the stable and unstable classes. The results presented in this paper are encouraging as the inferred models are able to identify the stable and unstable phases independent of any other operating condition.

Simultaneous optimization of discretization and memory estimation for model inference is a topic of future research. While the results obtained with Markov modeling for the combustion instability problem are inspiring, further investigation with transient data is required for better characterization of the process. More thorough comparison of the proposed models with HMM models of similar state-space size is also an important topic of future work.

ACKNOWLEDGMENTS

The authors would like to thank Professor Domenic Santavicca and Mr. Jihang Li of Center for Propulsion, Penn State for kindly providing the experimental data for combustion used in this work.

REFERENCES

- [1] Prognostic data repository: Bearing data set nsf i/ucrc center for intelligent maintenance systems, 2010. [Online]. Available: <http://ti.arc.nasa.gov/tech/dash/pcoe/prognostic-data-repository/>
- [2] H. Akaike, "A new look at the statistical model identification," *IEEE Transactions on Automatic Control*, vol. 19, no. 6, pp. 716–723, Dec 1974.
- [3] A. Banaszuk, P. G. Mehta, and G. Hagen, "The role of control in design: From fixing problems to the design of dynamics," *Control Engineering Practice*, vol. 15, no. 10, pp. 1292–1305, 2007.
- [4] A. Banaszuk, P. G. Mehta, C. A. Jacobson, and A. I. Khibnik, "Limits of achievable performance of controlled combustion processes," *Control Systems Technology, IEEE Transactions on*, vol. 14, no. 5, pp. 881–895, 2006.
- [5] C. M. Bishop, *Pattern recognition and machine learning*. Springer, 2006.
- [6] S. Candell, D. Durox, T. Schuller, J.-F. Bourgoign, and J. P. Moeck, "Dynamics of swirling flames," *Annual review of fluid mechanics*, vol. 46, pp. 147–173, 2014.
- [7] I. Chattopadhyay and H. Lipson, "Abductive learning of quantized stochastic processes with probabilistic finite automata," *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, vol. 371, no. 1984, p. 20110543, 2013.
- [8] F. Darema, "Dynamic data driven applications systems: New capabilities for application simulations and measurements," in *5th International Conference on Computational Science - ICCS 2005*, Atlanta, GA; United States, 2005.
- [9] B. C. Geiger, T. Petrov, G. Kubin, and H. Koepl, "Optimal kullback-leibler aggregation via information bottleneck," *Automatic Control, IEEE Transactions on*, vol. 60, no. 4, pp. 1010–1022, 2015.
- [10] R. M. Gray, *Entropy and information*. Springer, 1990.

- [11] Y. Huang and V. Yang, "Dynamics and stability of lean-premixed swirl-stabilized combustion," *Progress in Energy and Combustion Science*, vol. 35, no. 4, pp. 293–364, 2009.
- [12] D. K. Jha, A. Srivastav, and A. Ray, "Temporal learning in video data using deep learning and Gaussian processes," in *Workshop on Machine Learning for Prognostics and Health Management at 2016 KDD, San Francisco, CA*, 2016.
- [13] D. K. Jha, A. Srivastav, K. Mukherjee, and A. Ray, "Depth estimation in Markov models of time-series data via spectral analysis," in *American Control Conference (ACC)*, 2015. IEEE, 2015, pp. 5812–5817.
- [14] J. Lin, E. Keogh, L. Wei, and S. Lonardi, "Experiencing SAX: a novel symbolic representation of time series," *Data Mining and Knowledge Discovery*, vol. 15, no. 2, pp. 107–144, October 2007.
- [15] D. Lind and B. Marcus, *An introduction to symbolic dynamics and coding*. Cambridge University Press, 1995.
- [16] K. Marton, "Bounding d -distance by informational divergence: a method to prove measure concentration," *The Annals of Probability*, vol. 24, no. 2, pp. 857–866, 1996.
- [17] J. P. Moeck, J.-F. Bourgouin, D. Durox, T. Schuller, and S. Candel, "Nonlinear interaction between a precessing vortex core and acoustic oscillations in a turbulent swirling flame," *Combustion and Flame*, vol. 159, no. 8, pp. 2650–2668, 2012.
- [18] K. Mukherjee and A. Ray, "State splitting and merging in probabilistic finite state automata for signal representation and analysis," *Signal Processing*, vol. 104, pp. 105–119, 2014.
- [19] M. Murugesan and R. Sujith, "Combustion noise is scale-free: transition from scale-free to order at the onset of thermoacoustic instability," *Journal of Fluid Mechanics*, vol. 772, pp. 225–245, 2015.
- [20] V. Nair, G. Thampi, and R. Sujith, "Intermittency route to thermoacoustic instability in turbulent combustors," *Journal of Fluid Mechanics*, vol. 756, pp. 470–487, 2014.
- [21] J. O'Connor, V. Acharya, and T. Liewu, "Transverse combustion instabilities: Acoustic, fluid mechanic, and flame processes," *Progress in Energy and Combustion Science*, vol. 49, pp. 1–39, 2015.
- [22] H. Qiu, J. Lee, J. Lin, and G. Yu, "Wavelet filter-based weak signature detection method and its application on rolling element bearing prognostics," *Journal of sound and vibration*, vol. 289, no. 4, pp. 1066–1090, 2006.
- [23] V. Rajagopalan and A. Ray, "Symbolic time series analysis via wavelet-based partitioning," *Signal Processing*, vol. 86, no. 11, pp. 3309–3320, 2006.
- [24] A. Ray, "Symbolic dynamic analysis of complex systems for anomaly detection," *Signal Processing*, vol. 84, no. 7, pp. 1115–1130, July 2004.
- [25] S. Sarkar, S. R. Chakravarthy, V. Ramanan, and A. Ray, "Dynamic data-driven prediction of instability in a swirl-stabilized combustor," *International Journal of Spray and Combustion Dynamics*, p. 1756827716642091, 2016.
- [26] G. Schwarz, "Estimating the dimension of a model," *Ann. Statist.*, vol. 6, no. 2, pp. 461–464, 03 1978.
- [27] Sé, b. Ducruix, T. Schuller, D. Durox, Sé, and b. Candel, "Combustion dynamics and instabilities: Elementary coupling and driving mechanisms," *Journal of Propulsion and Power*, vol. 19, no. 5, pp. 722–734, 2003.
- [28] C. R. Shalizi and K. L. Shalizi, "Blind construction of optimal nonlinear recursive predictors for discrete sequences," in *Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence*, ser. UAI '04, 2004, pp. 504–511.
- [29] A. Srivastav, "Estimating the size of temporal memory for symbolic analysis of time-series data," *American Control Conference, Portland, OR, USA*, pp. 1126–1131, June 2014.
- [30] D. A. Tobon-Mejia, K. Medjaher, N. Zerhouni, and G. Tripot, "A data-driven failure prognostics method based on mixture of gaussians hidden Markov models," *IEEE Transactions on reliability*, vol. 61, no. 2, pp. 491–503, 2012.
- [31] E. Vidal, F. Thollard, C. De La Higuera, F. Casacuberta, and R. C. Carrasco, "Probabilistic finite-state machines-part i," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 27, no. 7, pp. 1013–1025, 2005.
- [32] M. Vidyasagar, "A metric between probability distributions on finite sets of different cardinalities and applications to order reduction," *Automatic Control, IEEE Transactions on*, vol. 57, no. 10, pp. 2464–2477, 2012.
- [33] N. Virani, D. K. Jha, and A. Ray, "Sequential hypothesis tests using Markov models of time series data," in *Workshop on Machine Learning for Prognostics and Health Management at 2016 KDD, San Francisco, CA*, 2016.
- [34] R. Xu and D. Wunsch, "Survey of clustering algorithms," *Neural Networks, IEEE Transactions on*, vol. 16, no. 3, pp. 645–678, 2005.
- [35] Y. Xu, S. M. Salapaka, and C. L. Beck, "Aggregation of graph models and Markov chains by deterministic annealing," *Automatic Control, IEEE Transactions on*, vol. 59, no. 10, pp. 2807–2812, 2014.