

1 Formalization Of Probabilistic Generators For Stochastic Dynamical Systems

Throughout this paper, Σ denotes a finite alphabet of symbols. The set of all finite but possibly unbounded strings on Σ is denoted by Σ^* , the Kleene $*$ operation [8]. The set of finite strings over Σ form a concatenative monoid, with the empty word λ as identity. Concatenation of two strings $x, y \in \Sigma^*$ is written as xy . Thus $xy = x\lambda y = xy\lambda = \lambda xy$. The set of strictly infinite strings on Σ is denoted as Σ^ω , where ω denotes the first transfinite cardinal. For a string $x \in \Sigma^*$, $|x|$ denotes the length of x and for a set A , $|A|$ denotes its cardinality.

1.1 Quantized Stochastic Processes

Definition 1 (QSP). A QSP \mathcal{H} is a discrete time Σ -valued strictly stationary, ergodic stochastic process, i.e.

$$\mathcal{H} \quad \{X_t \mid X_t \text{ is a } \Sigma\text{-valued random variable for } t \in \mathbb{N} \cup \{\circ\}\}$$

A stochastic process is ergodic if moments can be calculated from a single, sufficiently long realization, and strictly stationary if moments are not functions of time.

We next formalize the connection of QSPs to PFSA generators.

Definition 2 (σ -Algebra On Infinite Strings). For the set of infinite strings on Σ , we define \mathbb{B} to be the smallest σ -algebra generated by $\{x\Sigma^\omega \mid x \in \Sigma^*\}$.

Lemma 1. Any QSP induces a probability space $\Sigma^\omega, \mathbb{B}, \mu$.

Proof. Using stationarity of QSP \mathcal{H} we can construct a probability measure $\mu : \mathbb{B} \rightarrow \{0, 1\}$ by defining for any sequence $x \in \Sigma^* \setminus \{\lambda\}$, and a sufficiently large number of realizations N_R of \mathcal{H} , with fixed initial conditions:

$$\mu x \Sigma^\omega = \lim_{N_R \rightarrow \infty} \frac{(\# \text{ of initial occurrences of } x)}{(\# \text{ of initial occurrences of } x \text{ among all sequences of length } |x|)}$$

and extending the measure to elements of $\mathbb{B} \setminus B$ via at most countable sums. Note that $\mu \Sigma^\omega = \sum_{x \in \Sigma^*} \mu x \Sigma^\omega = 1$, and for the null word $\mu \lambda \Sigma^\omega = \mu \Sigma^\omega = 1$. \square

For notational brevity, we denote $\mu x \Sigma^\omega$ as $Pr x$.

Classically, states are induced via the Nerode equivalence, which defines two strings to be equivalent if and only if any finite extension of the strings is either both accepted, or both rejected [8] by the language under consideration. We use a probabilistic extension [5].

Definition 3 (Probabilistic Nerode Relation). $\Sigma^\omega, \mathbb{B}, \mu$ induces an equivalence relation \sim_N on the set of finite strings Σ^* as:

$$\begin{aligned} \forall x, y \in \Sigma^*, x \sim_N y \iff \forall z \in \Sigma^* ((Pr xz / Pr yz) \circ) \vee \\ |Pr xz / Pr x - Pr yz / Pr y| < \epsilon \end{aligned} \tag{1}$$

For $x \in \Sigma^*$, the equivalence class of x is denoted as $[x]$.

It follows that \sim_N is right invariant, i.e.

$$x \sim_N y \Rightarrow \forall z \in \Sigma^*, xz \sim_N yz \tag{2}$$

A right-invariant equivalence on Σ^* always induces an automaton structure.

Definition 4 (Initial-Marked PFSA). An Initial-Marked PFSA is a 5-tuple $(Q, \Sigma, \delta, \tilde{\pi}, q_0)$, where Q is a finite state set, Σ is the alphabet, $\delta : Q \times \Sigma \rightarrow Q$ is the state transition function, and $\tilde{\pi} : Q \times \Sigma \rightarrow \{0, 1\}$ specifies the conditional symbol-generation probabilities. δ and $\tilde{\pi}$ are recursively extended to arbitrary $y \sigma x \in \Sigma^*$ as $\delta q, \sigma x \delta q, \sigma, x$, and $\tilde{\pi} q, \sigma x \tilde{\pi} q, \sigma \tilde{\pi} \delta q, \sigma, x$. $q_0 \in Q$ is the initial state. If the next symbol is specified, our resultant state is fixed; similar to Probabilistic Deterministic Automata [7]. However, unlike the latter, we lack final states. Additionally, we assume our graphs to be strongly connected.

Definition 4 has no notion of a final state, and later we will remove initial state dependence using ergodicity. First we formalize how the notion of a PFSA arises from a QSP.

Lemma 2 (From QSP To A PFSA). *If the probabilistic Nerode relation has a finite index, then there exists an initial-marked PFSA generator.*

Proof. Every QSP represented as a probability space $\Sigma^\omega, \mathbb{B}, \mu$ induces a probabilistic automaton $Q, \Sigma, \delta, \tilde{\pi}, q_0$, where Q is the set of equivalence classes of the probabilistic Nerode relation (Definition 3), Σ is the alphabet, and:

$$\delta x, \sigma \quad x\sigma \tag{3}$$

$$\tilde{\pi}x, \sigma \quad \frac{Prx'\sigma}{Prx'} \text{ for any choice of } x' \in x \tag{4}$$

q_0 is identified with λ , and finite index of \sim_N implies $|Q| < \infty$. \square

The above construction yields a *minimal realization* unique up to state renaming.

Corollary 1 (To Lemma 2: Null-word Probability). *For the PFSA $Q, \Sigma, \delta, \tilde{\pi}$ induced from a QSP \mathcal{H} :*

$$\forall q \in Q, \tilde{\pi}q, \lambda = 1 \tag{5}$$

Proof. For $q \in Q$ let $x \in \Sigma^*$ such that $x \sim q$. From Eq. (4), we have:

$$\tilde{\pi}q, \lambda = \frac{Prx'\lambda}{Prx'} \text{ for } x' \in x \tag{6}$$

$$\frac{Prx'}{Prx'} = 1 \tag{7}$$

\square

While many reported approaches *define* the probability of the null-word to be unity, we can derive it from our formulation.

1.2 Canonical Representations

Next we define canonical representations to remove initial-state dependence. We use $\tilde{\Pi}$ to denote the matrix representation of $\tilde{\pi}$, i.e., $\tilde{\Pi}_{ij} = \tilde{\pi}q_i, \sigma_j$, $q_i \in Q, \sigma_j \in \Sigma$. We need the notion of transformation matrices Γ_σ .

Definition 5 (Transformation Matrices). *For an initial-marked PFSA $G = Q, \Sigma, \delta, \tilde{\pi}, q_0$, the symbol-specific transformation matrices $\Gamma_\sigma \in \{0, 1\}^{|Q| \times |Q|}$ are:*

$$\Gamma_\sigma|_{ij} = \begin{cases} \tilde{\pi}q_i, \sigma, & \text{if } \delta q_i, \sigma = q_j \\ 0, & \text{otherwise} \end{cases} \tag{8}$$

Transformation matrices have a single non-zero entry per row, reflecting our generation rule that given a state and a generated symbol, the next state is fixed. States in the canonical representation (denoted as ρ_x) are identified with probability distributions over states of the initial-marked PFSA. Here x denotes the string in Σ^* realizing this distribution, beginning from the stationary distribution on the states of the initial-marked representation. ρ_x is an equivalence class, and hence x is not unique.

Definition 6 (Canonical Representations). *An initial-marked PFSA $G = Q, \Sigma, \delta, \tilde{\pi}, q_0$ uniquely induces a canonical representation $Q^C, \Sigma, \delta^C, \tilde{\pi}^C$, with Q^C being the set of probability distributions over Q , $\delta^C : Q^C \times \Sigma \rightarrow Q^C$, and $\tilde{\pi}^C : Q^C \times \Sigma \rightarrow \{0, 1\}$, as follows:*

1. Construct the stationary distribution on Q using the transition probabilities of the Markov Chain induced by G , and include this as the first element ρ_λ of Q^C . The transition matrix for this induced chain is the row-stochastic matrix $M \in \{0, 1\}^{|Q| \times |Q|}$, with $M_{ij} = \sum_{\sigma \in \Sigma} \tilde{\pi}q_i, \sigma$.

2. Define δ^C and $\tilde{\pi}^C$ recursively:

$$\delta^C \varphi_x, \sigma = \frac{1}{\|\varphi_x \Gamma_\sigma\|_1} \varphi_x \Gamma_\sigma \triangleq \varphi_{x\sigma} \quad (9)$$

$$\tilde{\pi}^C \varphi_x, \sigma = \varphi_x \tilde{\Pi} \quad (10)$$

For a QSP \mathcal{H} , the canonical representation is denoted as $\mathcal{C}_{\mathcal{H}}$.

Ergodicity of QSPs, which makes φ_λ independent of the initial state in the initial-marked PFSA, implies that the canonical representation is initial state independent (See Figure ??), and subsumes the initial-marked representation in the following sense:

Let $\mathcal{E} = \{e^i \in \mathbb{O}^{|\Sigma|^{|Q|}}, i = 1, \dots, |\mathcal{Q}|\}$ denote the set of distributions satisfying:

$$e^i|_j = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

Then we note:

1. If we execute the canonical construction with an initial distribution from \mathcal{E} , then we get the initial-marked PFSA (with the initial marking missing).
2. If during the construction we encounter $\varphi_x \in \mathcal{E}$ for some x , then we stay within the graph of the initial-marked PFSA for all right extensions of x . This thus eliminates the need of knowing the initial state explicitly (See Figure ??), provided we find string x , which takes us within or close to \mathcal{E} .

Consequently, we denote the initial-marked PFSA induced by a QSP \mathcal{H} , with the initial marking removed, as $\mathcal{P}_{\mathcal{H}}$, and refer to it simply as a "PFSA" (dropping the qualifier "initial-marked"). States in $\mathcal{P}_{\mathcal{H}}$ are representable as states in $\mathcal{C}_{\mathcal{H}}$ as elements of \mathcal{E} . Next we establish a key result: we always encounter a state arbitrarily close to some element in \mathcal{E} in the canonical construction starting from the stationary distribution φ_λ .

Theorem 1 (ϵ -Synchronization of Probabilistic Automata). *For any QSP \mathcal{H} over Σ , the PFSA $\mathcal{P}_{\mathcal{H}}$ satisfies:*

$$\forall \epsilon' > 0, \exists x \in \Sigma^*, \exists \vartheta \in \mathcal{E}, \|\varphi_x - \vartheta\| \leq \epsilon' \quad (12)$$

where the norm used is unimportant.

Proof. We show that all PFSA are at least approximately synchronizable [3, 9], which is not true for deterministic automata. If the graph of $\mathcal{P}_{\mathcal{H}}$ (i.e., the deterministic automaton obtained by removing the arc probabilities) is synchronizable, then Eq. (12) trivially holds true for $\epsilon' = 0$ for any synchronizing string x . Thus, we assume the graph of $\mathcal{P}_{\mathcal{H}}$ to be non-synchronizable. From definition of non-synchronizability, it follows:

$$\forall q_i, q_j \in Q, \text{ with } q_i \neq q_j, \forall x \in \Sigma^*, \delta q_i, x \neq \delta q_j, x \quad (13)$$

If the PFSA has a single state, then every string satisfies the condition in Eq. (12). Hence, we assume that the PFSA has more than one state. Now if we have:

$$\forall x \in \Sigma^*, \frac{Prx'x}{Prx'} \neq \frac{Prx''x}{Prx''} \text{ where } x' \neq q_i, x'' \neq q_j \quad (14)$$

then, by the Definition 3, we have a contradiction $q_i \neq q_j$. Hence $\exists x_o$ such that

$$\frac{Prx'x_o}{Prx'} \neq \frac{Prx''x_o}{Prx''} \text{ where } x' \neq q_i, x'' \neq q_j \quad (15)$$

Since:

$$\sum_{x \in \Sigma^*} \frac{Prx'x}{Prx'} = 1, \text{ for any } x' \text{ where } x' \neq q_i \quad (16)$$

we conclude without loss of generality that $\forall q_i, q_j \in Q$, with $q_i \neq q_j$:

$$\exists x^{ij} \in \Sigma^*, \frac{Prx'x^{ij}}{Prx'} > \frac{Prx''x^{ij}}{Prx''} \text{ where } x' \neq q_i, x'' \neq q_j$$

It follows from induction that if we start with a distribution φ on Q such that $\varphi_i = \varphi_j = 0.5$, then for any $\epsilon' > 0$ we can construct a finite string x_o^{ij} such that if $\delta q_i, x_o^{ij} \rightarrow q_r, \delta q_j, x_o^{ij} \rightarrow q_s$, then for the new distribution φ' after execution of x_o^{ij} will satisfy $\varphi'_s > 1 - \epsilon'$. Recalling that \mathcal{P}_H is strongly connected, Now, for any $q_t \in Q$, there exists a string $y \in \Sigma^*$, such that $\delta q_s, y \rightarrow q_t$. Setting $x_*^{i,j \rightarrow t} = x_o^{ij}y$, we can ensure that the distribution φ'' obtained after execution of x_*^{ij} satisfies $\varphi''_t > 1 - \epsilon'$ for any q_t of our choice. For arbitrary initial distributions φ^A on Q , we consider contributions arising from simultaneously executing $x_*^{i,j \rightarrow t}$ from states other than just q_i and q_j . Nevertheless, it is easy to see that executing $x_*^{i,j \rightarrow t}$ implies that in the new distribution $\varphi^{A'}$, we have $\varphi_t^{A'} > \varphi_i^A \varphi_j^A - \epsilon'$. It immediately follows that executing of the string $x^{1,2 \rightarrow |Q|}x^{3,4 \rightarrow |Q|} \dots x^{n-1,n \rightarrow |Q|}$, where

$$n \begin{cases} |Q| & \text{if } |Q| \text{ is even} \\ |Q| - 1 & \text{otherwise} \end{cases} \quad (17)$$

would result in a final distribution $\varphi^{A''}$ which satisfies $\varphi_{|Q|}^{A''} > 1 - \frac{1}{2}n\epsilon'$. Appropriate scaling of ϵ' then completes the proof. \square

Theorem 1 induces the notion of ϵ -synchronizing strings, and guarantees their existence for arbitrary PFSA.

Definition 7 (ϵ -synchronizing Strings). *An ϵ -synchronizing string $x \in \Sigma^*$ for a PFSA is one that satisfies:*

$$\exists \vartheta \in \mathcal{E}, \|\varphi_x - \vartheta\| \leq \epsilon \quad (18)$$

The norm used is unimportant.

Theorem 1 does not yield an algorithm for computing synchronizing strings (See Theorem 3). It simply shows that one always exists. As a corollary, we estimate an asymptotic upper bound on the effort required to find it.

Corollary 2 (To Theorem 1). *At most $O(1/\epsilon)$ strings need to be analyzed to find an ϵ -synchronizing string.*

Proof. Theorem 1 works by multiplying entries from the $\tilde{\Pi}$ matrix, which cannot be all identical (otherwise the states would collapse). Let the minimum difference between two unequal entries be η . Then, following the construction in Theorem 1, the length ℓ of the synchronizing string, up to linear scaling, satisfies: $\eta^\ell \leq O\epsilon$, implying $\ell \leq O \log_{1/\eta} \epsilon^{-1}$. The number of strings to be analyzed therefore is at most all strings of length ℓ , which is given by

$$|\Sigma|^\ell \cdot |\Sigma|^{O \log_{1/\eta} \epsilon^{-1}} \leq O(1/\epsilon) \quad (19)$$

\square

Next, we describe the basic principle of our inference algorithm. PFSA states are not observable; we observe symbols generated from hidden states. This leads us to the notion of *symbolic derivatives*, which are computable from observations.

We denote the set of probability distributions over a set of cardinality k as \mathcal{D}^k . First, we specify a count function.

Definition 8 (Symbolic Count Function). *For a string s over Σ , the count function $\#^s : \Sigma^* \rightarrow \mathbb{N} \cup \{0\}$, counts the number of times a particular substring occurs in s . The count is overlapping, i.e., in a string $s = 0001$, we count the number of occurrences of 00 as 0001 and 0001, implying $\#^s 00 = 2$.*

Definition 9 (Symbolic Derivative). *For a string s generated by a QSP over Σ , the symbolic derivative is a function $\phi^s : \Sigma^* \rightarrow \mathcal{D}^{|\Sigma| - 1}$ as:*

$$\phi^s x|_i = \frac{\#^s x \sigma_i}{\sum_{\sigma_i \in \Sigma} \#^s x \sigma_i} \quad (20)$$

Thus, $\forall x \in \Sigma^*, \phi^s x$ is a probability distribution over Σ . $\phi^s x$ is referred to as the symbolic derivative at x .

For $q_i \in Q$, $\tilde{\pi}$ induces a distribution over Σ as $\tilde{\pi}q_i, \sigma_1, \dots, \tilde{\pi}q_i, \sigma_{|\Sigma|}$. We denote this as $\tilde{\pi}q_i, \cdot$. We show that the symbolic derivative at x can be used to estimate this distribution for $q_i \mid x$, provided x is ϵ -synchronizing.

Theorem 2 (ϵ -Convergence). *If $x \in \Sigma^*$ is ϵ -synchronizing, then:*

$$\forall \epsilon > 0, \lim_{|s| \rightarrow} \|\phi^s x - \tilde{\pi}x, \cdot\| \leq_{a.s.} \epsilon \quad (21)$$

Proof. The proof follows from the Glivenko-Cantelli theorem [15] on uniform convergence of empirical distributions. Since x is ϵ -synchronizing, we have:

$$\forall \epsilon > 0, \exists \vartheta \in \mathcal{E}, \|\rho_x - \vartheta\| \leq \epsilon \quad (22)$$

Let x ϵ -synchronize to q . Thus, every time we encounter x while reading s , we are guaranteed to be distributed over Q , and at most ϵ distance from the element of \mathcal{E} corresponding to q . Assuming that we encounter $n|s|$ occurrences of x within s , we note that $\phi^s x$ is an approximate empirical distribution for $\tilde{\pi}q, \cdot$. Denoting $F_{n|s|}$ as the perfect empirical distributions (*i.e.* ones that would be obtained for $\epsilon = 0$), we have:

$$\begin{aligned} & \lim_{|s| \rightarrow} \|\phi^s x - \tilde{\pi}q, \cdot\| \\ & \quad \lim_{|s| \rightarrow} \|\phi^s x - F_{n|s|} \circ F_{n|s|} - \tilde{\pi}q, \cdot\| \\ & \quad \stackrel{\text{a.s. } 0 \text{ by Glivenko-Cantelli}}{\leq} \lim_{|s| \rightarrow} \|\phi^s x - F_{n|s|}\| \overbrace{\lim_{|s| \rightarrow} \|F_{n|s|} - \tilde{\pi}q, \cdot\|}^{\text{a.s. } 0 \text{ by Glivenko-Cantelli}} \leq_{a.s.} \epsilon \end{aligned}$$

We use $n|s| \rightarrow$ as $|s| \rightarrow$, implied by the strong connectivity of our PFSA. \square

Next we describe identification of ϵ -synchronizing strings given a sufficiently long observed string s . Theorem 1 guarantees existence, and Corollary 2 establishes that O_1/ϵ substrings need to be analyzed till we encounter an ϵ -synchronizing string. These do not provide an executable algorithm, which arises from an inspection of the geometric structure of the set of probability vectors over Σ , obtained by constructing $\phi^s x$ for different choices of the candidate string x .

Definition 10 (Derivative Heap). *Given a string s generated by a QSP, a derivative heap $\mathcal{D}^s : 2^{\Sigma^*} \rightarrow \mathcal{D}^{|\Sigma|-1}$ is the set of probability distributions over Σ calculated for a given subset of strings $L \subset \Sigma^*$ as follows:*

$$\mathcal{D}^s L = \{\phi^s x \mid x \in L \subset \Sigma^*\} \quad (23)$$

Lemma 3 (Limiting Geometry). *Let $\mathcal{D} = \lim_{|s| \rightarrow} \lim_{L \rightarrow \Sigma^*} \mathcal{D}^s L$, and \mathcal{U} be the convex hull of \mathcal{D} . If u is a vertex of \mathcal{U} , then*

$$\exists q \in Q, \text{ such that } u = \tilde{\pi}q, \cdot \quad (24)$$

Proof. Recalling Theorem 2, the result follows from noting that any element of \mathcal{D} is a convex combination of elements from the set $\{\tilde{\pi}q_1, \cdot, \dots, \tilde{\pi}q_{|\mathcal{Q}|}, \cdot\}$. \square

Lemma 3 does not claim that the number of vertices of the convex hull of \mathbb{D} equals the number of states, but that every vertex corresponds to a state. We cannot generate \mathcal{D} in practice since we have a finite observed string s , and we can only calculate $\phi^s x$ for a finite number of x . Instead, we show that choosing a string corresponding to the vertex of the convex hull of the heap, constructed by considering O_1/ϵ strings, gives us an ϵ -synchronizing string with high probability.

Theorem 3 (Derivative Heap Approximation). *For s generated by a QSP, let $\mathcal{D}^s L$ be computed with $L \subset \Sigma^{O(\log_1/\epsilon)}$. If for some $x_0 \in \Sigma^{O(\log_1/\epsilon)}$, $\phi^s x_0$ is a vertex of the convex hull of $\mathcal{D}^s L$, then*

$$\text{Prob } x_0 \text{ is not } \epsilon\text{-synchronizing} \leq e^{-|s|\epsilon O_1} \quad (25)$$

Proof. The result follows from Sanov's Theorem [6] for convex set of probability distributions. Note that if $|s| \rightarrow \infty$, then x_o is guaranteed to be ϵ -synchronizing (Theorem 1, and Corollary 2). Denoting the number of times we encounter x_o in s as $n|s|$, and since \mathcal{D} is a convex set of distributions (allowing us to drop the polynomial factor in Sanov's bound), we apply Sanov's Theorem to the case of finite s :

$$\text{Prob}(KL(\phi^s x_o || \rho_{x_o} \tilde{\Pi}) > \epsilon) \leq e^{-n|s|\epsilon} \quad (26)$$

where $KL \cdot \cdot$ denotes the Kullback-Leibler divergence [12]. From the bound [16]:

$$\frac{1}{4} \|\phi^s x_o - \rho_{x_o} \tilde{\Pi}\|^2 \leq KL(\phi^s x_o || \rho_{x_o} \tilde{\Pi}) \quad (27)$$

and $n|s| \rightarrow |s|\alpha / |s|O_1$, where $\alpha > 0$ is the stationary probability of encountering x_o in s , we conclude:

$$\text{Prob}(\|\phi^s x_o - \rho_{x_o} \tilde{\Pi}\| > \epsilon) \leq 2e^{-\frac{1}{2}|s|\epsilon O_1} e^{-|s|\epsilon O_1} \quad (28)$$

□

Figure ?? illustrates PFSAs, derivative heaps, and ϵ -synchronizing strings. Next, we present our inference algorithm. Next, we use the preceding theoretical development to construct an effective procedure to infer $PFSA \mathcal{P}_H$ from a sufficiently long run from a QSP H , and a pre-specified $\epsilon > 0$.

2 Algorithm GenESeSS

2.1 Implementation Steps

We call our algorithm “Generator Extraction Using Self-similar Semantics”, or GenESeSS which for an observed sequence s , consists of three steps:

1. *Identification of ϵ -synchronizing string x_o :* Construct a derivative heap $\mathcal{D}^s L$ using the observed trace s (Definition 10), and set L consisting of all strings up to a sufficiently large, but finite, depth. We suggest as initial choice of L as $\log_{|\Sigma|} 1/\epsilon$. In L is sufficiently large, then the inferred model structure will not change for larger values. We then identify a vertex of the convex hull for \mathcal{D} , via any standard algorithm for computing the hull [2]. Choose x_o as the string mapping to this vertex.

2. *Identification of the structure of \mathcal{P}_H , i.e., transition function δ :* We generate δ as follows: For each state q , we associate a string identifier $x_q^{id} \in x_o \Sigma^*$, and a probability distribution h_q on Σ (which is an approximation of the $\tilde{\Pi}$ -row corresponding to state q). We extend the structure recursively:

- (a) Initialize the set Q as $Q = \{q_o\}$, and set $x_{q_o}^{id} = x_o$, $h_{q_o} = \phi^s x_o$.
- (b) For each state $q \in Q$, compute for each symbol $\sigma \in \Sigma$, find symbolic derivative $\phi^s x_q^{id} \sigma$. If $\|\phi^s x_q^{id} \sigma - h_{q'} \sigma\| \leq \epsilon$ for some $q' \in Q$, then define $\delta q, \sigma \rightarrow q'$. If, on the other hand, no such q' can be found in Q , then add a new state q' to Q , and define $x_{q'}^{id} = x_q^{id} \sigma$, $h_{q'} = \phi^s x_q^{id} \sigma$.

The process terminates when every $q \in Q$ has a target state, for each $\sigma \in \Sigma$. Then, if necessary, we ensure strong connectivity using [14].

3. *Identification of arc probabilities, i.e., function $\tilde{\Pi}$:*

- (a) Choose an arbitrary initial state $q \in Q$.
- (b) Run sequence s through the identified graph, as directed by δ , i.e., if current state is q , and the next symbol read from s is σ , then move to $\delta q, \sigma$. Count arc traversals, i.e, generate numbers N_j^i where $q_i \xrightarrow[N_j^i]{\sigma_j} q_k$.
- (c) Generate $\tilde{\Pi}$ by row normalization, i.e., $\tilde{\Pi}_{ij} = N_j^i / \sum_j N_j^i$

[7, 4] use similar recursive structure extension. However, with no notion of ϵ -synchronization, they are restricted to inferring only synchronizable or short-memory models, or large approximations for long-memory ones (See Figure ??).

2.2 Complexity Analysis & PAC Learnability

While h_q (in step 2) approximates \tilde{N} rows, we find the arc probabilities via normalization of traversal count. h_q only uses sequences in $x_0 \Sigma^*$, while traversal counting uses the entire sequence s , and is more accurate.

GenESeSS has no upper bound on the number of states; which is a function of the complexity of the process itself.

Theorem 4 (Time Complexity). *Asymptotic time complexity of GenESeSS is:*

$$\mathcal{T} = O((1/\epsilon |Q|) \times |s| \times |\Sigma|) \quad (29)$$

Proof. GenESeSS performs the following computations:

- C1 Computation of a derivative heap by computing $\phi^s x$ for $O(1/\epsilon)$ strings (Corollary 2), each of which involves reading the input s and normalization to distributions over Σ , thus contributing $O(1/\epsilon) \times |s| \times |\Sigma|$.
- C2 Finding a vertex of the convex hull of the heap, which, at worst, involves inspecting $O(\epsilon)$ points (encoded by strings generating the heap), contributing $O(1/\epsilon) \times |\Sigma|$, where each inspection is done in $O(|\Sigma|)$ time.
- C3 Finding δ , involving computing derivatives at string-identifiers (Step 2), thus contributing $O(|Q|) \times |s| \times |\Sigma|$. Strong connectivity [14] requires $O(|Q|) \times |\Sigma|$, and hence is unimportant.
- C4 Identification of arc probabilities using traversal counts and normalization, done in time linear in the number of arcs, i.e $O(|Q|) \times |\Sigma|$.

Summing the contributions, and using $|s| > |\Sigma|$,

$$\begin{aligned} \mathcal{T} &= O(1/\epsilon \times |s| \times |\Sigma|) + O(1/\epsilon \times |\Sigma|) + O(|Q| \times |s| \times |\Sigma|) + O(|Q| \times |\Sigma|) \\ &= O((1/\epsilon |Q|) \times |s| \times |\Sigma|) \end{aligned}$$

□

Theorem 4 shows that GenESeSS is polynomial in $O(1/\epsilon)$, size of input s , model size $|Q|$, and alphabet size $|\Sigma|$. In practice, $|Q| \ll 1/\epsilon$, implying that

$$\mathcal{T} = O\left(\frac{|s||\Sigma|}{\epsilon}\right) \quad (30)$$

An identification method is said to identify a target language L_* in the Probably Approximately Correct (PAC) sense [17, 1, 10], if it always halts and outputs L such that:

$$\exists \epsilon, \delta > 0, PdL_*, L \leq \epsilon \geq 1 - \delta \quad (31)$$

where $d(\cdot, \cdot)$ is a metric on the space of target languages. A class of languages is efficiently PAC-learnable if there exists an algorithm that PAC-identifies every language in the class, and runs in time polynomial in $1/\epsilon$, $1/\delta$, length of sample input, and inferred model size. We prove PAC-learnability of QSPs, by first establishing a metric on the space of probabilistic automata over Σ .

2.3 PAC Identifiability Of QSPs

Lemma 4 (Metric For Probabilistic Automata). *For two strongly connected PFSAs G_1, G_2 , denote the symbolic derivative at $x \in \Sigma^*$ as $\phi_{G_1}^s x$ and $\phi_{G_2}^s x$ respectively. Then,*

$$\Theta G_1, G_2 \sup_{x \in \Sigma^*} \left\{ Jx \lim_{|s_1| \rightarrow \infty} \lim_{|s_2| \rightarrow \infty} \|\phi_{G_1}^{s_1} x - \phi_{G_2}^{s_2} x\| \right\}$$

defines a metric on the space of probabilistic automata on Σ .

Proof. Non-negativity and symmetry follows immediately. Triangular inequality follows from noting that $\|\phi_{G_1}^{s_1} x - \phi_{G_2}^{s_2} x\|$ is upper bounded by 1, and therefore for any chosen order of the strings in Σ^* , we have two ℓ sequences, which would satisfy the triangular inequality under the sup norm. The metric is well-defined since for any sufficiently long s_1, s_2 , the symbolic derivatives at arbitrary x are uniformly convergent to some linear combination of the rows of the corresponding \tilde{H} matrices. \square

Theorem 5 (PAC-Learnability). *QSPs for which the probabilistic Nerode equivalence has a finite index is PAC-learnable using PFSAs, i.e., for $\epsilon, \eta > 0$, and for every sufficiently long sequence s generated by QSP \mathcal{H} , we can compute $\mathcal{P}'_{\mathcal{H}}$ as an estimate for $\mathcal{P}_{\mathcal{H}}$ such that:*

$$\text{Prob}(\Theta \mathcal{P}_{\mathcal{H}}, \mathcal{P}'_{\mathcal{H}} \leq \epsilon) \geq 1 - \eta \quad (32)$$

The algorithm runs in time polynomial in $1/\epsilon, 1/\eta$, input length $|s|$ and model size.

Proof. GenESeSS construction implies that, once the initial ϵ -synchronizing string x_0 is identified, there is no scope of the model error to be more than ϵ . Hence:

$$\begin{aligned} \text{Prob}(\Theta \mathcal{P}_{\mathcal{H}}, \mathcal{P}'_{\mathcal{H}} \leq \epsilon) &\geq 1 - \text{Prob}(\|\phi^s x_0 - \wp_{x_0} \tilde{H}\| > \epsilon) \\ &\Rightarrow \text{Prob}(\Theta \mathcal{P}_{\mathcal{H}}, \mathcal{P}'_{\mathcal{H}} \leq \epsilon) \geq 1 - e^{-|s|\epsilon O_1} \quad (\text{Using Eq. (28)}) \end{aligned}$$

Thus, for any $\eta > 0$, if we have $|s| = O(1/\epsilon \log 1/\eta)$, then the required condition of Eq. (32) is met. The polynomial runtimes are established in Theorem 4. \square

2.3.1 Remark On Kearns' Hardness Result

We are immune to Kearns' hardness result [11], since $\epsilon > 0$ enforces state distinguishability [13].

3 Summary & Conclusion

We establish the notion of causal generators for stationary ergodic quantized stochastic processes on a rigorous measure-theoretic foundation, and propose a new inference algorithm GenESeSS for identification of probabilistic automata models from sufficiently long traces. We show that our approach can learn processes with long range dependencies, which yield non-synchronizable automata. Additionally, we establish that GenESeSS is computationally efficient in the PAC sense. The results are validated on synthetic and real datasets. Future research will investigate the prediction problem in greater detail in the context of new applications.

References

- [1] Dana Angluin. Computational learning theory: survey and selected bibliography. In *Proceedings of the twenty-fourth annual ACM symposium on Theory of computing*, STOC '92, pages 351–369, New York, NY, USA, 1992. ACM.
- [2] C. Bradford Barber, David P. Dobkin, and Hannu Huhdanpaa. The quickhull algorithm for convex hulls. *ACM Trans. Math. Softw.*, 22(4):469–483, dec 1996.
- [3] S. Bogdanovic, B. Imreh, M. Ceric, and T. Petkovic. Directable automata and their generalizations - a survey. *Novi Sad Journal of Mathematics*, 29(2):31–74, 1999.
- [4] Jorge Castro and Ricard Gavaldà. Towards feasible pac-learning of probabilistic deterministic finite automata. In Alexander Clark, François Coste, and Laurent Miclet, editors, *ICGI*, volume 5278 of *Lecture Notes in Computer Science*, pages 163–174. Springer, 2008.
- [5] I. Chattopadhyay and A. Ray. Structural transformations of probabilistic finite state machines. *International Journal of Control*, 81(5):820–835, May 2008.

- [6] I. Csiszár. Sanov property, generalized I-projection and a conditional limit theorem. *Ann. Probab.*, 12:768–793, 1984.
- [7] Ricard Gavaldà, Philipp W. Keller, Joelle Pineau, and Doina Precup. Pac-learning of markov models with hidden state. In Johannes Fürnkranz, Tobias Scheffer, and Myra Spiliopoulou, editors, *ECML*, volume 4212 of *Lecture Notes in Computer Science*, pages 150–161. Springer, 2006.
- [8] J. E. Hopcroft, R. Motwani, and J. D. Ullman. *Introduction to Automata Theory, Languages, and Computation*, 2nd ed. Addison-Wesley, 2001.
- [9] M. Ito and Jürgen Duske. On cofinal and definite automata. *Acta Cybern.*, 6:181–189, 1984.
- [10] M. J. Kearns and U. V. Vazirani. *An Introduction to Computational Learning Theory*. MIT Press, Cambridge, MA, USA, 1994.
- [11] Michael Kearns, Yishay Mansour, Dana Ron, Ronitt Rubinfeld, Robert E. Schapire, and Linda Sellie. On the learnability of discrete distributions. In *Proceedings of the twenty-sixth annual ACM symposium on Theory of computing*, STOC ’94, pages 273–282, New York, NY, USA, 1994. ACM.
- [12] E. L. Lehmann and Joseph P. Romano. *Testing statistical hypotheses*. Springer Texts in Statistics. Springer, New York, third edition, 2005.
- [13] Dana Ron, Yoram Singer, and Naftali Tishby. Learning probabilistic automata with variable memory length. In *Computational Learing Theory*, pages 35–46, 1994.
- [14] Robert Tarjan. Depth-first search and linear graph algorithms. *SIAM Journal on Computing*, 1(2):146–160, 1972.
- [15] Flemming Topsøe. On the glivenko-cantelli theorem. *Probability Theory and Related Fields*, 14:239–250, 1970. 10.1007/BF01111419.
- [16] A.B. Tsybakov. *Introduction to nonparametric estimation*. Springer series in statistics. Springer, 2009.
- [17] L. G. Valiant. A theory of the learnable. *Commun. ACM*, 27:1134–1142, November 1984.