

A Non-Parametric Bayesian Approach Towards Online Sequence Learning

Stiven S. Dias
Embraer S.A.
R&D Team
Av. Brig. Faria Lima, 2170
S. J. dos Campos, Brazil

stiven.dias@embraer.com.br

Marcelo G. S. Bruno
Instituto Tecnológico de
Aeronáutica
Pç. Mar. Eduardo Gomes, 50
S. J. dos Campos, Brazil

bruno@ita.br

Alberto F. De Souza
Universidade Federal do
Espírito Santo
Av. Fernando Ferrari, 514
Vitória, Brazil

alberto@inf.ufes.br

ABSTRACT

Forecasting the next sample of an incoming sequence of noisy observations is a recurring problem in several real-world applications. In general, the streamed observations may originate from a possibly non-stationary stochastic process with non-linear dynamics. Very often though, to avoid devising explicit probabilistic models from scratch, we alternatively learn data-driven surrogate models from past observations to approximate complex dynamics. Time-varying systems in turn pose additional hurdles as such models must be continuously updated to accommodate changes in their underlying behavior. In this paper, we introduce a novel grid-based filter designed to inherently deal with the continuous learning problem. Specifically, Bayesian-grounded procedures are employed to both recursively predict observations and incrementally build a suitable non-parametric, probabilistic model of the indirectly observed phenomenon. Moreover, we also present a one-shot learning, memory-based implementation of the proposed filter, a shallow weightless neural network which is able to efficiently store and retrieve associative input-output pairs. Lastly, we demonstrate the effectiveness of our filter in performing simple anomaly detection tasks.

1. INTRODUCTION

Time-series modeling is paramount for many applications such as anomaly detection and prediction [1], stock market forecast [31] and even network attack detection [8]. Moreover, Bayesian inference – e.g. through the Bayes filter – provides a well-grounded theoretical framework relying on probabilistic models to forecast incoming observations of a stochastic phenomenon. However, in several engineering applications, designing an explicit probabilistic model of the time-series by hand is inadvisable as the data distribution can be prohibitively complex and often drifts, changing significantly over time. The ability to continuously learn suitable probabilistic models directly from available observations without supervision thus has the potential to further unlock the application of the Bayesian paradigm to time-series forecast.

1.1 Related Work

Online sequence learning narrows down the range of suitable methods to timely predict future observations, since they are also required to simultaneously update their inner models of

the temporal phenomenon using accumulated observations. Several methods exist for time-series prediction using different approaches [23, 25]. Few of them like [1, 30] are able to deal with the continuous learning (CL) problem, i.e. to build and improve a feasible model to accurately predict future observations as incoming observations arrive. Even fewer approaches though rely on non-parametric statistics [32, 28] to simultaneously learn models and predict observations [7, 9], but, to the best of our knowledge, none of them have adapted the Bayes filter to recursively learn model’s hyperparameters and jointly make predictions using a fully non-parametric, Bayesian inference approach as discussed in this paper.

Furthermore, shallow weightless neural networks (WNNs) [2] have been successfully employed to learn from data and solve challenging problems [13, 14, 12] as far as proper data encoding schemes [27] are employed. Despite their inherent ability for both one-shot learning and CL, so far, their few enthusiasts have not been able to effectively stack WNN layers into deeper networks to solve increasingly complex problems, as was the case for regular neural networks. On the other hand, although quite useful to successfully solve real-world, time-series forecasting problems [10, 3], deep neural models were not intrinsically designed for CL, as they often require an offline, expensive optimization procedure [5] with salting racks and workarounds [6] to timely / properly converge and adjust their multitude of parameters – e.g. weights and biases – using all accumulated training data so far; avoiding thus catastrophic forgetting issues [22].

An alternative, bio-inspired model which supports CL by design has been incrementally polished by Numenta researchers [19, 20]. However, despite their insightful findings [18] on how the underlying biological mechanisms – in special, biological neurons and layers within cortical cells – of the human neocortex interact with each other to create human intelligence as an emerging behavior, they did not provide so far a rigorous interpretation of their well-known hierarchical temporal memory (HTM) model [19] using e.g. the Bayesian inference framework.

In this paper, we model first the time-series CL and prediction problem using a fully Bayesian approach and propose therefore an explainable non-parametric, grid-based filter accompanied by an efficient, one-shot learning implementation based on a shallow WNN with two layers stacked in a meaningful, Bayesian interpretable way. Nevertheless, we also borrow some ideas from Cui et al. [11] and employ sparse representations over high-dimensional spaces to represent observed symbols.

1.2 Contributions and Paper Outline

The main contribution of this paper is twofold: it introduces i) a novel grid-based filter which allows one to perform inference on high-dimensional binary spaces and effectively predict incoming observations residing in continuous-valued spaces; and ii) an efficient grid-filter implementation based on virtual-generalized random access memory (VG-RAM) nodes [24] which enables time-series CL seamlessly and yields (auditable) outputs which are fully traceable to filter’s inputs. *Paper Outline:* The paper is divided into six sections. Sec. 1 is this Introduction. In Sec. 2, we state the online sequence learning and prediction problem under a Bayesian framework. In Sec. 3, we introduce our solution approach. Sec. 4 describes our associative random access memory (RAM)-based implementation of the solution proposed in Sec. 3. Simulations results for an application example are presented in Sec. 5. Finally, we offer our conclusions in Sec. 6.

1.3 Notation

We use uppercase letters, e.g. \mathbf{X} , to denote random vectors and lowercase letters, e.g. \mathbf{x} , to denote samples of random vectors. We also define the unconditional probability of an event A and the conditional probability of event A given the knowledge that an event B has occurred as $\Pr\{A\}$ and $\Pr\{A|B\}$, respectively. Moreover, we define the joint probability of events A and B occurring simultaneously as $\Pr\{A, B\}$.

We also employ the Iverson bracket – which takes as argument any logical proposition Λ – defined as $\llbracket \Lambda \rrbracket = 1$, if Λ is true, and $\llbracket \Lambda \rrbracket = 0$, otherwise. Furthermore, let

$$\mathbb{Z}_2^d \triangleq \{0, 1\} \times \dots \times \{0, 1\}$$

denote an d -dimensional binary space, we define then the Hamming distance between two vectors $\mathbf{a} \triangleq [a_1 \dots a_d]^\top$ and $\mathbf{b} \triangleq [b_1 \dots b_d]^\top$ in \mathbb{Z}_2^d as $d_H(\mathbf{a}, \mathbf{b}) \triangleq \sum_{i=1}^d \llbracket a_i \neq b_i \rrbracket$. The inner product between any two vectors \mathbf{a} and \mathbf{b} now in \mathbb{R}^d is defined as $\langle \mathbf{a}, \mathbf{b} \rangle \triangleq \sum_{i=1}^d a_i \cdot b_i$ and the L^p norm of an arbitrary vector $\mathbf{c} \triangleq [c_1 \dots c_d]^\top \in \mathbb{R}^d$ is defined as $\|\mathbf{c}\|_p \triangleq \left(\sum_{i=1}^d |c_i|^p \right)^{\frac{1}{p}}$. Lastly, we define the one-hot indicator vector as $\mathbf{1}_j \triangleq [\iota_{1,j} \dots \iota_{d,j}]^\top \in \mathbb{Z}_2^d$ with $\|\mathbf{1}_j\|_1 = 1$ such that $\iota_{i,j} = \llbracket i = j \rrbracket, \forall i \in \{1, \dots, d\}$.

Now, let \mathbf{X} be a d -dimensional random vector whose realizations $\mathbf{x} \sim \mathbf{X}$ assume values in a finite sampling space $\Omega \triangleq \{\boldsymbol{\mu}^{(1)}, \dots, \boldsymbol{\mu}^{(c)}\} \subseteq \mathbb{R}^d$ with cardinality $c = |\Omega|$. We denote a categorical distribution over Ω as $Cat(\Omega; \boldsymbol{\pi})$, in which the c -dimensional parameter vector $\boldsymbol{\pi} \triangleq [\pi_1 \dots \pi_c]^\top \in \mathbb{R}_+^c$, $\|\boldsymbol{\pi}\|_1 = 1$, collects the probabilities of the c possible categories such that its probability mass function (p.m.f.) is given by $P(\mathbf{x}) \triangleq \Pr\{\mathbf{X} = \mathbf{x}\} = \sum_{i=1}^c \pi_i \cdot \llbracket \mathbf{x} = \boldsymbol{\mu}^{(i)} \rrbracket$.

Moreover, let us further assume that $\boldsymbol{\pi}$, a valid p.m.f. over Ω , is a realization of a random vector $\boldsymbol{\Pi}$ in the $(c-1)$ -dimensional simplex \mathbb{S}_c distributed according to a Dirichlet distribution $Dir(c; \boldsymbol{\alpha})$ with hyperparameters collected by the vector $\boldsymbol{\alpha} \triangleq [\alpha_1 \dots \alpha_c]^\top$, $\alpha_i \in \mathbb{R}_+, \forall i \in \{1, \dots, c\}$. We define the probability density function (p.d.f.) of the Dirichlet distribution over \mathbb{S}_c as $p(\boldsymbol{\pi}) = (B(\boldsymbol{\alpha}))^{-1} \prod_{i=1}^c \pi_i^{\alpha_i - 1}$, where $B(\boldsymbol{\alpha}) \triangleq (\Gamma(\alpha_0))^{-1} \prod_{i=1}^c \Gamma(\alpha_i)$ is the Beta function with $\alpha_0 = \sum_{i=1}^c \alpha_i$ and $\Gamma(x) \triangleq \int_0^\infty t^{x-1} e^{-t} dt$ is the Gamma function with x in the set of non-negative real numbers \mathbb{R}_+ .

2. PROBLEM STATEMENT

Let $\check{\mathbf{y}}_{0:\infty} \triangleq (\check{\mathbf{y}}_0, \check{\mathbf{y}}_1, \check{\mathbf{y}}_2, \dots)$ be an *online* sequence of observations – e.g. a streamed sequence of dense feature vectors – in a continuous space $\Omega_y \subset \mathbb{R}^{d_y}$ at discrete time instants $t_0 < t_1 < t_2 < \dots$. Given the sub-sequence of observations $\check{\mathbf{y}}_{0:n} \triangleq (\check{\mathbf{y}}_0, \dots, \check{\mathbf{y}}_n)$ up to instant n , our ultimate goal is to predict, at instant n , the next observation $\check{\mathbf{y}}_{n+1} \triangleq [\check{y}_{n+1,1} \dots \check{y}_{n+1,d_y}]^\top \in \Omega_y$.

2.1 Bayesian formulation

Now, let $\mathbf{m}_{y \rightarrow z} : \Omega_y \mapsto \Omega_z$ denote a *known* encoding function which maps each observation $\check{\mathbf{y}}_n$ into a sparse representation $\check{\mathbf{z}}_n \triangleq [\check{z}_{n,1} \dots \check{z}_{n,d_z}]^\top$ within a high-dimensional binary space $\Omega_z \subset \mathbb{Z}_2^{d_z}$ such that $d_z \gg d_y$ and $d_z \gg \|\check{\mathbf{z}}_n\|_1$. Similarly, let $\mathbf{m}_{z \rightarrow y} : \Omega_z \mapsto \Omega_y$ define a *known* decoding function mapping arbitrary encoded observations $\mathbf{z}_n \in \Omega_z$ back to observations $\mathbf{y}_n \in \Omega_y$.

As illustrated in Fig. 1, let us further suppose that any particular sequence of encoded observations $\check{\mathbf{z}}_{0:n} \triangleq (\check{\mathbf{z}}_0, \dots, \check{\mathbf{z}}_n)$ up to instant n is emitted by a Markovian process with an underlying sequence of hidden states $\mathbf{s}_{0:n} \triangleq (\mathbf{s}_0, \dots, \mathbf{s}_n)$ within a high-dimensional binary space $\Omega_s \subset \mathbb{Z}_2^{d_s}, \forall n \geq 0$, $d_s \gg \|\mathbf{s}_n\|_1$, such that the online sequence of output observations can be written as

$$\check{\mathbf{y}}_{0:n} \equiv (\mathbf{m}_{z \rightarrow y}(\check{\mathbf{z}}_0), \dots, \mathbf{m}_{z \rightarrow y}(\check{\mathbf{z}}_n)). \quad (1)$$

Note that, as shown in Fig. 1, the sequences $\check{\mathbf{y}}_{0:n}$, $\check{\mathbf{z}}_{0:n}$ and $\mathbf{s}_{0:n}$ are realizations of the sequences of random vectors $\mathbf{Y}_{0:n} \triangleq (\mathbf{Y}_0, \dots, \mathbf{Y}_n)$, $\mathbf{Z}_{0:n} \triangleq (\mathbf{Z}_0, \dots, \mathbf{Z}_n)$ and $\mathbf{S}_{0:n} \triangleq (\mathbf{S}_0, \dots, \mathbf{S}_n)$, respectively.

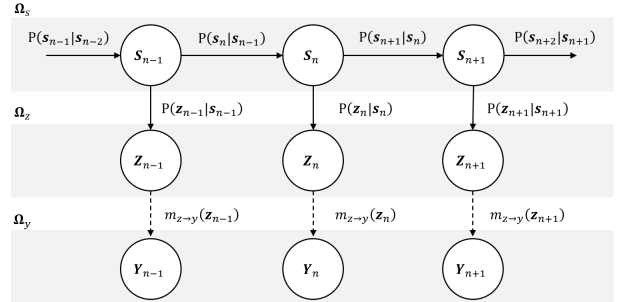


Figure 1: Hidden Markov chain representing the observed stochastic process.

Next, conditioned on the sequence of observations $\check{\mathbf{z}}_{0:n-1}$ up to instant $n-1$, the random state vector \mathbf{S}_n at instant n is distributed according to the conditional p.m.f.

$$P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n-1}) \triangleq \Pr\{\mathbf{S}_n = \mathbf{s}_n | \mathbf{Z}_0 = \check{\mathbf{z}}_0, \dots, \mathbf{Z}_{n-1} = \check{\mathbf{z}}_{n-1}\}. \quad (2)$$

Since, by construction, the observed process is first order Markovian and the observations are not correlated over time, the following usual assumptions hold respectively for the likelihood function and the state transition p.m.f.

$$P(\check{\mathbf{z}}_n | \mathbf{s}_{0:n}, \check{\mathbf{z}}_{0:n-1}) \equiv P(\check{\mathbf{z}}_n | \mathbf{s}_n) \quad (3)$$

$$P(\mathbf{s}_{n+1} | \mathbf{s}_{0:n}, \check{\mathbf{z}}_{0:n}) \equiv P(\mathbf{s}_{n+1} | \mathbf{s}_n). \quad (4)$$

Thus, upon the arrival of the encoded observation $\check{\mathbf{z}}_n = \mathbf{m}_{y \rightarrow z}(\check{\mathbf{y}}_n)$ at instant n , the predicted marginal posterior $P(\mathbf{s}_{n+1} | \check{\mathbf{z}}_{0:n})$ can be recursively computed from the predicted

marginal posterior $P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n-1})$ in (2) – built from the encoded observations $\check{\mathbf{z}}_{0:n-1}$ up to instant $n-1$ – as

$$P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n}) \propto P(\check{\mathbf{z}}_n | \mathbf{s}_n) P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n-1}) \quad (5)$$

$$P(\mathbf{s}_{n+1} | \check{\mathbf{z}}_{0:n}) = \sum_{\mathbf{s}_n \in \Omega_s} P(\mathbf{s}_{n+1} | \mathbf{s}_n) P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n}) \quad (6)$$

with proportionality constant in (5) given by the reciprocal of

$$P(\check{\mathbf{z}}_n | \check{\mathbf{z}}_{0:n-1}) = \sum_{\mathbf{s}_n \in \Omega_s} P(\check{\mathbf{z}}_n | \mathbf{s}_n) P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n-1}). \quad (7)$$

Eqs. (5) and (6) follow respectively by the straightforward application of the Bayes rule and the Chapman-Kolmogorov equation given assumptions (3) and (4).

Note also that, conditioned on the sequence of encoded observations $\check{\mathbf{z}}_{0:n}$ up to instant n , the random vector \mathbf{Z}_{n+1} at instant $n+1$ is distributed according to the p.m.f.

$$P(\mathbf{z}_{n+1} | \check{\mathbf{z}}_{0:n}) = \sum_{\mathbf{s}_{n+1} \in \Omega_s} P(\mathbf{z}_{n+1} | \mathbf{s}_{n+1}) P(\mathbf{s}_{n+1} | \check{\mathbf{z}}_{0:n}). \quad (8)$$

We can compute therefore the maximum a posteriori (MAP) estimate of the next encoded observation $\check{\mathbf{z}}_{n+1}$ given all encoded observations $\check{\mathbf{z}}_{0:n}$ up to instant n as

$$\hat{\mathbf{z}}_{n+1|n} = \arg \max_{\mathbf{z}_{n+1} \in \Omega_z} P(\mathbf{z}_{n+1} | \check{\mathbf{z}}_{0:n}) \quad (9)$$

and predict then the next output observation $\check{\mathbf{y}}_{n+1}$ at instant $n+1$ as

$$\hat{\mathbf{y}}_{n+1|n} \triangleq \mathbf{m}_{z \rightarrow y}(\hat{\mathbf{z}}_{n+1|n}). \quad (10)$$

In this paper, however, we consider that neither the observation model $P(\mathbf{z}_n | \mathbf{s}_n)$ nor the dynamic model $P(\mathbf{s}_n | \mathbf{s}_{n-1})$ is provided or known *a priori*. At first glance, this is an ill-defined problem to solve using the conventional, model-based Bayesian framework [21]. However, we introduce here a novel approach to incrementally build – as new, incoming observations are assimilated – non-parametric representations of $P(\mathbf{z}_n | \mathbf{s}_n)$ and $P(\mathbf{s}_n | \mathbf{s}_{n-1})$ which, as indicated in Eqs. (5) through (8), entail the required knowledge regarding the observed stochastic process to recursively compute the predicted posteriors $P(\mathbf{s}_{n+1} | \check{\mathbf{z}}_{0:n})$ and $P(\mathbf{z}_{n+1} | \check{\mathbf{z}}_{0:n})$ at instant n from the previous posterior $P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n-1})$ built at instant $n-1$.

2.2 Problem breakdown

Our ultimate goal can be finally unfolded into two tasks: i) *continuously* learn suitable, non-parametric probabilistic models for $P(\mathbf{s}_n | \mathbf{s}_{n-1})$ and $P(\mathbf{z}_n | \mathbf{s}_n)$ to be able to recursively approximate, at each time instant n , the marginal posterior $P(\mathbf{z}_{n+1} | \check{\mathbf{z}}_{0:n})$ from the available data $\check{\mathbf{z}}_{0:n}$ as indicated in Eqs. (5) through (8), and ii) use the predicted marginal posterior $P(\mathbf{z}_{n+1} | \check{\mathbf{z}}_{0:n})$ given the observations $\check{\mathbf{z}}_{0:n}$ up to instant n to predict the next output observation $\check{\mathbf{y}}_{n+1}$ at instant $n+1$ as shown in Eqs. (9) and (10).

3. SOLUTION APPROACH

We define first the required overall structure of a suitable encoding-decoding (codec) scheme. In the sequel, we derive a novel grid-based filter to perform the aforementioned tasks. In particular, we employ a non-parametric, Bayesian approach to recursively predict the next observation and learn the model hyperparameters.

3.1 Codec structure

Let $\Omega_z \triangleq \{\boldsymbol{\mu}_z^{(\ell)}\}$, $\ell \in \{1, \dots, c_z\}$, be a set of sparsely encoded, binary vectors $\boldsymbol{\mu}_z^{(\ell)} \in \mathbb{Z}_2^{d_z}$ with cardinality $c_z \triangleq |\Omega_z|$ such that

$$\begin{aligned} \forall \ell, \quad & \|\boldsymbol{\mu}_z^{(\ell)}\|_1 \leq a_z \\ \forall \ell_1 \neq \ell_2, \quad & b_z \leq d_H(\boldsymbol{\mu}_z^{(\ell_1)}, \boldsymbol{\mu}_z^{(\ell_2)}) \leq 2a_z, \end{aligned}$$

where $a_z \ll d_z$ is the maximum number of activated bits and b_z is the minimum Hamming distance between any two vectors in Ω_z . Thus, a_z controls the sparseness in Ω_z , whereas b_z determines the discriminability of its elements in terms of the Hamming distance. Moreover, let $\{\Omega_y^{(\ell)}\}$, $\ell \in \{1, \dots, c_z\}$, denote a set partition of the output observation space Ω_y such that

$$\bigcup_{\ell=1}^{c_z} \Omega_y^{(\ell)} = \Omega_y$$

and its subsets $\Omega_y^{(\ell)}$ are pairwise disjoint, i.e., $\forall \ell_1 \neq \ell_2$, $\Omega_y^{(\ell_1)} \cap \Omega_y^{(\ell_2)} = \emptyset$. We define then the following lossy, encoding function $\mathbf{m}_{y \rightarrow z} : \mathbb{R}^{d_y} \mapsto \{\boldsymbol{\mu}_z^{(\ell)}\} \subset \mathbb{Z}_2^{d_z}$

$$\begin{aligned} \mathbf{z} &= \mathbf{m}_{y \rightarrow z}(\mathbf{y}) \\ &\triangleq \sum_{\ell=1}^{c_z} \boldsymbol{\mu}_z^{(\ell)} \mathbb{I}[\mathbf{y} \in \Omega_y^{(\ell)}]. \end{aligned} \quad (11)$$

Note that $\{\boldsymbol{\mu}_z^{(\ell)}\}$, $\ell \in \{1, \dots, c_z\}$ is an encoded alphabet. In particular, $\forall \ell \in \{1, \dots, c_z\}$, $\boldsymbol{\mu}_z^{(\ell)}$ assigns a sparse, binary symbol to the set partition $\Omega_y^{(\ell)}$ with a_z activated bits and at least b_z bits far from the remaining $c_z - 1$ symbols in $\Omega_z \subset \mathbb{Z}_2^{d_z}$.

Note also that $\mathbf{m}_{y \rightarrow z}(\cdot)$ is by construction a deterministic, surjective function. We define therefore the following decoding function $\mathbf{m}_{z \rightarrow y} : \{\boldsymbol{\mu}_z^{(\ell)}\} \mapsto \mathbb{R}^{d_y}$

$$\begin{aligned} \mathbf{y} &= \mathbf{m}_{z \rightarrow y}(\mathbf{z}) \\ &\triangleq \sum_{\ell=1}^{c_z} \boldsymbol{\mu}_y^{(\ell)} \mathbb{I}[\mathbf{z} = \boldsymbol{\mu}_z^{(\ell)}], \end{aligned} \quad (12)$$

in which $\boldsymbol{\mu}_y^{(\ell)}$ can be e.g. the geometric centroid of the set partition $\Omega_y^{(\ell)}$. In this case, $\forall \ell \in \{1, \dots, c_z\}$, we can assume that $\Omega_y^{(\ell)}$ is convex by construction such that its centroid $\boldsymbol{\mu}_y^{(\ell)} \in \Omega_y^{(\ell)}$. Lastly, we follow the lead in [29] and further assume that close observations in Ω_y yield close symbols in Ω_z . In particular, if the Euclidean distance $\|\mathbf{y} - \mathbf{y}'\|_2 \leq \|\mathbf{y}'' - \mathbf{y}'''\|_2$ for $\mathbf{y}, \mathbf{y}', \mathbf{y}'', \mathbf{y}''' \in \Omega_y$, then the Hamming distance $d_H(\mathbf{m}_{y \rightarrow z}(\mathbf{y}), \mathbf{m}_{y \rightarrow z}(\mathbf{y}')) \leq d_H(\mathbf{m}_{y \rightarrow z}(\mathbf{y}''), \mathbf{m}_{y \rightarrow z}(\mathbf{y}'''))$.

3.2 Grid-based filter

Let the c_s -dimensional vector

$$\boldsymbol{\pi}_{n|n-1} \triangleq \left[\pi_{n|n-1}^{(1)} \quad \dots \quad \pi_{n|n-1}^{(\ell')} \quad \dots \quad \pi_{n|n-1}^{(c_s)} \right]^\top \quad (13)$$

collecting the conditional probabilities, $\forall \ell' \in \{1, \dots, c_s\}$,

$$\pi_{n|n-1}^{(\ell')} \triangleq Pr\{\mathbf{S}_n = \boldsymbol{\mu}_s^{(\ell')} | \mathbf{Z}_0 = \check{\mathbf{z}}_0, \dots, \mathbf{Z}_{n-1} = \check{\mathbf{z}}_{n-1}\}$$

represent the predicted marginal posterior $P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n-1})$. Similarly, let the c_z -dimensional vector

$$\tilde{\boldsymbol{\pi}}_{n+1|n} \triangleq \left[\tilde{\pi}_{n+1|n}^{(1)} \quad \dots \quad \tilde{\pi}_{n+1|n}^{(\ell)} \quad \dots \quad \tilde{\pi}_{n+1|n}^{(c_z)} \right]^\top \quad (14)$$

collecting the conditional probabilities, $\forall \ell \in \{1, \dots, c_z\}$,

$$\tilde{\pi}_{n+1|n}^{(\ell)} \triangleq Pr\{\mathbf{Z}_{n+1} = \boldsymbol{\mu}_z^{(\ell)} | \mathbf{Z}_0 = \check{\mathbf{z}}_0, \dots, \mathbf{Z}_{n-1} = \check{\mathbf{z}}_n\}$$

represent the predicted marginal posterior $P(\mathbf{z}_{n+1} | \check{\mathbf{z}}_{0:n})$. We can combine the update and prediction steps in Eqs. (5) and (6), respectively, into a single, one-step-ahead update rule as

$$P(\mathbf{s}_{n+1} | \check{\mathbf{z}}_{0:n}) \propto \sum_{\mathbf{s}_n \in \Omega_s} P(\mathbf{s}_{n+1} | \mathbf{s}_n) P(\check{\mathbf{z}}_n | \mathbf{s}_n) P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n-1}). \quad (15)$$

Thus, by assuming that $\check{\mathbf{z}}_n = \boldsymbol{\mu}_z^{(k)}$ for some $k \in \{1, \dots, c_z\}$, we can rewrite (15) in a compact notation using matrix-vector multiplication as

$$\boldsymbol{\pi}_{n+1|n} \propto \mathbf{F}_{n+1|n}^{(k)} \cdot \boldsymbol{\pi}_{n|n-1} \quad (16)$$

with proportionality constant such that $\|\boldsymbol{\pi}_{n+1|n}\|_1 = 1$. Moreover, the i -th element $f_{n+1|n}^{(k,j,i)}$ of the j -th row

$$\mathbf{f}_{n+1|n}^{(k,j)} \triangleq \begin{bmatrix} f_{n+1|n}^{(k,j,1)} & \dots & f_{n+1|n}^{(k,j,i)} & \dots & f_{n+1|n}^{(k,j,c_s)} \end{bmatrix}^\top$$

of the modified, $c_s \times c_s$ state transition matrix

$$\mathbf{F}_{n+1|n}^{(k)} \triangleq \begin{bmatrix} \mathbf{f}_{n+1|n}^{(k,1)} & \dots & \mathbf{f}_{n+1|n}^{(k,j)} & \dots & \mathbf{f}_{n+1|n}^{(k,c_s)} \end{bmatrix}^\top$$

– which depends on the observation $\check{\mathbf{z}}_n$ – is defined as

$$f_{n+1|n}^{(k,j,i)} \triangleq Pr\{\mathbf{S}_{n+1} = \boldsymbol{\mu}_s^{(j)} | \mathbf{S}_n = \boldsymbol{\mu}_s^{(i)}\} \\ \times Pr\{\mathbf{Z}_n = \boldsymbol{\mu}_z^{(k)} | \mathbf{S}_n = \boldsymbol{\mu}_s^{(i)}\}$$

such that $\tilde{\pi}_{n+1|n}^{(j)} \propto \langle \mathbf{f}_{n+1|n}^{(k,j)}, \boldsymbol{\pi}_{n|n-1} \rangle$, $\forall j \in \{1, \dots, c_s\}$.

Conversely, the predicted marginal posterior $P(\mathbf{z}_{n+1} | \check{\mathbf{z}}_{0:n})$ in (8), can be rewritten in compact notation as

$$\tilde{\boldsymbol{\pi}}_{n+1|n} \propto \mathbf{H}_{n+1} \cdot \boldsymbol{\pi}_{n+1|n} \quad (17)$$

with proportionality constant such that $\|\tilde{\boldsymbol{\pi}}_{n+1|n}\|_1 = 1$. Furthermore, the j -th element $h_{n+1}^{(\ell,j)}$ of the ℓ -th row

$$\mathbf{h}_{n+1}^{(\ell)} \triangleq \begin{bmatrix} h_{n+1}^{(\ell,1)} & \dots & h_{n+1}^{(\ell,j)} & \dots & h_{n+1}^{(\ell,c_s)} \end{bmatrix}^\top$$

of the $c_z \times c_s$ observation matrix

$$\mathbf{H}_{n+1} \triangleq \begin{bmatrix} \mathbf{h}_{n+1}^{(1)} & \dots & \mathbf{h}_{n+1}^{(\ell)} & \dots & \mathbf{h}_{n+1}^{(c_z)} \end{bmatrix}^\top$$

is defined as

$$h_{n+1}^{(\ell,j)} \triangleq Pr\{\mathbf{Z}_{n+1} = \boldsymbol{\mu}_z^{(\ell)} | \mathbf{S}_{n+1} = \boldsymbol{\mu}_s^{(j)}\}$$

such that $\tilde{\pi}_{n+1|n}^{(\ell)} \propto \langle \mathbf{h}_{n+1}^{(\ell)}, \boldsymbol{\pi}_{n+1|n} \rangle$, $\forall \ell \in \{1, \dots, c_z\}$.

Note though that we can not use Eqs. (16) and (17) to compute the posteriors $P(\mathbf{s}_{n+1} | \check{\mathbf{z}}_{0:n})$ and $P(\mathbf{z}_{n+1} | \check{\mathbf{z}}_{0:n})$ since the matrices $\mathbf{F}_{n+1|n}^{(k)}$ and \mathbf{H}_{n+1} are unknown. The Baum-Welch [4] algorithm is a widespread expectation-maximization (EM) procedure to estimate the state transition probabilities $P(\mathbf{s}_{n+1} | \mathbf{s}_n)$ and the symbol emitting probabilities $P(\check{\mathbf{z}}_n | \mathbf{s}_n)$ required by both $\mathbf{F}_{n+1|n}^{(k)}$ and \mathbf{H}_{n+1} . However, we need to continuously learn these models online, recursively as new observations arrive. Thus, we need to build models for them using a different approach.

3.3 Belief representation

The predicted marginal posterior $P(\mathbf{s}_n | \check{\mathbf{z}}_{0:n-1})$ can be represented by an alternative vector – referred to as sparsely encoded belief hereafter in the paper – in the high-dimensional binary space $\boldsymbol{\Psi}_s \triangleq \mathbb{Z}_2^{c_s+p_s}$

$$\mathbf{b}_{n|n-1} \triangleq \begin{bmatrix} b_{n|n-1}^{(1)} & \dots & b_{n|n-1}^{(c_s+p_s)} \end{bmatrix}^\top \in \boldsymbol{\Psi}_s \quad (18)$$

with up to p_s active bits such that $1 \leq \|\mathbf{b}_{n|n-1}\|_1 \leq p_s$. As indicated in Appendix A, we can employ a deterministic procedure to build

$$\mathbf{b}_{n|n-1} = Bel(\boldsymbol{\pi}_{n|n-1}; p_s) \quad (19)$$

at any $n \geq 0$. Moreover, we can approximately retrieve the original posterior $\boldsymbol{\pi}_{n|n-1}$ by decoding the sparse belief (see Appendix A for details)

$$\boldsymbol{\pi}_{n|n-1} \approx Bel^{-1}(\mathbf{b}_{n|n-1}; p_s). \quad (20)$$

Hence, the belief $\mathbf{b}_{n|n-1}$ can be seen as a sparse symbol in the high-dimensional binary space $\boldsymbol{\Psi}_s$ corresponding to a valid p.m.f. $\boldsymbol{\pi}_{n|n-1}$ collecting probabilities of symbols in Ω_s . As discussed in Appendix A, the proposed encoding procedure $Bel(\cdot)$ (see Algorithm 4) is lossy in the sense that multiple categorical distributions can map to the same sparse belief. Consequently, as indicated in Eq. (20), the decoding procedure $Bel^{-1}(\cdot)$ (see Algorithm 5) will not be able to fully recover the original belief. Furthermore, this may harm the performance of the proposed grid-based filter, since, as highlighted in the sequel (refer to Algorithm 1), at each time step n , the grid-based filter employs the encoding procedure to build the predicted belief $\mathbf{b}_{n+1|n} = Bel(\boldsymbol{\pi}_{n+1|n}; p_s)$ with up to p_s activated bits, a binary sparse representation of the predicted distribution $\boldsymbol{\pi}_{n+1|n}$.

3.4 Re-interpreting filtering steps

From Eqs. (16) and (20), we can see that, at instant n , the predicted marginal posterior $\boldsymbol{\pi}_{n+1|n}$ is approximately a function $\mathbf{f}_n : \Omega_z \times \boldsymbol{\Psi}_s \mapsto \mathbb{S}_{c_s}$ of the current observation $\check{\mathbf{z}}_n = \boldsymbol{\mu}_z^{(k)}$ and the belief $\mathbf{b}_{n|n-1}$ at instant $n-1$, i.e.

$$\boldsymbol{\pi}_{n+1|n} = C_{n+1|n}^{-1} \cdot \mathbf{F}_{n+1|n}^{(k)} \cdot \boldsymbol{\pi}_{n|n-1} \\ \approx C_{n+1|n}^{-1} \cdot \mathbf{F}_{n+1|n}^{(k)} \cdot Bel^{-1}(\mathbf{b}_{n|n-1}; p_s) \\ \equiv \mathbf{f}_n(\check{\mathbf{z}}_n, \mathbf{b}_{n|n-1}), \quad (21)$$

where the normalization constant $C_{n+1|n}$ is chosen such that $\|\boldsymbol{\pi}_{n+1|n}\|_1 = 1$. Similarly, by examining Eq. (17) and noting that $\boldsymbol{\pi}_{n+1|n} \approx Bel^{-1}(\mathbf{b}_{n+1|n}; p_s)$, we conclude that, at instant n , the predicted posterior $\tilde{\boldsymbol{\pi}}_{n+1|n}$ is also an approximate function $\mathbf{h}_{n+1} : \boldsymbol{\Psi}_s \mapsto \mathbb{S}_{c_z}$ of the predicted belief $\mathbf{b}_{n+1|n}$ at instant n , i.e.

$$\tilde{\boldsymbol{\pi}}_{n+1|n} = \tilde{C}_{n+1|n}^{-1} \cdot \mathbf{H}_{n+1} \cdot \boldsymbol{\pi}_{n+1|n} \\ \approx \tilde{C}_{n+1|n}^{-1} \cdot \mathbf{H}_{n+1} \cdot Bel^{-1}(\mathbf{b}_{n+1|n}; p_s) \\ \equiv \mathbf{h}_{n+1}(\mathbf{b}_{n+1|n}) \quad (22)$$

with $\mathbf{b}_{n+1|n} = Bel(\boldsymbol{\pi}_{n+1|n}; p_s)$ and $\tilde{C}_{n+1|n}$ selected such that $\|\tilde{\boldsymbol{\pi}}_{n+1|n}\|_1 = 1$. Thus, at each instant n , it suffices to propagate the sparse beliefs from $\mathbf{b}_{n|n-1}$ to $\mathbf{b}_{n+1|n}$ and plug them into Eqs. (21) and (22) to compute the predicted marginal posteriors $\boldsymbol{\pi}_{n+1|n}$ and $\tilde{\boldsymbol{\pi}}_{n+1|n}$, respectively.

3.5 Non-parametric estimation

We follow the lead of [17] and assume that, conditioned on a particular sequence of hidden states $\check{\mathbf{s}}_{0:n}$ up to instant n , the predicted posterior $\pi_{n+1|n}$ is a realization of a random vector $\tilde{\Pi}_{n+1|n}$ in the $(c_s - 1)$ -dimensional simplex \mathbb{S}_{c_s} following a Dirichlet distribution $Dir(c_s; \alpha_{n|n-1})$ with hyperparameters $\alpha_{n|n-1} \triangleq \left[\alpha_{n|n-1}^{(1)} \quad \dots \quad \alpha_{n|n-1}^{(c_s)} \right]^\top \in \mathbb{R}_+^{c_s}$, i.e.

$$\tilde{\Pi}_{n+1|n} | \check{\mathbf{s}}_{0:n} \sim Dir(c_s; \alpha_{n|n-1}). \quad (23)$$

We also assume that $\alpha_{n|n-1}$ is indexed by the observation $\check{\mathbf{z}}_n$ at instant n and the previous belief $\mathbf{b}_{n|n-1}$ as per Eq. (21). Thus, given the pair $(\check{\mathbf{z}}_n, \mathbf{b}_{n|n-1})$, we can access the corresponding hyperparameters as $\alpha_{n|n-1} \equiv \alpha_{n|n-1}(\check{\mathbf{z}}_n, \mathbf{b}_{n|n-1})$ and use (23) to draw the predicted posterior $\pi_{n+1|n}$. In the sequel, we employ Algorithm 4 (see Appendix A) to encode the predicted posterior $\pi_{n+1|n}$ into the sparse belief

$$\mathbf{b}_{n+1|n} = Bel(\pi_{n+1|n}; p_s). \quad (24)$$

Conditioned on the sequence of observations $\check{\mathbf{z}}_{0:n}$ up to instant n , we also assume that the predicted posterior $\tilde{\pi}_{n+1|n}$ is a realization of a Dirichlet random vector $\tilde{\Pi}_{n+1|n}$ in the $(c_z - 1)$ -dimensional simplex \mathbb{S}_{c_z} with distribution $Dir(c_z; \tilde{\alpha}_{n|n-1})$ and hyperparameters $\tilde{\alpha}_{n|n-1} \triangleq \left[\tilde{\alpha}_{n|n-1}^{(1)} \quad \dots \quad \tilde{\alpha}_{n|n-1}^{(c_z)} \right]^\top \in \mathbb{R}_+^{c_z}$, i.e.

$$\tilde{\Pi}_{n+1|n} | \check{\mathbf{z}}_{0:n} \sim Dir(c_z; \tilde{\alpha}_{n|n-1}), \quad (25)$$

in which $\tilde{\alpha}_{n|n-1}$ is yet again indexed by the predicted belief $\mathbf{b}_{n+1|n}$ computed as in (24) following Eq. (22). Thus, we access the hyperparameters as $\tilde{\alpha}_{n|n-1} \equiv \tilde{\alpha}_{n|n-1}(\mathbf{b}_{n+1|n})$ and use (25) to draw the predicted posterior $\tilde{\pi}_{n+1|n}$. Finally, at instant n , we can predict the next encoded observation $\check{\mathbf{z}}_{n+1}$ at instant n from the predicted belief $\tilde{\pi}_{n+1|n}$ as

$$\begin{aligned} \hat{\mathbf{z}}_{n+1|n} &= \boldsymbol{\mu}_z^{(\hat{k})} \in \Omega_z \\ \hat{k} &= \arg \max_{1 \leq \ell \leq c_z} \tilde{\pi}_{n+1|n}^{(\ell)}. \end{aligned} \quad (26)$$

The probability $\tilde{\pi}_{n+1|n}^{(\hat{k})}$ can be seen thus as a confidence score of the filter prediction $\hat{\mathbf{z}}_{n+1|n}$, whereas the probability $1 - \tilde{\pi}_{n+1|n}^{(\hat{k})}$ with $k \in \{1, \dots, c_z\}$ corresponding to true observed symbol $\check{\mathbf{z}}_{n+1} = \boldsymbol{\mu}_z^{(k)}$ at instant $n+1$ can be seen as a novelty (or anomaly) score of the observation $\check{\mathbf{z}}_{n+1}$.

Algorithm 1 summarizes a single iteration of the proposed non-parametric, grid-based filtering (nP-GbF) procedure at instant n considering some memory-based mechanism \mathcal{M} to store-retrieve the hyperparameters.

Algorithm 1 nP-GbF($\check{\mathbf{z}}_n, \mathbf{b}_{n|n-1}, \mathcal{M}$)

- 1: Retrieve $\alpha_{n|n-1} \equiv \alpha_{n|n-1}(\check{\mathbf{z}}_n, \mathbf{b}_{n|n-1})$ from \mathcal{M} .
 - 2: Draw $\pi_{n+1|n} \sim Dir(c_s; \alpha_{n|n-1})$ as in (23).
 - 3: Build the sparsely encoded belief $\mathbf{b}_{n+1|n}$ as in (24).
 - 4: Retrieve $\tilde{\alpha}_{n|n-1} \equiv \tilde{\alpha}_{n|n-1}(\mathbf{b}_{n+1|n})$ from \mathcal{M} .
 - 5: Draw $\tilde{\pi}_{n+1|n} \sim Dir(c_z; \tilde{\alpha}_{n|n-1})$ as in (25).
 - 6: Predict next observation $\hat{\mathbf{z}}_{n+1|n} = \boldsymbol{\mu}_z^{(\hat{k})}$ as in (26).
 - 7: **return** $(\hat{\mathbf{z}}_{n+1|n}, \hat{k}, \mathbf{b}_{n+1|n}, \tilde{\pi}_{n+1|n})$
-

Figures 2 and 3 show a running example of the nP-GbF procedure. In Step 1, we retrieve the pseudo-counts $\alpha_{n|n-1}$ from

\mathcal{M} given the observation $\check{\mathbf{z}}_n$ at instant n and the previous belief $\mathbf{b}_{n|n-1}$. In the sequel, we sample the predicted posterior $\pi_{n+1|n}$ from the Dirichlet distribution $Dir(c_s; \alpha_{n|n-1})$ in Step 2 and build the predicted belief $\mathbf{b}_{n+1|n}$ in Step 3. Note that the recovered posterior $Bel^{-1}(\mathbf{b}_{n+1|n}; p_s)$ shown in Fig. 2 does not match the sampled posterior $\pi_{n+1|n}$ in Step 2. Specifically, the probabilities $\pi_{n+1|n}^{(\ell'_2)}$ and $\pi_{n+1|n}^{(\ell'_3)}$ in the original posterior $\pi_{n+1|n}$ collapsed into a single probability at location ℓ'_2 , since the indexes $\ell'_2 = 12$ and $\ell'_3 = 14$ were too close, i.e. $\ell'_3 - \ell'_2 = 2 < p_s = 8$. On the other hand, note that the probability $\pi_{n+1|n}^{(\ell'_1)}$ at the location $\ell'_1 = 4$ remained virtually intact in the recovered posterior $Bel^{-1}(\mathbf{b}_{n+1|n}; p_s)$.

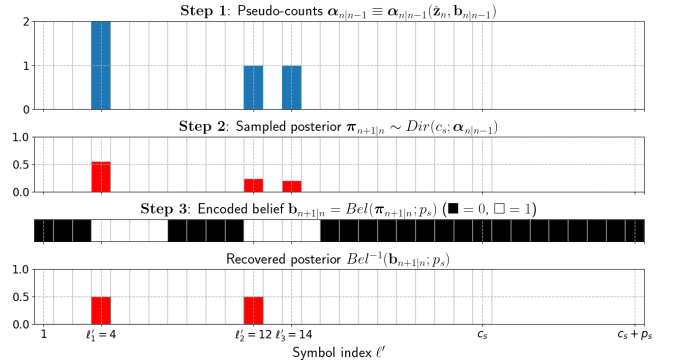


Figure 2: Running example illustrating Steps 1 through 3 of the nP-GbF procedure considering $c_s = 24$, $p_s = 8$ and $\alpha_{n|n-1}^{(\ell')} = 2^{-12}$, $\forall \ell'$, except for $\ell'_1 = 4$, $\ell'_2 = 12$ and $\ell'_3 = 14$.

Similarly, as show in Fig. 3, we retrieve the pseudo-counts $\tilde{\alpha}_{n|n-1}$ from \mathcal{M} in Step 4 using now the predicted belief $\mathbf{b}_{n+1|n}$ at instant n obtained in Step 3. Then, we sample the predicted posterior $\tilde{\pi}_{n+1|n}$ from the Dirichlet distribution $Dir(c_z; \tilde{\alpha}_{n|n-1})$ in Step 5 and, finally, predict the next observation $\hat{\mathbf{z}}_{n+1|n}$ in Step 6 as the symbol $\boldsymbol{\mu}_z^{(\hat{k})} \in \Omega_z$ with the highest probability $\tilde{\pi}_{n+1|n}^{(\hat{k})}$ in the predicted posterior $\tilde{\pi}_{n+1|n}$. In the running example, the highest probability $\tilde{\pi}_{n+1|n}^{(\ell_2)}$ occurs at the location ℓ_2 , i.e. $\hat{k} = \ell_2 = 16$.

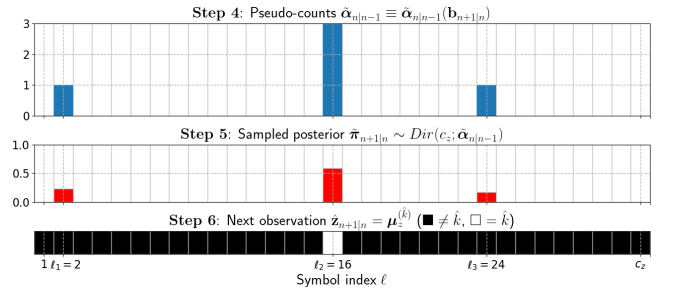


Figure 3: Running example illustrating Steps 3 through 6 of the nP-GbF procedure considering $c_z = 32$ and $\tilde{\alpha}_{n|n-1}^{(\ell)} = 2^{-12}$, $\forall \ell$, except for $\ell_1 = 2$, $\ell_2 = 16$ and $\ell_3 = 24$.

3.6 Hidden states formation

In general, the hidden states are not directly observable. As expected, the filtering steps (see Algorithm 1) do not require one to directly observe them. However, as shown in Appendix B, to be able to update the model hyperparameters

$\alpha_{n|n-1}$ using a frequentist approach, we need to pin the index j of the hidden state symbol $\mu_s^{(j)}$ at instant $n+1$ from which the observation \check{z}_{n+1} is emitted. To this end, we further assume that the hidden state \mathbf{s}_n represents the last ξ observations. Specifically, at instant $n+1$, upon the arrival of the new encoded observation $\check{z}_{n+1} = \mu_z^{(k)}$, $k \in \{1, \dots, c_z\}$, we select a unique index $j \in \{1, \dots, c_s\}$ and assign it to the sub-sequence of observations $\check{z}_{n+2-\xi:n+1}$, i.e. we make $\check{s}_{n+1} = \mu_s^{(j)} \in \Omega_s$. Lastly, in this paper, we increment the selected indexes from $\{1, \dots, c_s\}$ by p_s to make sure their belief representations will not overlap (see Appendix A). Then, we update the hyperparameters $\alpha_{n|n-1}$ and $\tilde{\alpha}_{n|n-1}$ as follows.

3.7 Non-parametric learning

By construction, conditioned on the predicted marginal posterior $\pi_{n+1|n}$ at instant n , the hidden state \mathbf{S}_{n+1} at instant $n+1$ is distributed according to a categorical distribution $Cat(\Omega_s; \pi_{n+1|n})$ over Ω_s , i.e.

$$\mathbf{S}_{n+1} | \pi_{n+1|n} \sim Cat(\Omega_s; \pi_{n+1|n}). \quad (27)$$

Since the Dirichlet distribution is the conjugate prior of the categorical distribution, it follows from Eqs. (23) and (27) that, conditioned on the built sequence of symbols $\check{s}_{0:n+1}$ up to instant $n+1$, the random vector $\mathbf{\Pi}_{n+1|n}$ is also Dirichlet

$$\mathbf{\Pi}_{n+1|n} | \check{s}_{n+1} = \mu_s^{(j)}, \check{s}_{0:n} \sim Dir(c_s; \alpha_{n+1|n}), \quad (28)$$

in which the new hyperparameters are updated using a frequentist approach as (proof details are in Appendix B)

$$\alpha_{n+1|n} = \alpha_{n|n-1} + \mathbf{1}_j. \quad (29)$$

Note that, conditioned on $\tilde{\pi}_{n+1|n}$, the next observation \mathbf{Z}_{n+1} at instant $n+1$ is also distributed according to a categorical distribution $Cat(\Omega_z; \tilde{\pi}_{n+1|n})$ over Ω_z , i.e.

$$\mathbf{Z}_{n+1} | \tilde{\pi}_{n+1|n} \sim Cat(\Omega_z; \tilde{\pi}_{n+1|n}). \quad (30)$$

Thus, it follows from Eqs. (25) and (30) that the same principle is applicable to $\tilde{\mathbf{\Pi}}_{n+1|n}$. Conditioned on the encoded sequence of observations $\check{z}_{0:n+1}$ up to instant $n+1$, we can write

$$\tilde{\mathbf{\Pi}}_{n+1|n} | \check{z}_{n+1} = \mu_z^{(k)}, \check{z}_{0:n} \sim Dir(c_z; \tilde{\alpha}_{n+1|n}), \quad (31)$$

where the new hyperparameters are given by (see Appendix B)

$$\tilde{\alpha}_{n+1|n} = \tilde{\alpha}_{n|n-1} + \mathbf{1}_k. \quad (32)$$

Even though learning happens at instant $n+1$, it is worth noting that the updated hyperparameters $\alpha_{n+1|n}$ must be properly stored indexed by both the observation \check{z}_n and the previous belief $\mathbf{b}_{n|n-1}$ from instant n such that $\alpha_{n+1|n} \equiv \alpha_{n+1|n}(\check{z}_n, \mathbf{b}_{n|n-1})$. Conversely, the updated hyperparameters $\tilde{\alpha}_{n+1|n}$ must be stored indexed by the predicted belief $\mathbf{b}_{n+1|n}$, i.e. $\tilde{\alpha}_{n+1|n} \equiv \tilde{\alpha}_{n+1|n}(\mathbf{b}_{n+1|n})$. Thus, in this context, we need to draw a new posterior $\pi_{n+1|n} \sim Dir(c_s; \alpha_{n+1|n})$ using (23) and rebuild the sparse belief $\mathbf{b}_{n+1|n}$ as in (24), before storing the updated hyperparameters $\tilde{\alpha}_{n+1|n}$ – in some preexisting storage \mathcal{M} – indexed by $\mathbf{b}_{n+1|n}$.

The non-parametric, memory-based learning (nP-MbL) procedure is outlined in Algorithm 2. We assume again that all hyperparameters are retrieved from and stored back to the same memory-based mechanism \mathcal{M} and that the observation and the hidden state at instant $n+1$ are respectively $\check{z}_{n+1} = \mu_z^{(k)}$ and $\check{s}_{n+1} = \mu_s^{(j)}$. Note that we weighted the increments in (29) and (32) by a learning rate $\eta \geq 1$ to speed-up sequence learning.

Algorithm 2 nP-MbL($j, k, \eta, \check{z}_n, \mathbf{b}_{n|n-1}, \mathbf{b}_{n+1|n}, \mathcal{M}$)

- 1: Retrieve $\alpha_{n|n-1} \equiv \alpha_{n|n-1}(\check{z}_n, \mathbf{b}_{n|n-1})$ from \mathcal{M} .
 - 2: Update $\alpha_{n+1|n} \leftarrow \alpha_{n|n-1} + \eta \cdot \mathbf{1}_j$ as in (29).
 - 3: Store $\alpha_{n+1|n}(\check{z}_n, \mathbf{b}_{n|n-1}) \equiv \alpha_{n+1|n}$ into \mathcal{M} .
 - 4: Draw $\pi_{n+1|n} \sim Dir(c_s; \alpha_{n+1|n})$ as in (23).
 - 5: Encode $\pi_{n+1|n}$ into $\mathbf{b}_{n+1|n}$ as in (24).
 - 6: Retrieve $\tilde{\alpha}_{n|n-1} \equiv \tilde{\alpha}_{n|n-1}(\mathbf{b}_{n+1|n})$ from \mathcal{M} .
 - 7: Update $\tilde{\alpha}_{n+1|n} \leftarrow \tilde{\alpha}_{n|n-1} + \eta \cdot \mathbf{1}_k$ as in (32).
 - 8: Store $\tilde{\alpha}_{n+1|n}(\mathbf{b}_{n+1|n}) \equiv \tilde{\alpha}_{n+1|n}$ into \mathcal{M} .
 - 9: **return** \mathcal{M}
-

4. RAM-BASED IMPLEMENTATION

In this section we propose an efficient implementation of the memory-based mechanism \mathcal{M} using a shallow, two-layer depth WNN. Specifically, each layer consists of an one-dimensional array of VG-RAM nodes allowing us thus to i) quickly retrieve input-output pairs indexed by binary input patterns and ii) perform one-shot learning by storing new associative input-output pairs readily.

A VG-RAM node stores a set \mathbf{M} of input-output pairs – a.k.a. lookup table – such that, $\forall (\beta', \alpha') \in \mathbf{M}$, $\beta' \in \mathbb{Z}_2^{d_\beta}$ is a d_β -dimensional binary input pattern and $\alpha' \in \mathbb{R}_+$ is the associated output, in our case, a non-negative hyperparameter. In particular, given an input pattern $\beta \in \mathbb{Z}_2^{d_\beta}$, the node finds the hyperparameter $\bar{\alpha}$ whose corresponding input pattern $\bar{\beta}$ has the closest Hamming distance to it. More precisely, the node retrieves the hyperparameter $\bar{\alpha}$ in $(\bar{\beta}, \bar{\alpha}) \leftarrow \text{FIND}(\mathbf{M}, \beta)$ from \mathbf{M} for a given β such that

$$(\bar{\beta}, \bar{\alpha}) = \arg \min_{(\beta', \alpha') \in \mathbf{M}} d_H(\beta', \beta).$$

On the other hand, to learn a new input-output pair (β, α) , the node makes $\mathbf{M} \leftarrow \text{INIT}(\mathbf{M}, \beta, \alpha) \triangleq \mathbf{M} \cup \{(\beta, \alpha)\}$. Lastly, to update the hyperparameter $\bar{\alpha}$ associated to an existing input pattern $\bar{\beta}$, replacing it by some $\tilde{\alpha} \in \mathbb{R}_+$, the node makes $\mathbf{M} \leftarrow \text{UPDT}(\mathbf{M}, \bar{\beta}, \tilde{\alpha}, \bar{\alpha}) \triangleq (\mathbf{M} \setminus \{(\bar{\beta}, \bar{\alpha})\}) \cup \{(\bar{\beta}, \tilde{\alpha})\}$.

4.1 Non-parametric estimation

Figure 4 graphically illustrates Steps 1 through 3 of Algorithm 1 and summarizes therefore how the function $\mathbf{b}_{n+1|n} \equiv \mathbf{f}_n(\check{z}_n, \mathbf{b}_{n|n-1})$ is implemented. In particular, the figure shows how we employ a c_s -dimensional array of VG-RAM nodes to retrieve the hyperparameters $\alpha_{n|n-1}$ based on the stored input-output pairs whose input patterns have the closest Hamming distance to the concatenated input vector $\check{z}_n \oplus \mathbf{b}_{n|n-1} \in \mathbb{Z}_2^{c_z+c_s+p_s}$. Figure 5 depicts Steps 4 through 6 of Algorithm 1, which in turn implements the function $\tilde{\mathbf{b}}_{n+1|n} \equiv \mathbf{h}_n(\mathbf{b}_{n+1|n})$. The figure indicates how an additional c_z -dimensional array of VG-RAM nodes is employed to retrieve $\tilde{\alpha}_{n|n-1}$ based on the stored associative pairs whose input patterns have the closest Hamming distance to the input $\mathbf{b}_{n+1|n} \in \mathbb{Z}_2^{c_s+p_s}$.

Note that, in filtering Steps 2 and 5, the filter assigns different probabilities $\pi_{n+1|n}^{(\ell')}$ and $\tilde{\pi}_{n+1|n}^{(\ell)}$ to symbols in Ω_s and Ω_z according to the retrieved hyperparameters $\alpha_{n+1|n}^{(\ell')}$ and $\tilde{\alpha}_{n+1|n}^{(\ell)}$ in Steps 1 and 4, respectively. In Step 3, the filter also assigns an indexable label (symbol in Ψ_s) to the predicted posterior $\pi_{n+1|n}$ by building the sparse belief $\mathbf{b}_{n+1|n} = \text{Bel}(\pi_{n+1|n})$,

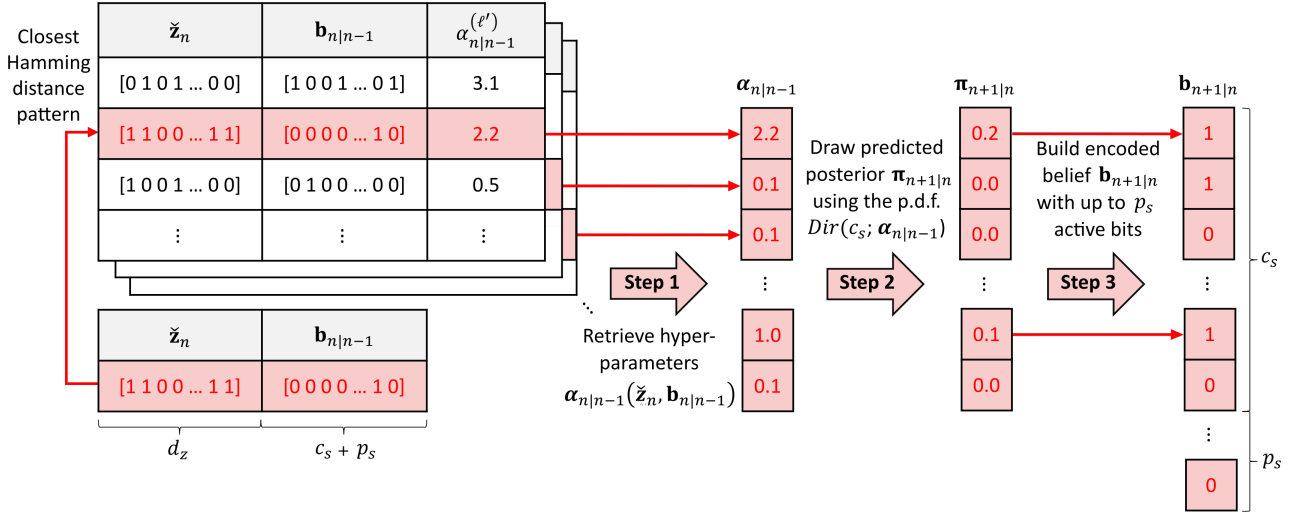


Figure 4: WNN layer \mathcal{L}_1 : one-step-ahead state belief prediction.

which is in turn employed to retrieve the hyperparameters $\tilde{\alpha}_{n+1|n}^{(\ell)}$ at Step 4. We select finally the predicted observation $\hat{\mathbf{z}}_{n+1|n}$ in Step 6 – based on the predicted marginal posterior $\tilde{\pi}_{n+1|n}$ – as the base vector $\boldsymbol{\mu}_z^{(\hat{k})}$ with the highest probability $\tilde{\pi}_{n+1|n}^{(\hat{k})}$ using (26).

4.2 Non-parametric learning

Let $\mathbf{M}_1^{(\ell')}$ and $\mathbf{M}_2^{(\ell)}$ denote respectively the lookup tables of the ℓ' -th and ℓ -th VG-RAM nodes at WNN layers \mathcal{L}_1 and \mathcal{L}_2 . At instant n , the ℓ' -th VG-RAM node in Fig. 4 retrieves

$$\left(\cdot, \alpha_{n|n-1}^{(\ell')}\right) \leftarrow \text{FIND} \left(\mathbf{M}_1^{(\ell')}, \tilde{\mathbf{z}}_n \oplus \mathbf{b}_{n|n-1} \right)$$

and, therefore, contributes indirectly to the activation of the ℓ' -th bit of the belief $\mathbf{b}_{n+1|n}$ through the Dirichlet-sampled probability $\pi_{n+1|n}^{(\ell')}$. By the same token, the ℓ -th node in Fig. 5 retrieves

$$\left(\cdot, \tilde{\alpha}_{n|n-1}^{(\ell)}\right) \leftarrow \text{FIND} \left(\mathbf{M}_2^{(\ell)}, \mathbf{b}_{n+1|n} \right)$$

at instant n driving therefore the odds of the ℓ -th symbol $\boldsymbol{\mu}_z^{(\ell)}$ in Ω_z being selected as the prediction $\hat{\mathbf{z}}_{n+1|n} = \boldsymbol{\mu}_z^{(\hat{k})}$, i.e. $\hat{k} = \ell$, by shaping the Dirichlet-sampled probability $\tilde{\pi}_{n+1|n}^{(\ell)}$. Then, at instant $n + 1$, the VG-RAM nodes finally learn the new hyperparameters $\alpha_{n+1|n}^{(\ell')}$ and $\tilde{\alpha}_{n+1|n}^{(\ell)}$ by storing new input-output association pairs

$$\mathbf{M}_1^{(\ell')} \leftarrow \text{INIT} \left(\mathbf{M}_1^{(\ell')}, \tilde{\mathbf{z}}_n \oplus \mathbf{b}_{n|n-1}, \alpha_{n+1|n}^{(\ell')} \right)$$

$$\mathbf{M}_2^{(\ell)} \leftarrow \text{INIT} \left(\mathbf{M}_2^{(\ell)}, \mathbf{b}_{n+1|n}, \tilde{\alpha}_{n+1|n}^{(\ell)} \right)$$

or by updating the existing ones as

$$\mathbf{M}_1^{(\ell')} \leftarrow \text{UPDT} \left(\mathbf{M}_1^{(\ell')}, \tilde{\mathbf{z}}_n \oplus \mathbf{b}_{n|n-1}, \alpha_{n|n-1}^{(\ell')}, \alpha_{n+1|n}^{(\ell')} \right)$$

$$\mathbf{M}_2^{(\ell)} \leftarrow \text{UPDT} \left(\mathbf{M}_2^{(\ell)}, \mathbf{b}_{n+1|n}, \tilde{\alpha}_{n|n-1}^{(\ell)}, \tilde{\alpha}_{n+1|n}^{(\ell)} \right).$$

Intuitively, the ℓ' -th VG-RAM node in layer \mathcal{L}_1 learns – in a frequentist fashion – how likely is to activate the correspond-

ing p_s bits $b_{n+1|n}^{(\ell')}, \dots, b_{n+1|n}^{(\ell'+p_s)}$ along the layer's output belief $\mathbf{b}_{n+1|n}$ when exposed to a given input pattern $\tilde{\mathbf{z}}_n \oplus \mathbf{b}_{n|n-1}$ such that the likelihood of activating those bits increases as the corresponding input pattern (or a similar one in terms of the Hamming distance) is repeatedly presented to the layer's input. The VG-RAM layer \mathcal{L}_1 thus learns how to predict its output belief given the current input pattern. Note that, as we are not estimating the hidden state \mathbf{s}_{n+1} , there is no need to explicitly build the set Ω_s . As shown in Appendix B, we just need the symbol indexes to update the hyperparameters $\alpha_{n|n-1}$. The set Ω_z , in turn, is required by both the WNN layer \mathcal{L}_2 to make the prediction $\hat{\mathbf{z}}_{n+1|n}$ in Step 6 and the codec scheme to encode the observation $\hat{\mathbf{y}}_{n+1}$ as in (11) and decode the prediction $\hat{\mathbf{z}}_{n+1|n}$ as in (12).

5. APPLICATION EXAMPLE

In this section, we illustrate the use of the proposed nP-GbF and nP-MbL procedures in an anomaly detection (AD) application. In particular, we build an anomaly detector to handle streaming, real-time data and assess its performance using four synthetic toy time-series, specifically designed to evaluate AD algorithms.

Algorithm 3 shows how we can employ Algorithms 1 and 2 to build a non-parametric, CL anomaly detection (nP-CLAD) procedure. In this procedure, we compute first the anomaly score Σ_{n+1}^{ad} as the complement of the probability $\tilde{\pi}_{n+1|n}^{(k)}$ in $\tilde{\pi}_{n+1|n}$ corresponding to the encoded observation $\tilde{\mathbf{z}}_{n+1} = \boldsymbol{\mu}_z^{(k)}$ at instant $n + 1$. Then, we verify if the corresponding anomaly score Σ_{n+1}^{ad} exceeds a threshold τ . Note therefore that this test takes into account the predicted probability $\tilde{\pi}_{n+1|n}^{(k)}$ of the observed symbol $\tilde{\mathbf{z}}_{n+1} = \boldsymbol{\mu}_z^{(k)}$ in $\tilde{\pi}_{n+1|n}$. Alternatively, we can modify this criterion to check if the encoded observation $\tilde{\mathbf{z}}_{n+1}$ mismatches the predicted observation $\hat{\mathbf{z}}_{n+1|n}$ – symbol in Ω_z with the highest assigned probability in the predicted (sampled) posterior $\tilde{\pi}_{n+1|n}$ – by a maximum number of bits δ , i.e. $d_H(\tilde{\mathbf{z}}_{n+1}, \hat{\mathbf{z}}_{n+1|n}) \geq \delta$.

For convenience, we summarize the hyperparameters employed by the sparse codec, the WNN layers – \mathcal{L}_1 and \mathcal{L}_2 – and the nP-CLAD algorithm in Table 1.

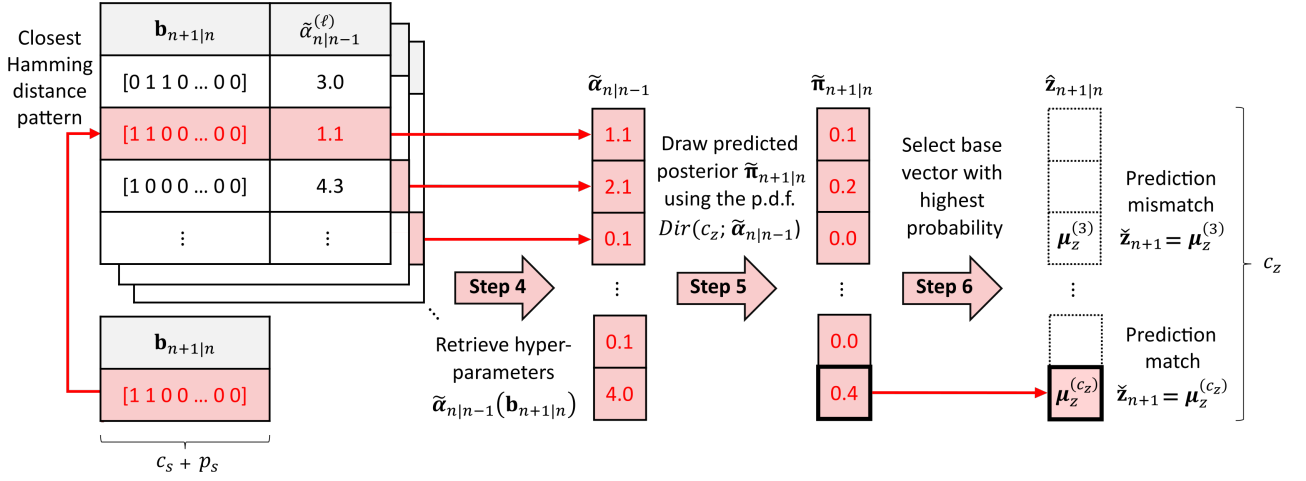


Figure 5: WNN layer \mathcal{L}_2 : next observation belief prediction & matching.

Algorithm 3 nP-CLAD($\check{\mathbf{y}}_{n+1}, \check{\mathbf{z}}_n, \mathbf{b}_{n|n-1}, \mathcal{M}$)

- 1: Encode observation $\check{\mathbf{y}}_{n+1}$ as $\check{\mathbf{z}}_{n+1} = \mu_z^{(k)}$ using (11).
- 2: Build the hidden state $\check{\mathbf{s}}_{n+1} = \mu_s^{(j)}$ from $\check{\mathbf{z}}_{n+2-\xi:n+1}$.
- 3: Predict the observation \mathbf{z}_{n+1} using Algorithm 1

$$\left(\check{\mathbf{z}}_{n+1|n}, \hat{k}, \mathbf{b}_{n+1|n}, \tilde{\pi}_{n+1|n} \right) \leftarrow \text{nP-GbF}(\check{\mathbf{z}}_n, \mathbf{b}_{n|n-1}, \mathcal{M}).$$

- 4: Compute the anomaly score as $\Sigma_{n+1}^{ad} \leftarrow 1 - \tilde{\pi}_{n+1|n}^{(k)}$.
- 5: Detect the anomaly by checking $\Lambda_{n+1}^{ad} \leftarrow \llbracket \Sigma_{n+1}^{ad} > \tau \rrbracket$.
- 6: Learn the non-parametric model using Algorithm 2

$$\mathcal{M} \leftarrow \text{nP-MbL}(j, k, \eta, \check{\mathbf{z}}_n, \mathbf{b}_{n|n-1}, \mathbf{b}_{n+1|n}, \mathcal{M}).$$

- 7: **return** $(\Lambda_{n+1}^{ad}, \Sigma_{n+1}^{ad}, \check{\mathbf{z}}_{n+1}, \mathbf{b}_{n+1|n}, \mathcal{M})$

Table 1: Summary of the WNN model hyperparameters.

Codec: map observations $\check{\mathbf{y}}_n \in \Omega_y$ into symbols $\check{\mathbf{z}}_n \in \Omega_z$	
a_z	Maximum number of activated bits in $\check{\mathbf{z}}_n$
b_z	Minimum Hamming distance of symbols in Ω_z
d_z	Total number of bits of $\check{\mathbf{z}}_n$
Layer \mathcal{L}_1 : learn how to predict the posterior $\pi_{n+1 n}$	
c_s	Number of symbols in Ω_s
$\forall \ell', \alpha_{\ell'}$	Initial Dirichlet $Dir(c_s; \alpha_{0 n-1})$ pseudo-counts
p_s	Number of activated bits in the belief $\mathbf{b}_{n n-1}$
Layer \mathcal{L}_2 : learn how to predict the posterior $\tilde{\pi}_{n+1 n}$	
c_z	Number of symbols in Ω_z
$\forall \ell, \tilde{\alpha}_{\ell}$	Initial Dirichlet $Dir(c_z; \tilde{\alpha}_{0 n-1})$ pseudo-counts
nP-CLAD: check $\check{\mathbf{y}}_{n+1}$ using its anomaly score Σ_{n+1}^{ad}	
ξ	Length of the sub-sequence $\check{\mathbf{z}}_{n+2-\xi:n+1}$
τ	Anomaly detection threshold
η	Learning rate employed by the WNN model

We initialize $\mathbf{b}_{0|n-1}$ as per nP-GbF filter Steps 2 and 3 by making $\alpha_{0|n-1} = \alpha_0$ in Algorithm 1. We encode the streamed observations $\check{\mathbf{y}}_{0:\infty}$ in Step 1 using a fixed-range (0 – 100)

numeric encoder with 961 buckets – unique symbols in Ω_z – and $d_z = 2^{10}$ bits to encode the time-series values as indicated in [29]. We also initialize $\check{\mathbf{z}}_{1-\xi:-1}$ by replicating $\check{\mathbf{z}}_0$ to be able to build $\check{\mathbf{s}}_n, \forall n \geq 0$, as required in Step 2 of Algorithm 3. There is no need to decode network predictions in this AD application.

Table 2 summarizes the hyperparameters employed to setup the WNN model and run all experiments. We should highlight that only 11 hyperparameters are required to initialize the proposed RAM-based model plus the nP-CLAD algorithm. The hyperparameter d_s is not required, since there is no need to build the set Ω_s .

Table 2: Hyperparameters employed by the codec, the WNN layers, and the nP-CLAD algorithm throughout experiments.

Codec	Layer \mathcal{L}_1	Layer \mathcal{L}_2	nP-CLAD				
a_z	2^6	c_s	2^{10}	c_z	2^{10}	ξ	2^3
b_z	2^2	$\forall \ell', \alpha_{\ell'}$	2^{-12}	$\forall \ell, \tilde{\alpha}_{\ell}$	2^{-12}	τ	0.75
d_z	2^{10}	p_s	2^3	–	–	η	2^1

Figures 6 through 9 show the AD results considering four artificial time-series with randomly anomalous spikes. Despite being quite simple, these AD experiments illustrate the filter ability to learn increasingly complex sequences online and detect point anomalies. We highlighted the true anomalies – intervals shaded in red – and indicated the detected anomalies – red circles – in the last 500 time instants (green area), since the filter was still learning the time-series in the first half of the simulation (gray area). The nP-CLAD was able to successfully detect all random spikes along the evaluated time-series using the anomaly scores provided by the non-parametric model. Nevertheless, it is worth noting that the filter still needs to be validated in real-world settings using e.g. Numenta Anomaly Benchmark (NAB) [1].

The experiments were conducted in an Intel(R) Core(TM) i7-9750H central processing unit (CPU) @ 2.60GHz with 32GB of RAM. Note that the inference time depends on the number of entries stored in the VGRAM nodes' memory, since the nodes memory is virtual and its size increases as the WNN model learns a new sequence.

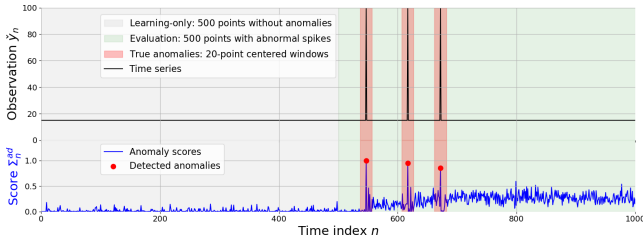


Figure 6: Constant time-series with anomalous spikes.

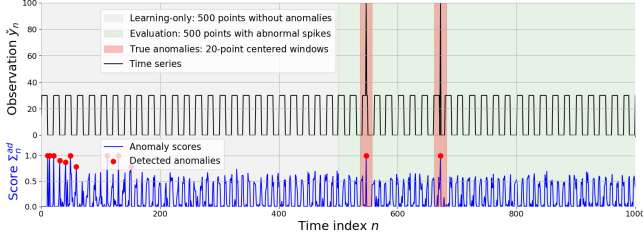


Figure 7: Squared time-series with anomalous spikes.

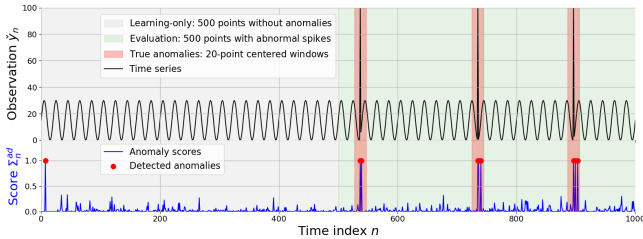


Figure 8: Sinusoidal time-series with anomalous spikes.

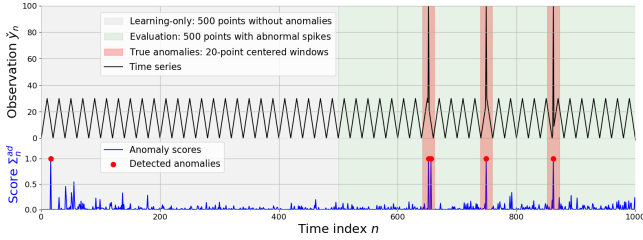


Figure 9: Sawtooth time-series with anomalous spikes.

Table 3 indicates the mean and maximum inference times in milliseconds – considering one iteration of Algorithm 3 –, and the accumulated memory size in kilobytes at the end of each experiment. As expected, the accumulated memory size increases according to the complexity of the learned time series. Finally, a non-optimized Python implementation of the proposed algorithms and reproducible experiments can be found in [15].

Table 3: Computational performance considering different time series. Inference times are in milliseconds (ms), whereas accumulated memory sizes are in kilobytes (kB).

Time series	Inference time		Accumulated memory size
	Average	Maximum	
Constant	10.74 ms	24.82 ms	58.45 kB
Squared	11.22 ms	32.03 ms	151.17 kB
Sinusoidal	11.72 ms	31.93 ms	168.30 kB
Sawtooth	11.50 ms	35.62 ms	178.40 kB

6. CONCLUSIONS

6.1 Summary of paper contributions

We presented in this paper a fully non-parametric, grid-based filter designed to intrinsically learn an interpretable probabilistic model of the underlying Markovian process (see Fig. 1) emitting the incoming sequence of observations, which allows one to make explainable and auditable predictions – relying on clear input-to-output causal chains connected to traceable memory entries – with promising, but still limited results as shown in toy AD examples. Moreover, the proposed RAM-based filter implementation performs one-shot, CL without requiring any special hardware to accomplish it. Indeed, VG-RAM nodes can be efficiently implemented through software – using e.g. C++ to really tame the machine – conceived to run on general-purpose, CPUs using ordinary, abundant RAM hardware, which is orders of magnitude cheaper than state-of-the-art graphics processing units (GPUs) specially designed for machine learning tasks.

6.2 Future research directions

As future research, we propose to adapt the sparse auto-encoder introduced in [26] to learn (offline) a high-dimensional binary representation of the continuous-valued observations in an embedded space. We also intend to employ mismatched predictions as learning signals as proposed by Cui et al. in [11]. Furthermore, we can limit the memory capacity of the VG-RAM nodes and individually force them to eventually replace associative input-output pairs with small hyperparameters by new ones, limiting thus the required network resources by forcing its nodes to gracefully forget in the long run input patterns which are not frequently seen, i.e. presented as input to their corresponding layers. Lastly, we may use sparse binary word vectors as proposed by Faruqi et al. in [16] to encode natural language data and check the suitability of the proposed filter to build a tiny language model which is able to continuously learn from new corpora.

7. ADDITIONAL AUTHORS

Additional authors: Thiago Oliveira-Santos and Claudine Badue (Applied Computational Intelligence Institute, UFES [todsantos, claudine]@inf.ufes.br).

8. REFERENCES

- [1] S. Ahmad, A. Lavin, S. Purdy, and Z. Agha. Unsupervised real-time anomaly detection for streaming data. *Neurocomputing*, 262:134–147, 2017.
- [2] I. Aleksander. From WISARD to MAGNUS: A family of weightless virtual neural machines. *RAM-Based Neural Networks*, pages 18–30, 1998.
- [3] S. Barra, S. M. Carta, A. Corrigan, A. S. Podda, and D. R. Recupero. Deep learning and time series-to-image encoding for financial forecasting. *IEEE/CAA Journal of Automatica Sinica*, 7(3):683–692, 2020.
- [4] L. E. Baum, T. Petrie, G. Soules, and N. Weiss. A maximization technique occurring in the statistical analysis of probabilistic functions of markov chains. *The Annals of Mathematical Statistics*, 41(1):164–171, 1970.

- [5] L. Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010: 19th International Conference on Computational Statistics*, pages 177–186. Springer, 2010. Paris France, August 22–27, 2010 Keynote, Invited and Contributed Papers.
- [6] L. Bottou. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade: Second Edition*, pages 421–436. Springer, 2012.
- [7] A. Canale and M. Ruggiero. Bayesian nonparametric forecasting of monotonic functional time series. *Electronic Journal of Statistics*, 10(2):3265–3286, 2016.
- [8] Y. Chen, S. Liu, Z. Wang, D. Wu, Y. Li, B. Xing, B. Guo, and L. Chen. Online parallel attack detection method for industrial control based on multi-bandpass filter. *IEEE Internet of Things Journal*, 11(1):880–888, 2024.
- [9] C. Cheng, A. Sa-Ngasoongsong, O. Beyca, T. Le, H. Yang, Z. Kong, and S. T. Bukkapatnam. Time series forecasting for nonlinear and non-stationary processes: a review and comparative study. *IIE Transactions*, 47(10):1053–1071, 2015.
- [10] V. K. R. Chimmula and L. Zhang. Time series forecasting of COVID-19 transmission in canada using lstm networks. *Chaos, Solitons & Fractals*, 135:109864, 2020.
- [11] Y. Cui, S. Ahmad, and J. Hawkins. Continuous online sequence learning with an unsupervised neural network model. *Neural Computation*, 28(11):2474–2504, 2016.
- [12] M. De Gregorio and M. Giordano. An experimental evaluation of weightless neural networks for multi-class classification. *Applied Soft Computing*, 72:338–354, 2018.
- [13] A. F. De Souza, C. Badue, F. Pedroni, E. Oliveira, S. S. Dias, H. Oliveira, and S. F. de Souza. Face recognition with VG-RAM weightless neural networks. In *Artificial Neural Networks-ICANN 2008: 18th International Conference, Prague, Czech Republic, September 3-6, 2008, Proceedings, Part I 18*, pages 951–960. Springer, 2008.
- [14] A. F. De Souza, F. Pedroni, E. Oliveira, P. M. Ciarelli, W. F. Henrique, L. Veronese, and C. Badue. Automated multi-label text categorization with VG-RAM weightless neural networks. *Neurocomputing*, 72(10–12):2209–2217, 2009.
- [15] S. S. Dias. Weightless neural network library (wnnlib), 2025. Available from <https://github.com/stivenschwanz/wnnlib/tree/safeai>.
- [16] M. Faruqui, Y. Tsvetkov, D. Yogatama, C. Dyer, and N. A. Smith. Sparse overcomplete word vector representations. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1491–1500, 2015.
- [17] Z. Ghahramani. Bayesian non-parametrics and the probabilistic approach to modelling. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 371(1984):20110553, 2013.
- [18] J. Hawkins and S. Ahmad. Why neurons have thousands of synapses, a theory of sequence memory in neocortex. *Frontiers in Neural Circuits*, 10(23):1–13, 2016.
- [19] J. Hawkins, S. Ahmad, D. Dubinsky, et al. Cortical learning algorithm and hierarchical temporal memory. *Numenta Whitepaper*, pages 1–68, 2011.
- [20] A. Iyer, K. Grewal, A. Velu, L. O. Souza, J. Forest, and S. Ahmad. Avoiding catastrophe: Active dendrites enable multi-task learning in dynamic environments. *Frontiers in Neurobotics*, 16(846219):1–23, 2022.
- [21] S. M. Kay. *Fundamentals of Statistical Signal Processing: Estimation Theory*. Prentice Hall PTR, 1993.
- [22] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017.
- [23] B. Lim and S. Zohren. Time-series forecasting with deep learning: a survey. *Philosophical Transactions of the Royal Society A*, 379(2194):20200209, 2021.
- [24] T. B. Ludermir, A. C. P. d. L. F. Carvalho, A. P. Braga, and M. C. De Souto. Weightless neural models: a review of current and past works. *Neural Computing Surveys*, 2:41–61, 1999.
- [25] G. Mahalakshmi, S. Sridevi, and S. Rajaram. A survey on forecasting of time series data. In *2016 International Conference on Computing Technologies and Intelligent Data Engineering (ICCTIDE'16)*, pages 1–8. IEEE, 2016.
- [26] A. Makhzani and B. Frey. K-sparse autoencoders. *arXiv preprint arXiv:1312.5663*, 2013.
- [27] R. J. Mitchell, J. Bishop, S. Box, and J. Hawker. Comparison of some methods for processing ‘grey level’ data in weightless networks. In *RAM-based Neural Networks*, pages 61–70. World Scientific, 1998.
- [28] P. Müller, F. A. Quintana, A. Jara, and T. Hanson. *Bayesian nonparametric data analysis*, volume 1. Springer, 2015.
- [29] S. Purdy. Encoding data for HTM systems. *arXiv preprint arXiv:1602.05925*, 2016.
- [30] T. Singh, R. Kalra, S. Mishra, Satakshi, and M. Kumar. An efficient real-time stock prediction exploiting incremental learning and deep learning. *Evolving Systems*, 14(6):919–937, 2023.
- [31] R. Sousa, T. Lima, A. Abelha, and J. Machado. Hierarchical temporal memory theory approach to stock market time series forecasting. *Electronics*, 10(14):1630, 2021.
- [32] J. Xuan, J. Lu, and G. Zhang. A survey on bayesian non-parametric learning. *ACM Computing Surveys (CSUR)*, 52(1):1–36, 2019.

APPENDIX

A. BELIEF ENCODING

Algorithm 4 indicates how to encode the categorical distribution $\boldsymbol{\pi}$ with c categories into a sparsely encoded belief \mathbf{b} in the binary space \mathbb{Z}_2^{c+p} with up to p activated bits. Note that this is a lossy encoding scheme. Besides the floor operation in Line 5, which rounds off the scaled probabilities $p \cdot \pi_i$, $\forall i \in \{1, \dots, c\}$, non null probabilities $\pi_i > 0$ and $\pi_j > 0$, $i \neq j$, in $\boldsymbol{\pi}$ with close indexes $|i - j| < p$ can generate a sequence of contiguous active bits in \mathbf{b} . Thus, as illustrated in the running example in Sec. 3 (inspect Fig. 2), depending on the probability values π_i and π_j , multiple valid categorical distributions can be mapped to the same belief by this algorithm.

Algorithm 4 $Bel(\boldsymbol{\pi} = [\pi_1 \dots \pi_c]^\top; p)$

```

1: Initialize  $\mathbf{b} = [b_1 \dots b_{c+p}]^\top$  such that,  $\forall i \in \{1, 2, \dots, c+p\}$ ,  $b_i \leftarrow 0$ .
2:  $acc \leftarrow 0$ 
3: for  $i \leftarrow 1$  to  $c+p$  do
4:   if  $i \leq c$  then
5:      $acc \leftarrow acc + \lfloor p \cdot \pi_i + \frac{1}{2} \rfloor$ 
6:   end if
7:   if  $acc > 0$  then
8:      $b_i \leftarrow 1$ 
9:      $acc \leftarrow acc - 1$ 
10:  end if
11: end for
12: return  $\mathbf{b}$ 

```

Algorithm 5 shows in turn how to decode the sparsely encoded belief \mathbf{b} into the posterior $\boldsymbol{\pi}$ within the $(c-1)$ -dimensional simplex \mathbb{S}_c . It is worth noting that this algorithm will assign a non null probability π_i to the first symbol belonging to a contiguous sequence $b_i, b_{i+1}, b_{i+2}, \dots$ of activate bits in \mathbf{b} . Thus, one should avoid assigning similar symbols to close locations in $\boldsymbol{\pi}$ such that it is unlikely to produce a single contiguous sequence of active bits within \mathbf{b} . Otherwise, the $Bel(\cdot)$ procedure in Algorithm 4 will activate overlapping sequences of contiguous bits in the binary space \mathbb{Z}_2^{c+p} for those symbols, therefore yielding an ambiguous representation which cannot be properly recovered by the $Bel^{-1}(\cdot)$ procedure in Algorithm 5.

Algorithm 5 $Bel^{-1}(\mathbf{b} = [b_1 \dots b_{c+p}]^\top; p)$

```

1: Initialize  $\boldsymbol{\pi} = [\pi_1 \dots \pi_c]^\top$  such that,  $\forall i \in \{1, 2, \dots, c\}$ ,  $\pi_i \leftarrow 0$ .
2:  $acc \leftarrow 0$ 
3: for  $i \leftarrow c+p$  to 1 do
4:    $acc \leftarrow acc + b_i$ 
5:   if  $i \leq c \wedge (i = 1 \vee b_{i-1} = 0)$  then
6:      $\pi_i \leftarrow acc$ 
7:      $acc \leftarrow 0$ 
8:   end if
9: end for
10: Normalize  $\boldsymbol{\pi}$  such that  $\sum_{i=1}^c \pi_i = 1$ .
11: return  $\boldsymbol{\pi}$ 

```

B. SEQUENCE LEARNING

According to assumption (30), $p(\mathbf{z}_{n+1} | \tilde{\boldsymbol{\pi}}_{n+1|n})$ is a categorical distribution. Thus, for a particular realization $\check{\mathbf{z}}_{n+1} = \boldsymbol{\mu}_z^{(k)}$ of \mathbf{Z}_{n+1} , the likelihood function becomes

$$\begin{aligned} p(\check{\mathbf{z}}_{n+1} | \tilde{\boldsymbol{\pi}}_{n+1|n}) &= p(\mathbf{z}_{n+1} | \tilde{\boldsymbol{\pi}}_{n+1|n}) \Big|_{\mathbf{z}_{n+1} = \check{\mathbf{z}}_{n+1}} \\ &= \sum_{\ell=1}^{c_z} \tilde{\pi}_{n+1|n}^{(\ell)} \cdot \mathbb{I}[\mathbf{z}_{n+1} = \boldsymbol{\mu}_z^{(\ell)}] \Big|_{\mathbf{z}_{n+1} = \boldsymbol{\mu}_z^{(k)}} \\ &= \tilde{\pi}_{n+1|n}^{(k)}. \end{aligned} \quad (33)$$

From Eq. (33) and assumption (25), it follows that we can write the updated posterior $p(\tilde{\boldsymbol{\pi}}_{n+1|n} | \check{\mathbf{z}}_{0:n+1})$ using the Bayes rule as

$$\begin{aligned} p(\tilde{\boldsymbol{\pi}}_{n+1|n} | \check{\mathbf{z}}_{0:n+1}) &\propto p(\check{\mathbf{z}}_{n+1} | \tilde{\boldsymbol{\pi}}_{n+1|n}) p(\tilde{\boldsymbol{\pi}}_{n+1|n} | \check{\mathbf{z}}_{0:n}) \\ &= \frac{\tilde{\pi}_{n+1|n}^{(k)}}{B(\tilde{\boldsymbol{\alpha}}_{n|n-1})} \prod_{\ell=1}^{c_z} \left(\tilde{\pi}_{n+1|n}^{(\ell)} \right)^{\tilde{\alpha}_{n|n-1}^{(\ell)} - 1} \\ &\propto \tilde{\pi}_{n+1|n}^{(k)} \prod_{\ell=1}^{c_z} \left(\tilde{\pi}_{n+1|n}^{(\ell)} \right)^{\tilde{\alpha}_{n|n-1}^{(\ell)} - 1}, \end{aligned} \quad (34)$$

since the Beta function $B(\tilde{\boldsymbol{\alpha}}_{n|n-1})$ is constant for a given $\tilde{\boldsymbol{\alpha}}_{n|n-1}$. We can rewrite thus the right-hand side of (34) as

$$\left(\tilde{\pi}_{n+1|n}^{(k)} \right)^{\left(\tilde{\alpha}_{n|n-1}^{(k)} + 1 \right) - 1} \prod_{\ell=1, \ell \neq k}^{c_z} \left(\tilde{\pi}_{n+1|n}^{(\ell)} \right)^{\tilde{\alpha}_{n|n-1}^{(\ell)} - 1}. \quad (35)$$

Thus, given $\check{\mathbf{z}}_{n+1} = \boldsymbol{\mu}_z^{(k)}$, the updated posterior can be rewritten as

$$p(\tilde{\boldsymbol{\pi}}_{n+1|n} | \check{\mathbf{z}}_{0:n+1}) = \frac{1}{B(\tilde{\boldsymbol{\alpha}}_{n+1|n})} \prod_{\ell=1}^{c_z} \left(\tilde{\pi}_{n+1|n}^{(\ell)} \right)^{\tilde{\alpha}_{n+1|n}^{(\ell)} - 1} \quad (36)$$

with new normalization constant given by the reciprocal of $B(\tilde{\boldsymbol{\alpha}}_{n+1|n})$ such that, conditioned on both $\check{\mathbf{z}}_{n+1} = \boldsymbol{\mu}_z^{(k)}$ and $\check{\mathbf{z}}_{0:n}$, $\tilde{\boldsymbol{\Pi}}_{n+1|n}$ is also Dirichlet

$$\tilde{\boldsymbol{\Pi}}_{n+1|n} | \check{\mathbf{z}}_{n+1} = \boldsymbol{\mu}_z^{(k)}, \check{\mathbf{z}}_{0:n} \sim \text{Dir}(c_z; \tilde{\boldsymbol{\alpha}}_{n+1|n})$$

with new hyperparameters $\tilde{\boldsymbol{\alpha}}_{n+1|n} = \tilde{\boldsymbol{\alpha}}_{n|n-1} + \mathbf{1}_k$.

Using a similar approach, given assumptions (23) and (27), we can show that, conditioned on a particular realization $\check{\mathbf{s}}_{n+1} = \boldsymbol{\mu}_s^{(j)}$ of \mathbf{S}_{n+1} , the Bayes rule application

$$\begin{aligned} p(\boldsymbol{\pi}_{n+1|n} | \check{\mathbf{s}}_{0:n+1}) &\propto p(\check{\mathbf{s}}_{n+1} | \boldsymbol{\pi}_{n+1|n}) p(\boldsymbol{\pi}_{n+1|n} | \check{\mathbf{s}}_{0:n}) \\ &= \frac{\pi_{n+1|n}^{(j)}}{B(\boldsymbol{\alpha}_{n|n-1})} \prod_{\ell'=1}^{c_s} \left(\pi_{n+1|n}^{(\ell')} \right)^{\alpha_{n|n-1}^{(\ell')} - 1} \\ &\propto \pi_{n+1|n}^{(j)} \prod_{\ell'=1}^{c_s} \left(\pi_{n+1|n}^{(\ell')} \right)^{\alpha_{n|n-1}^{(\ell')} - 1} \end{aligned} \quad (37)$$

yields an updated posterior

$$p(\boldsymbol{\pi}_{n+1|n} | \check{\mathbf{s}}_{0:n+1}) = \frac{1}{B(\boldsymbol{\alpha}_{n+1|n})} \prod_{\ell'=1}^{c_s} \left(\pi_{n+1|n}^{(\ell')} \right)^{\alpha_{n+1|n}^{(\ell')} - 1} \quad (38)$$

with new normalization constant given now by the reciprocal of $B(\boldsymbol{\alpha}_{n+1|n})$, which is also Dirichlet

$$\boldsymbol{\Pi}_{n+1|n} | \check{\mathbf{s}}_{n+1} = \boldsymbol{\mu}_s^{(j)}, \check{\mathbf{s}}_{0:n} \sim \text{Dir}(c_s; \boldsymbol{\alpha}_{n+1|n})$$

with new hyperparameters $\boldsymbol{\alpha}_{n+1|n} = \boldsymbol{\alpha}_{n|n-1} + \mathbf{1}_j$.