



Contents lists available at ScienceDirect

International Journal of Forecasting

journal homepage: www.elsevier.com/locate/ijforecast

Coupling LSTM neural networks and state-space models through analytically tractable inference

Van-Dai Vuong*, Luong-Ha Nguyen, James-A. Goulet

Department of Civil, Geologic and Mining Engineering, Polytechnique Montréal, 2500 chemin de Polytechnique, Montréal, Québec, H3T 1J4, Canada

ARTICLE INFO

Keywords:Bayesian inference
Probabilistic method
Long short-term memory
State-space models
Time series forecasting
Hybrid model

ABSTRACT

Long short-term memory (LSTM) neural networks and state-space models (SSMs) are effective tools for time series forecasting. Coupling these methods to exploit their advantages is not a trivial task because their respective inference procedures rely on different mechanisms. In this paper, we present formulations that allow for analytically tractable inference in Bayesian LSTMs and the probabilistic coupling between Bayesian LSTMs and SSMs. This is enabled by using analytical Gaussian inference as a single mechanism for inferring both the LSTM's parameters as well as the posterior for the SSM's hidden states. We show through several experimental comparisons that the resulting hybrid model retains the interpretability feature of SSMs, while exploiting the ability of LSTMs to learn complex seasonal patterns with minimal manual setups.

© 2024 The Author(s). Published by Elsevier B.V. on behalf of International Institute of Forecasters. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

1. Introduction

Time series forecasting plays a key role in various industries such as finance, transportation, energy, infrastructure and healthcare. Recurrent neural networks (RNNs) (Rumelhart et al., 1986) and especially the long short-term memory (LSTM) architecture (Hochreiter & Schmidhuber, 1997) have been shown to be effective tools for time series forecasting problems (Hewamalage et al., 2021). The main strength of neural networks (NNs) is their ability to automatically identify complex patterns while requiring minimal manual setups (Januschowski et al., 2020). Nevertheless, without making assumptions about the data, the results from NNs are typically difficult to interpret (Rangapuram et al., 2018). By contrast, state-space models (SSMs) (Durbin & Koopman, 2001) provide a probabilistic framework for decomposing time series into interpretable patterns such as level, trend and seasonality.

In this paper, we propose a new method that allows coupling Bayesian LSTMs with SSMs to obtain a hybrid

model that retains the interpretability feature of SSMs, while exploiting the ability of LSTMs to learn complex patterns automatically. Coupling these methods is not a trivial task because their respective inference procedures rely on different mechanisms. SSMs are probabilistic models which rely on Bayesian inference, whereas LSTMs typically optimize their parameters using backpropagation (Rumelhart et al., 1986) and gradient descent (GD). As a result, one needs to choose between the two mechanisms when coupling these models. Typically, existing hybrid models including those proposed by Smyl (2020) and Rangapuram et al. (2018) have relied on backpropagation and GD. In this paper, we explore the probabilistic coupling between LSTMs and SSMs, while using Bayesian inference for inferring both the LSTM's parameters as well as the posterior for SSM's hidden states.

Recently, the Tractable Approximate Gaussian Inference (TAGI) method (Goulet et al., 2021) has allowed analytical Bayesian inference in neural networks. The performance of TAGI has been demonstrated on feedforward

* Corresponding author.

E-mail address: van-dai.vuong@polymtl.ca (V.-D. Vuong).

<https://doi.org/10.1016/j.ijforecast.2024.04.002>

0169-2070/© 2024 The Author(s). Published by Elsevier B.V. on behalf of International Institute of Forecasters. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

neural networks (FNNs) (Goulet et al., 2021), convolutional neural networks (CNNs), generative adversarial networks (GANs) (Nguyen & Goulet, 2021), and reinforcement learning (RL) (Nguyen & Goulet, 2022), yet it has never been tested for RNNs due to the lack of a mathematical framework to do so. In this paper, we present the mathematical formulations which enable using the TAGI method with the LSTM architecture for obtaining analytically the posterior mean vectors and diagonal covariance matrices for the latent variables and model parameters. Because all variables in this TAGI-LSTM are now Gaussians, this new framework can thus be directly coupled with SSMs.

The paper is organized as follows: Section 2 reviews Bayesian RNNs along with hybrid models which couple RNNs with SSMs. Section 3 reviews the TAGI method along with the LSTM architecture which are the foundations for this work. Section 4 first presents the mathematical formulations that have been developed to apply TAGI to the LSTM architecture, and then introduces the probabilistic coupling between TAGI-LSTM and SSMs. Finally, Section 5 compares our methods with other benchmark models on both simulated and real datasets.

The following notation is used throughout the manuscript; x : lowercases denote variables; \mathbf{x} : bold lowercases denote vectors; X : italic uppercases denote random variables; \mathbf{X} : bold italic uppercases denote vectors or matrices of random variables; \mathbf{X} : bold upright uppercases denote deterministic matrices; $\mathbf{y}_{1:t}$: denote observations from 1 to t ; $\boldsymbol{\mu}$ and $\mathbb{E}[\cdot]$: mean vectors; $\boldsymbol{\Sigma}$: covariance matrices; $\mathbf{X}_{t|t-1} \sim \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{t|t-1}^{\mathbf{X}}, \boldsymbol{\Sigma}_{t|t-1}^{\mathbf{X}}) \equiv \mathcal{N}(\boldsymbol{\mu}_{t|t-1}^{\mathbf{X}}, \boldsymbol{\Sigma}_{t|t-1}^{\mathbf{X}})$: prior probability density functions (PDFs) for the random variables \mathbf{X}_t given data $\mathbf{y}_{1:t-1}$ with $\boldsymbol{\mu}_{t|t-1}^{\mathbf{X}} \equiv \mathbb{E}[\mathbf{X}_t | \mathbf{y}_{1:t-1}]$ and $\boldsymbol{\Sigma}_{t|t-1}^{\mathbf{X}} \equiv \text{cov}(\mathbf{X}_t | \mathbf{y}_{1:t-1})$; $\mathbf{X}_{t|t} \sim \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{t|t}^{\mathbf{X}}, \boldsymbol{\Sigma}_{t|t}^{\mathbf{X}})$: posterior PDFs for the random variables \mathbf{X}_t given data $\mathbf{y}_{1:t}$; and $\mathbf{X}_{t|T} \sim \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{t|T}^{\mathbf{X}}, \boldsymbol{\Sigma}_{t|T}^{\mathbf{X}})$: smoothed estimates for the random variables \mathbf{X}_t given data $\mathbf{y}_{1:T}$ where T is the last training time step.

2. Related works

In this section, we review Bayesian recurrent neural networks (RNNs) along with hybrid models combining RNNs with state-space models (SSMs).

2.1. Bayesian recurrent neural networks

Bayesian NNs place a prior distribution over the parameters $f(\boldsymbol{\theta})$, and the goal is to find the posterior distribution $f(\boldsymbol{\theta} | \mathbf{x}, \mathbf{y})$, where \mathbf{x} are the input covariates and \mathbf{y} are the observed responses. This posterior is widely acknowledged as being intractable for NNs (Goodfellow et al., 2016), thus methods such as the Laplace approximation (MacKay, 1992) and variational inference (VI) (Hernandez-Lobato & Adams, 2015; Sun et al., 2017; Wu et al., 2019) have been proposed to approximate it. Chien and Ku (2016) developed an approximation for the Hessian matrix in order to use the Laplace approximation in Bayesian RNNs. Fortunato et al. (2019) applied the variational scheme Bayes by Backprop (Blundell et al., 2015)

to RNNs for estimating the parameters' posterior distribution. In a general context, VI has been shown to have a limited scalability to large datasets (Hernandez-Lobato & Adams, 2015). A part of the reason for it is that multiple passes through data are required for optimization, and there are typically more parameters in the VI-based RNNs than the deterministic ones sharing the same architecture.

Gal and Ghahramani (2016a) have approached Bayesian RNNs in a different way. They showed that performing Monte Carlo dropout (MC-dropout) and weight decay in a deterministic RNN trained with backpropagation and GD enables providing accurate predictive uncertainties. Creating Bayesian RNNs by applying MC-dropout to a deterministic RNN is straightforward, yet it cannot quantify the posterior for the parameters. In this approach, dropout can either be applied to only the non-recurrent connections (Pham et al., 2014; Zaremba et al., 2014), to the LSTM cell states (Moon et al., 2015), or to all of the input, output and recurrent connections (Gal & Ghahramani, 2016a).

The abovementioned Bayesian RNN methods rely on backpropagation and GD to optimize their parameters. In this paper, we provide a mathematical framework that allows applying the TAGI method to the LSTM architecture for quantifying the parameters' posterior distribution analytically using Bayesian inference.

2.2. Hybrid models

Some studies have established hybrid models combining RNNs and SSMs by using the former ones to either model nonlinear SSM's transition and/or observation equations, or define the SSM's parameters. Krishnan et al. (2017), Zheng et al. (2017) and Zaheer et al. (2017) have used RNNs to model the nonlinear transition of the hidden states, allowing their models to be non-Markovian. However, the common limitation among methods using NNs to parametrize the SSM's nonlinear dynamic models is that the marginal log-likelihood function is intractable so that they rely on either Monte Carlo (MC) methods (Zheng et al., 2017) or variational inference (Fraccaro et al., 2016; Krishnan et al., 2017) to approximate this function in order to update their neural network parameters. Fraccaro et al. (2016) and Rangapuram et al. (2018) have kept the SSM's transition and observation models linear which allows performing exact inference. The DeepState method (Rangapuram et al., 2018) uses a global LSTM to learn from the whole dataset and a local linear SSM for each time series. The global LSTM outputs the SSM's time-varying parameters defining the model matrices for each local SSM. In a different approach, Smyl (2020) used the deterministic exponential smoothing equations as a data processing tool for a LSTM. The level and seasonality components are estimated for each time series, and are used to normalize and deseasonalize data on-the-fly. However, the use of the deterministic form of exponential smoothing leads to point estimates for the baseline level and trend components.

The existing methods mentioned above have already attempted to couple SSMs with the deterministic RNNs. In this paper, we explore another direction by providing a probabilistic coupling between Bayesian LSTMs and SSMs.



Fig. 1. Graphical representation of a feedforward neural network. Black arrows represent the network forward connections; red and cyan arrows represent the inference directions for the hidden states and for the parameters, respectively.

3. Background

In this section, we review the theoretical basics of the LSTM architecture and the fundamentals of the TAGI method.

3.1. Long short-term memory neural network

In the context of univariate time series, dependencies are the relations between the current observation and the past ones. RNNs typically experience vanishing and exploding gradients which limit their capacity to learn long-term dependencies (Goodfellow et al., 2016). The long short-term memory (LSTM) network (Hochreiter & Schmidhuber, 1997) is a special type of RNN designed to address this problem. LSTMs use gating mechanisms to automatically select and store the dependency information in their memory. There are two kinds of memory in a LSTM cell, the hidden states \mathbf{h} encode the short-term dependency information, whereas the cell states \mathbf{c} store the long-term one. A LSTM cell consists in four gates including the forget, input, output and candidate gates, $\{\mathbf{f}, \mathbf{i}, \mathbf{o}, \tilde{\mathbf{c}}\} \in \mathbb{R}^N$, the cell states $\mathbf{c} \in \mathbb{R}^N$, and the hidden states $\mathbf{h} \in \mathbb{R}^N$, where N is the number of hidden units. All operations in a LSTM cell can be summarized by the following LSTM recursive equations

$$\mathbf{f}_t = \sigma(\mathbf{W}^f \mathbf{x}_t + \mathbf{U}^f \mathbf{h}_{t-1} + \mathbf{b}^f), \quad (1a)$$

$$\mathbf{i}_t = \sigma(\mathbf{W}^i \mathbf{x}_t + \mathbf{U}^i \mathbf{h}_{t-1} + \mathbf{b}^i), \quad (1b)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}^o \mathbf{x}_t + \mathbf{U}^o \mathbf{h}_{t-1} + \mathbf{b}^o), \quad (1c)$$

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}^c \mathbf{x}_t + \mathbf{U}^c \mathbf{h}_{t-1} + \mathbf{b}^c), \quad (1d)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t, \quad (1e)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t), \quad (1f)$$

where $\mathbf{x}_t \in \mathbb{R}^M$ is the input vector, $\{\mathbf{W}^f, \mathbf{W}^i, \mathbf{W}^o, \mathbf{W}^c\} \in \mathbb{R}^{M \times N}$ are the weight matrices for the input, $\{\mathbf{U}^f, \mathbf{U}^i, \mathbf{U}^o, \mathbf{U}^c\} \in \mathbb{R}^{N \times N}$ are the weight matrices for the hidden states, $\{\mathbf{b}^f, \mathbf{b}^i, \mathbf{b}^o, \mathbf{b}^c\} \in \mathbb{R}^N$ are the bias vectors, with M being the input vector's length, $\sigma(\cdot)$ is the logistic sigmoid function and $\tanh(\cdot)$ is the hyperbolic tangent function, and \odot denotes the element-wise multiplication operation.

A LSTM cell takes the input \mathbf{x}_t and the hidden states \mathbf{h}_{t-1} from the previous time step as inputs (Eqs. (1a)–(1d)), and outputs the hidden states \mathbf{h}_t (Eq. (1f)) which are then used to make final predictions. The cell and the hidden states, \mathbf{c}_t and \mathbf{h}_t , are updated at every time step in order to retain relevant information which is useful for current and future predictions. The updating process for the cell states \mathbf{c}_t consists of two steps, and is done using Eq. (1e). The first step is to discard irrelevant information from the previous cell states \mathbf{c}_{t-1} , by performing the element-wise product $\mathbf{f}_t \odot \mathbf{c}_{t-1}$. The second step consists in including the new information into the cell states \mathbf{c}_t by adding the element-wise product $\mathbf{i}_t \odot \tilde{\mathbf{c}}_t$. After updating the cell states \mathbf{c}_t , the hidden states \mathbf{h}_t are updated using Eq. (1f).

3.2. Tractable approximate Gaussian inference

Tractable Approximate Gaussian Inference (TAGI) (Goulet et al., 2021) is an analytically tractable framework for Bayesian neural networks consisting in two main steps: forward and backward. The forward step propagates the uncertainties from the input layer and the network parameters up to the output layer. The backward step employs a layer-wise recursive procedure to infer analytically the parameters $\boldsymbol{\theta}$ and the hidden states \mathbf{Z} from the observations \mathbf{y} .

Fig. 1 presents a compact graphical representation for a feedforward neural network where \mathbf{x} is the input vector; $\mathbf{z}^{(j)}$ and $\boldsymbol{\theta}^{(j)} = \{\mathbf{W}^{(j)}, \mathbf{B}^{(j)}\}$, $\forall j = 1 : L$, are the hidden states and parameters of the j^{th} hidden layer; $\mathbf{z}^{(0)}$ is the output vector; \mathbf{y} are the observations; and L is the number of hidden layers. TAGI considers the network's input, output, hidden states as well as parameters as Gaussian random variables such that $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X)$, $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}_Z, \boldsymbol{\Sigma}_Z)$, and $\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta)$.

The j^{th} hidden layer contains the hidden units $\mathbf{z}^{(j)}$ and their activation units $\mathbf{a}^{(j)}$ such that $\mathbf{a}^{(j)} = \phi(\mathbf{z}^{(j)})$, where $\phi(\cdot)$ is the activation function. The forward step starts by passing information from the input layer to the first hidden layer using

$$\mathbf{Z}^{(1)} = \mathbf{W}^{(0)} \mathbf{X} + \mathbf{B}^{(0)}, \quad (2)$$

$$\mathbf{A}^{(1)} = \phi(\mathbf{Z}^{(1)}). \quad (3)$$

Eq. (2) involves two types of operations including additions and multiplications of Gaussian random variables. Because TAGI assumes that the hidden states $\mathbf{Z}^{(1)}$ are Gaussians, these two operations must also result in Gaussian random variables, even if the multiplication operation does not lead to exact results. In order to overcome this problem, TAGI relies on the *Gaussian multiplication approximation* (GMA) to model the product of two Gaussian random variables by a Gaussian distribution whose exact mean and variance are calculated analytically as presented in Appendix A. In addition, the output moments for the nonlinear activation function $\phi(\cdot)$ in Eq. (3) cannot be evaluated analytically; TAGI uses the local linearization of activation functions $\phi(\cdot)$ to obtain the output moments as detailed in Appendix B. Using these approximations allows calculating analytically the mean vectors and covariance matrices defining the PDFs for the hidden and activation units. Similarly, going from a j^{th} layer to the subsequent $j + 1^{\text{th}}$ is done using the same Eqs. (2) and (3), while replacing the input \mathbf{X} by the activation units $\mathbf{A}^{(j)} \sim \mathcal{N}(\boldsymbol{\mu}_A^{(j)}, \boldsymbol{\Sigma}_A^{(j)})$, and $\boldsymbol{\theta}^{(0)}$ by the parameters $\boldsymbol{\theta}^{(j)} \sim \mathcal{N}(\boldsymbol{\mu}_\theta^{(j)}, \boldsymbol{\Sigma}_\theta^{(j)})$ of the j^{th} layer. This allows propagating the uncertainties from the input \mathbf{X} and the parameters $\boldsymbol{\theta}$ through the network up to the output layer $\mathbf{Z}^{(0)}$.

The relation between the output layer $\mathbf{z}^{(0)}$ and the observations \mathbf{y} is defined by the observation equation as

$$\mathbf{y} = \mathbf{z}^{(0)} + \mathbf{v}, \quad \mathbf{v} : \mathbf{V} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{v}}), \quad (4)$$

where \mathbf{v} denotes the observation error. The moments of the predictive distribution $\mathbf{Y} \sim \mathcal{N}(\mu_{\mathbf{Y}}, \Sigma_{\mathbf{Y}})$ can be obtained following $\mu_{\mathbf{Y}} = \mu_{\mathbf{Z}^{(0)}}$, $\Sigma_{\mathbf{Y}} = \Sigma_{\mathbf{Z}^{(0)}} + \Sigma_{\mathbf{v}}$.

In the context of TAGI, inferring the hidden states and parameters means estimating the conditional probability distributions $f(\mathbf{z}|\mathbf{y}) = \mathcal{N}(\mu_{\mathbf{z}|\mathbf{y}}, \Sigma_{\mathbf{z}|\mathbf{y}})$ and $f(\theta|\mathbf{y}) = \mathcal{N}(\mu_{\theta|\mathbf{y}}, \Sigma_{\theta|\mathbf{y}})$. The backward step first computes the posterior for the output $\mathbf{Z}^{(0)}$ using the Gaussian conditional equations

$$\begin{aligned} f(\mathbf{z}^{(0)} | \mathbf{y}) &= \mathcal{N}(\mu_{\mathbf{Z}^{(0)}|\mathbf{y}}, \Sigma_{\mathbf{Z}^{(0)}|\mathbf{y}}), \\ \mu_{\mathbf{Z}^{(0)}|\mathbf{y}} &= \mu_{\mathbf{Z}^{(0)}} + \Sigma_{\mathbf{Y}^{(0)}}^{-1} \Sigma_{\mathbf{Y}^{(0)}|\mathbf{Z}^{(0)}} (\mathbf{y} - \mu_{\mathbf{Y}}), \\ \Sigma_{\mathbf{Z}^{(0)}|\mathbf{y}} &= \Sigma_{\mathbf{Z}^{(0)}} - \Sigma_{\mathbf{Y}^{(0)}|\mathbf{Z}^{(0)}} \Sigma_{\mathbf{Y}^{(0)}}^{-1} \Sigma_{\mathbf{Y}^{(0)}|\mathbf{Z}^{(0)}}, \end{aligned} \quad (5)$$

where $\Sigma_{\mathbf{Y}^{(0)}|\mathbf{Z}^{(0)}} = \Sigma_{\mathbf{Z}^{(0)}}$. This is depicted by the red arrow from \mathbf{y} to $\mathbf{z}^{(0)}$ in Fig. 1. The layer-wise recursive inference procedure is then recursively applied to infer the hidden states and parameters of each layer from the last to the first layer as presented by the red and cyan arrows in Fig. 1. For maintaining the linear complexity with respect to the number of parameters in the network, TAGI assumes diagonal covariance matrices for the hidden states and the parameters, and relies on the conditional independence assumptions of hidden units between layers, i.e., $\mathbf{Z}^{(j-1)} \perp\!\!\!\perp \mathbf{Z}^{(j+1)} | \mathbf{z}^{(j)}$. The posterior distribution for the hidden states of each layer $\mathbf{Z}^{(j)}$ is estimated analytically using

$$\begin{aligned} f(\mathbf{z}^{(j)} | \mathbf{y}) &= \mathcal{N}(\mu_{\mathbf{z}^{(j)}|\mathbf{y}}, \Sigma_{\mathbf{z}^{(j)}|\mathbf{y}}), \\ \mu_{\mathbf{z}^{(j)}|\mathbf{y}} &= \mu_{\mathbf{z}^{(j)}} + \mathbf{J}_{\mathbf{z}^{(j)}} (\mu_{\mathbf{z}^{(j+1)}} - \mu_{\mathbf{z}^{(j+1)}}), \\ \Sigma_{\mathbf{z}^{(j)}|\mathbf{y}} &= \Sigma_{\mathbf{z}^{(j)}} + \mathbf{J}_{\mathbf{z}^{(j)}} (\Sigma_{\mathbf{z}^{(j+1)}} - \Sigma_{\mathbf{z}^{(j+1)}}) \mathbf{J}_{\mathbf{z}^{(j)}}^T, \\ \mathbf{J}_{\mathbf{z}^{(j)}} &= \text{cov}(\mathbf{Z}^{(j)}, \mathbf{Z}^{(j+1)}) (\Sigma_{\mathbf{z}^{(j+1)}})^{-1}, \end{aligned} \quad (6)$$

where $\text{cov}(\mathbf{Z}^{(j)}, \mathbf{Z}^{(j+1)})$ is the covariance between the hidden states of the j^{th} and $j+1^{\text{th}}$ layers. In parallel, the parameters $\theta^{(j)}$ can be inferred using the same Eq. (6) where the variable $\mathbf{Z}^{(j)}$ is replaced by $\theta^{(j)}$. The calculations for obtaining the covariances $\text{cov}(\mathbf{Z}^{(j)}, \mathbf{Z}^{(j+1)})$ and $\text{cov}(\theta^{(j)}, \mathbf{Z}^{(j+1)})$, as well as other quantities in Eq. (6) have been detailed by Goulet et al. (2021).

A weakly informative prior is employed when initializing the hyper-parameters $\{\mu_{\theta}^{(0)}, \Sigma_{\theta}^{(0)}\}$. In order to learn the parameters efficiently, TAGI repeats the inference over multiple epochs, analogously to the empirical Bayes approach (Efron, 2010). This means that the posterior hyper-parameters $\{\mu_{\theta|\mathbf{y}}^{(i)}, \Sigma_{\theta|\mathbf{y}}^{(i)}\}$ at the i^{th} epoch are used as the prior hyper-parameters for the next $i+1^{\text{th}}$ epoch.

TAGI can theoretically be used with any existing NN architecture to create the corresponding Bayesian neural network. However, this has never been done for RNNs. In this paper, we provide the mathematical formulations to use TAGI with the LSTM architecture.

4. Methods

In this section, we first present the analytically tractable TAGI-LSTM neural network, then introduce the probabilistic coupling between it and SSMS to create TAGI-LSTM/SSM hybrid models.

4.1. TAGI-LSTM

This section introduces the mathematical formulations for the TAGI-LSTM model where the LSTM's parameters and hidden states presented in Section 3.1 are inferred analytically using the TAGI method presented in Section 3.2. In order to apply the TAGI method to LSTM, we employ the same architecture formulations presented in Eq. (1) to calculate the gates $\{\mathbf{f}, \mathbf{i}, \tilde{\mathbf{c}}, \mathbf{o}\}$, the cell states \mathbf{c} , and the hidden states \mathbf{h} , and we consider them, along with the network parameters θ , as Gaussian random variables. In order to maintain the linear computational complexity of the TAGI method, we need to rely on the same independence assumption employed by Goulet et al. (2021), that is considering a diagonal covariance structure for all LSTM's gates, hidden states, cell states, and parameters. Fig. 2a shows a graphical representation of a LSTM cell, whereas Fig. 2b presents the graph for an example of stacked TAGI-LSTM network having an input layer containing the covariates \mathbf{x} , two LSTM layers, and a fully connected output layer. The observations \mathbf{y} are related to the outputs $\mathbf{z}^{(0)}$ through Eq. (4).

4.1.1. Forward step

At a time step t , we denote the marginal prior knowledge for the hidden states of a LSTM layer given the past data $\mathbf{y}_{1:t-1}$ by the Gaussian PDF $\mathbf{H}_{t|t-1} \sim \mathcal{N}(\mu_{\mathbf{H}_{t|t-1}}^{\mathbf{H}}, \Sigma_{\mathbf{H}_{t|t-1}}^{\mathbf{H}})$, where $\mu_{\mathbf{H}_{t|t-1}}^{\mathbf{H}} \equiv \mathbb{E}[\mathbf{H}_t | \mathbf{y}_{1:t-1}]$, and $\Sigma_{\mathbf{H}_{t|t-1}}^{\mathbf{H}} \equiv \text{cov}(\mathbf{H}_t | \mathbf{y}_{1:t-1})$. In the forward step, we want to pass information from the input covariates \mathbf{X}_t through the LSTM layers, up to the output layer. This corresponds to estimating the prior knowledge for the hidden and cell states of each LSTM layer, $\mathbf{H}_{t|t-1}^{(j)}$ and $\mathbf{C}_{t|t-1}^{(j)}$, as well as the prior $\mathbf{Z}_{t|t-1}^{(0)} \sim \mathcal{N}(\mu_{\mathbf{Z}_{t|t-1}^{(0)}}, \Sigma_{\mathbf{Z}_{t|t-1}^{(0)}})$ for the hidden states of the output layer.

In order to pass information through a LSTM layer, we first need to obtain the prior knowledge for the four LSTM gates. In the probabilistic context where all quantities are modelled by Gaussian random variables, we define the hidden states for the forget gate as

$$\mathbf{Z}_t^f = \mathbf{W}_t^f \mathbf{X}_t + \mathbf{U}_t^f \mathbf{H}_{t-1} + \mathbf{B}_t^f.$$

Each hidden state $\mathbf{Z}_{i,t}^f$ (at time step t , and cell i) for the forget gate can be obtained as

$$\mathbf{Z}_{i,t}^f = \mathbf{W}_{i,t}^f \mathbf{X}_t + \mathbf{U}_{i,t}^f \mathbf{H}_{t-1} + \mathbf{B}_{i,t}^f, \quad (7)$$

where $i = 1 : N$, $\mathbf{W}_{i,t}^f \in \mathbb{R}^{1 \times M}$, $\mathbf{U}_{i,t}^f \in \mathbb{R}^{1 \times N}$ are row weight matrices, N being the LSTM cell's size, and M is the input's size. We employ diagonal covariance structures for the PDFs of the input covariates \mathbf{X}_t and the hidden states \mathbf{H}_{t-1} . Under this assumption, $\mathbf{Z}_{i,t}^f$ is the sum of $M + N$ independent products and a bias term as shown in Eq. (7). Goulet et al. (2021) have shown that when adding a large number of independent product terms, the

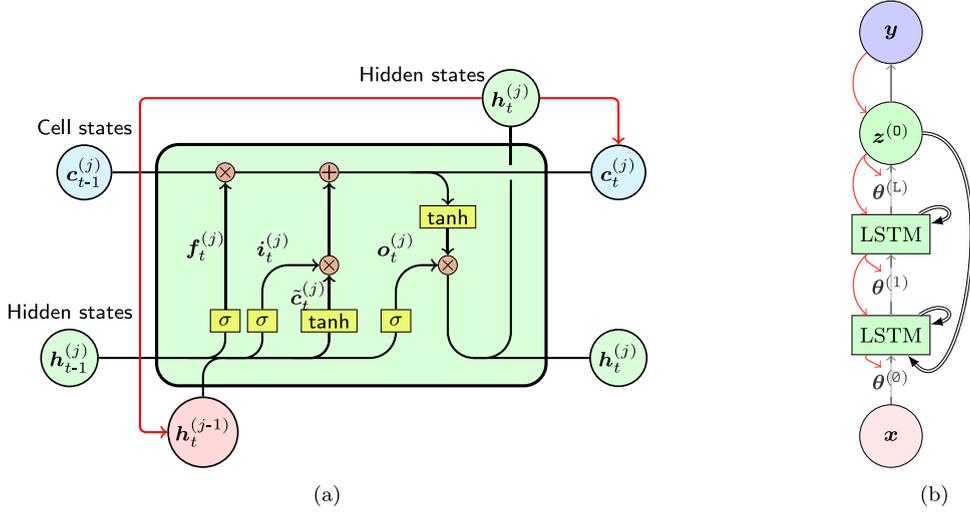


Fig. 2. (a) Graphical representation of a LSTM cell. Red arrows represent the inference procedure to update the j -th LSTM layer from the subsequent j th LSTM layer. (b) Graphical representation of a stacked TAGI-LSTM network. Black arrows represent the network's forward connections, red arrows represent the layer-wise inference paths, and double arrows represent the recurrent connections.

correlation between the resulting pairs of output hidden states tends to zero such that assuming $Z_{i,t}^f \perp\!\!\!\perp Z_{j,t}^f$ is valid. Therefore, the hidden states for the forget gate Z_t^f are modelled by a Gaussian PDF having a diagonal covariance matrix, and the first two moments for its units are computed from

$$\begin{aligned} \mathbb{E}[Z_{i,t|t-1}^f] &= \mathbb{E}[\mathbf{W}_{i,t|t-1}^f \mathbf{X}_{t|t-1}] \\ &\quad + \mathbb{E}[\mathbf{U}_{i,t|t-1}^f \mathbf{H}_{t-1|t-1}] + \mathbb{E}[\mathbf{B}_{i,t|t-1}^f], \\ \text{var}(Z_{i,t|t-1}^f) &= \text{var}(\mathbf{W}_{i,t|t-1}^f \mathbf{X}_{t|t-1}) \\ &\quad + \text{var}(\mathbf{U}_{i,t|t-1}^f \mathbf{H}_{t-1|t-1}) + \text{var}(\mathbf{B}_{i,t|t-1}^f), \end{aligned}$$

where the means and variances of the product terms $\mathbf{W}_{i,t|t-1}^f \mathbf{X}_{t|t-1}$ and $\mathbf{U}_{i,t|t-1}^f \mathbf{H}_{t-1|t-1}$ are calculated exactly using the GMA equations given in Appendix A.

From Eq. (1a), we apply the locally linearized sigmoid activation function $\tilde{\sigma}(\cdot)$ to the hidden states of the forget gate to estimate its output values $F_t = \tilde{\sigma}(Z_t^f)$. As a result, the forget gate F_t also has a diagonal covariance matrix such that $F_{i,t} \perp\!\!\!\perp F_{j,t}$. The equations for obtaining its mean vector and covariance matrix are presented in Appendix B. Similarly, the prior knowledge for the other LSTM gates can be estimated using the same procedure. Because the equations used to calculate other gates $\{I, \tilde{C}, O\}$ involve the same operations, i.e., the sum of several independent product terms, we can extend the independence assumption not only to these gates but also between all LSTM gates such that $I_{i,t} \perp\!\!\!\perp I_{j,t}$, $\tilde{C}_{i,t} \perp\!\!\!\perp \tilde{C}_{j,t}$, $O_{i,t} \perp\!\!\!\perp O_{j,t}$, and $F_t \perp\!\!\!\perp I_t \perp\!\!\!\perp \tilde{C}_t \perp\!\!\!\perp O_t$. As presented in Eq. (1e), the calculations of the cell states only involve element-wise operations from independent components. Therefore, all the cell states $C_{i,t}$ can be considered as independent and are obtained by

$$C_{i,t} = F_{i,t} C_{i,t-1} + I_{i,t} \tilde{C}_{i,t}.$$

Because Eqs. (1a)–(1d) which are used to estimate the LSTM gates $\{F_t, I_t, \tilde{C}_t, O_t\}$ do not involve the cell states

C_{t-1} , there is no direct information path between these gates and the cell states C_{t-1} other than through the hidden states H_{t-1} . Therefore, the LSTM gates at time step t and the cell states at time step $t-1$ are conditionally independent given the hidden states h_{t-1} , that is, $F_t \perp\!\!\!\perp I_t \perp\!\!\!\perp \tilde{C}_t \perp\!\!\!\perp O_t \perp\!\!\!\perp C_{t-1} | h_{t-1}$. Under this conditional independence assumption, we can apply the GMA equations to obtain the mean and variance for each cell state following

$$\begin{aligned} \mathbb{E}[C_{i,t|t-1}] &= \mathbb{E}[F_{i,t|t-1}] \cdot \mathbb{E}[C_{i,t-1|t-1}] \\ &\quad + \mathbb{E}[I_{i,t|t-1}] \cdot \mathbb{E}[\tilde{C}_{i,t|t-1}], \\ \text{var}(C_{i,t|t-1}) &= \text{var}(F_{i,t|t-1}) \cdot \text{var}(C_{i,t-1|t-1}) \\ &\quad + \text{var}(F_{i,t|t-1}) \cdot \mathbb{E}[C_{i,t-1|t-1}]^2 \\ &\quad + \text{var}(C_{i,t-1|t-1}) \cdot \mathbb{E}[F_{i,t|t-1}]^2 \\ &\quad + \text{var}(I_{i,t|t-1}) \cdot \text{var}(\tilde{C}_{i,t|t-1}) \\ &\quad + \text{var}(I_{i,t|t-1}) \cdot \mathbb{E}[\tilde{C}_{i,t|t-1}]^2 \\ &\quad + \text{var}(\tilde{C}_{i,t|t-1}) \cdot \mathbb{E}[I_{i,t|t-1}]^2. \end{aligned}$$

From Eq. (1f), the LSTM hidden states are obtained from $H_{i,t} = O_{i,t} \cdot \tanh(C_{i,t})$, where $\tanh(\cdot)$ is the locally linearized hyperbolic tangent activation function. Following the same reasoning used for the cell states, the hidden states $H_{i,t}$ are also considered as independent. Their mean and variance are estimated by

$$\begin{aligned} \mathbb{E}[H_{i,t|t-1}] &= \mathbb{E}[O_{i,t|t-1}] \cdot \mathbb{E}[\tanh(C_{i,t|t-1})], \\ \text{var}(H_{i,t|t-1}) &= \text{var}(O_{i,t|t-1}) \cdot \text{var}(\tanh(C_{i,t|t-1})) \\ &\quad + \text{var}(O_{i,t|t-1}) \cdot \mathbb{E}[\tanh(C_{i,t|t-1})]^2 \\ &\quad + \text{var}(\tanh(C_{i,t|t-1})) \cdot \mathbb{E}[O_{i,t|t-1}]^2. \end{aligned}$$

The hidden states for the fully-connected output layer are obtained from the hidden states of the last (L^{th}) LSTM layer using

$$Z_{i,t}^{(0)} = \mathbf{W}_{i,t}^{(L)} \mathbf{H}_t^{(L)} + \mathbf{B}_{i,t}^{(L)}. \quad (8)$$

Under the independence assumption, the mean and variance of each hidden state are given by

$$\begin{aligned}\mathbb{E}[Z_{i,t|t-1}^{(0)}] &= \mathbb{E}[\mathbf{W}_{i,t|t-1}^{(L)}\mathbf{H}_{t|t-1}^{(L)}] + \mathbb{E}[\mathbf{B}_{i,t|t-1}^{(L)}], \\ \text{var}(Z_{i,t|t-1}^{(0)}) &= \text{var}(\mathbf{W}_{i,t|t-1}^{(L)}\mathbf{H}_{t|t-1}^{(L)}) + \text{var}(\mathbf{B}_{i,t|t-1}^{(L)}).\end{aligned}$$

4.1.2. Backward step

The forward step presented in Section 4.1.1 can be regarded as sending information from the input layer to the output layer. In the backward step, we want to send the information in the opposite direction, from the output layer back to the input layer, in order to update the prior knowledge that has been obtained during the forward pass. Therefore, the objective of this step is to estimate the posterior PDFs for the hidden states and parameters of each layer in the network. For that, we rely on the conditional independence assumption between the hidden states of different layers in order to apply the layer-wise inference procedure, allowing to update the hidden states and parameters simultaneously within a same layer. This is essential for maintaining the computational tractability of the TAGI method. The different steps of the inference are depicted by the red arrows in Fig. 2b.

For the output layer, we compute the posterior $\mathbf{Z}_{t|t}^{(0)} \sim \mathcal{N}(\boldsymbol{\mu}_{t|t}^{(0)}, \boldsymbol{\Sigma}_{t|t}^{(0)})$ using the Gaussian conditional equations given in Eq. (5). This is analogous to updating the output $\mathbf{Z}^{(0)}$ of a TAGI-FNN as presented in Section 3.2. For the j^{th} LSTM layer, the posterior $\mathbf{H}_{t|t}^{(j)} \sim \mathcal{N}(\boldsymbol{\mu}_{t|t}^{(j)}, \boldsymbol{\Sigma}_{t|t}^{(j)})$ for its hidden states is estimated following

$$\begin{aligned}f(\mathbf{h}_{t|t}^{(j)}) &= \mathcal{N}(\boldsymbol{\mu}_{t|t}^{(j)}, \boldsymbol{\Sigma}_{t|t}^{(j)}), \\ \boldsymbol{\mu}_{t|t}^{(j)} &= \boldsymbol{\mu}_{t|t-1}^{(j)} + \mathbf{J}_{\mathbf{H}}(\boldsymbol{\mu}_{t|t}^{(j+1)} - \boldsymbol{\mu}_{t|t-1}^{(j+1)}), \\ \boldsymbol{\Sigma}_{t|t}^{(j)} &= \boldsymbol{\Sigma}_{t|t-1}^{(j)} + \mathbf{J}_{\mathbf{H}}(\boldsymbol{\Sigma}_{t|t}^{(j+1)} - \boldsymbol{\Sigma}_{t|t-1}^{(j+1)})\mathbf{J}_{\mathbf{H}}^{\top}, \\ \mathbf{J}_{\mathbf{H}} &= \text{cov}(\mathbf{H}_{t|t-1}^{(j)}, \mathbf{H}_{t|t-1}^{(j+1)})(\boldsymbol{\Sigma}_{t|t-1}^{(j+1)})^{-1},\end{aligned}\quad (9)$$

where $\text{cov}(\mathbf{H}_{t|t-1}^{(j)}, \mathbf{H}_{t|t-1}^{(j+1)})$ is the covariance between the hidden states of the j^{th} layer and the $j+1^{\text{th}}$. Note that $\mathbf{H}_{t|t}^{(j)}$ is indirectly related to $\mathbf{H}_{t|t}^{(j+1)}$ through the LSTM gates $\{\mathbf{F}_t^{(j+1)}, \mathbf{I}_t^{(j+1)}, \tilde{\mathbf{C}}_t^{(j+1)}, \mathbf{O}_t^{(j+1)}\}$ as presented in Eqs. (1a)–(1d). To obtain $\mathbf{H}_{t|t}^{(j)}$, we directly calculate the covariance matrix $\text{cov}(\mathbf{H}_{t|t-1}^{(j)}, \mathbf{H}_{t|t-1}^{(j+1)})$, bypassing the LSTM gates as shown by the red arrow from $\mathbf{h}_t^{(j)}$ to $\mathbf{h}_t^{(j-1)}$ shown in Fig. 2a. The posterior $\boldsymbol{\theta}_{t|t}^{(j)} \sim \mathcal{N}(\boldsymbol{\mu}_{t|t}^{(j)}, \boldsymbol{\Sigma}_{t|t}^{(j)})$ for the parameters can be estimated using the same Eq. (9), while replacing the variable $\mathbf{H}^{(j)}$ by $\boldsymbol{\theta}^{(j)}$. Note that the last LSTM layer (L^{th}) is connected with the output layer so that when applying Eq. (9) for inferring its hidden states and parameters, we need to calculate the covariances $\text{cov}(\mathbf{H}_{t|t-1}^{(L)}, \mathbf{Z}_{t|t-1}^{(0)})$ and $\text{cov}(\boldsymbol{\theta}_{t|t-1}^{(L)}, \mathbf{Z}_{t|t-1}^{(0)})$. The detailed formulations for the covariances $\text{cov}(\mathbf{H}_{t|t-1}^{(j)}, \mathbf{H}_{t|t-1}^{(j+1)})$ and $\text{cov}(\boldsymbol{\theta}_{t|t-1}^{(j)}, \mathbf{H}_{t|t-1}^{(j+1)})$ are given in Appendix C, whereas the ones for $\text{cov}(\mathbf{H}_{t|t-1}^{(L)}, \mathbf{Z}_{t|t-1}^{(0)})$ and $\text{cov}(\boldsymbol{\theta}_{t|t-1}^{(L)}, \mathbf{Z}_{t|t-1}^{(0)})$ have been detailed by Goulet et al. (2021).

When performing inference for a LSTM layer, we need to update the cell states in addition to updating the hidden states and parameters. To this end, we estimate the posterior $\mathbf{C}_{t|t}^{(j)} \sim \mathcal{N}(\boldsymbol{\mu}_{t|t}^{(j)}, \boldsymbol{\Sigma}_{t|t}^{(j)})$ for the cell states based

on the knowledge from the hidden states of a same LSTM layer using the same Rauch-Tung-Striebel (RTS) procedure (Rauch et al., 1965) used in Eq. (9) where the variable $\mathbf{H}^{(j)}$ is replaced by $\mathbf{C}^{(j)}$, and $\mathbf{H}^{(j+1)}$ is replaced by $\mathbf{H}^{(j)}$. This step is depicted by the red arrow from $\mathbf{h}_t^{(j)}$ to $\mathbf{c}_t^{(j)}$ shown in Fig. 2a. The diagonal cross-covariance matrix $\text{cov}(\mathbf{C}_{t|t-1}^{(j)}, \mathbf{H}_{t|t-1}^{(j)})$ between the cell and hidden states of the same LSTM layer that is required for estimating $\mathbf{C}_{t|t}^{(j)}$ is obtained following

$$\begin{aligned}\text{cov}(\mathbf{C}_{t|t-1}^{(j)}, \mathbf{H}_{t|t-1}^{(j)}) &= \text{var}(\mathbf{C}_{t|t-1}^{(j)}) \cdot \nabla_C \tilde{\text{tanh}}(\mathbb{E}[\mathbf{C}_{t|t-1}^{(j)}]) \\ &\quad \cdot \mathbb{E}[\mathbf{O}_{t|t-1}^{(j)}],\end{aligned}$$

where $\nabla_C \tilde{\text{tanh}}(\mathbb{E}[\mathbf{C}_{t|t-1}^{(j)}])$ is the gradient of the function $\tilde{\text{tanh}}(\mathbf{C})$ with respect to \mathbf{C} evaluated at the mean $\mathbb{E}[\mathbf{C}_{t|t-1}^{(j)}]$ for the j^{th} layer, cell i , at time t .

4.1.3. Smoothing for TAGI-LSTM

In SSMs, the transition model, $\mathbf{z}_t = f(\mathbf{z}_{t-1}) + \mathbf{w}_t$, describes the relationship between the hidden variables \mathbf{z} at two consecutive time steps where $f(\cdot)$ is the transition function and \mathbf{w}_t is a realization from the independent error process $\mathbf{W}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$. In other words, there is a flow of information through time from the first time step $t=0$ to the last time step $t=T$. The Kalman smoother (Rauch et al., 1965) leverages this connection in order to send information backward, and update the knowledge for past hidden variables from future information. Analogously, LSTM also has through-time connections as shown in Eqs. (1a)–(1e) where the hidden and cell states at time step $t-1$ are connected to those at t such that

$$\begin{aligned}\mathbf{c}_t &= \mathbf{g}(\mathbf{h}_{t-1}, \mathbf{c}_{t-1}, \mathbf{x}_t, \boldsymbol{\theta}), \\ \mathbf{h}_t &= \mathbf{k}(\mathbf{h}_{t-1}, \mathbf{c}_{t-1}, \mathbf{x}_t, \boldsymbol{\theta}),\end{aligned}$$

where the functions $\mathbf{g}(\cdot)$ and $\mathbf{k}(\cdot)$ are defined by Eqs. (1a)–(1e). Therefore, we can leverage these connections in order to perform smoothing for TAGI-LSTM. The backward step in Section 4.1.2 uses the smoothing equations at a single time step in order to update the parameters and hidden states through the architecture; we now use the same approach to update backward through time following the inference path depicted by the red arrows in Fig. 3.

The posterior knowledge $\mathbf{H}_{t|t}^{(j)}$ for the hidden states of the j^{th} LSTM layer obtained from the backward step only contains the past and present information, i.e., $\mathbf{y}_{1:t}$. With the smoothing procedure, we want to estimate $\mathbf{H}_{t|T}^{(j)} \sim \mathcal{N}(\boldsymbol{\mu}_{t|T}^{(j)}, \boldsymbol{\Sigma}_{t|T}^{(j)})$ containing both past and future information, i.e., the whole sequence of observations $\mathbf{y}_{1:T}$. For each LSTM layer, we go backward from the last to the first time step, and apply the RTS smoother procedure recursively following

$$\begin{aligned}\boldsymbol{\mu}_{t|T}^{(j)} &= \boldsymbol{\mu}_{t|t}^{(j)} + \mathbf{J}_{\mathbf{H}}(\boldsymbol{\mu}_{t+1|T}^{(j)} - \boldsymbol{\mu}_{t+1|t}^{(j)}), \\ \boldsymbol{\Sigma}_{t|T}^{(j)} &= \boldsymbol{\Sigma}_{t|t}^{(j)} + \mathbf{J}_{\mathbf{H}}(\boldsymbol{\Sigma}_{t+1|T}^{(j)} - \boldsymbol{\Sigma}_{t+1|t}^{(j)})\mathbf{J}_{\mathbf{H}}^{\top}, \\ \mathbf{J}_{\mathbf{H}} &= \text{cov}(\mathbf{H}_{t|t}^{(j)}, \mathbf{H}_{t+1|t}^{(j)})(\boldsymbol{\Sigma}_{t+1|t}^{(j)})^{-1},\end{aligned}\quad (10)$$

where $\text{cov}(\mathbf{H}_{t|t}^{(j)}, \mathbf{H}_{t+1|t}^{(j)})$ is the covariance between the hidden states of the j^{th} LSTM layer at time t and $t+1$. We

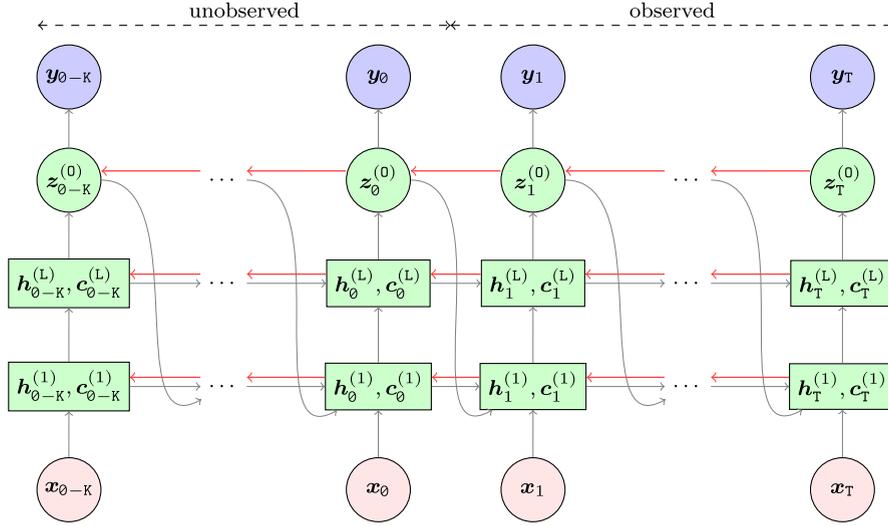


Fig. 3. Time-unfolded representation of a TAGI-LSTM network with explicit connections between times steps. Black arrows represent the network forward connections, red arrows represent the smoothing procedure. The observed data includes training data, whereas the unobserved one is data K time steps before the first training time step. This smoothing procedure allows to infer not only the smoothed estimates for hidden and cell states, and output during training time steps, but also the past ones before the training history.

can also use Eq. (10) to obtain the smoothed estimates $\mathbf{C}_{t|T}^{(j)}$ for the cell states, and $\mathbf{Z}_{t|T}^{(0)}$ for the hidden states of the output layer, by replacing the hidden states \mathbf{H} by the relevant variable. Appendix D presents the derivations of the covariances $\text{cov}(\mathbf{H}_{t|t}, \mathbf{H}_{t+1|t})$, $\text{cov}(\mathbf{C}_{t|t}, \mathbf{C}_{t+1|t})$, and $\text{cov}(\mathbf{Z}_{t|t}, \mathbf{Z}_{t+1|t})$ which are required to estimate $\mathbf{H}_{t|T}^{(j)}$, $\mathbf{C}_{t|T}^{(j)}$ and $\mathbf{Z}_{t|T}^{(0)}$. Note that because the network parameters θ are assumed to be constant through time, the smoother has no effect on them.

With the smoothing procedure, we start from the last time step T and either stop the procedure at $t = 0$ to infer the initial hidden and cell states, $\mathbf{H}_{0|T}$ and $\mathbf{C}_{0|T}$, or we can even go back K time steps before the first training time step for obtaining the smoothed estimates for the hidden and cell states as presented in Fig. 3. We then can use them to estimate the unobserved past observations $\mathbf{y}_{0-K:0}$. This means that we can use a single TAGI-LSTM model to estimate both future observations after the last training time and past observations before the first training time.

4.2. Coupling TAGI-LSTM and SSMs

In this section, we present the methodology to couple probabilistically the TAGI-LSTM presented in Section 4.1 with SSMs to create the TAGI-LSTM/SSM hybrid model. We consider a structured state-space model with additive errors where the hidden state vector \mathbf{z} consists in two sets of hidden states \mathbf{z}^B and $z^{(0)}$, $\mathbf{z} = [\mathbf{z}^B \ z^{(0)}]^T$. The baseline hidden state vector \mathbf{z}^B models the time series' baseline patterns such as level and trend, and is associated with the linear dynamical system

$$\mathbf{z}_t^B = \mathbf{A} \mathbf{z}_{t-1}^B + \mathbf{w}_t, \quad (11)$$

where \mathbf{A} is the transition matrix for the baseline hidden states, \mathbf{w}_t is a realization from the independent error process $\mathbf{w}_t : \mathbf{W} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$. The hidden state $z^{(0)}$ is

employed to model the repeated patterns such as seasonalities that are present in the data. Contrarily to the baseline hidden states which follow the transition model described in Eq. (11), the pattern hidden state is predicted using the output node of a TAGI-LSTM network

$$z_t^{(0)} = \text{TAGI-LSTM}(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}, \theta), \quad (12)$$

where \mathbf{x}_t is a vector of input covariates, \mathbf{h}_{t-1} are the TAGI-LSTM hidden states at the previous time step, and θ are the neural network parameters. Note that similarly to the baseline hidden states \mathbf{Z}^B , the output $Z^{(0)}$ is also Gaussian.

At a time t , the joint prior knowledge for the baseline hidden states given all the past data $\mathbf{y}_{1:t-1}$ is described by a Gaussian random vector $\mathbf{Z}_{t|t-1}^B \sim \mathcal{N}(\boldsymbol{\mu}_{t|t-1}^B, \boldsymbol{\Sigma}_{t|t-1}^B)$ which is obtained by propagating the uncertainty associated with the posterior at $t - 1$, i.e., $\mathbf{Z}_{t-1|t-1}^B$ through the transition model described in Eq. (11). The marginal prior at t for the pattern hidden state is a Gaussian random variable $Z_{t|t-1}^{(0)} \sim \mathcal{N}(\mu_{t|t-1}^{(0)}, \sigma_{t|t-1}^{(0)})$ obtained using the TAGI-LSTM one-step-ahead prediction presented in Section 4.1. In order to obtain the joint prior knowledge at t for both the baseline and the pattern hidden states, we need their cross-covariance $\text{cov}(\mathbf{Z}_{t|t-1}^B, Z_{t|t-1}^{(0)})$ for which no exact closed-form analytical expression is available. In practice, these covariance terms are non-zero, yet they are typically small such that the correlation coefficients between hidden states $\rho < 0.01$. Therefore, we rely on the simplifying hypothesis that $\mathbf{Z}_{t|t-1}^B$ and $Z_{t|t-1}^{(0)}$ are independent so that $\text{cov}(\mathbf{Z}_{t|t-1}^B, Z_{t|t-1}^{(0)}) = \mathbf{0}$. The validity of that hypothesis is investigated through experiments provided in Appendix E.

We can form the joint prior knowledge between the hidden states $\mathbf{Z}_{t|t-1} = [\mathbf{Z}_{t|t-1}^B \ Z_{t|t-1}^{(0)}]^T$ and the observation Y_t using the linear observation model

$$y_t = \mathbf{C} \mathbf{z}_t + v_t, \quad (13)$$

where v_t is a realization from the independent and identically distributed (i.i.d.) error $v_t : V \sim \mathcal{N}(0, \sigma_v^2)$. In the observation model presented in Eq. (13), the observation vector indicates which hidden state is observable such that $\mathbf{C} = [\mathbf{c}^B \ c^{(0)}]$, where $c^{(0)} = 1$ and the composition of \mathbf{c}^B depends on the dynamic system chosen for the baseline. With that joint prior information, we can use the Gaussian conditional equations in order to obtain the joint posterior $\mathbf{Z}_{t|t}$ which can be broken down into the posteriors $\mathbf{Z}_{t|t}^B$ and $Z_{t|t}^{(0)}$. The later is then used by the TAGI-LSTM method in order to infer both the Gaussian posterior for their model parameters $\theta \sim \mathcal{N}(\boldsymbol{\mu}_{t|t}^\theta, \mathbf{I}\sigma_{t|t}^\theta)$ and their hidden states $\mathbf{H}_t \sim \mathcal{N}(\boldsymbol{\mu}_{t|t}^H, \mathbf{I}\sigma_{t|t}^H)$.

The covariance $\text{cov}(\mathbf{Z}_{t|t}, \mathbf{Z}_{t+1|t})$ between the hidden states at two consecutive time steps is calculated by

$$\text{cov}(\mathbf{Z}_{t|t}, \mathbf{Z}_{t+1|t}) = \begin{bmatrix} \boldsymbol{\Sigma}_{t|t}^B \mathbf{A}^\top & \mathbf{0} \\ \mathbf{0} & \text{cov}(Z_{t|t}^{(0)}, Z_{t+1|t}^{(0)}) \end{bmatrix},$$

where $\text{cov}(Z_{t|t}^{(0)}, Z_{t+1|t}^{(0)})$ is calculated using the smoothing procedure presented in Section 4.1.3. Given this covariance, the Kalman smoother can be performed analytically for the TAGI-LSTM/SSM model. As it is the case for the TAGI model presented in Section 3.2, we train the model and learn the network parameters θ over multiple epochs. For each epoch, we perform a forward pass over time and a backward one using the Kalman smoother.

4.2.1. Exponential smoothing component

Using the local linear trend component (Durbin & Koopman, 2001) for the baseline hidden states \mathbf{z}^B , the TAGI-LSTM/SSM model presented in Section 4.2 can only model data which contains a linear trend. In this section, we formulate a parameter-free exponential smoothing component which can be used with the TAGI-LSTM/SSM to automatically detrend data containing complex long-term patterns.

Recall the state-space form of a simple exponential smoothing model (Hyndman & Athanasopoulos, 2018), but now we consider the smoothing parameter z^α and the observation error z^V as hidden states. The model is defined by the following equations

$$\begin{aligned} \text{Observation equation:} \quad & y_t = z_t^E + z_t^V, & (14) \\ \text{Transition equations:} \quad & z_t^E = z_{t-1}^E + \bar{z}_{t-1}^\alpha z_{t-1}^V, \\ & z_t^\alpha = z_{t-1}^\alpha, \\ & z_t^V = v_t, & (15) \end{aligned}$$

where z_t^E is the exponential smoothing hidden state, v_t is a realization from the i.i.d. error $v_t : V \sim \mathcal{N}(0, \sigma_v^2)$, and y_t is the observation. Because the smoothing parameter takes a value between zero and one, we apply the locally linearized sigmoid function that is already used for TAGI in Section 3.2 to the hidden state z_t^α . The transformed variable $\bar{z}_t^\alpha = \tilde{\sigma}(z_t^\alpha)$ is assumed to follow a Gaussian PDF where its moments are obtained by applying the equations presented in Appendix B. The product $\bar{z}_{t-1}^\alpha z_{t-1}^V$ uses the previous step's error z_{t-1}^V to adjust the predicted exponential smoothing hidden state z_t^E . We rely on the covariance $\text{cov}(\bar{Z}_{t-1|t-1}^\alpha \bar{Z}_{t-1|t-1}^V, Z_{t|t-1}^\alpha)$ in order to update the hidden state z_t^α . Note that if instead

we would use the product $\bar{z}_{t-1}^\alpha z_t^V$ to adjust z_t^E , the covariance $\text{cov}(\bar{Z}_{t-1|t-1}^\alpha \bar{Z}_{t-1|t-1}^V, Z_{t|t-1}^\alpha)$ would be zero, and thus we could not update z_t^α .

The transition model as given in Eq. (15) is nonlinear, however, Deka et al. (2022) showed that this nonlinear model can be reformulated as a linear dynamic model by creating a new hidden state $z_t^N = \bar{z}_{t-1}^\alpha z_{t-1}^V$ which is assumed to follow a Gaussian PDF where its moments are calculated exactly using the GMA equations given in Appendix A. The baseline hidden state vector $\mathbf{z}^B = [z^E \ z^\alpha \ z^V]^\top$ at time $t - 1$ is thus assumed to be Gaussian. The augmented baseline hidden state vector $\tilde{\mathbf{z}}^B = [\mathbf{z}^B \ z^N]^\top$ is defined by

$$\tilde{\mathbf{z}}_{t-1|t-1}^B \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_{t-1|t-1}^B, \tilde{\boldsymbol{\Sigma}}_{t-1|t-1}^B),$$

where the components of $\tilde{\boldsymbol{\mu}}_{t-1|t-1}^B$ and $\tilde{\boldsymbol{\Sigma}}_{t-1|t-1}^B$ are presented in Appendix F.

The nonlinear Eq. (15) now can be written in a linear form as

$$\tilde{\mathbf{z}}_t^B = \mathbf{A} \tilde{\mathbf{z}}_{t-1}^B + \mathbf{w}_t, \quad \mathbf{w}_t : \mathbf{W} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}), \quad (16)$$

where \mathbf{A} is the transition matrix and \mathbf{Q} is the covariance matrix. The prior for the augmented baseline hidden state vector is given by

$$\tilde{\mathbf{z}}_{t|t-1}^B \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_{t|t-1}^B, \tilde{\boldsymbol{\Sigma}}_{t|t-1}^B), \quad (17)$$

where $\tilde{\boldsymbol{\mu}}_{t|t-1}^B = \mathbf{A} \tilde{\boldsymbol{\mu}}_{t-1|t-1}^B$ and $\tilde{\boldsymbol{\Sigma}}_{t|t-1}^B = \mathbf{A} \tilde{\boldsymbol{\Sigma}}_{t-1|t-1}^B \mathbf{A}^\top + \mathbf{Q}$. The observation model given in Eq. (14) is written in a matrix form as

$$y_t = \mathbf{C} \tilde{\mathbf{z}}_t^B, \quad (18)$$

where \mathbf{C} is the observation matrix. The model matrices \mathbf{A} , \mathbf{Q} and \mathbf{C} are given in Appendix F. Using Eqs. (16)–(18), we can now apply the Kalman filter to have a parameter-free component that can capture complex long-term patterns in data. When we couple this new exponential smoothing component with TAGI-LSTM, it enables to automatically detrend the data as shown by the experiments in the next section.

5. Experiments

In this section, we first conduct a qualitative experiment using synthetic time series to verify the TAGI-LSTM model's capability to accurately retrieve the ground truth. Next, a quantitative analysis is conducted to compare the performance of the same LSTM architecture trained with either TAGI or gradient-based approaches. The Electricity and Traffic datasets (Yu et al., 2016) datasets are chosen for this comparison because they are stationary, i.e., each time series displays a constant baseline so that the LSTM models can analyze them directly without requiring data preprocessing to remove trends. Finally, after verifying TAGI-LSTM's capabilities on simulated and stationary datasets, we compare the performance of the TAGI-LSTM/SSM hybrid model with other benchmark methods on the non-stationary Tourism (monthly and quarterly) and M4 (hourly) datasets, displaying either linear or complex nonlinear trends. These datasets

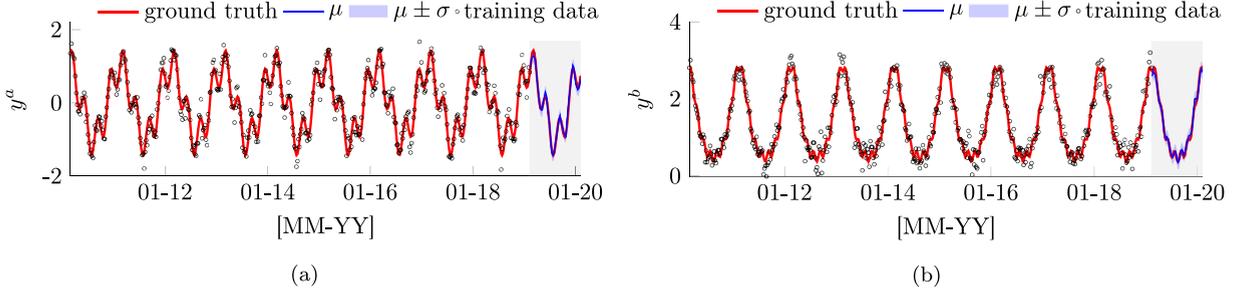


Fig. 4. Test predictions from TAGI-LSTM models for synthetic data generated using (a) Eq. (19a), and (b) Eq. (19b). The grey shaded area presents the forecast period, red line presents the ground truth, blue line presents test predictions along with $\pm\sigma$ confidence intervals (shade), black dots present training data.

are intended to demonstrate that the hybrid TAGI-LSTM/SSM model is well-suited to handle trended data without the need for preprocessing. Note that one cannot use plain LSTMs on these benchmarks without preprocessing the data beforehand to remove trends.

5.1. Verification of TAGI-LSTM on synthetic data

For the verification experiment, we generate ten years of weekly data (520 data points) from the following functions

$$y_t^a = \sin\left(\frac{2\pi t}{365.22}\right) + 0.5\sin\left(\frac{2\pi t}{365.22/4}\right) + v_t, \quad (19a)$$

$$y_t^b = \exp\left[\sin\left(\frac{2\pi t}{365.22}\right)\right] + \left[0.5\sin\left(\frac{2\pi t}{365.22/4}\right)\right]^2 + v_t, \quad (19b)$$

as depicted in Fig. 4, where $v_t : V \sim \mathcal{N}(0, 0.2^2)$ and t is the timestamp. We train TAGI-LSTM models using the first 9 years of data, and make multi-step-ahead predictions for the last year. Fig. 4 presents the ground truth which does not contain the error term v_t , the training data, as well as the test predictions and their $\pm\sigma$ uncertainties from our models. The test set's predictions show that the TAGI-LSTM can accurately retrieve the ground truth for these synthetic data.

5.2. Comparison of TAGI-LSTM with gradient-based LSTM networks on stationary data

In this section, we perform a quantitative experiment to compare the predictive performance of the TAGI-LSTM in comparison to two gradient-based LSTMs, i.e., deterministic and variational ones (Gal & Ghahramani, 2016b), on the same network architecture. We perform the comparison on the Electricity dataset (Yu et al., 2016) containing 370 hourly time series of electricity consumption, and the Traffic one (Yu et al., 2016) including 963 hourly time series of occupancy rate for different car lanes in the San Francisco bay area. These datasets are specifically selected because they are stationary which allows the LSTM models to analyze them directly after standardization without requiring ad-hoc procedures to

remove trends. As these datasets experience hourly and daily seasonal patterns, we include hour-of-the-day and day-of-the-week as time features in the covariate vector \mathbf{x}_t . We apply the moving window approach proposed by Bandara et al. (2020) so that the input covariate vector includes a lookback window of posterior values for $\mathbf{z}_{t-W:t-1}^{(0)}$ where W is the window's length. Detailed explanation for this approach can be found in Hewamalage et al. (2021). The task consists in providing predictions for seven consecutive days using the rolling window operation described in Yu et al. (2016). Using this operation, the model makes multi-step-ahead predictions for a 24-hour-window. Then, the observations for this window are made available so that they can be used to predict the next one. There are seven windows required to cover the test period. In order to make multi-step-ahead predictions, we apply the single-output forecasting procedure and recursively make one-step-ahead predictions.

The data is divided into training and test sets using three different splits (Oreshkin et al., 2020) as detailed in Appendix G. A window of data points at the end of the training data is reserved for validation, and it is not used for training. We train our models with 50 epochs, then identify the optimal epoch which maximizes the log-likelihood for the validation set, and obtain the network parameters at this optimal epoch. We build a separate LSTM model for each time series where all of the models use the same network architecture as used by Salinas et al. (2020). The standard deviation for the error σ_V is a hyper-parameter in our TAGI-LSTM model that is common to all time series. For variational LSTM, we perform 50 dropout passes and calculate the predictive mean and variance. The architecture and hyper-parameters used in this experiment are given in Appendix H.

The ρ -quantile loss, $\rho \in (0, 1)$, measures the accuracy of the predictive distribution at a specific quantile (Rangapuram et al., 2018), and is defined by

$$QL_\rho(\mathbf{y}, \hat{\mathbf{y}}^\rho) = 2 \frac{\sum_{i,t} [\rho \cdot \max(y_{i,t} - \hat{y}_{i,t}^\rho, 0) + (1 - \rho) \cdot \max(\hat{y}_{i,t}^\rho - y_{i,t}, 0)]}{\sum_{i,t} |y_{i,t}|},$$

where $y_{i,t}$ is the observation at time t of the i th time series and $\hat{y}_{i,t}^\rho$ is the corresponding predicted ρ -quantile. Note

Table 1

Comparison between TAGI-LSTM, deterministic and variational LSTMs on Electricity and Traffic datasets. $p50$ and $p90$ -loss metrics for test sets are calculated using the rolling window operation. Results are obtained by averaging the forecasts over five independent runs. Deterministic LSTM provides only point forecasts so that we report only the $p50$ -loss. Three train/test splits are used following Oreshkin et al. (2020). Bold fonts indicate the best results.

Split	Electricity						Traffic					
	2014-03-31		2014-09-01		last 7 days		2008-01-14		2008-06-15		last 7 days	
	$p50$	$p90$										
TAGI-LSTM	0.080	0.058	0.066	0.053	0.152	0.095	0.337	0.276	0.169	0.158	0.102	0.130
LSTM	0.086	–	0.077	–	0.159	–	0.319	–	0.158	–	0.098	–
Variational LSTM	0.080	0.056	0.064	0.052	0.160	0.100	0.323	0.261	0.162	0.154	0.123	0.133

that the $p50$ -loss is equal to the Normalized deviation (ND) metric reported in Oreshkin et al. (2020), Salinas et al. (2020), Wang et al. (2019), Yu et al. (2016), and the ND metric was original used to compare these datasets by Yu et al. (2016). The results reported in Table 1 are obtained by averaging the predictions from five independent runs with different initial seeds. Both TAGI-LSTM and variational LSTM provide predictions along with the associated uncertainties. Therefore, we report the $p50$ and $p90$ -losses for these methods. The deterministic LSTM provides point estimates so that we report only the $p50$ -loss. The results from Table 1 show that our method matches the performance of the two LSTMs trained with gradient-based optimization. The deterministic LSTM provides the best $p50$ -loss for the Traffic dataset on all splits which means that it provides the most accurate point forecasts. However, unlike TAGI-LSTM and variational LSTM, this method cannot provide predictive uncertainties. The variational LSTM provides slightly better $p90$ -losses compared to TAGI-LSTM in four out of six scenarios. However, the result of this method varies with respect to the number of dropout passes during test time. By contrast, TAGI-LSTM can analytically estimate both the predictions and the associated uncertainties without repeatedly performing dropout.

Robustness verifications for the TAGI-LSTM method are presented in Appendix I, where we compare the predictive performances of four different TAGI-LSTM models, each with a different number of hidden layers, ranging from one to four. The root mean square error (RMSE) and mean absolute scaled error (MASE) metrics for this experiment are presented in Appendix J.

5.3. Comparison of TAGI-LSTM/SSM with state-of-the-art methods on non-stationary data

After verifying TAGI-LSTM's capabilities on simulated and stationary datasets, this section validates the capacity of the hybrid TAGI-LSTM/SSM model to match the performance of state-of-the-art methods on the non-stationary Tourism (monthly and quarterly) (Athanasopoulos et al., 2011) and M4 (hourly) (Makridakis et al., 2018) datasets, containing either linear or complex nonlinear trends. We excluded the plain LSTM models from this experiment because they would have required further ad-hoc preprocessing in order to remove the trends from data before analyzing. The monthly Tourism dataset contains 366 time

series, the quarterly Tourism dataset contains 427 time series, and the hourly M4 dataset has 414 time series.

The baseline hidden state vector contains a level (L), a trend (T), and an exponential smoothing component, $\mathbf{z}^B = [z^L \ z^T \ z^E \ z^\alpha \ z^V]^T$, and the pattern hidden state $z^{(0)}$ models the seasonality. The model matrices are given in Appendix K. For the time features included in the covariate vector \mathbf{x}_t , we include hour-of-the-day and day-of-the-week for the hourly dataset, month-of-the-year for the monthly dataset, and quarter-of-the-year for the quarterly dataset. We build a separate model for each time series where all models share the same network architecture which is detailed in Appendix H. The standard deviation for the noise σ_V is learnt using the AGVI method (Deka, 2022; Deka & Goulet, 2023). We train our models on multiple epochs where the initial hidden states $\mu_0^{(i+1)}$ and $\Sigma_0^{(i+1)}$ at the $i+1^{\text{th}}$ epoch are the smoothed estimates $\mu_{0|T}^{(i)}$ and $\Sigma_{0|T}^{(i)}$ at the i^{th} epoch with T being the last training time.

We compare our method with the statistical methods ARIMA (Box et al., 1994) and ETS (Hyndman et al., 2008); the hybrid methods DeepState (Rangapuram et al., 2018) and ES-RNN (Smyl, 2020); and the pure neural networks DeepAR (Salinas et al., 2020) and N-Beats (Oreshkin et al., 2020). The two hybrid methods are similar to our model since all of them combine SSMs or exponential smoothing and LSTMs, and ES-RNN is the state-of-the-art hybrid method for the M4 dataset (Makridakis et al., 2020). DeepAR is a pure LSTM-based method, whereas N-Beats is the state-of-the-art neural network method for these datasets. Among them, ARIMA, ETS and our method adopt a local setup such that a separate model is built for each time series, and there is no information shared between models; DeepAR and N-Beats adopt a global setup, fitting a single model for multiple time series; whereas DeepState and ES-RNN use a mixed setup, using a global LSTM to learn across multiple time series and a local SSM or exponential smoothing model for each time series. Table 2 reports the $p50$ and $p90$ metrics for our TAGI-LSTM/SSM along with benchmark models. For a fair comparison, we report the values from the original papers (Oreshkin et al., 2020; Rangapuram et al., 2018; Skepetari, 2020). The results show that our method exceeds the performance of the ARIMA, ETS, DeepAR, and DeepState in both $p50$ and $p90$ metrics for all datasets except for the $p50$ loss of the monthly Tourism one. This shows the data-efficiency of our method such that we can adopt a local setup, that is building a separate

Table 2

$p50$ and $p90$ -loss performances on the Tourism (monthly and quarterly) and M4 (hourly) test sets. The results for N-Beats are obtained from Oreshkin et al. (2020), the one for ES-RNN is calculated by us from the submission downloaded from the M4 GitHub repository (Skepetari, 2020), and the results for other baseline methods are obtained from Rangapuram et al. (2018). The results for TAGI-LSTM/SSM are obtained by averaging the forecasts over three independent runs with different initial seeds. N-Beats provides point forecasts so that only the $p50$ -loss is reported. Bold fonts indicate the best results.

	Tourism (Monthly)		Tourism (Quarterly)		M4 (Hourly)	
	$p50$	$p90$	$p50$	$p90$	$p50$	$p90$
ARIMA	0.100	0.058	0.124	0.062	0.052	0.035
ETS	0.093	0.054	0.105	0.055	0.054	0.027
DeepState	0.138	0.067	0.098	0.047	0.044	0.027
ES-RNN	–	–	–	–	0.039	–
DeepAR	0.107	0.059	0.110	0.062	0.090	0.030
N-Beats	0.097	–	0.077	–	0.023	–
TAGI-LSTM/SSM	0.102	0.053	0.073	0.041	0.042	0.021

model for each time series, and provide better results than other models using local, global and mixed setups. N-Beats provides point forecasts so that we report only the $p50$ -loss. The ES-RNN and N-Beats methods provide superior point forecasts compared to our method. Note that these methods use more advanced neural network architecture such as dilated LSTMs with attention mechanism and residual connections. In addition, unlike our method, ES-RNN and N-Beats can provide interpretable components such as level, trend and seasonality, but are unable to provide their associated uncertainties.

In order to extend the comparisons, we re-obtained the results for the DeepState, DeepAR and N-Beats methods using the *gluonTS* library (Alex et al., 2020), for ARIMA and ETS models using the *forecast* library (Hyndman & Khandakar, 2008), and for the ES-RNN method using the *ESRNN* one (Gutierrez et al., 2021). Appendix L presents the comparisons on RMSE and MASE metrics, as well as the multiple comparisons with the best (MCB) test (Konig et al., 2005). These additional results further confirm that the TAGI-LSTM/SSM method proposed has a performance that is competitive with other state-of-the-art methods.

Fig. 5 shows an example of predictions on the test set, the hidden states, as well as their associated uncertainties for the time series #30 of the Tourism (monthly) dataset. These results confirm that our model can separate the long-term linear pattern which is captured by the level hidden state z^L , the long-term nonlinear pattern which is captured by the exponential smoothing hidden state z^E , and the seasonality modelled by the TAGI-LSTM. Additional examples for other time series are presented in Appendix M.

6. Conclusion

The new mathematical formulations proposed in this paper enable using the Tractable Approximate Gaussian Inference (TAGI) method with the LSTM neural network architecture. The approach allows estimating analytically the posterior mean vectors and diagonal covariance matrices for the hidden states and model parameters using approximate Bayesian inference. The experiments performed showed that for the same network architecture,

our TAGI-LSTM models provide on-par performance compared to the deterministic and variational LSTMs trained with backpropagation and gradient descent. With the method presented in this paper, one can now couple LSTM neural networks with the existing components borrowed from SSMs. We have demonstrated, for example, how we can couple TAGI-LSTM with the existing SSM's level and trend components, as well as the parameter-free exponential smoothing one in order to automatically detrend time series while providing interpretable results. The experimental results on the quarterly and monthly Tourism, and hourly M4 datasets have shown that our hybrid model is competitive with state-of-the-art methods.

CRedit authorship contribution statement

Van-Dai Vuong: Conceptualization, Methodology, Software, Investigation, Writing - Original Draft, Writing - Review & Editing. **Luong-Ha Nguyen:** Conceptualization, Software. **James-A. Goulet:** Conceptualization, Writing - Review & Editing, Supervision, Project administration, Funding acquisition.

Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Van-Dai Vuong, Luong-Ha Nguyen, and James-A. Goulet report financial support was provided by Natural Sciences and Engineering Research Council of Canada (NSERC) and Hydro-Québec.

Acknowledgments

This project is funded by Natural Sciences and Engineering Research Council of Canada (NSERC) and Hydro Québec under grant #ALLRP 578556-22.

Appendix A. Supplementary material

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.ijforecast.2024.04.002>.

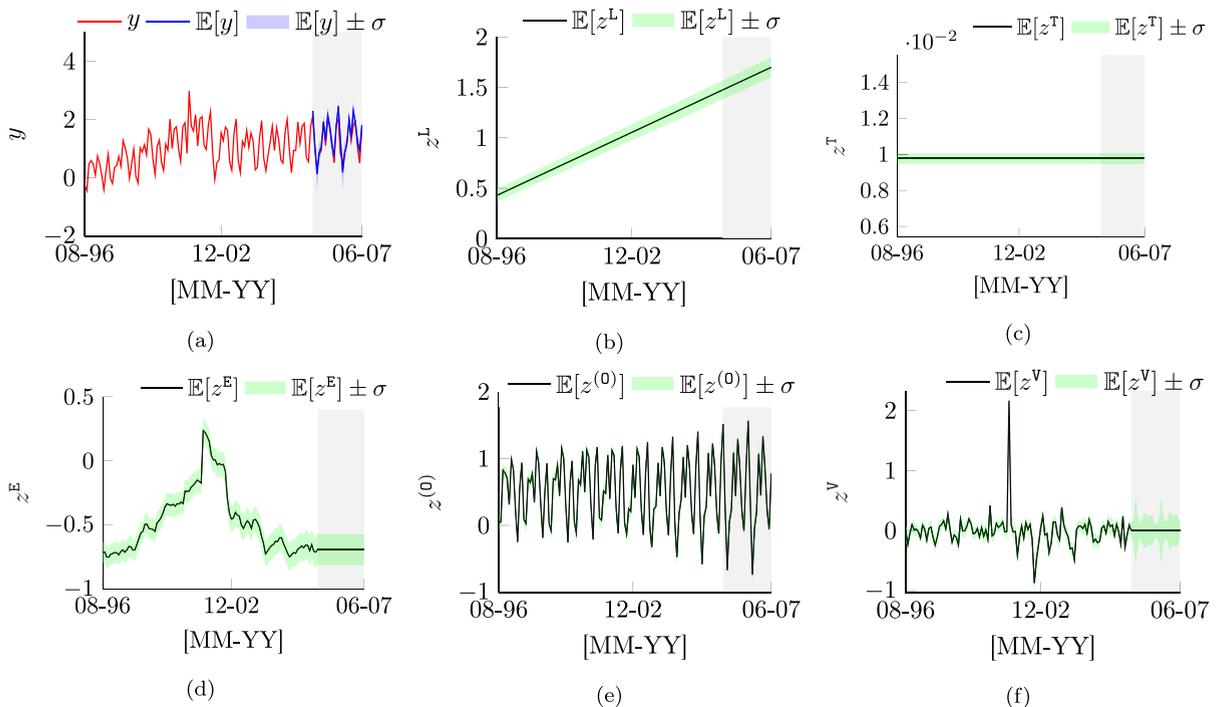


Fig. 5. An example of interpretable results provided by TAGI-LSTM/SSM using the proposed exponential smoothing component for the time series #30 of the Tourism (monthly) dataset. Values are presented in standardized space. Red line presents observations, blue line presents test predictions, black lines present components along with 68% confidence intervals (shade). The grey shaded areas present the forecast period. (a) Predictions on test set (b) level, (c) trend, (d) exponential smoothing, (e) seasonality modelled by TAGI-LSTM, and (f) error hidden state. The results show that the model can separate the linear long-term pattern which is captured by the level hidden state and the nonlinear long-term one which is captured by the exponential smoothing component.

References

- Alexandrov, A., Benidis, K., Bohlke-Schneider, M., Flunkert, V., Gasthaus, J., Januschowski, T., Maddix, D. C., Rangapuram, S., Salinas, D., Schulz, J., Stella, L., Türkmen, A. C., & Wang, Y. (2020). GluonTS: Probabilistic and neural time series modeling in Python. *Journal of Machine Learning Research*, 21(116), 1–6.
- Athanasopoulos, G., Hyndman, R. J., Song, H., & Wu, D. C. (2011). The tourism forecasting competition. *International Journal of Forecasting*, 27(3), 822–844.
- Bandara, K., Bergmeir, C., & Smyl, S. (2020). Forecasting across time series databases using recurrent neural networks on groups of similar series: A clustering approach. *Expert Systems with Applications*, 140, 112896.
- Blundell, C., Cornebise, J., Kavukcuoglu, K., & Wierstra, D. (2015). Weight uncertainty in neural network. In *Proceedings of the 32nd international conference on machine learning*, vol. 37 (pp. 1613–1622).
- Box, G. E. P., Reinsel, G. C., & Jenkins, G. M. (1994). *Time series analysis: forecasting and control*. Prentice-Hall.
- Chien, J.-T., & Ku, Y.-C. (2016). Bayesian recurrent neural network for language modeling. *IEEE Transactions on Neural Networks and Learning Systems*, 27(2), 361–374.
- Deka, B. (2022). *Analytical Bayesian parameter inference for probabilistic models with engineering applications* (Ph.D. thesis), Canada: Polytechnique Montréal.
- Deka, B., & Goulet, J.-A. (2023). Approximate Gaussian variance inference for state-space models. *International Journal of Adaptive Control and Signal Processing*, 37(11), 2934–2962.
- Deka, B., Nguyen, L. H., Amiri, S., & Goulet, J.-A. (2022). The Gaussian multiplicative approximation for state-space models. *Structural Control and Health Monitoring*, 29(3), e2904.
- Durbin, J., & Koopman, S. (2001). *Time series analysis by state space methods*. Oxford University Press.
- Efron, B. (2010). *Large-scale inference: Empirical Bayes methods for estimation, testing, and prediction*. Cambridge University Press.
- Fortunato, M., Blundell, C., & Vinyals, O. (2019). Bayesian recurrent neural networks. arXiv arXiv:1704.02798.
- Fracaro, M., Sønderby, S. K., Paquet, U., & Winther, O. (2016). Sequential neural models with stochastic layers. In *Proceedings of the 30th international conference on neural information processing systems* (pp. 2207–2215).
- Gal, Y., & Ghahramani, Z. (2016a). Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of the 33rd international conference on international conference on machine learning*, vol. 48 (pp. 1050–1059).
- Gal, Y., & Ghahramani, Z. (2016b). A theoretically grounded application of dropout in recurrent neural networks. In *Proceedings of the 30th international conference on neural information processing systems*, vol. 29 (pp. 1027–1035).
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- Goulet, J.-A., Nguyen, L. H., & Amiri, S. (2021). Tractable approximate Gaussian inference for Bayesian neural networks. *Journal of Machine Learning Research*, 22, 1–23.
- Gutierrez, K., Challu, C., Garza, F., & Mergenthaler, M. (2021). Pytorch implementation of the ES-RNN algorithm. GitHub repository. https://github.com/kdgutier/esrnn_torch. (Accessed 17 March 2023).
- Hernandez-Lobato, J. M., & Adams, R. (2015). Probabilistic backpropagation for scalable learning of Bayesian neural networks. In *Proceedings of the 32nd international conference on machine learning*, vol. 37 (pp. 1861–1869).
- Hewamalage, H., Bergmeir, C., & Bandara, K. (2021). Recurrent neural networks for time series forecasting: Current status and future directions. *International Journal of Forecasting*, 37(1), 388–427.

- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9, 1735–1780.
- Hyndman, R., & Athanasopoulos, G. (2018). *Forecasting: Principles and practice* (2nd ed.). OTexts.
- Hyndman, R. J., & Khandakar, Y. (2008). Automatic time series forecasting: the forecast package for R. *Journal of Statistical Software*, 26(3), 1–22.
- Hyndman, R., Koehler, A., Ord, J., & Snyder, R. (2008). *Forecasting with exponential smoothing: The state space approach*. Springer.
- Januschowski, T., Gasthaus, J., Wang, Y., Salinas, D., Flunkert, V., Bohlke-Schneider, M., & Callot, L. (2020). Criteria for classifying forecasting methods. *International Journal of Forecasting*, 36(1), 167–177.
- Koning, A. J., Franses, P. H., Hibon, M., & Stekler, H. (2005). The M3 competition: Statistical tests of the results. *International Journal of Forecasting*, 21(3), 397–409.
- Krishnan, R. G., Shalit, U., & Sontag, D. (2017). Structured inference networks for nonlinear state space models. In *Proceedings of the thirty-first AAAI conference on artificial intelligence* (pp. 2101–2109).
- MacKay, D. J. C. (1992). A practical Bayesian framework for backpropagation networks. *Neural Computation*, 4(3), 448–472.
- Makridakis, S., Spiliotis, E., & Assimakopoulos, V. (2018). The M4 competition: Results, findings, conclusion and way forward. *International Journal of Forecasting*, 34(4), 802–808.
- Makridakis, S., Spiliotis, E., & Assimakopoulos, V. (2020). The M4 Competition: 100,000 time series and 61 forecasting methods. *International Journal of Forecasting*, 36(1), 54–74.
- Moon, T., Choi, H., Lee, H., & Song, I. (2015). RNNDROP: A novel dropout for RNNs in ASR. In *2015 IEEE workshop on automatic speech recognition and understanding* (pp. 65–70).
- Nguyen, L.-H., & Goulet, J.-A. (2021). Analytically tractable inference in deep neural networks. In *The 4th workshop on tractable probabilistic modeling*.
- Nguyen, L.-H., & Goulet, J.-A. (2022). Analytically tractable hidden-states inference in Bayesian neural networks. *Journal of Machine Learning Research*, 23(50), 1–33.
- Oreshkin, B. N., Carpov, D., Chapados, N., & Bengio, Y. (2020). N-BEATS: Neural basis expansion analysis for interpretable time series forecasting. In *International conference on learning representations 2020*.
- Pham, V., Bluche, T., Kermorvant, C., & Louradour, J. (2014). Dropout improves recurrent neural networks for handwriting recognition. In *14th international conference on frontiers in handwriting recognition* (pp. 285–290).
- Rangapuram, S. S., Seeger, M. W., Gasthaus, J., Stella, L., Wang, Y., & Januschowski, T. (2018). Deep state space models for time series forecasting. In *Advances in neural information processing systems: vol. 31*.
- Rauch, H. E., Tung, F., & Striebel, C. T. (1965). Maximum likelihood estimates of linear dynamic systems. *American Institute of Aeronautics and Astronautics*, 3(8), 1445–1450.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323, 533–536.
- Salinas, D., Flunkert, V., Gasthaus, J., & Januschowski, T. (2020). DeepAR: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3), 1181–1191.
- Skeptari, E. (2020). Data, benchmarks, and methods submitted to the M4 forecasting competition. <https://github.com/Mcompetitions/M4-methods>. (Accessed 27 February 2023).
- Smyl, S. (2020). A hybrid method of exponential smoothing and recurrent neural networks for time series forecasting. *International Journal of Forecasting*, 36(1), 75–85.
- Sun, S., Chen, C., & Carin, L. (2017). Learning Structured Weight Uncertainty in Bayesian Neural Networks. In *Proceedings of the 20th international conference on artificial intelligence and statistics, vol. 54* (pp. 1283–1292).
- Wang, Y., Smola, A., Maddix, D., Gasthaus, J., Foster, D., & Januschowski, T. (2019). Deep factors for forecasting. In *Proceedings of the 36th international conference on machine learning, vol. 97* (pp. 6607–6617).
- Wu, A., Nowozin, S., Meeds, E., Turner, R. E., Hernandez-Lobato, J. M., & Gaunt, A. L. (2019). Deterministic variational inference for robust Bayesian neural networks. In *International conference on learning representations 2019*.
- Yu, H.-F., Rao, N., & Dhillon, I. S. (2016). Temporal regularized matrix factorization for high-dimensional time series prediction. In *Advances in neural information processing systems: vol. 29*.
- Zaheer, M., Ahmed, A., & Smola, A. J. (2017). Latent LSTM allocation: Joint clustering and non-linear dynamic modeling of sequence data. In *Proceedings of the 34th international conference on machine learning, vol. 70* (pp. 3967–3976).
- Zaremba, W., Sutskever, I., & Vinyals, O. (2014). Recurrent neural network regularization. arXiv arXiv:1409.2329.
- Zheng, X., Zaheer, M., Ahmed, A., Wang, Y., Xing, E. P., & Smola, A. (2017). State space LSTM models with particle MCMC inference. arXiv arXiv:1711.11179.