

Titre: Title:	Approximate Gaussian variance inference for state-space models
Auteurs: Authors:	Bhargob Deka, & James Alexandre Goulet
Date:	2023
Type:	Article de revue / Article
Référence: Citation:	Deka, B., & Goulet, J. A. (2023). Approximate Gaussian variance inference for state-space models. International Journal of Adaptive Control and Signal Processing, 29 pages. https://doi.org/10.1002/acs.3667

 **Document en libre accès dans PolyPublie**
Open Access document in PolyPublie

URL de PolyPublie: PolyPublie URL:	https://publications.polymtl.ca/54786/
Version:	Version officielle de l'éditeur / Published version Révisé par les pairs / Refereed
Conditions d'utilisation: Terms of Use:	CC BY-NC-ND

 **Document publié chez l'éditeur officiel**
Document issued by the official publisher

Titre de la revue: Journal Title:	International Journal of Adaptive Control and Signal Processing
Maison d'édition: Publisher:	Wiley
URL officiel: Official URL:	https://doi.org/10.1002/acs.3667
Mention légale: Legal notice:	This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

Approximate Gaussian variance inference for state-space models

Bhargob Deka  | James-A. Goulet

Department of Civil, Geologic and Mining Engineering, Polytechnique Montréal, Montréal, Québec, Canada

Correspondence

Bhargob Deka, Department of Civil, Geologic and Mining Engineering, Polytechnique Montréal, Montréal, Québec, Canada.

Email: bhargob.deka@polymtl.ca

Funding information

Hydro-Québec; Institut de Valorisation des Données; Natural Sciences and Engineering Research Council of Canada

Summary

State-space models require an accurate knowledge of the process error (\mathbf{Q}) and measurement error (\mathbf{R}) covariance matrices for exact state estimation. Even though the matrix \mathbf{R} can be, in many situations, considered to be known from the measuring instrument specifications, it is still a challenge to infer the \mathbf{Q} matrix online while providing reliable estimates along with a low computational cost. In this article, we propose an analytically tractable online Bayesian inference method for inferring the \mathbf{Q} matrix in state-space models. We refer to this method as *approximate Gaussian variance inference* (AGVI) using which we are able to treat the error variance and covariance terms in the full \mathbf{Q} matrix as Gaussian hidden states and infer them simultaneously with the other hidden states in a closed-form manner. The two case studies show that the method is able to provide statistically consistent estimates for the mean and uncertainties of the error variance terms for univariate and multivariate cases. The method also exceeds the performance of the existing adaptive Kalman filter methods both in terms of accuracy and computational efficiency.

KEYWORDS

Bayesian inference, closed-form inference, Gaussian multiplicative approximation, online parameter estimation, process error covariance matrix, state-space models

1 | INTRODUCTION

For linear dynamic systems, the Kalman filter is an exact state estimator if the process error (\mathbf{Q}) and the measurement error (\mathbf{R}) covariance matrices are known.¹ In most practical situations, the deterministic part of the model which includes the transition and the observation models is formulated based on known system dynamics. In contrast, the stochastic part representing the \mathbf{Q} and \mathbf{R} matrices is either unknown or only approximately known.¹⁻³ Previous studies have also shown that using incorrect error covariance matrices may result in large estimation errors or even cause divergence.^{1,4,5} Even though in many situations, the matrix \mathbf{R} can be considered to be known from measuring instrument specification, the \mathbf{Q} matrix is often unknown. Hence, an accurate estimation of the matrix \mathbf{Q} is necessary for the exact state estimation.^{4,6}

This article provides an analytical Bayesian inference method called the *approximate Gaussian variance inference* (AGVI) for performing closed-form online estimation of the error variance and covariance terms in the full \mathbf{Q} matrix. By definition, the expected value of the square of the univariate process error W^2 is equal to the error variance parameter, that is, $\mathbb{E}[W^2] = \sigma_W^2$, given that W has a zero mean. With the approximation that W^2 is Gaussian such that

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2023 The Authors. *International Journal of Adaptive Control and Signal Processing* published by John Wiley & Sons Ltd.

$W^2 \sim \mathcal{N}(w^2; \mathbb{E}[W^2], \text{var}(W^2))$, the error variance parameter is the same as the mean parameter for the probability density function (PDF) of W^2 . Subsequently, considering that this mean parameter $\mathbb{E}[W^2]$ is a random variable itself, inferring its posterior becomes analogous to computing the posterior for the error variance term. The AGVI method utilizes this definition and formulates the relationship between the process error W , the square of process error W^2 , and $\mathbb{E}[W^2]$ by leveraging the Gaussian multiplicative approximation (GMA)^{7,8} that provides the exact moments for W^2 . Thereafter, the Gaussian conjugate prior^{9,10} is used to analytically infer the unknown mean parameter for W^2 , that is, $\mathbb{E}[W^2] = \sigma_W^2$, using closed-form equations. The methodology is also extended to the multivariate observation model case where one error variance term σ_W^2 is inferred for each observation equation along with the covariance between each pair of process error terms. The article also provides a closed-form square-root filtering technique using the Cholesky decomposition such that the estimated \mathbf{Q} matrix is always positive semi-definite (PSD).

The layout of this article is as follows: Section 2 presents related works; Section 3 provides the AGVI method for estimating the univariate process error variance, and Section 4 extends the methodology to the multivariate case for inferring the full \mathbf{Q} matrix. Section 5 presents two case studies showing the application of the method; Section 5.1 presents the online estimation for the univariate error variance term and Section 5.2 shows the applicability of AGVI in inferring the full \mathbf{Q} matrix for multiple time series.

1.1 | Notations

The following notations are used throughout the manuscript; x : lowercase denotes a variable, X : slanted uppercase denotes a random variable, \mathbf{X} : bold upright uppercase denotes a deterministic matrix, \mathbf{x} : bold lowercase denotes a vector, \mathbf{X} : bold slanted uppercase denotes vector of random variables, $\mathbf{y}_{1:t}$: denote observations from 1 to t , $\boldsymbol{\mu}$: mean vector, $\boldsymbol{\Sigma}$: covariance matrix, W : random process error, W^2 : square of process error W , $\overline{W^2}$: random variable representing the expected value for W^2 , $W^i W^j$: product of any i th and j th process error, $\overline{W^i W^j}$: random variable for the expected value of $W^i W^j$, \mathbf{W}^p : random vector of all the product terms $W^i W^j$, $\overline{\mathbf{W}^p}$: random vector of all expected values of \mathbf{W}^p , $\boldsymbol{\Sigma}^{XW}$: covariance matrix between X and W , \mathbf{L}^W : upper triangular random matrix, $\overline{\mathbf{L}^W}$: random vector of all elements in \mathbf{L}^W .

Throughout the article, consistency is maintained in the notations as we transition from the univariate to the multivariate case study. For instance, we use W^2 to represent the random variable denoting the square of the univariate random process error W . Similarly, $W^i W^j$ represents the random variable for the product of W^i and W^j , where the superscripts indicate the i th and j th error terms. To denote the random expected value for $W^i W^j$, we use a bar on top, denoted by $\overline{W^i W^j}$. In the case of the univariate scenario, this becomes $\overline{W^2}$, while for any i th process error, it is expressed as $(\overline{W^i})^2$. The expected value and variance for $\overline{W^i W^j}$ are denoted as $\mu^{\overline{W^i W^j}}$ and $(\sigma^{\overline{W^i W^j}})^2$, respectively, with the corresponding random variable provided in the superscript.

2 | RELATED WORKS

The adaptive Kalman filters (AKF) were developed to estimate both the states and the error covariance matrices together by adaptively adjusting the Kalman filter to the measured data such that the estimation errors can be either bounded or reduced.⁴ The AKFs are broadly grouped as follows: (1) *correlation methods*,^{1,11-16} (2) *covariance-matching methods* (CMM),¹⁷⁻¹⁹ (3) *maximum likelihood methods*,^{20,21} and (4) *Bayesian methods*.^{3,22-25} One such AKF is the *innovation correlation method* (ICM)^{1,16} that uses the auto-correlation function of the innovations to form a system of linear equations involving the unknown covariance matrices. A least-square method is then used to solve these equations simultaneously to obtain the estimates for the \mathbf{Q} and \mathbf{R} matrices. The correlation methods provide unbiased estimates for linear time-invariant (LTI) systems and only asymptotically unbiased estimates for linear time-varying (LTV) systems.²⁶ A recent correlation approach called the *measurement difference method* (MDM)^{2,14,26-28} was proposed that is capable of providing unbiased and weakly consistent estimates for LTV systems, even for datasets with small number of measurements.^{26,28} MDM does not require any prior knowledge on the error covariance matrices and can also be performed online using recursive least-square methods.^{28,29} On the other hand, there are CMM methods such as the *adaptive limited memory filter* (ALMF)¹⁷ that computes sample covariance matrices at each time step for both the state prediction error and the innovation sequence using either the entire past data or over a moving window. However, such methods have shown to produce biased estimates for the covariance matrices and often fail to ensure the positive-definiteness of matrices when the sample size of the data is small.^{13,30} Shumway and Stoffer²⁰ provided a framework that uses the *expectation-maximization*

(EM) algorithm³¹ to obtain both the states and the error covariance matrices even when the data is irregularly spaced, but this can only be applied offline.³⁰ An extensive amount of literature exists for the AKF methods under the Bayesian category. The Bayesian methods include *state augmentation methods* primarily relying on nonlinear estimation techniques such as the *extended Kalman filter* (EKF),³⁰ the *unscented Kalman filter* (UKF),³² or the *particle filters*^{33,34} for the joint estimation of both the states and the error covariance matrices (ECM). While most methods in this category identify the error variances offline, Kontoroupi and Smyth³² provided an online estimation method by employing an approximation of the inverse gamma conjugacy. Online Bayesian inference methods also exist that rely on Markov chain Monte Carlo moves within a particle filter but they suffer from the particle degeneracy problem.^{23,33,35} Another Bayesian method is the *interactive multiple models* (IMM)³⁶ that defines multiple models each having a separate dynamic model with its own ECM as well as the transitional probabilities between one model i and another model j at any given time step t . A set of several Kalman filters are run in parallel to evaluate the state estimates for each model simultaneously. However, exact estimates using IMM can be obtained only when an infinitely large number of models are considered.²²

The *variational Bayes* (VB) methods have been proposed to approximate the intractable joint posterior PDF of the states and the covariance matrices at a comparatively lower computational cost than using the particle filters or the multiple model methods.^{2,22} Sarkaa and Hartikainen³⁷ proposed the VB-AKF method to obtain the full \mathbf{R} matrix using an inverse Wishart conjugate prior. However, the same methodology could not be applied to obtain the \mathbf{Q} matrix, since it does not appear in simple conjugate prior form, as opposed to the \mathbf{R} matrix.^{37,38} Ardeshiri et al.³⁹ proposed a VB based RTS smoother to obtain both the \mathbf{Q} and \mathbf{R} matrices, but it can only evaluate the error covariance matrices offline.^{38,40} Huang et al.⁴⁰ proposed an online VB-AKF method, referred to as VBAKF-PR, to directly estimate the joint distribution of the states, the state prediction error covariance matrix, and the \mathbf{R} matrix by using the conjugacy of the inverse Wishart prior for the covariance matrices. However, the method requires an accurate nominal \mathbf{Q} matrix based on problem-specific expertise without which the performance degrades drastically. Moreover, the method has additional parameters such as the tuning parameter, the forgetting factor, and the number of fixed-point iterations per time step that needs to be tuned. The *sliding window variational adaptive Kalman filter* (SWVAKF) method overcomes the limitation of VBAKF-PR as it is robust to the initialization of the nominal \mathbf{Q} matrix and proved to be computationally more efficient by avoiding the fixed-point iteration step. Hence, there are several methods that have been proposed to estimate the error covariance matrices. However, most methods are either offline in nature,^{1,20} restricted to linear dynamic systems^{24,40} or are computationally demanding.^{31,33,36} Furthermore, there is no closed-form method to obtain these matrices and none of the available methods have demonstrated the capacity to estimate a high-dimensional full \mathbf{Q} matrix. Hence, there is still the challenge to develop a method that performs closed-form online estimation of the matrix \mathbf{Q} and that is scalable to high-dimensional domains.

3 | UNIVARIATE PROCESS ERROR

This section presents the mathematical formulation of the AGVI method for inferring the variance parameter σ_W^2 associated with the univariate process error $W \sim \mathcal{N}(w; 0, \sigma_W^2)$ in the context of state-space models.

3.1 | Problem formulation

Let us consider an N -dimensional hidden state vector at time $t-1$, $\mathbf{x}_{t-1} = [x_1 \ x_2 \ \dots \ x_N]^T$, having a Gaussian PDF such that $\mathbf{X}_{t-1|t-1} \sim \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{t-1|t-1}, \boldsymbol{\Sigma}_{t-1|t-1})$ where $\boldsymbol{\mu}_{t-1|t-1} = \mathbb{E}[\mathbf{X}_{t-1} | \mathbf{y}_{1:t-1}]$ is the prior mean and $\boldsymbol{\Sigma}_{t-1|t-1} = \text{var}(\mathbf{X}_{t-1} | \mathbf{y}_{1:t-1})$ is the prior covariance matrix. Note that for brevity, the notation $\mathbf{X}_{t-1|t-1}$ is used as a shorthand for $\mathbf{X}_{t-1} | \mathbf{y}_{1:t-1}$. The transition and the observation equations for the linear Gaussian state-space models,⁴¹ are given by

$$\begin{aligned} \mathbf{x}_t &= \mathbf{A}\mathbf{x}_{t-1} + \mathbf{w}_t, & \mathbf{w} : \mathbf{W} &\sim \mathcal{N}(\mathbf{w}; 0, \mathbf{Q}), \\ y_t &= \mathbf{C}\mathbf{x}_t + v_t, & v : V &\sim \mathcal{N}(v; 0, \mathbf{R}), \end{aligned} \quad (1)$$

where \mathbf{A} is the transition matrix, $\mathbf{w}_t = [w_1 \ w_2 \ \dots \ w_N]^T$ is a vector of process error terms for which \mathbf{Q} is the process error covariance matrix, y_t is the observation, \mathbf{C} is the observation matrix, and v_t is the observation error for which the observation error variance is $\mathbf{R} = \sigma_V^2$. The \mathbf{A} and \mathbf{Q} matrices are constructed by assembling S specific components given by

$$\begin{aligned} \mathbf{A} &= \text{blkdiag}(\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_S), \\ \mathbf{Q} &= \text{blkdiag}(\mathbf{Q}_1, \mathbf{Q}_2, \dots, \mathbf{Q}_S), \end{aligned} \quad (2)$$

where $\text{blkdiag}(\cdot, \cdot)$ refers to block diagonal assembly of the individual components. The \mathbf{Q} matrix in Equation (2) can be further described by

$$\mathbf{Q} = \text{blkdiag}(\mathbf{Q}_1(\sigma_{W_1}^2, \Delta t), \mathbf{Q}_2(\sigma_{W_2}^2, \Delta t), \dots, \mathbf{Q}_s(\sigma_{W_s}^2, \Delta t)), \quad (3)$$

where each component $\mathbf{Q}_i(\sigma_{W_i}^2, \Delta t)$ can be represented as a function of the error variance parameter $\sigma_{W_i}^2$ and the time difference between two successive observations Δt . For example, consider that the model comprises two generic components, namely the local trend (LT) and the autoregressive (AR), for which the global \mathbf{Q} matrix is

$$\mathbf{Q} = \begin{bmatrix} \sigma_{\text{LT}}^2 \cdot \begin{bmatrix} \frac{\Delta t^4}{4} & \frac{\Delta t^3}{2} \\ \frac{\Delta t^3}{2} & \Delta t^2 \end{bmatrix} & \mathbf{0} \\ \mathbf{0} & \sigma_{\text{AR}}^2 \end{bmatrix}, \quad (4)$$

where the \mathbf{Q} matrices associated with the local trend \mathbf{Q}_1 and the autoregressive component \mathbf{Q}_2 are

$$\mathbf{Q}_1 = \sigma_{\text{LT}}^2 \cdot \begin{bmatrix} \frac{\Delta t^4}{4} & \frac{\Delta t^3}{2} \\ \frac{\Delta t^3}{2} & \Delta t^2 \end{bmatrix}, \quad \mathbf{Q}_2 = \sigma_{\text{AR}}^2,$$

where σ_{LT}^2 and σ_{AR}^2 are the process error variance terms associated with the LT and the AR components, respectively.⁴¹ Both matrices \mathbf{Q}_1 and \mathbf{Q}_2 are assembled in a block diagonal arrangement to get the \mathbf{Q} matrix as shown by Equation (4). Moreover, for a time series consisting of a single observation variable y_t , it is only possible to infer $\sigma_{W_i}^2$ for one component \mathbf{Q}_i , while all other should be either known or 0. This is because only a single unknown variable can be uniquely solved per equation. Hence, for each time series there is one unique process error variance that can be inferred. The next section describes the various steps for performing AGVI in order to obtain the posterior PDF for a single error variance parameter.

3.2 | Methodology

The proposed method considers the expected value of W^2 as a Gaussian random variable represented by $\overline{W^2}$ such that

$$\overline{W^2} \sim \mathcal{N}(\overline{w^2}; \mu^{\overline{W^2}}, (\sigma^{\overline{W^2}})^2), \quad (5)$$

where $\mu^{\overline{W^2}}$ and $(\sigma^{\overline{W^2}})^2$ are the hyperprior mean and variance for $\overline{W^2}$. Using Equation (5), the PDF of W can be re-written as

$$W \sim \mathcal{N}(w; 0, \overline{w^2}). \quad (6)$$

Hence, the first objective is to obtain the marginal PDF of W such that the random variance $\overline{W^2}$ can be marginalized out. The following lemmas are invoked to show that the marginal PDF of W can be obtained using the marginal PDF of W^2 . The subsequent proposition uses these lemmas to provide the prior predictive PDF for W at a time t such that $W_{t|t-1} \sim \mathcal{N}(w_t; 0, \mu_{t-1|t-1}^{\overline{W^2}})$.

Lemma 1. *Given that W is Gaussian with a zero mean and W^2 is approximated as a Gaussian random variable given by $W^2 \sim \mathcal{N}(w^2; \mu^{W^2}, (\sigma^{W^2})^2)$ for which the exact moments are provided by GMA (see Appendix A), it can be shown that the PDF of W^2 is dependent only on the mean parameter μ^{W^2} so that*

$$W^2 \sim \mathcal{N}(w^2; \mu^{W^2}, 2(\mu^{W^2})^2),$$

where the variance term $(\sigma^{W^2})^2$ is equal to $2(\mu^{W^2})^2$. As a result, the PDF $f(w^2 | \mu^{W^2}, (\sigma^{W^2})^2)$ can be shown by $f(w^2 | \mu^{W^2})$.

Proof. See Appendix B.1. ■

Lemma 2. Given that the parameter μ^{W^2} in $f(w^2|\mu^{W^2})$ is considered as a Gaussian random variable $\overline{W^2} \sim \mathcal{N}(\overline{w^2}; \mu^{\overline{W^2}}, (\sigma^{\overline{W^2}})^2)$, the mean and variance of the prior predictive PDF of $W_{t|t-1}^2$ are given by

$$\begin{aligned}\mu_{t|t-1}^{W^2} &= \mu_{t-1|t-1}^{\overline{W^2}}, \\ (\sigma_{t|t-1}^{W^2})^2 &= 3(\sigma_{t-1|t-1}^{\overline{W^2}})^2 + 2(\mu_{t-1|t-1}^{\overline{W^2}})^2,\end{aligned}$$

where $\mu_{t-1|t-1}^{\overline{W^2}}$ and $(\sigma_{t-1|t-1}^{\overline{W^2}})^2$ are the prior moments for $\overline{W^2}_{t-1|t-1}$.

Proof. See Appendix B.2. ■

Proposition 1. Considering that the mean parameter μ^{W^2} is itself a random variable $\overline{W^2}$ so that

$$\overline{W^2}_{t-1|t-1} \sim \mathcal{N}(\overline{w^2}_{t-1}; \mu_{t-1|t-1}^{\overline{W^2}}, (\sigma_{t-1|t-1}^{\overline{W^2}})^2),$$

where $\mu_{t-1|t-1}^{\overline{W^2}}$ and $(\sigma_{t-1|t-1}^{\overline{W^2}})^2$ are the hyperprior mean and variance for $\overline{W^2}_{t-1|t-1}$, the error variance σ_W^2 can be made equal to

$$\sigma_W^2 = \mu_{t-1|t-1}^{\overline{W^2}}. \quad (7)$$

Proof. Using Lemmas 1 and 2, and considering the one-to-one relationship between the moments of W and W^2 , the prior predictive PDF of $W_{t|t-1}$ can be formulated as

$$f(w_t) = \mathcal{N}(w_t; 0, \mu_{t-1|t-1}^{\overline{W^2}}),$$

where by Lemma 2, the variance of $W_{t|t-1}$ is $\sigma_W^2 = \mathbb{E}[W_{t|t-1}^2] = \mu_{t-1|t-1}^{\overline{W^2}}$. ■

The next objective is to perform the prediction step in the filtering procedure using the model matrices \mathbf{A} , \mathbf{C} , \mathbf{Q} , and \mathbf{R} defined in Section 3.1 and the prior knowledge of σ_W^2 . The transition model for $\overline{w^2}$ is $\overline{w^2}_t = \overline{w^2}_{t-1}$, where the hidden state $\overline{w^2}$ is assumed to be constant from $t-1$ to t . Using the prior knowledge for W , the augmented state vector \mathbf{h}_{t-1} at any time $t-1$ is given by

$$\mathbf{h}_{t-1} = [\mathbf{x} \ w]_{t-1}^\top. \quad (8)$$

The prior predictive PDF of the hidden states $\mathbf{H}_{t|t-1}$ is given by

$$\mathbf{H}_{t|t-1} \sim \mathcal{N}(\mathbf{h}_t; \mu_{t|t-1}^{\mathbf{H}}, \Sigma_{t|t-1}^{\mathbf{H}}),$$

where, using Equations (1) and (7), the mean vector and the covariance matrix are given by

$$\begin{aligned}\mu_{t|t-1}^{\mathbf{H}} &= \begin{bmatrix} \mathbf{A}\mu_{t-1|t-1} \\ 0 \end{bmatrix}_{t|t-1}, \\ \Sigma_{t|t-1}^{\mathbf{H}} &= \begin{bmatrix} \mathbf{A}\Sigma_{t-1|t-1}\mathbf{A}^\top + \mathbf{Q} & \Sigma^{XW} \\ (\Sigma^{XW})^\top & \mu_{t-1|t-1}^{\overline{W^2}} \end{bmatrix}_{t|t-1}.\end{aligned} \quad (9)$$

The covariance term $\Sigma_{t|t-1}^{XW}$ between $\mathbf{X}_{t|t-1}$ and $W_{t|t-1}$ in Equation (9) is formulated as

$$\begin{aligned}\Sigma_{t|t-1}^{XW} &= \text{cov}(\mathbf{A}\mathbf{X}_{t-1|t-1} + \mathbf{W}_{t|t-1}, W_{t|t-1}), \\ &= \text{cov}(\mathbf{W}, W)_{t|t-1},\end{aligned} \quad (10)$$

where $\mathbf{W}_{t|t-1}$ is a vector of random variables representing the process error terms in the state vector \mathbf{h} . Moreover the hidden states $\mathbf{X}_{t-1|t-1}$ and the process error $W_{t|t-1}$ are assumed to be independent of each other. The mean and variance of $Y_{t|t-1} \sim \mathcal{N}(y_t; \mu_Y, \sigma_Y^2)$ are given by

$$\begin{aligned}\mu_Y &= \mathbf{C}\boldsymbol{\mu}_{t|t-1} + \mu_V, \\ \sigma_Y^2 &= \mathbf{C}\boldsymbol{\Sigma}_{t|t-1}\mathbf{C}^\top + \sigma_V^2,\end{aligned}$$

given that \mathbf{X} and V are assumed to be independent of each other. The covariance term $\boldsymbol{\Sigma}_{HY}$ between $\mathbf{H}_{t|t-1}$ and $Y_{t|t-1}$ is

$$\boldsymbol{\Sigma}_{HY} = \boldsymbol{\Sigma}_{t|t-1}^H \mathbf{F}_t^\top,$$

where the observation matrix is $\mathbf{F}_t = [\mathbf{C} \ 0]$.

The inference for the parameter σ_W^2 requires two update steps; In the first step, the posterior PDF $f(\mathbf{h}_t | \mathbf{y}_{1:t})$ is estimated using the observation model defined in Equation (1) so that

$$f(\mathbf{h}_t | \mathbf{y}_{1:t}) = \frac{f(\mathbf{h}_t, y_t | \mathbf{y}_{1:t-1})}{f(y_t | \mathbf{y}_{1:t-1})} \approx \mathcal{N}(\mathbf{h}_t; \boldsymbol{\mu}_{t|t}^H, \boldsymbol{\Sigma}_{t|t}^H), \quad (11)$$

which we approximate by a Gaussian distribution with a posterior mean vector $\boldsymbol{\mu}_{t|t}^H$ and a covariance matrix $\boldsymbol{\Sigma}_{t|t}^H$ that are obtained using the predicted moments provided in Equation (9) and the Gaussian conditional equations^{9,10} such that

$$\begin{aligned}\boldsymbol{\mu}_{t|t}^H &= \boldsymbol{\mu}_{t|t-1}^H + \frac{\boldsymbol{\Sigma}_{HY}}{\sigma_Y^2}(y_t - \mu_Y), \\ \boldsymbol{\Sigma}_{t|t}^H &= \boldsymbol{\Sigma}_{t|t-1}^H - \frac{\boldsymbol{\Sigma}_{HY} \cdot \boldsymbol{\Sigma}_{HY}^\top}{\sigma_Y^2}.\end{aligned}$$

Now that we have the posterior PDF $f(w_t | \mathbf{y}_{1:t})$ from Equation (11), we move to the second update step where we use this new information of W at time t to update our current knowledge of $\overline{W^2}$. Figure 1 shows a graphical model representing the relationship between the random variables $Y_{t|t-1}$, $\mathbf{X}_{t|t-1}$, $W_{t|t-1}$, $W_{t|t-1}^2$, and $\overline{W^2}_{t|t-1}$. Note that while considering $\mu_W = 0$, the first moment of W^2 is equal to the second moment of W (under Lemma 1). In this case, the knowledge of W is fully defined by the knowledge of W^2 , which is denoted in Figure 1 by an undirected solid line between the nodes W^2 and W . Following the structure depicted in Figure 1, the subsequent lemmas are provided for obtaining the posterior knowledge of $\overline{W^2}_{t|t}$.

Lemma 3. *Considering the joint PDF of the random variables $Y_{t|t-1}$, $W_{t|t-1}^2$, and $\overline{W^2}_{t|t-1}$, and marginalizing out W^2 from the joint PDF, the posterior PDF of W^2 can be obtained by the following integral*

$$f(\overline{w^2}_t | \mathbf{y}_{1:t}) = \int f(w_t^2 | \mathbf{y}_{1:t}) \cdot f(\overline{w^2}_t | w_t^2, \mathbf{y}_{1:t-1}) dw_t^2.$$

Proof. See Appendix B.3. ■

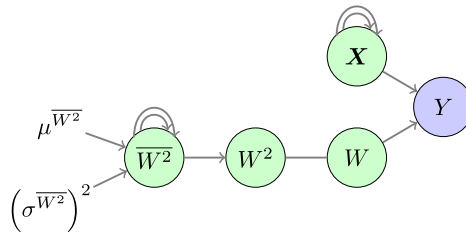


FIGURE 1 Illustration showing the graphical model for the online inference of the error variance parameter. The hidden and observed state variables are denoted by green and violet nodes. The double arrows on the nodes \mathbf{X} and $\overline{W^2}$ represent that these variables are learnt recursively over time. For brevity, the subscript $t|t-1$ is dropped from each of the variables.

Lemma 4. *The posterior mean and variance of $W_{t|t}^2$ are*

$$\begin{aligned}\mu_{t|t}^{W^2} &= (\mu_{t|t}^W)^2 + (\sigma_{t|t}^W)^2, \\ (\sigma_{t|t}^{W^2})^2 &= 2(\sigma_{t|t}^W)^4 + 4(\sigma_{t|t}^W)^2(\mu_{t|t}^W)^2.\end{aligned}$$

Proof. See Appendix B.4. ■

The Lemmas 3 and 4 are used for proving the following proposition.

Proposition 2. *The posterior mean and variance of $\overline{W}_{t|t}^2 \sim \mathcal{N}(\mu_{t|t}^{\overline{W}^2}, (\sigma_{t|t}^{\overline{W}^2})^2)$ are given by*

$$\begin{aligned}\mu_{t|t}^{\overline{W}^2} &= \mu_{t|t-1}^{\overline{W}^2} + k_t(\mu_{t|t}^{W^2} - \mu_{t|t-1}^{W^2}), \\ (\sigma_{t|t}^{\overline{W}^2})^2 &= (\sigma_{t|t-1}^{\overline{W}^2})^2 + k_t^2((\sigma_{t|t}^{W^2})^2 - (\sigma_{t|t-1}^{W^2})^2), \\ k &= \frac{(\sigma_{t-1|t-1}^{\overline{W}^2})^2}{(\sigma_{t|t-1}^{W^2})^2}.\end{aligned}$$

Proof. See Appendix B.5. ■

Both the update steps 1 and 2 are employed recursively as observations are collected in order to first estimate the posterior knowledge of W and then use this to update our knowledge of the expected value of W^2 , that is, \overline{W}^2 , which is a variable that is equal to σ_W^2 , the parameter we seek to infer. All the steps performed in a particular time step t are summarized in Algorithm 1 as provided in Appendix D.

4 | MULTIVARIATE PROCESS ERROR

This section extends the mathematical formulation of the AGVI method for inferring error variance and covariance terms comprising a full \mathbf{Q} matrix.

4.1 | Problem formulation

Let us consider D observed time series for which the global state vector is $\mathbf{x}_t = [\mathbf{x}_t^1 \ \mathbf{x}_t^2 \ \dots \ \mathbf{x}_t^D]^\top$, where $\mathbf{x}_t^j, \forall j \in \{1, 2, \dots, D\}$ refers to the concatenation of all S_j generic components for the j th time series. Similarly, the vector of correlated process errors is assembled following $\mathbf{w}_t = [\mathbf{w}_t^1 \ \mathbf{w}_t^2 \ \dots \ \mathbf{w}_t^D]^\top$. The global transition, observation, process error covariance, and observation error covariance matrices are assembled block diagonally as

$$\begin{aligned}\mathbf{A} &= \text{blkdiag}[\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_D], \\ \mathbf{C} &= \text{blkdiag}[\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_D], \\ \mathbf{Q} &= \text{blkdiag}[\mathbf{Q}_1, \mathbf{Q}_2, \dots, \mathbf{Q}_D], \\ \mathbf{R} &= \text{blkdiag}[\mathbf{R}_1, \mathbf{R}_2, \dots, \mathbf{R}_D].\end{aligned}$$

Covariance matrices $\text{cov}(\mathbf{W}^k, \mathbf{W}^n)$ exists between the process errors \mathbf{W}^k and \mathbf{W}^n of the k th and n th time series respectively, where $k, n \in \{1, 2, \dots, D\}$. The process error covariance matrix \mathbf{Q} can be reformulated as follows

$$\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{Q}_{1,2} & \dots & \mathbf{Q}_{1,D} \\ \vdots & \mathbf{Q}_2 & \dots & \mathbf{Q}_{2,D} \\ \vdots & \dots & \ddots & \vdots \\ \text{sym.} & \dots & \dots & \mathbf{Q}_D \end{bmatrix}, \quad (12)$$

where the covariance term $\text{cov}(\mathbf{W}^k, \mathbf{W}^n)$ is represented by $\mathbf{Q}_{k,n}$. The sub-matrices within the matrix $\mathbf{Q}_{k,n}$ in Equation (12) are themselves represented by $\text{cov}(\mathbf{W}^{jk}, \mathbf{W}^{mn}) = \mathbf{Q}_{jk,mn}$, where $j \in \{1, 2, \dots, S_j\}$ and $m \in \{1, 2, \dots, S_m\}$ are the j th and m th component of the k th and n th time series, respectively. As described in Section 3.1, each of the sub-matrices $\mathbf{Q}_{jk,mn}(\sigma_{jk}^2, \sigma_{mn}^2, \Delta t)$ can be represented as a function of the error variance parameters σ_{jk}^2 , σ_{mn}^2 , and Δt . Moreover, each of the elements within the sub-matrix $\mathbf{Q}_{jk,mn}$ is given by $\text{cov}(W^{ijk}, W^{lmn})$, which provides the covariance between the i th process error term of the j th component in the k th time series, W^{ijk} , and the l th process error term of the m th component in the n th time series, W^{lmn} . For example, let us consider two time series each modeled using a local trend (LT) component as described in Section 3.1. The global \mathbf{Q} matrix is assembled block diagonally such that

$$\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{Q}_{1,2} \\ \mathbf{Q}_{2,1} & \mathbf{Q}_2 \end{bmatrix},$$

where \mathbf{Q}_1 and \mathbf{Q}_2 are the process error covariance matrices associated with their individual local trend components and $\mathbf{Q}_{1,2}$ is the cross-covariance matrix between the process errors of the two time series. Each of these covariance matrices are defined as follows: $\mathbf{Q}_1 = \sigma_{\text{LT}_1}^2 \cdot \mathbf{J}$, $\mathbf{Q}_2 = \sigma_{\text{LT}_2}^2 \cdot \mathbf{J}$, and $\mathbf{Q}_{1,2} = \sigma_{\text{LT}_{12}} \cdot \mathbf{J}$, where $\sigma_{\text{LT}_1}^2$, and $\sigma_{\text{LT}_2}^2$ are the two error variance terms for each of the LT component, $\sigma_{\text{LT}_{12}}$ is the covariance term between the two process error random variables W^{LT_1} and W^{LT_2} . For a constant acceleration kinematic model,^{42,43} the matrix \mathbf{J} is defined such that

$$\mathbf{J} = \begin{bmatrix} \frac{\Delta t^4}{4} & \frac{\Delta t^3}{2} \\ \frac{\Delta t^3}{2} & \Delta t^2 \end{bmatrix}.$$

Hence, for this case the terms to be inferred are: $\sigma_{\text{LT}_1}^2$, $\sigma_{\text{LT}_2}^2$, and $\sigma_{\text{LT}_{12}}$. Similarly, for multiple time series, the goal is to infer one error variance term per time series along with the covariance terms for each pair of process error terms.

4.2 | Methodology

Let us consider the multivariate process error term $\mathbf{w} = [w^1 \ w^2 \ \dots \ w^i \ \dots \ w^D]^T$, where w^i , $\forall i \in \{1, 2, \dots, D\}$ represents the process error term for the i th time series for which the variance term has to be inferred. Given that the expected value of W is zero, the covariance term between the i th and j th process error is

$$\text{cov}(W^i, W^j) = \mathbb{E}[W^i W^j] - \mathbb{E}[W^i] \mathbb{E}[W^j] = \mathbb{E}[W^i W^j],$$

the covariance matrix $\Sigma^{\mathbf{W}}$ is given by

$$\Sigma^{\mathbf{W}} = \begin{bmatrix} \mathbb{E}[(W^1)^2] & \mathbb{E}[W^1 W^2] & \dots & \mathbb{E}[W^1 W^D] \\ \vdots & \mathbb{E}[(W^2)^2] & \dots & \mathbb{E}[W^2 W^D] \\ \vdots & \dots & \ddots & \vdots \\ \text{sym.} & \dots & \dots & \mathbb{E}[(W^D)^2] \end{bmatrix}, \quad (13)$$

where $\text{var}(W^i) = \mathbb{E}[(W^i)^2]$ is the error variance for the i th time series, and $\text{cov}(W^i, W^j) = \mathbb{E}[W^i W^j]$ is the covariance term between the two process errors for the i th and j th time series. Similarly to the univariate process error, let us consider the approximation that each of the product terms $W^i W^j$ is a Gaussian random variable such that

$$W^i W^j \sim \mathcal{N}(w^i w^j; \mu^{W^i W^j}, (\sigma^{W^i W^j})^2), \quad (14)$$

where $\mathbb{E}[W^i W^j] = \mu^{W^i W^j}$ is the mean parameter and $\text{var}(W^i W^j) = (\sigma^{W^i W^j})^2$ is the variance. For D time series, there are a total of $\frac{D(D+1)}{2}$ product terms which are represented by the random vector $\mathbf{w}^{\mathbf{P}} = [(w^1)^2 \ (w^2)^2 \ \dots \ w^i w^j \ \dots \ w^D w^{D-1}]^T$ such that

$$\mathbf{W}^{\mathbf{P}} \sim \mathcal{N}(\mathbf{w}^{\mathbf{P}}; \mu^{\mathbf{W}^{\mathbf{P}}}, \Sigma^{\mathbf{W}^{\mathbf{P}}}), \quad (15)$$

where using Equation (14), the mean vector of \mathbf{W}^p is given by

$$\boldsymbol{\mu}^{\mathbf{W}^p} = \left[\mu^{(W^1)^2} \mu^{(W^2)^2} \dots \mu^{(W^D)^2} \mu^{W^1 W^2} \dots \mu^{W^{D-1} W^D} \right]_{k \times 1}^T. \quad (16)$$

Similarly to Lemma 1, the covariance matrix $\boldsymbol{\Sigma}^{\mathbf{W}^p}$ can be obtained in terms of the mean parameters in $\boldsymbol{\mu}^{\mathbf{W}^p}$ such that

$$\boldsymbol{\Sigma}^{\mathbf{W}^p} = \begin{bmatrix} 2(\mu^{(W^1)^2})^2 & 2(\mu^{W^1 W^2})^2 & \dots & 2\mu^{W^1 W^{D-1}} \mu^{W^1 W^D} \\ \vdots & 2(\mu^{(W^2)^2})^2 & \dots & 2\mu^{W^2 W^{D-1}} \mu^{W^1 W^D} \\ \vdots & \dots & \dots & \vdots \\ \vdots & \dots & \ddots & \vdots \\ \text{sym.} & \dots & \dots & \mu^{(W^{D-1})^2} \mu^{(W^D)^2} + (\mu^{W^{D-1} W^D})^2 \end{bmatrix}_{k \times k},$$

where the variance $\text{var}(W^i W^j)$ and the covariance $\text{cov}(W^i W^j, W^k W^m)$ terms for the product of the errors $\forall i, j, k, m \in \{1, 2, \dots, D\}$ are obtained using the GMA equations presented in Appendix A. The mean vector $\boldsymbol{\mu}^{\mathbf{W}^p}$ defined in Equation (16) is considered to be random and denoted by $\overline{\mathbf{W}^p}$ having a Gaussian PDF given by

$$\overline{\mathbf{W}^p} \sim \mathcal{N}(\overline{\mathbf{w}^p}; \boldsymbol{\mu}^{\overline{\mathbf{W}^p}}, \boldsymbol{\Sigma}^{\overline{\mathbf{W}^p}}), \quad (17)$$

where the vector $\overline{\mathbf{w}^p} = \left[\overline{(w^1)^2} \ \overline{(w^2)^2} \ \dots \ \overline{(w^D)^2} \ \overline{w^1 w^2} \ \dots \ \overline{w^{D-1} w^D} \right]^T$. The mean vector and the covariance matrix of $\overline{\mathbf{W}^p}$ are

$$\boldsymbol{\mu}^{\overline{\mathbf{W}^p}} = \left[\mu^{\overline{(W^1)^2}} \ \mu^{\overline{(W^2)^2}} \ \dots \ \mu^{\overline{W^{D-1} W^D}} \right]_{k \times 1}^T,$$

$$\boldsymbol{\Sigma}^{\overline{\mathbf{W}^p}} = \begin{bmatrix} (\sigma^{\overline{(W^1)^2}})^2 & 0 & \dots & 0 \\ \vdots & (\sigma^{\overline{(W^2)^2}})^2 & \dots & 0 \\ \vdots & \dots & \dots & 0 \\ \vdots & \dots & \ddots & \vdots \\ \text{sym.} & \dots & \dots & (\sigma^{\overline{W^{D-1} W^D}})^2 \end{bmatrix}_{k \times k},$$

where the random variables in $\overline{\mathbf{W}^p}$ are assumed to be independent from each other as shown by the covariance matrix $\boldsymbol{\Sigma}^{\overline{\mathbf{W}^p}}$ where the off-diagonal terms are zero.

Using the hyperprior $\overline{\mathbf{W}^p}$ defined in Equation (17), the first objective is to obtain the covariance matrix $\boldsymbol{\Sigma}^{\mathbf{W}}$ defined in Equation (13) by obtaining the prior predictive PDF of $\mathbf{W}^p_{t|t-1}$ as provided by the following lemma and proposition.

Lemma 5. Using the transition model $\mathbf{w}^p_t = \mathbf{w}^p_{t-1}$, the prior predictive PDF of $\mathbf{W}^p_{t|t-1}$ is given by

$$\mathbf{W}^p_{t|t-1} \sim \mathcal{N}(\boldsymbol{\mu}^{\mathbf{W}^p_{t|t-1}}, \boldsymbol{\Sigma}^{\mathbf{W}^p_{t|t-1}}),$$

where the mean terms in $\boldsymbol{\mu}^{\mathbf{W}^p_{t|t-1}}$, and the variance and covariance terms in $\boldsymbol{\Sigma}^{\mathbf{W}^p_{t|t-1}}$ are given by

$$\begin{aligned} \mathbb{E}[W^i W^j] &= \mu^{\overline{W^i W^j}}, \\ \text{var}((W^i)^2) &= 3(\sigma^{\overline{(W^i)^2}})^2 + 2(\mu^{\overline{(W^i)^2}})^2, \\ \text{var}(W^i W^j) &= (\sigma^{\overline{W^i W^j}})^2 \\ &\quad + \frac{(\mu^{\overline{W^i W^j}})^2}{\mu^{\overline{(W^i)^2}} \mu^{\overline{(W^j)^2}} + (\mu^{\overline{W^i W^j}})^2} \cdot (\sigma^{\overline{W^i W^j}})^2 \\ &\quad + \mu^{\overline{(W^i)^2}} \mu^{\overline{(W^j)^2}} + (\mu^{\overline{W^i W^j}})^2, \\ \text{cov}(W^i W^j, W^l W^m) &= \mu^{\overline{W^i W^l}} \mu^{\overline{W^j W^m}} + \mu^{\overline{W^i W^m}} \mu^{\overline{W^j W^l}}. \end{aligned}$$

Proof. See Appendix C.1. ■

Proposition 3. The prior predictive PDF of \mathbf{W} has a zero mean vector and covariance matrix $\Sigma_{t|t-1}^{\mathbf{W}}$ defined by

$$\Sigma_{t|t-1}^{\mathbf{W}} = \begin{bmatrix} \mu^{(W^1)^2} & \mu^{\overline{W^1 W^2}} & \cdots & \mu^{\overline{W^1 W^D}} \\ \vdots & \mu^{(W^2)^2} & \cdots & \mu^{\overline{W^2 W^D}} \\ \vdots & \cdots & \ddots & \vdots \\ \text{sym.} & \cdots & \cdots & \mu^{(W^D)^2} \end{bmatrix}_{t|t-1}. \quad (18)$$

Proof. Using Lemma 5, the covariance matrix $\Sigma^{\mathbf{W}}$ for the prior predictive PDF of \mathbf{W} is obtained by substituting the terms $\mathbb{E}[W^i W^j]$ in Equation (13) by the mean parameters of $\overline{\mathbf{W}^{\mathcal{P}}}$, that is, $\mu^{\overline{W^i W^j}}$. ■

In order to maintain positive semi-definiteness of $\Sigma_{t|t-1}^{\mathbf{W}}$ shown by Equation (18), the prior information is built from a random vector $\overrightarrow{\mathbf{L}^{\mathbf{W}}}$ using Cholesky decomposition as shown by the following lemma.

Lemma 6. Any ij th element of $\Sigma^{\mathbf{W}}$ is obtained such that

$$\mu^{\overline{W^i W^j}} = \mathbb{E} \left[\sum_{k=1}^D L_{jk} L_{ki} \right],$$

where all elements of $\overrightarrow{\mathbf{L}^{\mathbf{W}}}$ are assumed to be Gaussian, $L_{ij} \sim \mathcal{N}(\mu_{L_{ij}}, \sigma_{L_{ij}}^2)$, and the expectation of the product terms are obtained using the GMA equations. Moreover, any covariance term between the random vectors $\overrightarrow{\mathbf{L}^{\mathbf{W}}}$ and $\overline{\mathbf{W}^{\mathcal{P}}}$ given by $\Sigma_{t|t-1}^{\overrightarrow{\mathbf{L}^{\mathbf{W}}} \overline{\mathbf{W}^{\mathcal{P}}}}$, can be shown as

$$\text{cov}(L_{ij}, \overline{W^i W^j}) = \text{cov} \left(L_{ij}, \sum_{k=1}^D L_{jk} L_{ki} \right).$$

Proof. See Appendix C.2. ■

Using the prior predictive PDF of \mathbf{W} , the next objective is to perform the prediction step. Let us consider the augmented vector of hidden states $\mathbf{h}_{t-1} = [\mathbf{x}_{t-1}^{\top} \mathbf{w}_{t-1}^{\top}]^{\top}$ such that the PDF of $\mathbf{H}_{t|t-1} \sim \mathcal{N}(\mathbf{h}_t, \mu_{t|t-1}^{\mathbf{H}}, \Sigma_{t|t-1}^{\mathbf{H}})$ has a mean vector $\mu_{t|t-1}^{\mathbf{H}}$ and a covariance matrix $\Sigma_{t|t-1}^{\mathbf{H}}$ defined by

$$\mu_{t|t-1}^{\mathbf{H}} = \begin{bmatrix} \mu_{t|t-1}^{\top} \\ \mathbf{0} \end{bmatrix}^{\top}, \quad (19)$$

$$\Sigma_{t|t-1}^{\mathbf{H}} = \begin{bmatrix} \mathbf{A} \Sigma_{t-1|t-1} \mathbf{A}^{\top} + \mathbf{Q} & \Sigma^{\mathbf{XW}} \\ (\Sigma^{\mathbf{XW}})^{\top} & \Sigma^{\mathbf{W}} \end{bmatrix}_{t|t-1}, \quad (20)$$

where the covariance matrix $\Sigma^{\mathbf{W}}$ defined in Equation (18) is obtained using the prior knowledge of $\overrightarrow{\mathbf{L}^{\mathbf{W}}}$ defined in the Cholesky space as stated in Lemma 6. Similarly to Equation (10), the covariance matrix between \mathbf{X} and \mathbf{W} is given by

$$\begin{aligned} \Sigma_{t|t-1}^{\mathbf{XW}} &= \text{cov}(\mathbf{X}, \mathbf{W})_{t|t-1} \\ &= \text{cov}(\mathbf{X}, [W^1 \ W^2 \ \cdots \ W^D]^{\top})_{t|t-1}, \\ &= \text{cov}(\mathbf{A} \mathbf{X}_{t-1|t-1} + \mathbf{W}_{t|t-1}, [W^1 \ W^2 \ \cdots \ W^D]^{\top}_{t|t-1}), \\ &= \text{cov}(\mathbf{W}, [W^1 \ W^2 \ \cdots \ W^D]^{\top})_{t|t-1}, \end{aligned}$$

where $\mathbf{W}_{t|t-1} = [W^1 \ W^2 \ \cdots \ W^D]^{\top}_{t|t-1}$ is a vector of random variables that includes one process error term W from each of the D time series.

The inference for the covariance matrix Σ^W requires two update steps. Similarly to the univariate case shown in Section 3, the Gaussian conditional equations are used to perform the first update step to obtain the posterior PDF of \mathbf{H} shown by

$$\mathbf{H}_{t|t} \sim \mathcal{N}(\mathbf{h}_t, \mu_{t|t}^H, \Sigma_{t|t}^H). \quad (21)$$

We now move to the second update step where we use the posterior PDF $f(\mathbf{w}_t | \mathbf{y}_{1:t})$ obtained from Equation (21), and the GMA equations to obtain the posterior PDF $f(\mathbf{w}^p_t | \mathbf{y}_{1:t})$ such that

$$f(\mathbf{w}^p_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{w}^p_t; \mu_{t|t}^{W^p}, \Sigma_{t|t}^{W^p}).$$

The posterior PDF of $\overline{\mathbf{W}^p}$ is defined using the following lemma.

Lemma 7. *The posterior mean, variance and covariance terms of $\overline{\mathbf{W}^p}$ are*

$$\begin{aligned} \overline{\mu}_{t|t}^{W^p} &= \overline{\mu}_{t|t-1}^{W^p} + \mathbf{K}_t (\mu_{t|t}^{W^p} - \mu_{t|t-1}^{W^p}), \\ \overline{\Sigma}_{t|t}^{W^p} &= \overline{\Sigma}_{t|t-1}^{W^p} + \mathbf{K}_t (\Sigma_{t|t}^{W^p} - \Sigma_{t|t-1}^{W^p}) \mathbf{K}_t^\top, \\ \mathbf{K}_t &= \Sigma_{t|t-1}^{W^p} \overline{\Sigma}_{t|t-1}^{W^p} (\Sigma_{t|t-1}^{W^p})^{-1}, \\ \overline{\Sigma}_{t|t-1}^{W^p} &= \overline{\Sigma}_{t|t-1}^{W^p}. \end{aligned}$$

Proof. See Appendix C.3. ■

Using the updated knowledge of $\overline{\mathbf{W}^p}$ in Lemma 7, the posterior moments for $\overrightarrow{\mathbf{L}^W}$ in the Cholesky space is defined using the following proposition.

Proposition 4. *The posterior moments of $\overrightarrow{\mathbf{L}^W}$ are*

$$\begin{aligned} \overline{\mu}_{t|t}^{\overrightarrow{\mathbf{L}^W}} &= \overline{\mu}_{t|t-1}^{\overrightarrow{\mathbf{L}^W}} + \mathbf{K}_t^L (\mu_{t|t}^{\overrightarrow{\mathbf{L}^W}} - \mu_{t|t-1}^{\overrightarrow{\mathbf{L}^W}}), \\ \overline{\Sigma}_{t|t}^{\overrightarrow{\mathbf{L}^W}} &= \overline{\Sigma}_{t|t-1}^{\overrightarrow{\mathbf{L}^W}} + \mathbf{K}_t^L (\Sigma_{t|t}^{\overrightarrow{\mathbf{L}^W}} - \Sigma_{t|t-1}^{\overrightarrow{\mathbf{L}^W}}) (\mathbf{K}_t^L)^\top, \\ \mathbf{K}_t^L &= \Sigma_{t|t-1}^{\overrightarrow{\mathbf{L}^W}} \overline{\Sigma}_{t|t-1}^{\overrightarrow{\mathbf{L}^W}} (\Sigma_{t|t-1}^{\overrightarrow{\mathbf{L}^W}})^{-1}. \end{aligned}$$

Proof. Proposition 4 is derived using the Lemmas 5–7. ■

Both steps are employed recursively in order to estimate the elements of the covariance matrix Σ^W and then use this to update our knowledge of the mean vector of \mathbf{W}^p , that is, $\overline{\mathbf{W}^p}$. All the steps performed in a particular time step t are summarized in Algorithm 2 as provided in Appendix D.

5 | APPLIED EXAMPLES

In this section, two case studies are presented to demonstrate the use of AGVI for online inference of error variance terms. Case study 1 focuses on the online estimation of error variance for a linear time-varying (LTV) model. The study includes statistical consistency tests to showcase the filter's optimality, t -tests to prove that the estimates are unbiased, empirical validation of uncertainty associated with error variance estimates, and analysis of the impact of $\frac{Q}{R}$ ratio on the posterior mean estimate $\mu_{T|T}$ of the error variance term. The AGVI method is also compared to two adaptive Kalman filtering (AKF) approaches: the *sliding window variational adaptive Kalman filter* (SWVAKF)²⁴ and the *measurement difference method* (MDM).^{2,14,27,28} These AKF methods fall under different categories, with MDM being a correlation method and SWVAKF being a Bayesian method. Case study 2 presents a simulated multivariate random walk model with a full process error covariance matrix \mathbf{Q} and a comparison of the AGVI method to the AKF methods.

5.1 | Case Study 1: Univariate linear time-varying

For this case study, a linear time-varying (LTV) dynamic model is considered given by

$$\begin{aligned} x_t &= A_t x_{t-1} + w_t \quad w : W \sim \mathcal{N}(w; 0, \sigma_W^2), \\ y_t &= C_t x_t + v_t \quad v : V \sim \mathcal{N}(v; 0, \sigma_V^2), \end{aligned} \quad (22)$$

where the transition (\mathbf{A}_t) and observation (\mathbf{C}_t) equations at a time t are

$$\begin{aligned} A_t &= 0.8 - 0.1 \sin\left(\frac{7\pi t}{T}\right), \\ C_t &= 1 - 0.99 \sin\left(\frac{100\pi t}{T}\right), \end{aligned}$$

in which T represents the total number of time steps. The process error (w) and observation error (v) are Gaussian and are assumed to be independent at any time t . The process error variance σ_W^2 is unknown and needs to be inferred. The observation error variance σ_V^2 is assumed to be known and is equal to the true σ_W^2 which is randomly selected from the prior PDF of $\overline{W^2}_{0|0}$ such that

$$\overline{W^2}_{0|0} \sim \mathcal{N}\left(\overline{w^2}_0; \mu_{0|0}^{\overline{W^2}}, (\sigma_{0|0}^{\overline{W^2}})^2\right).$$

Three different true values for σ_W^2 are generated by considering different prior initialization for the pair of $\{\mu_{0|0}^{\overline{W^2}}, (\sigma_{0|0}^{\overline{W^2}})^2\}$ such that the three cases are (a) $\{\mu_{0|0}^{\overline{W^2}} = 0.2, (\sigma_{0|0}^{\overline{W^2}})^2 = 0.01\}$, (b) $\{\mu_{0|0}^{\overline{W^2}} = 2, (\sigma_{0|0}^{\overline{W^2}})^2 = 1\}$, and (c) $\{\mu_{0|0}^{\overline{W^2}} = 20, (\sigma_{0|0}^{\overline{W^2}})^2 = 100\}$.

5.1.1 | Online estimation of error variance

Data was simulated using the model specified in Equation (22) with a total of $T = 1000$ observations. Figure 2 shows the online state estimation of the error variance term for each of the three cases. These results confirm that the method is able to perform online inference for different magnitudes of the error variance starting from arbitrary initial estimates.

Two-tailed t -tests were carried out to test the null hypothesis that the error variance estimates obtained using AGVI are equal to the true values. This empirical test, if accepted, proves that the estimates are unbiased. A 95% significance level is chosen for which the critical t -value is 1.96, given that the degrees of freedom are $T - 1 = 999$. The t -test is carried out five times for each three cases. Table 1 shows the t -values where five different runs were carried out for each case. The results show that the computed t -values are within the bounds of the critical t -value, that is, $-1.96 < t\text{-value} < 1.96$, proving that the estimates are unbiased.

5.1.2 | Statistical consistency

The optimality of the filter is evaluated using two chi-square (χ^2) tests that rely on the *normalised estimation error squared* (NEES) and the *normalised innovation error squared* (NIS) values.⁴² These tests are conducted using 50 random

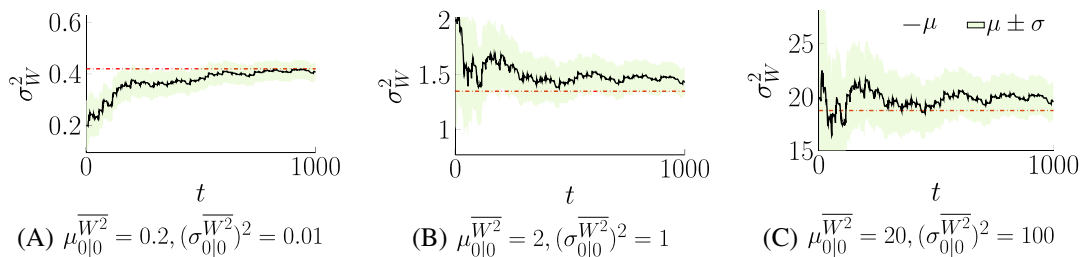


FIGURE 2 Online estimation of the error variance term for each of the three cases for which the different prior initializations are (A) $\mu_{0|0}^{\overline{W^2}} = 0.2, (\sigma_{0|0}^{\overline{W^2}})^2 = 0.01$, (B) $\mu_{0|0}^{\overline{W^2}} = 2, (\sigma_{0|0}^{\overline{W^2}})^2 = 1$, and (C) $\mu_{0|0}^{\overline{W^2}} = 20, (\sigma_{0|0}^{\overline{W^2}})^2 = 100$. The true σ_W^2 value in each case is shown in red dashed line, while the estimated values and their $\pm 1\sigma$ uncertainty bound are shown in black and green shaded area.

TABLE 1 t -values computed in all three cases, that is, (a) $\{\mu_{0|0}^{\overline{W^2}} = 0.2, (\sigma_{0|0}^{\overline{W^2}})^2 = 0.01\}$, (b) $\{\mu_{0|0}^{\overline{W^2}} = 2, (\sigma_{0|0}^{\overline{W^2}})^2 = 1\}$, and (c) $\{\mu_{0|0}^{\overline{W^2}} = 20, (\sigma_{0|0}^{\overline{W^2}})^2 = 100\}$.

t-values		
case (a)	case (b)	case (c)
-1.67	0.46	1.53
-1.29	-0.94	1.51
-0.87	1.26	-1.19
-0.79	1.21	1.03
-1.88	-0.64	-1.87

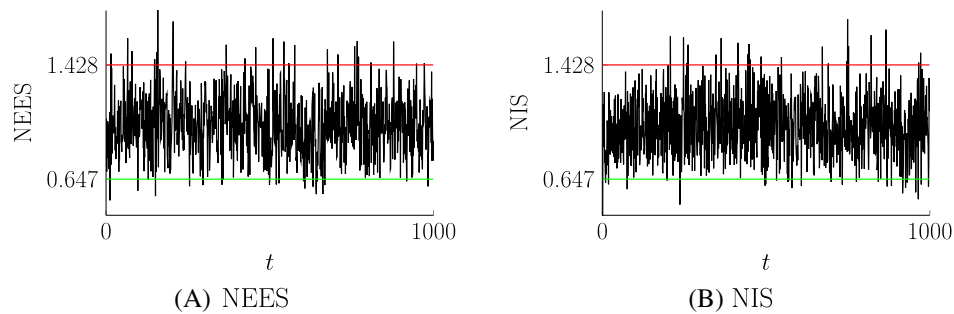


FIGURE 3 Illustration showing the average normalized state estimation error squared (NEES) and the average normalized innovation squared (NIS) for the case study (A), that is, $\{\mu_{0|0}^{\overline{W^2}} = 0.2, (\sigma_{0|0}^{\overline{W^2}})^2 = 0.01\}$, with its 95% probability region given by [0.647, 1.428] is marked by the green and red lines.

TABLE 2 Average number of points outside the 95% probability region for the NEES and NIS values in all the three cases, that is, (a) $\{\mu_{0|0}^{\overline{W^2}} = 0.2, (\sigma_{0|0}^{\overline{W^2}})^2 = 0.01\}$, (b) $\{\mu_{0|0}^{\overline{W^2}} = 2, (\sigma_{0|0}^{\overline{W^2}})^2 = 1\}$, and (c) $\{\mu_{0|0}^{\overline{W^2}} = 20, (\sigma_{0|0}^{\overline{W^2}})^2 = 100\}$.

Criteria	Average number of points		
	Case (a)	Case (b)	Case (c)
NEES	51	50.8	51.2
NIS	53	51.2	49.8

simulations. Considering a 95% confidence interval (C.I.) and the degrees of freedom $\nu = X = Y = 1$, that is, the size of the state and observation vector, the two-sided probability region is given by $[\chi_{50}^2(0.025), \chi_{50}^2(0.975)] = [32.3, 71.4]$. By dividing the range by 50, we obtain the probability region for the average NEES and NIS values [0.647, 1.428]. Figure 3 illustrates an example of the 95% region marked by the green and red lines for both the average NEES and NIS values in case (a). From the definition of the test, there should be approximately 5% of the total number of points outside the 95% region. The length of the time series for the case study is 1000 and hence, approximately 50 points should be outside the region. Table 2 presents the average number of points outside the probability region for both the NEES and NIS tests in all three cases where each of the 50 runs are carried out five times in order to compute the average value. The results verify that the filter is optimal and provide consistent estimates for the error variance term, given that the number of points outside the 95% probability region are in accordance to the theoretical results.

5.1.3 | Statistical consistency for the variance of the error variance

In order to check the statistical consistency for the variance of the error variance term, we created 1000 simulated time series where the true values of the error variance in each time series is generated from the prior knowledge of $\overline{W^2}_{0|0}$. Figure 4 presents, for each time step, the percentage of realizations (γ) where the true value lies within the confidence

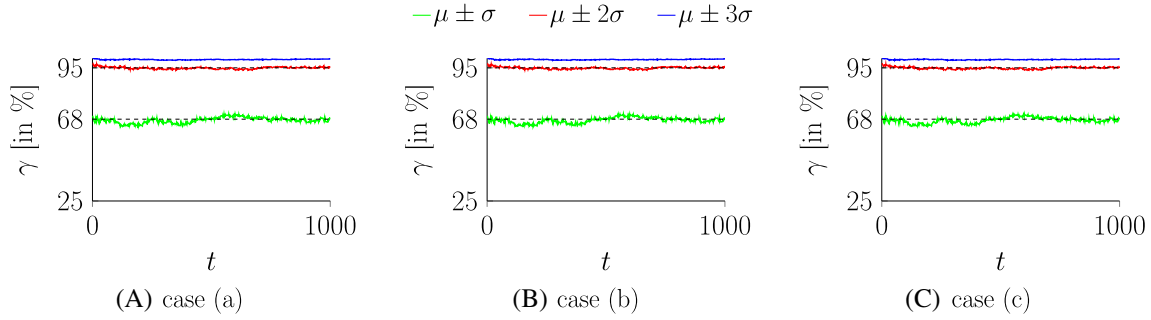


FIGURE 4 Empirical consistency check for the variance of the error variance estimate, where γ is the percentage of realizations where the true value lies within the three C.I. for the cases (A) $\{\mu_{0|0}^{\overline{W}^2} = 0.2, (\sigma_{0|0}^{\overline{W}^2})^2 = 0.01\}$, (B) $\{\mu_{0|0}^{\overline{W}^2} = 2, (\sigma_{0|0}^{\overline{W}^2})^2 = 1\}$, and (C) $\{\mu_{0|0}^{\overline{W}^2} = 20, (\sigma_{0|0}^{\overline{W}^2})^2 = 100\}$.

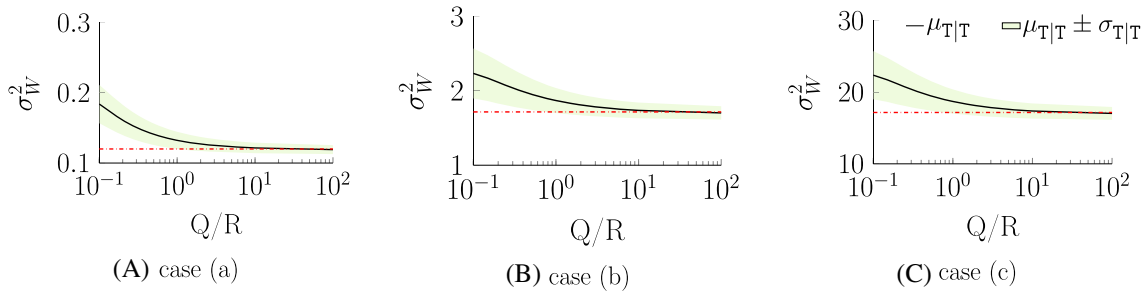


FIGURE 5 The posterior mean estimate and C.I. of the error variance for different values of $\frac{Q}{R}$ for the cases (A) $\{\mu_{0|0}^{\overline{W}^2} = 0.2, (\sigma_{0|0}^{\overline{W}^2})^2 = 0.01\}$, (B) $\{\mu_{0|0}^{\overline{W}^2} = 2, (\sigma_{0|0}^{\overline{W}^2})^2 = 1\}$, and (C) $\{\mu_{0|0}^{\overline{W}^2} = 20, (\sigma_{0|0}^{\overline{W}^2})^2 = 100\}$. Note that the x-axis is in log-scale.

interval (C.I.) for 1, 2, and 3 standard deviations from the mean estimate in each of the three case studies. The results in Figure 4 show that the γ values match the theoretical C.I. quantities, that is, {68, 95, 99}%, for the Gaussian distribution supporting the hypothesis that the Gaussian PDF for the error variance is adequate at each time step.

5.1.4 | Impact of the $\frac{Q}{R}$ ratio

Also, we noticed the effect of the $\frac{Q}{R} = \frac{\sigma_W^2}{\sigma_V^2}$ ratio on the estimation accuracy. Figure 5 shows the posterior mean estimate $\mu_{T|T}$ and the confidence interval $\mu_{T|T} \pm \sigma_{T|T}$ for the error variance after T time steps with respect to different $\frac{Q}{R}$ values for the three cases, where $T = 1000$ is the total length of the time series. The results validate that the AGVI method is accurate for $\frac{Q}{R} \geq 10$. For $1 < \frac{Q}{R} < 10$, the estimates are accurate with negligible biases in comparison to the true values, whereas for $\frac{Q}{R} < 1$, the estimates are inaccurate with large biases. This phenomenon is explained by the fact that, given that the system is observable,^{43,44} the Kalman gain has a higher value with an increase in the $\frac{Q}{R}$ ratio. As a result, the Kalman filter puts more weight on the measurements. Hence, we obtain a better mean and variance estimate of W by learning from each measurement which in turn provide better estimates for \overline{W}^2 .

5.1.5 | Comparison of AGVI with adaptive Kalman filters

The AGVI method is compared with two existing adaptive Kalman filter (AKF) methods, namely SWVAKF⁴⁵ and MDM.^{27,28} Three cases are created for comparison, where the true values for the process error variances are (a) $\sigma_W^2 = 0.42$, (b) $\sigma_W^2 = 1.35$, and (c) $\sigma_W^2 = 18.75$, using the same dynamic model shown in Equation (22). In this case study, we use the publicly available codes for MDM⁴⁶ and SWVAKF⁴⁵ that only provides point estimates of \mathbf{Q} and \mathbf{R} matrices. As the codes are not available for obtaining the uncertainty associated with the error variance estimates, our study is restricted to analyzing the mean estimates, as made available by the authors. The MDM method involves no prior knowledge of the hidden states or the error variances, and the user-defined lag parameter is set to 1. For SWVAKF, the same values are used

TABLE 3 Comparison of the average RMSE values and the computational time (s) obtained from each method in all three cases where the true values are (a) $\sigma_W^2 = 0.42$, (b) $\sigma_W^2 = 1.35$, and (c) $\sigma_W^2 = 18.75$. The results are averaged over five independent runs. Both the methods are picked from different AKF categories where AGVI and SWVAKF are Bayesian methods, MDM is a correlation method.

Category	Methods	RMSE			Time (s)
		$\sigma_W^2 = 0.42$	$\sigma_W^2 = 1.35$	$\sigma_W^2 = 18.75$	
Bayesian	AGVI	0.043	0.083	2.06	0.004
Bayesian	SWVAKF	0.067	0.138	3.18	2.638
Correlation	MDM	0.053	0.094	2.65	0.069

for filtering parameters as provided in the method's implementation code.²⁴ The prior knowledge for the hidden states in both AGVI and SWVAKF are set to $\mu_{0|0} = 0$ and $\sigma_{0|0}^2 = 100$.

Table 3 compares the average root mean square error (RMSE) values and computational time obtained using all three methods for the three cases, where the true values are (a) $\sigma_W^2 = 0.42$, (b) $\sigma_W^2 = 1.35$, and (c) $\sigma_W^2 = 18.75$. The results show that AGVI outperforms all methods in terms of predictive capacity. In comparison to SWVAKF, which is a Bayesian method, AGVI is more than two orders of magnitude faster, and compared to MDM, which is a correlation method, it is an order of magnitude faster. Thus, we conclude that AGVI provides unbiased and consistent estimates for the process error variance given that the observation error variance is known. The comparative study shows that AGVI provides better predictive capacity and computational speed than both AKF methods.

5.2 | Case Study 2: Multivariate random walk model

This case study is conducted using five simulated datasets of 1000 time steps with a transition process error having a full covariance matrix \mathbf{Q} . The vector of hidden states \mathbf{x}_t associated with the five time series is given by

$$\mathbf{x}_t = [x_t^{L1} x_t^{L2} x_t^{L3} x_t^{L4} x_t^{L5}]^T.$$

The state transition matrix \mathbf{A} and the observation matrix \mathbf{C} are defined as $\mathbf{A} = \mathbf{I}_5$, and $\mathbf{C} = \mathbf{I}_5$, The \mathbf{Q} and \mathbf{R} matrices are defined as

$$\mathbf{Q} = \begin{bmatrix} 1 & -0.3 & -0.2 & -0.1 & 0.25 \\ -0.3 & 3 & 0.35 & 0.4 & 0.45 \\ -0.2 & 0.35 & 4 & 0.5 & 0.55 \\ -0.1 & 0.4 & 0.5 & 0.8 & 0.6 \\ 0.25 & 0.45 & 0.55 & 0.6 & 2 \end{bmatrix},$$

$$\mathbf{R} = 0.1 \cdot \mathbf{I}_5,$$

where the off-diagonal covariance terms in the \mathbf{Q} matrix are selected arbitrarily such that it is symmetric and semi positive-definite.

5.2.1 | Online estimation of full \mathbf{Q} matrix

For AGVI, the prior knowledge for the augmented hidden states $\tilde{\boldsymbol{\mu}}_{0|0} = [\boldsymbol{\mu}_{0|0}; \overrightarrow{\boldsymbol{\mu}}_{0|0}^L]$ and $\tilde{\boldsymbol{\Sigma}}_{0|0} = \text{blkdiag}(\boldsymbol{\Sigma}_{0|0}, \overrightarrow{\boldsymbol{\Sigma}}_{0|0}^L)$ are initialized by

$$\begin{aligned} \tilde{\boldsymbol{\mu}}_{0|0} &= [\mathbf{0}_5^T \ 2 \cdot \mathbf{1}_5^T \ 0.8 \cdot \mathbf{1}_{10}^T]^T, \\ \tilde{\boldsymbol{\Sigma}}_{0|0} &= \text{diag}([\mathbf{1}_5^T \ 0.5 \cdot \mathbf{1}_5^T \ 0.5 \cdot \mathbf{1}_{10}^T]), \end{aligned} \quad (23)$$

where $\mathbf{0}$ and $\mathbf{1}$ represent vector of zeros and ones, respectively. The mean vector and the covariance matrix for

$$\overrightarrow{\mathbf{L}}_{0|0}^L = [L_{11} \ L_{22} \ \cdots \ L_{55} \ L_{12} \ \cdots \ L_{45}]_{0|0}^T,$$

are given by

$$\begin{aligned}\vec{\mu}_{0|0}^{\mathbf{L}\vec{W}} &= [2 \cdot \mathbf{1}_5^T \ 0.8 \cdot \mathbf{1}_{10}^T]^T, \\ \vec{\Sigma}_{0|0}^{\mathbf{L}\vec{W}} &= \text{diag}([0.5 \cdot \mathbf{1}_5^T \ 0.5 \cdot \mathbf{1}_{10}^T]).\end{aligned}\quad (24)$$

Note that the matrix \mathbf{R} is assumed to be known and only the matrix \mathbf{Q} is inferred. Figure 6 shows the estimates obtained using AGVI for four elements of the \mathbf{Q} matrix, namely σ_{11}^2 , σ_{22}^2 , σ_{12} , and σ_{13} , chosen arbitrarily. The plots for the remaining elements are provided in Figure E1 (Appendix E).

The two-tailed t -test was carried out to test the null hypothesis that the error variance estimates obtained using AGVI are equal to the true values and hence, unbiased. A 95% significance level was chosen for which the critical t -value is 1.96, given that the degree of freedom is $T - 1 = 999$. Table 4 shows the average t -values for all the variance terms. The results show that the computed t -values are within the bounds of the critical t -value, that is, $-1.96 < t\text{-value} < 1.96$, proving that the estimates are unbiased.

5.2.2 | Statistical consistency

The normalized innovation square (NIS) metric is used to check the consistency of the estimator. The two-sided probability region for a 95% C.I. having five degrees of freedom, that is, the size of the observation vector \mathbf{Y} , is [0.831 12.833]. Considering that the total length of the training set is 1000, the theoretical 5% value for the number of acceptable points outside the 95% C.I. is 50. The different prior initialization are chosen such that $\vec{\mu}_{0|0}^{\mathbf{L}\vec{W}} = [\alpha \cdot \mathbf{1}_5^T \ \beta \cdot \mathbf{1}_{10}^T]^T$, where $\alpha = \{1.5 : 0.1 : 2\}$ and $\beta = \{0.5 : 0.1 : 1\}$ while considering the same covariance matrix as defined in Equation (24). Table 5 presents the average number of points outside the probability region for the different prior initialization of $\vec{\mu}_{0|0}^{\mathbf{L}\vec{W}}$ where the average value is computed using five simulated datasets for each combination of $\{\alpha, \beta\}$. The results show that, on average, there are ≈ 54 points that lie outside the 95% probability region, which is comparable to the theoretical value of 50. This verifies that the filter is optimal and provides consistent estimates for the error variance and covariance terms of the full \mathbf{Q} matrix.

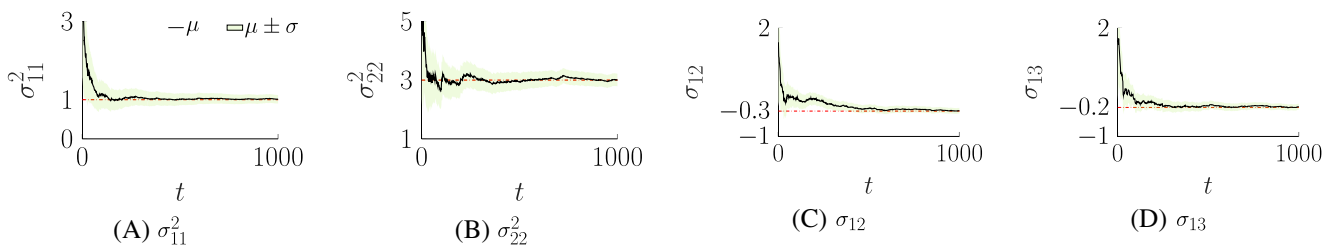


FIGURE 6 Online estimation of the error variance terms (A) σ_{11}^2 and (B) σ_{22}^2 and the covariance terms (C) σ_{12} and (D) σ_{13} from the full \mathbf{Q} matrix compared to their true values marked by the dashed red line. The estimated values are shown by the black solid line and their $\pm 1\sigma$ uncertainty bound is shown using the green shaded region.

TABLE 4 Average t -values for all the variance terms in the full \mathbf{Q} matrix. Five independent runs were carried out. The variance terms and the covariance terms are represented by σ_{ii}^2 and σ_{ij} , $\forall i, j \in 1, \dots, D$, respectively.

Variance terms															
σ_{11}^2	σ_{12}	σ_{13}	σ_{14}	σ_{15}	σ_{22}^2	σ_{23}	σ_{24}	σ_{25}	σ_{33}^2	σ_{34}	σ_{35}	σ_{44}^2	σ_{45}	σ_{55}^2	
t -values	0.6002	-0.0736	0.1006	0.1380	-0.2828	0.3387	0.2157	0.2527	-0.1177	-0.0423	-0.0981	-0.2424	0.0692	-0.0855	-0.0287

TABLE 5 Average number of points outside the 95% probability region for the different prior initialization of $\vec{\mu}_{0|0}^{\mathbf{L}\vec{W}}$. Each column presents the average value computed using the five simulated datasets for one combination of $\{\alpha, \beta\}$.

	{1.5, 0.5}	{1.6, 0.6}	{1.7, 0.7}	{1.8, 0.8}	{1.9, 0.9}	{2, 1}	Mean
NIS	54.8	54.6	53.8	53.8	53.6	53.6	54.03

TABLE 6 Comparison of the average RMSE values and the computational time for each method. Each of the methods are picked from different AKF categories where AGVI and SWVAKF are Bayesian methods, MDM is a correlation method. The variance terms and the covariance terms are represented by σ_{ii}^2 and σ_{ij} , $\forall i, j \in 1, \dots, D$, respectively.

Variance terms															Time (s)	
Methods	σ_{11}^2	σ_{12}	σ_{13}	σ_{14}	σ_{15}	σ_{22}^2	σ_{23}	σ_{24}	σ_{25}	σ_{33}^2	σ_{34}	σ_{35}	σ_{44}^2	σ_{45}		σ_{55}^2
AGVI	0.0766	0.0409	0.0160	0.0476	0.0139	0.0468	0.2097	0.0004	0.0578	0.1653	0.0528	0.0123	0.0251	0.0115	0.0422	1.4
SWVAKF	0.0246	0.0789	0.1364	0.0857	0.1023	0.5607	0.0754	0.0995	0.2539	1.1391	0.1661	0.1366	0.0130	0.0265	0.2304	9.4
MDM	0.0034	0.0197	0.0675	0.0437	0.0629	0.4431	0.5396	0.0581	0.1306	0.0347	0.0009	0.0736	0.0413	0.0609	0.1412	0.03

5.2.3 | Comparison with adaptive Kalman filters

For this case study, the AGVI method is compared to SWVAKF and MDM. For the SWVAKF method, the hidden states are initialized similarly to Equation (23), where the mean vector is $\mu_{0|0} = \mathbf{0}_5$ and the covariance matrix is $\Sigma_{0|0} = \mathbf{I}_5$. The same values are used for filtering parameters as provided in the method's implementation code.²⁴ The MDM method involves no prior knowledge of the hidden states or error variances, and the user-defined lag parameter is set to 1. Table 6 shows the average RMSE values for estimating some of the elements chosen arbitrarily from the \mathbf{Q} matrix as well as the average computational time for each method. The results show that AGVI outperforms other methods in terms of predictive capacity for most of the variance and covariance terms. In comparison to SWVAKF which is a Bayesian method, it is more than an order of magnitude faster. The MDM is an order of magnitude faster than AGVI but the current implementation code only provides point estimates, whereas AGVI estimates the mean as well as the variance of the error variance terms. The results for both case studies were obtained using the AGVI MATLAB library.⁴⁷

6 | CONCLUSION

The approximate Gaussian variance inference (AGVI) method proposed in this article is an analytically tractable Bayesian inference method that provides many advantages. First, it allows for online inference of the process error variance and covariance terms in the full \mathbf{Q} matrix as Gaussian random variables. Second, it provides accurate, unbiased, and statistically consistent estimates of the mean and uncertainties associated with the error variance terms at each time step. Third, it employs a closed-form square-root filtering technique using the Cholesky decomposition to maintain the estimated \mathbf{Q} matrix as positive semi-definite.

The case study 1 shows the application of the AGVI method for a linear time-varying model where the univariate process error variance was inferred starting from different prior initializations. The t -test proves that the estimates are unbiased. The statistical consistency tests verify that the filter is optimal and that the AGVI method provide consistent estimates for the mean as well as the uncertainties associated with the error variance term. In comparison to the existing adaptive Kalman filter (AKF) methods, the AGVI method provides a better predictive capacity in all cases. In comparison to SWVAKF, which is a Bayesian method, AGVI is more than two orders of magnitude faster, and compared to MDM, which is a correlation method, it is an order of magnitude faster. The case study 2 shows the application of AGVI for a multivariate random walk model with a full \mathbf{Q} matrix and compares its performance with the two AKF methods. The results show that AGVI outperforms all methods in terms of predictive capacity for most of the variance and covariance terms, and yields statistically consistent estimates. Hence, the proposed method is capable of online estimation of error variances in regard to state-space models involving multiple time series.

ACKNOWLEDGMENTS

This project was funded by the Natural Sciences and Engineering Research Council of Canada (NSERC), Hydro-Québec (HQ), Hydro-Québec's Research Institute (IREQ), Institute For Data Valorization (IVADO). The authors thank Luong Ha Nguyen, Post doc, Department of Civil, Geologic and Mining Engineering, Polytechnique Montréal for his help in the project.

ORCID

Bhargob Deka  <https://orcid.org/0000-0002-8585-0738>

REFERENCES

1. Mehra R. On the identification of variances and adaptive Kalman filtering. *IEEE Trans Autom Control*. 1970;15(2):175-184.
2. Duník J, Straka O, Kost O, Havlík J. Noise covariance matrices in state-space models: A survey and comparison of estimation methods—Part I. *Int J Adapt Control Signal Process*. 2017;31(11):1505-1543.
3. Matisko P, Havlena V. Noise covariance estimation for Kalman filter tuning using Bayesian approach and Monte Carlo. *Int J Adapt Control Signal Process*. 2013;27(11):957-973.
4. Mehra R. Approaches to adaptive filtering. *IEEE Trans Autom Control*. 1972;17(5):693-698.
5. Wang J. Stochastic Modeling for Real-Time Kinematic GPS/GLONASS Positioning. *Navigation*. 1999;46(4):297-305.
6. Mehra R. On-line identification of linear dynamic systems with applications to Kalman filtering. *IEEE Trans Autom Control*. 1971;16(1):12-21.

7. Goulet JA, Nguyen LH, Amiri S. Tractable approximate Gaussian inference for Bayesian neural networks. *J Mach Learn Res.* 2021;22(251):1-23.
8. Deka B, Ha Nguyen L, Amiri S, Goulet JA. The Gaussian multiplicative approximation for state-space models. *Struct Control Health Monit.* 2021;29:e2904.
9. Gelman A, Carlin B, Stern S, Dunson B, Vehtari A, Rubin B. *Bayesian Data Analysis.* Chapman and Hall/CRC; 2013.
10. Murphy KP. *Machine Learning: A Probabilistic Perspective.* MIT Press; 2012.
11. Belanger PR. Estimation of noise covariance matrices for a linear time-varying stochastic process. *Automatica.* 1974;10(3):267-275.
12. Zhou J, Luecke R. Estimation of the covariances of the process noise and measurement noise for a linear discrete dynamic system. *Comput Chem Eng.* 1995;19(2):187-195.
13. Odelson BJ, Rajamani MR, Rawlings JB. A new autocovariance least-squares method for estimating noise covariances. *Automatica.* 2006;42(2):303-308.
14. Duník J, Straka O, Kost O. Measurement difference autocovariance method for noise covariance matrices estimation. *2016 IEEE 55th Conference on Decision and Control (CDC).* IEEE; 2016:365-370.
15. Mussot V, Mercère G, Dairay T, Arvis V, Vayssettes J. Noise covariance matrix estimation with subspace model identification for Kalman filtering. *Int J Adapt Control Signal Process.* 2021;35(4):591-611.
16. Brumana A, Piroddi L. A multi-tone central divided difference frequency tracker with adaptive process noise covariance tuning. *Int J Adapt Control Signal Process.* 2020;34(7):877-900.
17. Myers K, Tapley B. Adaptive sequential estimation with unknown noise statistics. *IEEE Trans Autom Control.* 1976;21(4):520-523.
18. Sage P, Husa W. Algorithms for sequential adaptive estimation of prior statistics. *1969 IEEE Symposium on Adaptive Processes (8th) Decision and Control.* IEEE; 1969:61-61.
19. Assa A, Plataniotis KN. Adaptive Kalman filtering by covariance sampling. *IEEE Signal Process Lett.* 2017;24(9):1288-1292.
20. Shumway RH, Stoffer DS, Stoffer DS. *Time Series Analysis and its Applications.* Springer; 2000.
21. Kashyap R. Maximum likelihood identification of stochastic linear systems. *IEEE Trans Autom Control.* 1970;15(1):25-34.
22. Sarkka S, Nummenmaa A. Recursive noise adaptive Kalman filtering by variational Bayesian approximations. *IEEE Trans Autom Control.* 2009;54(3):596-600.
23. Kantas N, Doucet A, Singh SS, Maciejowski J, Chopin N. On particle methods for parameter estimation in state-space models. *Stat Sci.* 2015;30(3):328-351.
24. Huang Y, Zhu F, Jia G, Zhang Y. A slide window variational adaptive Kalman filter. *IEEE Trans Circuits Syst II Exp Briefs.* 2020;67(12):3552-3556.
25. Pavelková L, Kárný M. State and parameter estimation of state-space model with entry-wise correlated uniform noise. *Int J Adapt Control Signal Process.* 2014;28(11):1189-1205.
26. Duník J, Kost O, Straka O. Design of measurement difference autocovariance method for estimation of process and measurement noise covariances. *Automatica.* 2018;90:16-24.
27. Duník J, Kost O, Straka O, Blasch E. Covariance estimation and Gaussianity assessment for state and measurement noise. *J Guid Control Dyn.* 2020;43(1):132-139.
28. Kost O, Duník J, Straka O. Measurement difference method: A universal tool for noise identification. *IEEE Trans Autom Control.* 2022;68(3):1792-1799.
29. Kay SM. *Fundamentals of Statistical Signal Processing: Estimation Theory.* Prentice-Hall, Inc.; 1993.
30. Bavdekar VA, Deshpande AP, Patwardhan SC. Identification of process and measurement noise covariance for state and parameter estimation using extended Kalman filter. *J Process Control.* 2011;21(4):585-601.
31. Särkkä S. *Bayesian filtering and smoothing.* Cambridge University Press; 2013.
32. Kontoroupi T, Smyth AW. Online noise identification for joint state and parameter estimation of nonlinear systems. *ASCE-ASME J Risk Uncertain Eng Syst A Civil Eng.* 2016;2(3):B4015006.
33. Storvik G. Particle filters for state-space models with the presence of unknown static parameters. *IEEE Trans Signal Process.* 2002;50(2):281-289.
34. Nguyen H, Goulet JA. Real-time anomaly detection with Bayesian dynamic linear models. *Struct Control Health Monit.* 2019;26:e2404.
35. Fearnhead P. Markov chain Monte Carlo, sufficient statistics, and particle filters. *J Comput Graph Stat.* 2002;11(4):848-862.
36. Li XR, Bar-Shalom Y. A recursive multiple model approach to noise identification. *IEEE Trans Aerosp Electron Syst.* 1994;30(3):671-684.
37. Särkkä S, Hartikainen J. Non-linear noise adaptive Kalman filtering via variational Bayes. *2013 IEEE International Workshop on Machine Learning for Signal Processing.* IEEE; 2013:1-6.
38. Ma J, Lan H, Wang Z, Wang X, Pan Q, Moran B. Improved adaptive Kalman filter with unknown process noise covariance. *2018 21st International Conference on Information Fusion.* IEEE; 2018:1-5.
39. Ardeshiri T, Özkan E, Orguner U, Gustafsson F. Approximate Bayesian smoothing with unknown process and measurement noise covariances. *IEEE Signal Process Lett.* 2015;22(12):2450-2454.
40. Huang Y, Zhang Y, Wu Z, Li N, Chambers J. A novel adaptive Kalman filter with inaccurate process and measurement noise covariance matrices. *IEEE Trans Autom Control.* 2017;63(2):594-601.
41. Goulet JA. *Probabilistic Machine Learning for Civil Engineers.* MIT Press; 2020.
42. Bar-Shalom Y, Li XR, Kirubarajan T. *Estimation with Applications to Tracking and Navigation: Theory Algorithms and Software.* John Wiley & Sons; 2004.

43. Simon D. *Optimal State Estimation: Kalman, H Infinity, and Nonlinear Approaches*. John Wiley & Sons; 2006.
44. Bryson AE. Applied optimal control: Optimization. *Estimization Control*. 1975;2:2.
45. Huang Y, Zhu F, Jia G, Zhang Y. *Implementation codes for the paper A Slide Window Variational Adaptive Kalman Filter*. IEEE. https://www.researchgate.net/publication/342466130_Implementation_codes_for_the_paper_A_Slide_Window_Variational_Adaptive_Kalman_Filter; 2020
46. Dunik J, Straka O, Kost O. The identification and decision making research group (IDM). <https://idm.kky.zcu.cz/sw.html>; 2015.
47. Deka B, Goulet JA. AGVI. <https://github.com/CivML-PolyMtl/AGVI>; 2023.

How to cite this article: Deka B, Goulet J-A. Approximate Gaussian variance inference for state-space models. *Int J Adapt Control Signal Process*. 2023;1-29. doi: 10.1002/acs.3667

APPENDIX A. GAUSSIAN MULTIPLICATION APPROXIMATION (GMA)

Consider $\mathbf{X} = [X_1 X_2 X_3 X_4]^T$, a vector of Gaussian random variables such that $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Using the Gaussian moment generating function or 2^{nd} order Taylor series expansion, the following equations hold for the product of any two Gaussian random variables such that

$$\mathbb{E}[X_1 X_2] = \mu_1 \mu_2 + \text{cov}(X_1, X_2), \quad (\text{A1})$$

$$\begin{aligned} \text{var}(X_1 X_2) &= \sigma_1^2 \sigma_2^2 + \text{cov}(X_1, X_2)^2 \\ &\quad + 2\text{cov}(X_1, X_2) \mu_1 \mu_2 \\ &\quad + \sigma_1^2 \mu_2^2 + \sigma_2^2 \mu_1^2, \end{aligned} \quad (\text{A2})$$

$$\text{cov}(X_3, X_1 X_2) = \text{cov}(X_1, X_3) \mu_2 + \text{cov}(X_2, X_3) \mu_1, \quad (\text{A3})$$

$$\begin{aligned} \text{cov}(X_1 X_2, X_3 X_4) &= \text{cov}(X_1, X_3) \text{cov}(X_2, X_4) \\ &\quad + \text{cov}(X_1, X_4) \text{cov}(X_2, X_3) \\ &\quad + \text{cov}(X_1, X_3) \mu_2 \mu_4 \\ &\quad + \text{cov}(X_1, X_4) \mu_2 \mu_3 + \text{cov}(X_2, X_3) \mu_1 \mu_4 \\ &\quad + \text{cov}(X_2, X_4) \mu_1 \mu_3. \end{aligned} \quad (\text{A4})$$

These equations can be obtained either using moment generating functions⁷ or second-order Taylor series expansion.⁸

APPENDIX B. PROOFS FOR UNIVARIATE PROCESS ERROR

B.1 Proof for Lemma 1

Proof. Given that W is Gaussian and has a zero mean, the moments of W can be derived using a Gaussian moment generating function so that

$$\begin{aligned} \mu_W &= \mathbb{E}[W] = 0, \\ \sigma_W^2 &= \mathbb{E}[(W - \mu_W)^2] = \mathbb{E}[W^2] - \mathbb{E}[W]^2, \\ &= \mathbb{E}[W^2], \end{aligned} \quad (\text{B1})$$

$$\mathbb{E}[W^4] = 3\mathbb{E}[W^2]^2, \quad (\text{B2})$$

where using Equations (B1) and (B2), we can define the variance of W^2 such that

$$\text{var}(W^2) = \mathbb{E}[(W^4)] - \mathbb{E}[W^2]^2 = 2\mathbb{E}[W^2]^2. \quad (\text{B3})$$

If we make the approximation that $W^2 \sim \mathcal{N}(w^2; \mu^{W^2}, (\sigma^{W^2})^2)$ is a Gaussian random variable, then the PDF can be fully defined by its mean and variance,

$$\begin{aligned}\mu^{W^2} &= \mathbb{E}[W^2], \\ (\sigma^{W^2})^2 &= \text{var}(W^2) = 2\mathbb{E}[W^2]^2,\end{aligned}$$

where by using Equation (B3), the variance $\text{var}(W^2)$ can also be expressed in terms of the expected value $\mathbb{E}[W^2]$. Hence, the PDF of W^2 only depends on the unknown hyper parameter μ^{W^2} such that

$$\begin{aligned}f(w^2 | \mu^{W^2}, (\sigma^{W^2})^2) &\equiv f(w^2 | \mu^{W^2}), \\ &= \mathcal{N}(w^2, \mu^{W^2}, 2(\mu^{W^2})^2).\end{aligned}\tag{B4}$$

B.2 Proof for Lemma 2

Proof. Let us consider that the mean parameter μ^{W^2} is described by the random variable $\overline{W^2} : \overline{w^2} \in (0, \infty)$ for which

$$f(\overline{w^2}) \sim \mathcal{N}(\overline{w^2}; \mu^{\overline{W^2}}, (\sigma^{\overline{W^2}})^2).\tag{B5}$$

Using (B4) and (B5), we can rewrite the PDF of W^2 as

$$f(w^2 | \overline{w^2}) = \mathcal{N}(w^2; \overline{w^2}, 2(\overline{w^2})^2).\tag{B6}$$

Using the acyclic graph in Figure 1, the joint PDF of $Y_{t|t-1}$, $\mathbf{X}_{t|t-1}$, $W_{t|t-1}^2$, and $\overline{W}_{t|t-1}^2$ is shown by

$$\begin{aligned}f(y_t, x_t, w_t^2, \overline{w}_t^2 | \mathbf{y}_{1:t-1}) &= f(y_t | x_t, w_t^2) \cdot f(x_t | \mathbf{y}_{1:t-1}) \\ &\quad \cdot f(w_t^2 | \overline{w}_t^2) \cdot f(\overline{w}_t^2 | \mathbf{y}_{1:t-1}).\end{aligned}\tag{B7}$$

Using Equation (B6) and marginalizing $Y_{t|t-1}$, $\mathbf{X}_{t|t-1}$, and $\overline{W}_{t|t-1}^2$ from the joint PDF defined in Equation (B7), the prior predictive PDF of $W_{t|t-1}^2$ is

$$\begin{aligned}f(w_t^2 | \mathbf{y}_{1:t-1}) &= \int f(w_t^2 | \overline{w}_t^2) \cdot f(\overline{w}_t^2 | \mathbf{y}_{1:t-1}) d\overline{w}_t^2, \\ &= \int \mathcal{N}(w_t^2; \overline{w}_t^2, 2(\overline{w}_t^2)^2) \cdot \mathcal{N}(\overline{w}_t^2; \mu_{t-1|t-1}^{\overline{W^2}}, (\sigma_{t-1|t-1}^{\overline{W^2}})^2) d\overline{w}_t^2.\end{aligned}\tag{B8}$$

The integration in Equation (B8) can be solved in closed-form. The equivalent formulation to obtain this is to represent the Gaussian random variables W^2 and $\overline{W^2}$ in terms of the standard Gaussian variable ϵ and ζ shown by

$$W^2 = \overline{W^2} + \sqrt{2\overline{W^2}}\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)\tag{B9}$$

$$\overline{W^2} = \mu^{\overline{W^2}} + \sigma^{\overline{W^2}}\zeta, \quad \zeta \sim \mathcal{N}(0, 1).\tag{B10}$$

Using Equations (B9) and (B10), the mean and variance of the prior predictive PDF of $W_{t|t-1}^2$ are given by

$$\begin{aligned}\mathbb{E}[W_{t|t-1}^2] &= \mathbb{E}[\overline{W}_{t|t-1}^2] + \sqrt{2\mathbb{E}[\overline{W}_{t|t-1}^2]} \mathbb{E}[\epsilon], \\ &= \mu_{t-1|t-1}^{\overline{W^2}},\end{aligned}\tag{B11}$$

$$\begin{aligned}\text{var}(W_{t|t-1}^2) &= \text{var}(\overline{W}_{t|t-1}^2) + 2 \text{var}(\overline{W}_{t|t-1}^2 \epsilon), \\ &= (\sigma_{t-1|t-1}^{\overline{W^2}})^2 + 2(\text{var}(\overline{W}_{t|t-1}^2) \cdot \text{var}(\epsilon))\end{aligned}$$

$$\begin{aligned}
& + \text{var}(\epsilon) \cdot \mathbb{E}[\overline{W^2}_{t|t-1}]^2, \\
& = 3(\sigma_{t-1|t-1}^{\overline{W^2}})^2 + 2(\mu_{t-1|t-1}^{\overline{W^2}})^2,
\end{aligned} \tag{B12}$$

where the term $\text{var}(\overline{W^2}_{t|t-1}\epsilon)$ in Equation (B12) is obtained using the GMA equations,

$$\text{var}(\overline{W^2}_{t|t-1}\epsilon) = \text{var}(\overline{W^2}_{t|t-1}) \cdot \text{var}(\epsilon) + \text{var}(\epsilon) \cdot \mathbb{E}[\overline{W^2}_{t|t-1}]^2.$$

Using Equations (B1) and (B11), the error variance term is given by

$$\sigma_W^2 = \mu_{t-1|t-1}^{\overline{W^2}}. \tag{B13}$$

B.3 Proof for Lemma 3

Proof. Let us consider the joint PDF of $Y_{t|t-1}$, $W_{t|t-1}^2$, and $\overline{W^2}_{t|t-1}$ as shown by Figure 1,

$$\begin{aligned}
f(y_t, w_t^2, \overline{w^2}_t | \mathbf{y}_{1:t-1}) &= f(y_t | w_t^2) \cdot f(w_t^2 | \overline{w^2}_t) \cdot f(\overline{w^2}_t | \mathbf{y}_{1:t-1}). \\
&= f(y_t | w_t^2) \cdot \frac{f(w_t^2 | \overline{w^2}_t) \cdot f(\overline{w^2}_t | \mathbf{y}_{1:t-1})}{f(w_t^2 | \mathbf{y}_{1:t-1})} \cdot f(w_t^2 | \mathbf{y}_{1:t-1}), \\
&= f(y_t | w_t^2) \cdot \frac{f(w_t^2, \overline{w^2}_t | \mathbf{y}_{1:t-1})}{f(w_t^2 | \mathbf{y}_{1:t-1})} \cdot f(w_t^2 | \mathbf{y}_{1:t-1}), \\
&= f(y_t | w_t^2) \cdot f(\overline{w^2}_t | w_t^2, \mathbf{y}_{1:t-1}) \cdot f(w_t^2 | \mathbf{y}_{1:t-1}).
\end{aligned} \tag{B14}$$

By dividing both sides in Equation (B14) by $f(y_t | \mathbf{y}_{1:t-1})$ we obtain

$$\begin{aligned}
\frac{f(y_t, w_t^2, \overline{w^2}_t | \mathbf{y}_{1:t-1})}{f(y_t | \mathbf{y}_{1:t-1})} &= \frac{f(y_t | w_t^2) \cdot f(w_t^2 | \mathbf{y}_{1:t-1})}{f(y_t | \mathbf{y}_{1:t-1})} \cdot f(\overline{w^2}_t | w_t^2, \mathbf{y}_{1:t-1}), \\
f(w_t^2, \overline{w^2}_t | \mathbf{y}_{1:t}) &= f(w_t^2 | \mathbf{y}_{1:t}) \cdot f(\overline{w^2}_t | w_t^2, \mathbf{y}_{1:t-1}),
\end{aligned} \tag{B15}$$

By marginalizing out W^2 from the joint PDF defined in Equation (B15), the posterior PDF of $\overline{W^2}_{t|t}$ is obtained such that

$$f(\overline{w^2}_t | \mathbf{y}_{1:t}) = \int f(w_t^2 | \mathbf{y}_{1:t}) \cdot f(\overline{w^2}_t | w_t^2, \mathbf{y}_{1:t-1}) dw_t^2. \tag{B16}$$

B.4 Proof for Lemma 4

Proof. Using the GMA equations in Section A, the posterior PDF of W^2 can be shown by

$$f(w_t^2 | \mathbf{y}_{1:t}) \sim \mathcal{N}(w_t^2; \mu_{t|t}^{W^2}, (\sigma_{t|t}^{W^2})^2), \tag{B17}$$

where the mean and variance of $W_{t|t}^2$ are

$$\begin{aligned}
\mu_{t|t}^{W^2} &= (\mu_{t|t}^W)^2 + (\sigma_{t|t}^W)^2, \\
(\sigma_{t|t}^{W^2})^2 &= 2(\sigma_{t|t}^W)^4 + 4(\sigma_{t|t}^W)^2(\mu_{t|t}^W)^2.
\end{aligned}$$

B.5 Proof for Proposition 2

Proof. Given that the prior predictive PDF of both W^2 and $\overline{W^2}$ are Gaussian, the joint multivariate Gaussian PDF $f(\overline{w^2}_t, w^2_t | \mathbf{y}_{1:t-1})$ is shown by

$$f(\overline{w^2}_t, w^2_t | \mathbf{y}_{1:t-1}) = \mathcal{N} \left(\begin{pmatrix} \overline{w^2}_t \\ w^2_t \end{pmatrix}; \boldsymbol{\mu}_{t|t-1}^{\overline{W^2}, W^2}, \boldsymbol{\Sigma}_{t|t-1}^{\overline{W^2}, W^2} \right), \quad (\text{B18})$$

having a mean vector $\boldsymbol{\mu}_{t|t-1}^{\overline{W^2}, W^2}$ and a covariance matrix $\boldsymbol{\Sigma}_{t|t-1}^{\overline{W^2}, W^2}$ defined as

$$\begin{aligned} \boldsymbol{\mu}_{t|t-1}^{\overline{W^2}, W^2} &= \begin{bmatrix} \overline{\mu}_{t|t-1}^{\overline{W^2}} \\ \mu_{t|t-1}^{W^2} \end{bmatrix}^T, \\ \boldsymbol{\Sigma}_{t|t-1}^{\overline{W^2}, W^2} &= \begin{bmatrix} (\sigma_{t|t-1}^{\overline{W^2}})^2 & \text{cov}(\overline{W^2}, W^2)_{t|t-1} \\ \text{cov}(W^2, \overline{W^2})_{t|t-1} & (\sigma_{t|t-1}^{W^2})^2 \end{bmatrix}, \end{aligned} \quad (\text{B19})$$

and where using the transition model $w^2_t = w^2_{t-1}$, the mean and the variance of $W^2_{t|t-1} = W^2_{t-1|t-1}$ are given by Equations (B11) and (B12). The covariance term $\text{cov}(W^2, \overline{W^2})_{t|t-1}$ between $W^2_{t|t-1}$ and $\overline{W^2}_{t|t-1}$ in Equation (B19) is obtained using Equations (B5) and (B9), and the GMA equations from Section A so that

$$\begin{aligned} \text{cov}(W^2_{t|t-1}, \overline{W^2}_{t|t-1}) &= \text{cov}(W^2, \overline{W^2})_{t|t-1}, \\ &= \text{cov}(\overline{W^2} + \sqrt{2} \overline{W^2} \epsilon, \overline{W^2})_{t|t-1}, \\ &= \text{var}(\overline{W^2})_{t|t-1} + \sqrt{2} \text{cov}(\overline{W^2} \epsilon, \overline{W^2})_{t|t-1}, \\ &= \text{var}(\overline{W^2})_{t|t-1} + \sqrt{2} (\text{cov}(\overline{W^2}, \overline{W^2}) \mathbb{E}[\epsilon] + \text{cov}(\epsilon, \overline{W^2}) \mathbb{E}[\overline{W^2}]), \\ &= (\sigma_{t-1|t-1}^{\overline{W^2}})^2. \end{aligned}$$

Given that the joint Gaussian PDF is defined as shown by Equation (B18), the Gaussian conditional properties are used to obtain the conditional PDF $f(\overline{w^2}_t | w^2_t, \mathbf{y}_{1:t-1})$ which is part of the integrand shown in Equation (B16),

$$f(\overline{w^2}_t | w^2_t, \mathbf{y}_{1:t-1}) = \mathcal{N}(\overline{w^2}_t; \mu_{t|t-1}^{\overline{W^2} | W^2}, (\sigma_{t|t-1}^{\overline{W^2} | W^2})^2), \quad (\text{B20})$$

for which the conditional mean and variance are

$$\mu_{t|t-1}^{\overline{W^2} | W^2} = \mu_{t|t-1}^{\overline{W^2}} + k_t (w^2_t - \mu_{t|t-1}^{W^2}), \quad (\text{B21})$$

$$(\sigma_{t|t-1}^{\overline{W^2} | W^2})^2 = (\sigma_{t|t-1}^{\overline{W^2}})^2 - k_t^2 (\sigma_{t|t-1}^{W^2})^2, \quad (\text{B22})$$

$$\begin{aligned} k_t &= \frac{\text{cov}(W^2_{t|t-1}, \overline{W^2}_{t|t-1})}{(\sigma_{t|t-1}^{W^2})^2}, \\ &= \frac{(\sigma_{t-1|t-1}^{\overline{W^2}})^2}{(\sigma_{t|t-1}^{W^2})^2}. \end{aligned} \quad (\text{B23})$$

Using Equations (B20) and (B17), Equation (B16) is rewritten as

$$f(\overline{w}^2_t | \mathbf{y}_{1:t}) = \int \mathcal{N}(\overline{w}^2_t; \mu_{t|t-1}^{\overline{w}^2 | W^2}, (\sigma_{t|t-1}^{\overline{w}^2 | W^2})^2) \cdot \mathcal{N}(w_t^2; \mu_{t|t}^{W^2}, (\sigma_{t|t}^{W^2})^2) dw_t^2. \quad (\text{B24})$$

Equation (B24) can be solved in closed-form having a Gaussian PDF with a random mean, that is, $\mu_{t|t-1}^{\overline{w}^2 | W^2}$, and a constant variance, that is, $(\sigma_{t|t-1}^{\overline{w}^2 | W^2})^2$, shown by Equations (B21) and (B22). Hence, the PDF $f(\overline{w}^2_t | \mathbf{y}_{1:t})$ is also Gaussian such that

$$f(\overline{w}^2_t | \mathbf{y}_{1:t}) = \mathcal{N}(\overline{w}^2_t; \mu_{t|t}^{\overline{w}^2}, (\sigma_{t|t}^{\overline{w}^2})^2),$$

for which the posterior mean and the variance can be computed using the Kalman gain defined in Equation (B23) as

$$\begin{aligned} \mu_{t|t}^{\overline{w}^2} &= \mathbb{E} \left[\mu_{t|t-1}^{\overline{w}^2} + k_t (W_{t|t}^2 - \mu_{t|t-1}^{W^2}) \right], \\ &= \mu_{t|t-1}^{\overline{w}^2} + k_t (\mu_{t|t}^{W^2} - \mu_{t|t-1}^{W^2}), \\ (\sigma_{t|t}^{\overline{w}^2})^2 &= (\sigma_{t|t-1}^{\overline{w}^2})^2 - k_t^2 (\sigma_{t|t-1}^{W^2})^2 + k_t^2 \text{var}(W_{t|t}^2), \\ &= (\sigma_{t|t-1}^{\overline{w}^2})^2 + k_t^2 ((\sigma_{t|t}^{W^2})^2 - (\sigma_{t|t-1}^{W^2})^2). \end{aligned} \quad (\text{B25})$$

APPENDIX C. PROOFS FOR MULTIVARIATE PROCESS ERROR

C.1 Proof for Lemma 5

As described by Lemma B.2 for the univariate process error, the expected value $\mathbb{E}[(W^i)^2]$ and the variance terms $\text{var}((W^i)^2), \forall i \in \{1, 2, \dots, D\}$ for the prior predictive PDF of $\mathbf{W}^{\mathbf{P}}$ are given by

$$\begin{aligned} \mathbb{E}[(W^i)^2] &= (\mu^{\overline{(W^i)^2}}), \\ \text{var}((W^i)^2) &= 3(\sigma^{\overline{(W^i)^2}})^2 + 2(\mu^{\overline{(W^i)^2}})^2. \end{aligned}$$

Using the GMA equations, the mean and variance term of $W^i W^j$ are

$$\mathbb{E}[W^i W^j] = \text{cov}(W^i, W^j) = \overline{w^i w^j}, \quad (\text{C1})$$

$$\begin{aligned} \text{var}(W^i W^j) &= \text{var}(W^i) \text{var}(W^j) + \text{cov}(W^i, W^j)^2, \\ &= \mu^{\overline{(W^i)^2}} \mu^{\overline{(W^j)^2}} + (\overline{w^i w^j})^2, \end{aligned} \quad (\text{C2})$$

where using Equation (B13) from Proof B.2, $\text{var}(W^i) = \mu^{\overline{(W^i)^2}}$. Using Equations (C1) and (C2), the Gaussian random variable $W^i W^j \sim \mathcal{N}(w^i w^j; \overline{w^i w^j}, \mu^{\overline{(W^i)^2}} \mu^{\overline{(W^j)^2}} + (\overline{w^i w^j})^2)$ can be represented in terms of its standard Gaussian variable ϵ by

$$w^i w^j = \overline{w^i w^j} + \sqrt{\mu^{\overline{(W^i)^2}} \cdot \mu^{\overline{(W^j)^2}} + (\overline{w^i w^j})^2} \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1).$$

The moments of the prior predictive PDF of $W^i W^j$ are given by

$$\begin{aligned} \mathbb{E}[W^i W^j] &= \mathbb{E}[\overline{W^i W^j}] + \mathbb{E} \left[\sqrt{\mu^{\overline{(W^i)^2}} \mu^{\overline{(W^j)^2}} + (\overline{W^i W^j})^2} \cdot \epsilon \right], \\ &= \mu^{\overline{W^i W^j}}, \end{aligned} \quad (\text{C3})$$

$$\begin{aligned}
\text{var}(W^i W^j) &= \text{var}(\overline{W^i W^j}) + \text{var}\left(\sqrt{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\overline{W^i W^j})^2} \cdot \epsilon\right), \\
&= \text{var}(\overline{W^i W^j}) + \text{var}\left(\sqrt{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\overline{W^i W^j})^2}\right), \\
&\quad + \mathbb{E}\left[\sqrt{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\overline{W^i W^j})^2}\right]^2,
\end{aligned} \tag{C4}$$

where using GMA equations the term $\text{var}\left(\sqrt{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\overline{W^i W^j})^2} \cdot \epsilon\right)$ is obtained by

$$\text{var}\left(\sqrt{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\overline{W^i W^j})^2} \cdot \epsilon\right) = \text{var}\left(\sqrt{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\overline{W^i W^j})^2}\right) + \mathbb{E}\left[\sqrt{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\overline{W^i W^j})^2}\right]^2.$$

In order to simplify the notation in Equation (C4), let us consider $u = \overline{w^i w^j}$ and $t(u) = \sqrt{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\overline{w^i w^j})^2}$, so that using 1st order Taylor series expansion we get

$$\mathbb{E}[t(u)]^2 = t(\mathbb{E}[U])^2 = \mu^{(W^i)^2} \mu^{(W^j)^2} + (\mu^{\overline{w^i w^j}})^2, \tag{C5}$$

$$\begin{aligned}
\text{var}(t(u)) &= (t'(\mathbb{E}[U]))^2 \cdot \text{var}(U), \\
&= \frac{(\mu^{\overline{w^i w^j}})^2}{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\mu^{\overline{w^i w^j}})^2} \cdot (\sigma^{\overline{w^i w^j}})^2.
\end{aligned} \tag{C6}$$

Hence, combining Equations (C4), (C5), and (C6) we get

$$\begin{aligned}
\text{var}(W^i W^j) &= (\sigma^{\overline{w^i w^j}})^2 + \frac{(\mu^{\overline{w^i w^j}})^2}{\mu^{(W^i)^2} \mu^{(W^j)^2} + (\mu^{\overline{w^i w^j}})^2} \cdot (\sigma^{\overline{w^i w^j}})^2 \\
&\quad + \mu^{(W^i)^2} \mu^{(W^j)^2} + (\mu^{\overline{w^i w^j}})^2.
\end{aligned}$$

Using the GMA equations, and Equation (C3), the covariance between the product terms $W^i W^j$ and $W^l W^m$, $i, j, l, m \in \{1, 2, \dots, D\}$, is given by

$$\begin{aligned}
\text{cov}(W^i W^j, W^l W^m) &= \text{cov}(W^i, W^l) \text{cov}(W^j, W^m) + \text{cov}(W^i, W^m) \text{cov}(W^j, W^l), \\
&= \mathbb{E}[W^i W^l] \mathbb{E}[W^j W^m] + \mathbb{E}[W^i W^m] \mathbb{E}[W^j W^l], \\
&= \mu^{\overline{w^i w^l}} \mu^{\overline{w^j w^m}} + \mu^{\overline{w^i w^m}} \mu^{\overline{w^j w^l}}.
\end{aligned}$$

C.2 Proof for Lemma 6

Proof. The covariance matrix Σ^W in Equation (13) can be reformulated in terms of the random variables in \mathbf{W}^P given by

$$\Sigma^W = \begin{bmatrix} \overline{(W^1)^2} & \overline{W^1 W^2} & \dots & \overline{W^1 W^D} \\ \vdots & \overline{(W^2)^2} & \dots & \overline{W^2 W^D} \\ \vdots & \dots & \ddots & \vdots \\ \text{sym.} & \dots & \dots & \overline{(W^D)^2} \end{bmatrix}_{t|t-1}, \tag{C7}$$

where using Equation (C1), $\mathbb{E}[W^i W^j] = \overline{W^i W^j}$, $\forall i, j \in \{1, 2, \dots, D\}$. Let us consider \mathbf{L}^W is an upper triangular random matrix such that

$$\mathbf{L}^W = \begin{bmatrix} L_{11} & L_{12} & \cdots & L_{1D} \\ 0 & L_{22} & \cdots & L_{2D} \\ \vdots & \cdots & \ddots & \vdots \\ 0 & \cdots & 0 & L_{DD} \end{bmatrix}, \quad (\text{C8})$$

where each of the term is assumed to be a Gaussian random variable given by $L_{ij} \sim \mathcal{N}(\mu_{L_{ij}}, \sigma_{L_{ij}}^2)$. The elements of \mathbf{L}^W can be arranged in a random vector,

$$\overrightarrow{\mathbf{L}^W} = [L_{11} L_{22} L_{DD} L_{12} \cdots L_{ij} \cdots L_{D-1D}]^T,$$

such that $\overrightarrow{\mathbf{L}^W}$ is a Gaussian random vector given by

$$\overrightarrow{\mathbf{L}^W} \sim \mathcal{N}\left(\mu^{\overrightarrow{\mathbf{L}^W}}, \Sigma^{\overrightarrow{\mathbf{L}^W}}\right), \quad (\text{C9})$$

where $\mu^{\overrightarrow{\mathbf{L}^W}}$ and $\Sigma^{\overrightarrow{\mathbf{L}^W}}$ are the mean vector and the covariance matrix of $\overrightarrow{\mathbf{L}^W}$. Let us reproduce Σ^W using Equation (C8) such that

$$\Sigma^W = (\mathbf{L}^W)^T \mathbf{L}^W,$$

where each element $\overline{W^i W^j}$ of Σ^W defined in Equation (C7) is obtained using matrix multiplication so that

$$\overline{W^i W^j} = \sum_{k=1}^D L_{jk} L_{ki}, \quad \forall i, j \in \{1, \dots, D\},$$

where using Equation (C9) and the GMA equations we can determine the expected value, the variance, and the covariance terms of any element $\overline{W^i W^j}$ as follows,

$$\mathbb{E}[\overline{W^i W^j}] = \mathbb{E}\left[\sum_{k=1}^D L_{jk} L_{ki}\right], \quad \text{var}(\overline{W^i W^j}) = \text{var}\left(\sum_{k=1}^D L_{jk} L_{ki}\right). \quad (\text{C10})$$

Using Equation (C10), the elements of the prior predictive PDF of \mathbf{W} defined in Proposition 3 can be computed as

$$\mu^{\overline{W^i W^j}} = \mathbb{E}\left[\sum_{k=1}^D L_{jk} L_{ki}\right].$$

Similarly, the covariance between the random matrices, Σ^W and \mathbf{L}^W , is equivalent to finding the covariance between the random vectors $\overrightarrow{\mathbf{L}^W}$ and $\overrightarrow{\mathbf{W}^p}$ given by $\Sigma_{t|t-1}^{\overrightarrow{\mathbf{L}^W \mathbf{W}^p}}$, where any covariance term is obtained by

$$\text{cov}(L_{ij}, \overline{W^i W^j}) = \text{cov}\left(L_{ij}, \sum_{k=1}^D L_{jk} L_{ki}\right). \quad \blacksquare$$

C.3 Proof for Lemma 7

Proof. The prior knowledge of $\overline{\mathbf{W}^p}$ is updated by employing the prior predictive $\mathbf{W}^p_{t|t-1}$ and the posterior PDF $\mathbf{W}^p_{t|t}$ such that

$$f(\overline{\mathbf{w}^p}_t | \mathbf{y}_{1:t}) = \mathcal{N}(\overline{\mathbf{w}^p}_t; \boldsymbol{\mu}_{t|t}^{\overline{\mathbf{W}^p}}, \boldsymbol{\Sigma}_{t|t}^{\overline{\mathbf{W}^p}}),$$

where using Equations (B23)–(B25), the posterior mean, variance and covariance terms of $\overline{\mathbf{W}^p}$ are

$$\begin{aligned} \boldsymbol{\mu}_{t|t}^{\overline{\mathbf{W}^p}} &= \boldsymbol{\mu}_{t|t-1}^{\overline{\mathbf{W}^p}} + \mathbf{K}_t (\boldsymbol{\mu}_{t|t}^{\mathbf{W}^p} - \boldsymbol{\mu}_{t|t-1}^{\mathbf{W}^p}), \\ \boldsymbol{\Sigma}_{t|t}^{\overline{\mathbf{W}^p}} &= \boldsymbol{\Sigma}_{t|t-1}^{\overline{\mathbf{W}^p}} + \mathbf{K}_t (\boldsymbol{\Sigma}_{t|t}^{\mathbf{W}^p} - \boldsymbol{\Sigma}_{t|t-1}^{\mathbf{W}^p}) \mathbf{K}_t^\top, \\ \mathbf{K}_t &= \boldsymbol{\Sigma}_{t|t-1}^{\overline{\mathbf{W}^p}} (\boldsymbol{\Sigma}_{t|t-1}^{\mathbf{W}^p})^{-1}, \\ \boldsymbol{\Sigma}_{t|t-1}^{\overline{\mathbf{W}^p}} &= \boldsymbol{\Sigma}_{t|t-1}^{\overline{\mathbf{W}^p}}. \end{aligned}$$

■

APPENDIX D. ALGORITHMS FOR UNIVARIATE AND MULTIVARIATE PROCESS ERRORS

Algorithm 1. One-time step of the proposed AGVI method for univariate process error

Input: $\boldsymbol{\mu}_{t-1|t-1}$, $\boldsymbol{\Sigma}_{t-1|t-1}$, $\mu_{t-1|t-1}^{\overline{W^2}}$, $(\sigma_{t-1|t-1}^{\overline{W^2}})^2$, y_t , \mathbf{A} , \mathbf{C} , \mathbf{Q} , and σ_V^2

Prior knowledge for the error variance parameter:

1: $\sigma_W^2 = \mu_{t-1|t-1}^{\overline{W^2}}$

Prediction Step:

2: $\boldsymbol{\mu}_{t|t-1}^H = \begin{bmatrix} \mathbf{A} \boldsymbol{\mu}_{t-1|t-1} \\ 0 \end{bmatrix}_{t|t-1}$, $\boldsymbol{\Sigma}_{t|t-1}^H = \begin{bmatrix} \mathbf{A} \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{A}^\top + \mathbf{Q} & \boldsymbol{\Sigma}^{XW} \\ (\boldsymbol{\Sigma}^{XW})^\top & \mu_{t-1|t-1}^{\overline{W^2}} \end{bmatrix}_{t|t-1}$,

$\mu_Y = \mathbf{C} \boldsymbol{\mu}_{t|t-1}$, $\sigma_Y^2 = \mathbf{C} \boldsymbol{\Sigma}_{t|t-1} \mathbf{C}^\top + \sigma_V^2$, $\boldsymbol{\Sigma}_{HY} = \boldsymbol{\Sigma}_{t|t-1}^H \mathbf{F}_t^\top$

1st Update Step:

3: $\boldsymbol{\mu}_{t|t}^H = \boldsymbol{\mu}_{t|t-1}^H + \frac{\boldsymbol{\Sigma}_{HY}}{\sigma_Y^2} (y_t - \mu_Y)$, $\boldsymbol{\Sigma}_{t|t}^H = \boldsymbol{\Sigma}_{t|t-1}^H - \frac{\boldsymbol{\Sigma}_{HY} \boldsymbol{\Sigma}_{HY}^\top}{\sigma_Y^2}$

Posterior Moments for W^2 :

4: $\mu_{t|t}^{W^2} = (\mu_{t|t}^W)^2 + (\sigma_{t|t}^W)^2$,

$(\sigma_{t|t}^{W^2})^2 = 2(\sigma_{t|t}^W)^4 + 4(\sigma_{t|t}^W)^2 (\mu_{t|t}^W)^2$

2nd Update Step:

5: $\mu_{t|t}^{\overline{W^2}} = \mu_{t|t-1}^{\overline{W^2}} + k_t (\mu_{t|t}^{W^2} - \mu_{t|t-1}^{W^2})$, $(\sigma_{t|t}^{\overline{W^2}})^2 = (\sigma_{t|t-1}^{\overline{W^2}})^2 + k_t^2 ((\sigma_{t|t}^{W^2})^2 - (\sigma_{t|t-1}^{W^2})^2)$,

$k_t = \frac{(\sigma_{t-1|t-1}^{\overline{W^2}})^2}{(\sigma_{t|t-1}^{W^2})^2}$

6: **return** $\boldsymbol{\mu}_{t|t}$, $\boldsymbol{\Sigma}_{t|t}$, $\mu_{t|t}^{\overline{W^2}}$, and $(\sigma_{t|t}^{\overline{W^2}})^2$

Algorithm 2. One-time step of the AGVI method for multivariate process errors

Input: $\mu_{t-1|t-1}$, $\Sigma_{t-1|t-1}$, $\mu_{t-1|t-1}^{\overline{L^W}}$, $\Sigma_{t-1|t-1}^{\overline{L^W}}$, \mathbf{y}_t , \mathbf{A} , \mathbf{C} , \mathbf{Q} , and \mathbf{R}

Prior Predictive PDF of $\mathbf{W}_{t|t-1} \sim \mathcal{N}(\mathbf{w}_t; \mathbf{0}_{t|t-1}, \Sigma_{t|t-1}^W)$:

- 1: Any ij th element of $\Sigma_{t|t-1}^W$ is obtained using $\mu^{\overline{W^i W^j}} = \mathbb{E} [\sum_{k=1}^D L_{jk} L_{ki}]$

Prediction Step:

$$2: \mu_{t|t-1}^H = \begin{bmatrix} \mathbf{A} \mu_{t-1|t-1} \\ \mathbf{0} \end{bmatrix}_{t|t-1}, \quad \Sigma_{t|t-1}^H = \begin{bmatrix} \mathbf{A} \Sigma_{t-1|t-1} \mathbf{A}^\top + \mathbf{Q} & \Sigma^{XW} \\ (\Sigma^{XW})^\top & \Sigma^W \end{bmatrix}_{t|t-1},$$

$$\mu_Y = \mathbf{C} \mu_{t|t-1}, \quad \Sigma_Y = \mathbf{C} \Sigma_{t|t-1} \mathbf{C}^\top + \mathbf{R}, \quad \Sigma_{HY} = \Sigma_{t|t-1}^H \mathbf{F}^\top, \text{ where } \mathbf{F} = [\mathbf{C} \ \mathbf{0}]$$

1st Update Step:

$$3: \mu_{t|t}^H = \mu_{t|t-1}^H + \Sigma_{HY} \Sigma_Y^{-1} (\mathbf{y}_t - \mu_Y), \quad \Sigma_{t|t}^H = \Sigma_{t|t-1}^H - \Sigma_{HY} \Sigma_Y^{-1} \Sigma_{HY}^\top$$

- 4: Obtain the posterior PDF, $f(\mathbf{w}^p_t | \mathbf{y}_{1:t}) = \mathcal{N}(\mathbf{w}^p_t; \mu_{t|t}^{W^p}, \Sigma_{t|t}^{W^p})$

2nd Update Step:

$$5: \mu_{t|t}^{\overline{W^p}} = \mu_{t|t-1}^{\overline{W^p}} + \mathbf{K}_t (\mu_{t|t}^{W^p} - \mu_{t|t-1}^{W^p}), \quad \Sigma_{t|t}^{\overline{W^p}} = \Sigma_{t|t-1}^{\overline{W^p}} + \mathbf{K}_t (\Sigma_{t|t}^{W^p} - \Sigma_{t|t-1}^{W^p}) \mathbf{K}_t^\top,$$

$$\mathbf{K}_t = \Sigma_{t|t-1}^{\overline{W^p W^p}} (\Sigma_{t|t-1}^{W^p})^{-1}, \quad \Sigma_{t|t-1}^{\overline{W^p W^p}} = \Sigma_{t|t-1}^{\overline{W^p}}$$

Posterior moments of $\overline{L^W}$:

$$6: \mu_{t|t}^{\overline{L^W}} = \mu_{t|t-1}^{\overline{L^W}} + \mathbf{K}_t^L (\mu_{t|t}^{\overline{W^p}} - \mu_{t|t-1}^{\overline{W^p}}), \quad \Sigma_{t|t}^{\overline{L^W}} = \Sigma_{t|t-1}^{\overline{L^W}} + \mathbf{K}_t^L (\Sigma_{t|t}^{\overline{W^p}} - \Sigma_{t|t-1}^{\overline{W^p}}) (\mathbf{K}_t^L)^\top,$$

$$\mathbf{K}_t^L = \Sigma_{t|t-1}^{\overline{L^W W^p}} (\Sigma_{t|t-1}^{\overline{W^p}})^{-1}$$

- 7: **return** $\mu_{t|t}$, $\Sigma_{t|t}$, $\mu_{t|t}^{\overline{L^W}}$, and $\Sigma_{t|t}^{\overline{L^W}}$

APPENDIX E. ADDITIONAL RESULTS FOR CASE STUDY 2

E.1 Online inference of the variance and covariance terms in the full Q matrix

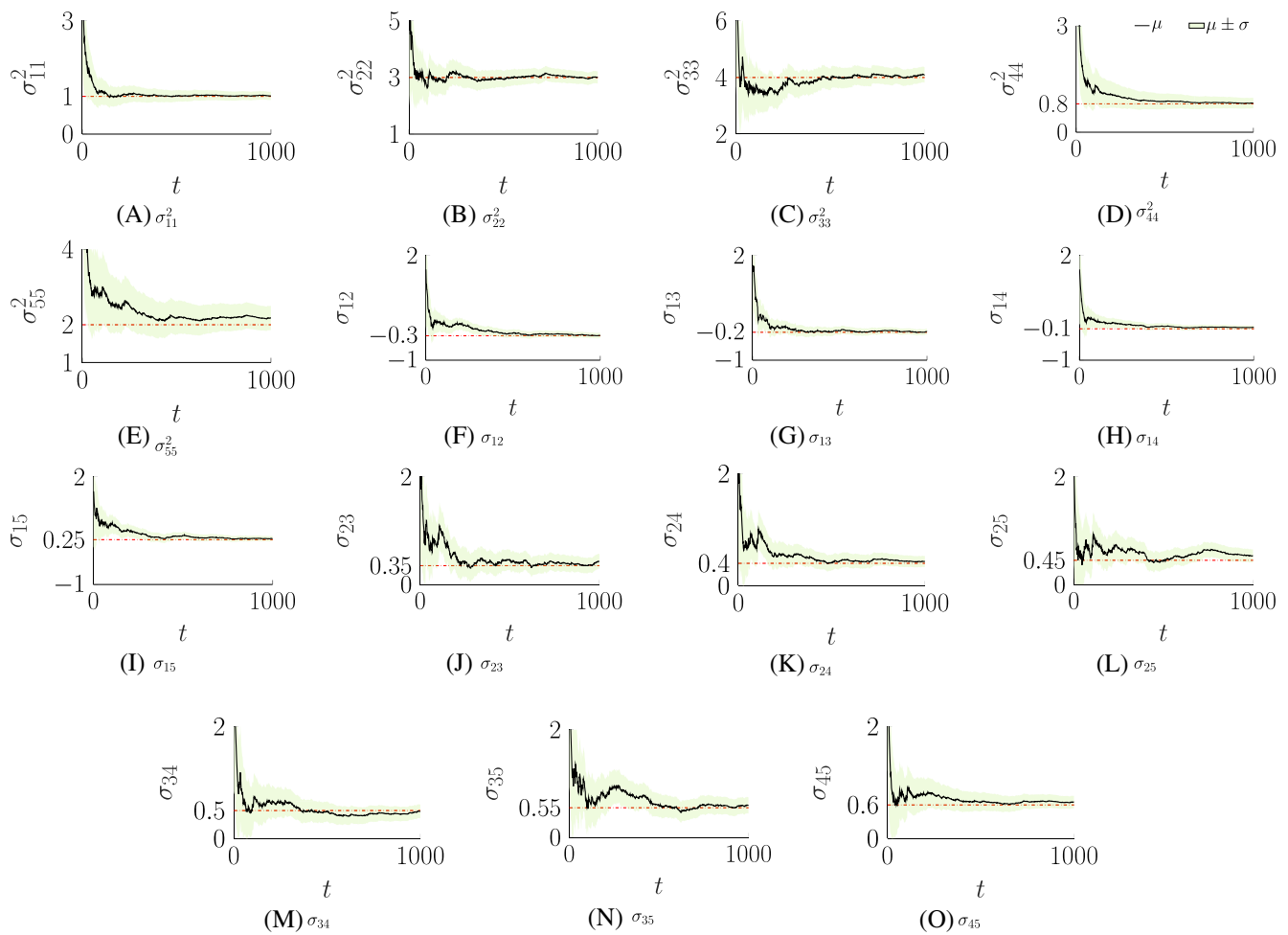


FIGURE E1 Online estimation of the error variance term and the covariance terms from the full Q matrix compared to their true values marked by the dashed red line. The estimated values are shown by the black solid line and their $\pm 1\sigma$ uncertainty bound is shown using the green shaded region.