

TI 2020-052/III  
Tinbergen Institute Discussion Paper

# Bellman filtering and smoothing for state-space models

**Revision: December 2023**

*Rutger-Jan Lange<sup>1</sup>*

Tinbergen Institute is the graduate school and research institute in economics of Erasmus University Rotterdam, the University of Amsterdam and Vrije Universiteit Amsterdam.

Contact: [discussionpapers@tinbergen.nl](mailto:discussionpapers@tinbergen.nl)

More TI discussion papers can be downloaded at <https://www.tinbergen.nl>

Tinbergen Institute has two locations:

Tinbergen Institute Amsterdam  
Gustav Mahlerplein 117  
1082 MS Amsterdam  
The Netherlands  
Tel.: +31(0)20 598 4580

Tinbergen Institute Rotterdam  
Burg. Oudlaan 50  
3062 PA Rotterdam  
The Netherlands  
Tel.: +31(0)10 408 8900

# Bellman filtering and smoothing for state-space models

Rutger-Jan Lange\*

Econometric Institute, Erasmus School of Economics, Rotterdam, The Netherlands

1 November 2023

## Abstract

This paper presents a new filter for state-space models based on Bellman’s dynamic-programming principle, allowing for nonlinearity, non-Gaussianity and degeneracy in the observation and/or state-transition equations. The resulting Bellman filter is a direct generalisation of the (iterated and extended) Kalman filter, enabling scalability to higher dimensions while remaining computationally inexpensive. It can also be extended to enable smoothing. Under suitable conditions, the Bellman-filtered states are stable over time and contractive towards a region around the true state at every time step. Static (hyper)parameters are estimated by maximising a filter-implied pseudo log-likelihood decomposition. In univariate simulation studies, the Bellman filter performs on par with state-of-the-art simulation-based techniques at a fraction of the computational cost. In two empirical applications, involving up to 150 spatial dimensions or highly degenerate/nonlinear state dynamics, the Bellman filter outperforms competing methods in both accuracy and speed.

JEL Classification Codes: C32, C53, C61

Keywords: dynamic programming, posterior mode, Kalman filter, particle filter

## 1 Introduction

### 1.1 State-space models

State-space models allow observations to be affected by an unobserved state that changes stochastically over time. For discrete times  $t = 1, 2, \dots, n$ , the observation  $\mathbf{y}_t \in \mathbb{R}^l$  is drawn from a conditional distribution,  $p(\mathbf{y}_t|\boldsymbol{\alpha}_t)$ , while the latent state  $\boldsymbol{\alpha}_t \in \mathbb{R}^m$  follows a first-order Markov process with a state-transition density,  $p(\boldsymbol{\alpha}_{t+1}|\boldsymbol{\alpha}_t)$ , and some initial condition,  $p(\boldsymbol{\alpha}_1)$ , i.e.

$$\mathbf{y}_t \sim p(\mathbf{y}_t|\boldsymbol{\alpha}_t), \quad \boldsymbol{\alpha}_{t+1} \sim p(\boldsymbol{\alpha}_{t+1}|\boldsymbol{\alpha}_t), \quad \boldsymbol{\alpha}_1 \sim p(\boldsymbol{\alpha}_1). \quad (1)$$

In a slight abuse of notation,  $p(\cdot|\cdot)$  and  $p(\cdot)$  denote *generic* conditional and marginal densities; i.e. any two  $p$ ’s need not denote the same probability density function (e.g. Durbin and Koopman, 2000, p. 6). For a given model, the functional form of all  $p$ ’s is considered known. These densities may further depend on a static (hyper)parameter  $\boldsymbol{\psi}$ , which for notational simplicity is suppressed. They may also depend on lags of  $\mathbf{y}_t$  or, more generally, any  $\mathcal{F}_{t-1}$ -measurable variables, where  $\mathcal{F}_{t-1}$  denotes the information set at time  $t-1$ .

---

\*Email: [lange@ese.eur.nl](mailto:lange@ese.eur.nl). Postal address: P.O. Box 1738, 3000 DR, Rotterdam, the Netherlands.

This potential dependence on  $\mathcal{F}_{t-1}$  is likewise suppressed for the sake of readability. Both the observation and state-transition densities may involve non-Gaussianity, nonlinearity and degeneracy.

Observations  $\mathbf{y}_t$  may take either continuous or discrete values in  $\mathbb{R}^l$ ; in the case of discrete observations,  $p(\mathbf{y}_t|\boldsymbol{\alpha}_t)$  is interpreted as a probability rather than a density. Latent states are assumed to take continuous values in  $\mathbb{R}^m$ ; hence, the state space can be viewed as ‘infinite dimensional’ even as  $m$  remains finite. This is in contrast with Markov-switching models (also known as hidden Markov models; see e.g. Künsch, 2001, p. 109 and Fuh, 2006, p. 2026), in which the state takes a finite number of (discrete) values.

Myriad examples of model (1) can be found in engineering, biology, geological physics, economics and mathematical finance (for a comprehensive overview, see Künsch, 2001, or Doucet et al., 2001). Examples in financial econometrics with continuous state spaces include models for count data (Singh and Roberts, 1992, Frühwirth-Schnatter and Wagner, 2006), intensity (Bauwens and Hautsch, 2006), duration (Bauwens and Veredas, 2004), volatility (Harvey et al., 1994, Ghysels et al., 1996, Jacquier et al., 2002, Taylor, 2008) and dependence structure (Hafner and Manner, 2012).

Model (1) presents researchers and practitioners with three important problems: (a) filtering, (b) smoothing and (c) parameter estimation. The filtering problem concerns the real-time estimation of the current state  $\boldsymbol{\alpha}_t$  conditional on the real-time data  $\mathbf{y}_1, \dots, \mathbf{y}_t$ , where the static parameter  $\boldsymbol{\psi}$  is considered known. The smoothing problem concerns the ex-post estimation of all latent states  $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n$  conditional on the full sample  $\mathbf{y}_1, \dots, \mathbf{y}_n$ , still assuming that  $\boldsymbol{\psi}$  is known. The parameter-estimation problem entails determining the parameter  $\boldsymbol{\psi}$ , where both this parameter and the latent states are assumed to be unknown.

The filtering and smoothing problems can be solved in closed form when model (1) is linear and Gaussian. Kalman’s (1960) filter then computes the real-time expectation of the state (i.e. the mean) and the most likely state (i.e. the mode), which are identical for these models (see Table 1). The Rauch, Tung and Striebel (RTS, 1965) smoother, colloquially known as the ‘Kalman smoother’, computes ex-post state estimates by complementing the (forward) Kalman filter with a subsequent backward recursion. Parameter estimation is typically performed by numerically maximising the log-likelihood function, which is known in closed form through the standard prediction-error decomposition (e.g Harvey, 1990, p. 126).

For the majority of state-space models, however, no exact methods are available for filtering, smoothing or likelihood computation. Here I present an approximate filter and smoother for the general state-space model (1), followed by an approximate parameter-estimation method. This paper thus addresses all three problems mentioned above.

## 1.2 Primary contribution: Filtering and smoothing using Bellman’s equation

This article develops an approximate filter and smoother that are generally applicable and computationally efficient even in higher dimensions. My point of departure is the view that optimisation may be computationally more attractive than integration—especially in higher dimensions. For this reason, I consider a filter and smoother based not on the mean but on the mode, which is also known as the *maximum a posteriori* (MAP) estimate (e.g. Koyama and Paninski, 2010, Liu and Ihler, 2013) or the posterior mode (e.g. Fahrmeir, 1992, Durbin and Koopman, 1997, Jungbacker and Koopman, 2007). In line with the literature, this approach relies on the assumption that the mode exists and is unique. This assumption is not overly restrictive in practice, although it is possible to formulate models for which it does not hold.<sup>1</sup>

<sup>1</sup>E.g. when the observation equation reads  $y_t = \alpha_t^2 + \varepsilon_t$  with  $\varepsilon_t \sim N(0, \sigma_\varepsilon^2)$ .

Table 1: Categorisation of filtering methods.

	<b>Discrete states</b>	<b>Continuously varying states</b>	
		Linear & Gaussian	Nonlinear and/or non-Gaussian
	Exact filters	Exact filters	Approximate filters
<b>Mean</b>	Baum and Petrie (1966) Hamilton (1989)	Kalman (1960)	Iterated extended KF (e.g. Anderson and Moore, 2012) Unscented KF (Julier and Uhlmann, 1997) Masreliez (1975) filter Numerical integration filter (Kitagawa, 1987) Discretisation filter (Farmer, 2021)
<b>Mode</b>	Viterbi (1967)	Kalman (1960)	Bellman filter (BF, this article) Special cases of BF: Fahrmeir’s (1992) mode estimator and Koyama et al.’s (2010) Laplace Gaussian filter

Note: The table should be considered indicative rather than exhaustive, and, for brevity, excludes simulation-based approaches. KF = Kalman filter. BF = Bellman filter.

Computing the mode in real time using plain-vanilla optimisation methods is, however, computationally cumbersome. A naive approach would be to re-estimate, at each time step  $t$ , all previous states of dimension  $m$ , requiring us to continually solve  $m \times t$  dimensional optimisation problems. Computing times per time step then scale as  $O(m^3 t^3)$ , implying a cumulative computing effort, up to time  $t$ , of  $O(m^3 t^4)$ . This escalating complexity over time may explain why the mode estimator has to date received scant attention as a potential filtering method.

My proposed solution to this drawback is to apply Bellman’s (1957) dynamic-programming principle, which yields a forward recursion in function space. The solution to this recursion at any time step is referred to as the *value function*, which maps the state space  $\mathbb{R}^m$  to values in  $\mathbb{R}$  and summarises the researcher’s knowledge of the state at time  $t$ . First, the argmax of the value function represents the most likely state at time  $t$  conditional on  $\mathbf{y}_1, \dots, \mathbf{y}_t$ ; hence, it acts as our filtered state estimate. Second, the negative Hessian matrix evaluated at the peak is indicative of the precision of this state estimate: a ‘sharper’ peak corresponds to a more precise state estimate. Recursively solving Bellman’s equation thus yields a feasible filtering method, producing at each time step both a filtered state and an associated measure of uncertainty.

Importantly for the present purpose, computing the argmax of the value function entails maximisation over a *single* state of dimension  $m$  for each time step. The required computing cost per time step remains constant at  $O(m^3)$ . The resulting cumulative computational complexity over  $t$  time steps then amounts to  $O(m^3 t)$ , which is identical to that of the (information form of the) Kalman filter. On the one hand, the computational complexity of  $O(t)$  means the Bellman filter can be classed as a filter in the strict sense of the term. On the other, the complexity of  $O(m^3)$  offers full scalability to higher dimensional state spaces; e.g. up to 150 dimensions in the application in section 9.

The price we pay for this reduced computational complexity is that Bellman’s recursion generally lacks an analytic solution; hence, we must resort to approximation, which can be viewed as a form of approximate dynamic programming (e.g. Bertsekas, 2012). One possibility is to discretise the (continuous) state space  $\mathbb{R}^m$ , forcing the state to take a finite number of (discrete) values. Bellman’s equation can then be solved exactly, yielding Viterbi’s (1967) algorithm (see Table 1), which has proven highly successful in engineering. However, this approach quickly becomes infeasible due to the curse of dimensionality (Künsch, 2001, p. 125, Liu, 2008, p. 29), as it requires the computation and storage of  $N^m$  values for each time step, where  $N$  is the number of gridpoints in each of  $m$  spatial directions (e.g.  $N = 100$  and  $m = 5$  is infeasible).

Instead, I take inspiration from another exact solution to Bellman’s forward recursion. As it turns out, Bellman’s recursion allows an exact solution if the entire model (1) is linear and Gaussian, yielding Kalman’s (1960) filter. The solution to Bellman’s equation is then a *function*, rather than a finite-dimensional object as in Viterbi’s case. This value function has a particularly simple form: it is multivariate quadratic at every time step, with a unique argmax that corresponds to Kalman’s filtered state. Moreover, its negative Hessian matrix equals the inverse of the usual Kalman-filtered covariance matrix. Hence, the Kalman filter represents an *exact function-space solution* to Bellman’s equation. This was long recognised in the engineering literature (e.g. Whittle, 1996, ch. 12; Whittle, 2004) before finding its way into the econometrics literature (Hansen and Sargent, 2013, ch. 8). Perhaps less widely known is the fact that the RTS (1965) smoother similarly corresponds to an exact—also multivariate quadratic—solution to a combination of Bellman’s forward and backward recursions (see section 6).

The basic premise of this article is that Bellman’s forward and backward recursions remain valid in the context of the general state-space model (1). Motivated by the exact solutions leading to the Kalman filter and RTS smoother, I deviate from the literature in exploring *function-space approximations* of value functions rather than discretising. Computing at every time step some parametric approximation of the value function yields a new class of (Bellman) filters and smoothers. Within the class of function-space approximations, I employ arguably the simplest non-trivial option: a multivariate quadratic function. This quadratic approximation is exact for linear Gaussian models and—given that value functions in filtering applications are typically smooth and possess global maxima—broadly applicable. The approximation can also be viewed as a second-order Taylor expansion of a generic smooth value function. This simple approximation approach yields immediate and novel extensions of the Kalman filter and smoother. The main contribution of this article is the insight that using function-space rather than discrete approximations allows us to avoid the curse of dimensionality, leading to a new class of filters and smoothers that are computationally frugal and turn out to be remarkably accurate.

### 1.3 Secondary contribution: Parameter estimation using likelihood approximation

To address the parameter-estimation problem, I deviate from the literature that relies on simulation-based approaches (e.g. Malik and Pitt, 2011, Koopman et al., 2015, Koopman et al., 2016) by presenting a deterministic and computationally efficient—albeit approximate—method based on the output of the Bellman filter. While no formal guarantees are offered, an extensive simulation study (section 8) demonstrates that the proposed estimator is no less accurate or efficient than (asymptotically exact) simulation-based methods, while requiring a fraction of the computational cost. Establishing the asymptotic properties of the estimator remains an open question.

Specifically, I propose to maximise an approximate version of the log-likelihood function that is immediately computable from the output of the Bellman filter. First, the (exact) log-likelihood function is decomposed into (a) the ‘fit’ of the Bellman-filtered states in view of the data, minus (b) the realised Kullback-Leibler (KL, see Kullback and Leibler, 1951) divergence between filtered and predicted state densities. While the former is known in closed form, the latter typically is not—except in the case of linear Gaussian state-space models, in which case it is multivariate quadratic. Second, I approximate this KL divergence term using a multivariate quadratic term computed from the output of the Bellman filter. The resulting pseudo log-likelihood function remains exact in the case of linear Gaussian models; more generally, it can be viewed as a second-order approximation of the log-likelihood function. It can be optimised using

standard gradient-based numerical optimisers, making approximate parameter estimation for the general state-space model (1) as simple and fast as maximum-likelihood estimation of the Kalman filter.

## 1.4 Limitations of existing methods

Existing approaches to filtering, smoothing and parameter estimation can be classified as either approximation- or simulation-based, each with their own disadvantages. First, approximate filtering methods (see Table 1) tend to be specialised in their applications. The extended and unscented Kalman filters account for nonlinearity, but assume additive noise and maintain the normality assumption. Conversely, West (1981) relaxes the normality assumption, while maintaining the linearity assumption. Masreliez’s (1975) filter is robust in the case of heavy-tailed observation noise but, due to the need to approximate integrals, computationally inefficient in higher dimensions. Similarly, numerical integration (Kitagawa, 1987) and other discretisation methods (Farmer, 2021) are flexible in theory, but restricted in practice by the curse of dimensionality. Fahrmeir’s (1992) method applies to observations drawn from an exponential distribution. Durbin and Koopman (2000) and Koyama et al. (2010) mostly rely on a linear Gaussian state equation. Müller and Petalas (2010) assume that deviations of the latent state from its equilibrium value are small. In the literature, no approximate filters seem to be available at the level of generality of model (1). Moreover, the aforementioned approaches tend to neglect the smoothing and parameter-estimation problems.

Second, simulation-based methods such as particle filters are widely applicable and easy to implement (for a textbook treatment, see e.g. Chopin and Papaspiliopoulos, 2020). However, the curse of dimensionality means that particle filters may struggle with high-dimensional state spaces (Surace et al., 2019). For the same reason, the importance-sampling method by Koopman et al. (2015, 2016, 2017) has not been applied in situations in which the state-space dimension exceeds two. Particle smoothing (as opposed to filtering) tends to be even more computationally expensive, as the computational cost scales with the number of particles squared (Kantas et al., 2015). Particle filters have also been applied to the parameter-estimation problem, but this remains challenging (Liu and West, 2001, Kantas et al., 2015); e.g. Malik and Pitt’s (2011) method applies only when the state space is one dimensional.

## 2 Main idea: Filtering using Bellman’s principle

The state-space model under consideration is given in equation (1). A realised path is denoted by  $(\mathbf{y}_1, \dots, \mathbf{y}_t)(\omega)$  for every event  $\omega \in \Omega$ , where  $\Omega$  denotes the event space of the underlying complete probability space of interest, denoted  $(\Omega, \mathcal{F}, \mathbb{P})$ . The logarithm of joint and conditional densities are written using generic notation as  $\ell(\cdot, \cdot) := \log p(\cdot, \cdot)$  and  $\ell(\cdot|\cdot) := \log p(\cdot|\cdot)$ , respectively, for potentially different  $p$ ’s. This section considers the filtering problem; any dependence on  $\boldsymbol{\psi}$  is suppressed.

The joint log-likelihood function of the states and the data is written as  $L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t) : \Omega \times \mathbb{R}^m \times \dots \times \mathbb{R}^m \rightarrow \mathbb{R}$ . Here, the data  $\mathbf{y}_1, \dots, \mathbf{y}_t$  are considered fixed and known, as indicated by the subscript, while the states  $\mathbf{a}_1, \dots, \mathbf{a}_t$  in Roman font are considered variables to be evaluated along any path. The true states  $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_t$  in Greek font remain unknown. For the state-space model (1), the joint log likelihood of the data and the states follows from the ‘probability chain rule’ (Godsill et al., 2004, p. 156):

$$L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t) = \sum_{i=1}^t \ell(\mathbf{y}_i|\mathbf{a}_i) + \sum_{i=2}^t \ell(\mathbf{a}_i|\mathbf{a}_{i-1}) + \ell(\mathbf{a}_1), \quad t \leq n. \quad (2)$$

This joint log likelihood is, *a priori*, a random function of the observations  $\mathbf{y}_1, \dots, \mathbf{y}_t$ , even though the data are considered known and fixed *ex post*. For clarity, I formalise the assumption that for some sufficiently large  $t$ , there exists a unique sequence of states, denoted  $\mathbf{a}_{1|t}, \dots, \mathbf{a}_{t|t}$ , that maximise equation (2).

**Assumption E (Existence of the mode)** *There exists some  $t_0 \geq 1$ , such that for all  $t \geq t_0$ , the mode  $(\mathbf{a}_{1|t}, \mathbf{a}_{2|t}, \dots, \mathbf{a}_{t|t})$  exists and is unique, where*

$$(\mathbf{a}_{1|t}, \mathbf{a}_{2|t}, \dots, \mathbf{a}_{t|t}) := \arg \max_{(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_t) \in \mathbb{R}^{m \times t}} L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t). \quad (3)$$

This assumption is labelled ‘E’ for existence, because it is required to underpin the main idea; later, Assumption 1-3 (in section 5) are used to derive the theoretical properties of the filter.

As equation (3) illustrates, elements of the mode at time  $t$  are denoted by  $\mathbf{a}_{i|t}$  for  $i \leq t$ , where  $i$  denotes the state that is estimated,  $t$  the information set used. The entire solution is a collection of  $t$  vectors, each of length  $m$ . Iterative solution methods for solving (3) were proposed in Durbin and Koopman (2000) and So (2003). When the mode (3) is computed for *each* time step  $t \geq t_0$ , we can extract a sequence of real-time state estimates  $\{\mathbf{a}_{t|t}\}_{t \geq t_0}$ , where each estimate  $\mathbf{a}_{t|t}$  is extracted from a *different* mode (3).

As time progresses, however, the computation of filtered states  $\{\mathbf{a}_{t|t}\}_t$  becomes ever more complicated—note that optimisation problem (3) involves  $m \times t$  optimisation variables at each time  $t$ . Indeed, solving problem (3) may become practically infeasible for large  $t$ . This raises the question whether it is possible to proceed in real time without solving an optimisation problem of ever-increasing complexity. As shown next, this can be achieved using Bellman’s dynamic-programming principle. To this end, I define the *value function* by maximising the joint log-likelihood function (2) with respect to all states apart from the most recent state  $\mathbf{a}_t \in \mathbb{R}^m$ ; such functions are also known as ‘profile’ log-likelihood functions (Murphy and Van der Vaart, 2000) in statistics and ‘stress’ functions in engineering (Whittle, 1981, p. 769).

**Definition 1 (Value function)** *Let Assumption E hold. For  $t \geq t_0$ , the value function  $V_t : \Omega \times \mathbb{R}^m \rightarrow \mathbb{R}$  is*

$$V_t(\mathbf{a}_t) := \max_{(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{t-1}) \in \mathbb{R}^{m \times (t-1)}} L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t), \quad \mathbf{a}_t \in \mathbb{R}^m. \quad (4)$$

The value function  $V_t(\cdot)$  encodes our knowledge of the state at time  $t$ , as indicated by the subscript, and depends on past and current data  $\mathbf{y}_1, \dots, \mathbf{y}_t$ , which are considered fixed, as well as on its argument  $\mathbf{a}_t$ , which is a continuous variable in  $\mathbb{R}^m$ . Naturally,  $\mathbf{a}_{t|t} = \arg \max_{\mathbf{a}_t} V_t(\mathbf{a}_t)$ , such that the last element of the mode (3) can be recovered from the value function. Usefully, the value function (4) satisfies a forward recursive equation, known as Bellman’s equation, which can be used for the purpose of filtering.

**Proposition 1 (Filtering using Bellman’s equation)** *Let Assumption E hold. The value function (4) satisfies Bellman’s forward recursion:*

$$V_t(\mathbf{a}_t) = \ell(\mathbf{y}_t | \mathbf{a}_t) + \max_{\mathbf{a}_{t-1} \in \mathbb{R}^m} \left\{ \ell(\mathbf{a}_t | \mathbf{a}_{t-1}) + V_{t-1}(\mathbf{a}_{t-1}) \right\}, \quad \mathbf{a}_t \in \mathbb{R}^m, \quad (5)$$

for all  $t_0 < t \leq n$ . Further,

$$\mathbf{a}_{t|t} := \arg \max_{\mathbf{a}_t \in \mathbb{R}^m} V_t(\mathbf{a}_t), \quad t_0 \leq t \leq n. \quad (6)$$

Bellman’s equation (5) is a forward recursion that relates the value function  $V_t(\mathbf{a}_t)$  to the (previous) value function  $V_{t-1}(\mathbf{a}_{t-1})$  by adding one term reflecting the state transition,  $\ell(\mathbf{a}_t|\mathbf{a}_{t-1})$ ; one term reflecting the observation density,  $\ell(\mathbf{y}_t|\mathbf{a}_t)$ ; and a subsequent maximisation over a single state variable,  $\mathbf{a}_{t-1} \in \mathbb{R}^m$ . The value function  $V_t(\mathbf{a}_t)$  still depends on the data  $\mathbf{y}_1, \dots, \mathbf{y}_{t-1}$ , but only indirectly, i.e. through the previous value function  $V_{t-1}(\mathbf{a}_{t-1})$ . Apart from assuming the existence of the mode, no (additional) assumptions are imposed on the log densities  $\ell(\mathbf{y}_t|\mathbf{a}_t)$  and  $\ell(\mathbf{a}_t|\mathbf{a}_{t-1})$ ; the proof in Supplement A uses only standard dynamic-programming arguments. As such, Bellman’s equation (5) is of quite general applicability. As the researcher receives the data  $\mathbf{y}_1$  through  $\mathbf{y}_t$ , she can iteratively compute a sequence of value functions (5), which imply a sequence of filtered state estimates via the respective maximisers (6).

**Remark 1** *For Markov-switching models, in which the latent state takes a finite number of (discrete) values, Bellman’s equation (5) can be solved exactly for all time steps, yielding Viterbi’s (1967) algorithm. Exact solubility of (5) tends to be lost when the states take continuous values.*

When latent states take values in a continuum, as in the present article, the solution to Bellman’s equation (5) is a *function* rather than a (finite-dimensional) vector as in Viterbi’s algorithm. While the value function cannot generally be found exactly, there is an exception to this rule, as highlighted next.

**Corollary 1 (Kalman filter as a special case)** *Take a linear Gaussian state-space model with observation equation  $\mathbf{y}_t = \mathbf{d} + \mathbf{Z}\boldsymbol{\alpha}_t + \boldsymbol{\varepsilon}_t$ , where  $\boldsymbol{\varepsilon}_t \sim i.i.d. \mathbf{N}(\mathbf{0}, \mathbf{H})$ , and state-transition equation  $\boldsymbol{\alpha}_t = \mathbf{c} + \mathbf{T}\boldsymbol{\alpha}_{t-1} + \boldsymbol{\eta}_t$ , where  $\boldsymbol{\eta}_t \sim i.i.d. \mathbf{N}(\mathbf{0}, \mathbf{Q})$  with a positive semidefinite covariance matrix  $\mathbf{Q}$ , such that Kalman’s (1960) filter applies. Assume the Kalman-filtered covariance matrices, denoted  $\{\mathbf{P}_{t|t}\}$ , are positive definite. Then (a) the value function is exactly multivariate quadratic at every time step, (b) the Bellman-filtered states are identical to the Kalman-filtered states, and (c) the negative Hessian matrix of the value function equals  $\mathbf{P}_{t|t}^{-1}$  at every time step.*

The proof of Corollary 1 is contained in section 4, where I treat the case of a linear Gaussian state equation but a general observation density. As is well known in engineering (e.g. Whittle, 1996, ch. 12), the exact solubility of Bellman’s equation in the case of linear Gaussian models is attributable to the quadratic nature of all terms appearing on its right-hand side. The left-hand side turns out to be quadratic as well, preserving exact solubility over time.

A key contribution of this article is the insight that Bellman’s equation continues to hold for state-space models that are not necessarily linear and Gaussian, even if analytic solubility is lost. In this case, I deviate from the literature in considering function-space approximations in solving Bellman’s recursion (5). I consider a particularly simple approximation—the multivariate quadratic function—which happens to be exact for linear Gaussian state-space models. A different class of Bellman filters, not explored here, would be obtained by using non-parametric approximations.

### 3 Bellman filter for general state-space models

#### 3.1 Non-degenerate case

This section develops the Bellman filter for the general state-space model (1) by approximating the value function, at every time step, by a multivariate quadratic function. I assume here that the observation and state-transition densities are non-degenerate; an extension to the degenerate case is set out below.

The Bellman-filtered state (6) requires a maximisation with respect to the current state,  $\mathbf{a}_t$ , while Bellman's equation (5) additionally contains a maximisation with respect to the lagged state,  $\mathbf{a}_{t-1}$ . Merging both steps generates a joint optimisation problem in both state variables:

$$\begin{bmatrix} \mathbf{a}_{t|t} \\ \mathbf{a}_{t-1|t} \end{bmatrix} := \underset{\begin{bmatrix} \mathbf{a}_t \\ \mathbf{a}_{t-1} \end{bmatrix} \in \mathbb{R}^{2m}}{\arg \max} \left\{ \ell(\mathbf{y}_t|\mathbf{a}_t) + \ell(\mathbf{a}_t|\mathbf{a}_{t-1}) + V_{t-1}(\mathbf{a}_{t-1}) \right\}. \quad (7)$$

The left-hand side features the filtered state,  $\mathbf{a}_{t|t}$ , as well as the revised estimate of the previous state, denoted  $\mathbf{a}_{t-1|t}$ . The computation of the latter, while not our main focus, is inherent to Bellman's equation and cannot be avoided. The right-hand side features two log densities denoted  $\ell(\cdot|\cdot) := \log p(\cdot|\cdot)$ , which are given in closed form by the state-space model (1).

While the lagged value function  $V_{t-1}(\cdot)$  on the right-hand side of optimisation (7) is typically unavailable in closed form, the shape around its peak turns out to be most relevant in the determination of the filtered state  $\mathbf{a}_{t|t}$ . I thus propose to approximate  $V_{t-1}(\mathbf{a}_{t-1})$  by a multivariate quadratic function that is parametrised by its argmax, denoted  $\mathbf{a}_{t-1|t-1} \in \mathbb{R}^m$ , and the negative Hessian matrix, denoted  $\mathbf{I}_{t-1|t-1} \in \mathbb{R}^{m \times m}$ , which is assumed positive definite and can be interpreted as an information (or 'precision') matrix. The approximation thus reads

$$V_{t-1}(\mathbf{a}_{t-1}) = -\frac{1}{2}(\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1})' \mathbf{I}_{t-1|t-1} (\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1}) + \text{constants}, \quad \mathbf{a}_{t-1} \in \mathbb{R}^m, \quad (8)$$

which for simplicity is written with equality. Constants can be ignored in the context of optimisation (7). Substituting the quadratic approximation (8) into maximisation (7) yields a viable function-space algorithm. For linear Gaussian state-space models, approximation (8) is exact and the bivariate optimisation (7) can be performed analytically, leading to (the information form of) the Kalman filter.

While optimisation (7) does not generally allow closed-form solutions, it is typically straightforward to write out analytically the steps of e.g. Newton's method (Nocedal and Wright, 2006):

$$\begin{bmatrix} \mathbf{a}_t \\ \mathbf{a}_{t-1} \end{bmatrix} \leftarrow \begin{bmatrix} \mathbf{a}_t \\ \mathbf{a}_{t-1} \end{bmatrix} + \begin{bmatrix} \mathbf{J}_t^{11} - \frac{d^2\ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} & \mathbf{J}_t^{12} \\ \mathbf{J}_t^{21} & \mathbf{I}_{t-1|t-1} + \mathbf{J}_t^{22} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{J}_t^1 + \frac{d\ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t} \\ \mathbf{J}_t^2 - \mathbf{I}_{t-1|t-1}(\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1}) \end{bmatrix}, \quad (9)$$

where, for notational simplicity, I use the assignment symbol; this allows the iterates (which appear on both the left- and right-hand sides) to be denoted by  $\mathbf{a}_t$  and  $\mathbf{a}_{t-1}$ . In Newton's step (9), derivatives related to the state-transition density are

$$\begin{bmatrix} \mathbf{J}_t^1 \\ \mathbf{J}_t^2 \end{bmatrix} := \begin{bmatrix} \frac{d\ell(\mathbf{a}_t|\mathbf{a}_{t-1})}{d\mathbf{a}_t} \\ \frac{d\ell(\mathbf{a}_t|\mathbf{a}_{t-1})}{d\mathbf{a}_{t-1}} \end{bmatrix}, \quad \begin{bmatrix} \mathbf{J}_t^{11} & \mathbf{J}_t^{12} \\ \mathbf{J}_t^{21} & \mathbf{J}_t^{22} \end{bmatrix} := - \begin{bmatrix} \frac{d^2\ell(\mathbf{a}_t|\mathbf{a}_{t-1})}{d\mathbf{a}_t d\mathbf{a}_t'} & \frac{d^2\ell(\mathbf{a}_t|\mathbf{a}_{t-1})}{d\mathbf{a}_t d\mathbf{a}_{t-1}'} \\ \frac{d^2\ell(\mathbf{a}_t|\mathbf{a}_{t-1})}{d\mathbf{a}_{t-1} d\mathbf{a}_t'} & \frac{d^2\ell(\mathbf{a}_t|\mathbf{a}_{t-1})}{d\mathbf{a}_{t-1} d\mathbf{a}_{t-1}'} \end{bmatrix}. \quad (10)$$

Fisher's optimisation method is obtained by replacing  $d^2\ell(\mathbf{y}_t|\mathbf{a}_t)/(d\mathbf{a}_t d\mathbf{a}_t')$  in equation (9) with its expectation conditional on  $\mathbf{a}_t$ . When the observation and state-transition densities in model (1) are given, it is straightforward (if tedious) to compute all required derivatives. As  $\mathbf{I}_{t-1|t-1}$  is assumed to be invertible, analytic block-matrix inversion can be used for each Newton step (9), reducing the size of matrices to be numerically inverted from  $2m \times 2m$  to  $m \times m$  (see Supplement B for details). The resulting algorithm is shown under step 4 in Table 2. Alternatively, black-box numerical optimisers may be used to solve (7),

Table 2: Bellman filter for model (1).

Step	Method	Computation
1. Initialise		Set $\mathbf{a}_{0 0}$ equal to the unconditional mean of the latent state (or treat it as a static parameter to be estimated) and set $\mathbf{I}_{0 0}$ equal to some sufficiently large multiple of the identity matrix. Set $t = 1$ .
2. Predict		$\mathbf{a}_{t t-1} = \arg \max_{\mathbf{a}_t \in \mathbb{R}^m} \ell(\mathbf{a}_t   \mathbf{a}_{t-1 t-1})$ $\mathbf{I}_{t t-1} = \mathbf{J}_t^{11} - \mathbf{J}_t^{12}(\mathbf{I}_{t-1 t-1} + \mathbf{J}_t^{22})^{-1} \mathbf{J}_t^{21} \Big _{\mathbf{a}_t = \mathbf{a}_{t t-1}, \mathbf{a}_{t-1} = \mathbf{a}_{t-1 t-1}}$
3. Start		Set $\mathbf{a}_t \leftarrow \mathbf{a}_{t t-1}$ and $\mathbf{a}_{t-1} \leftarrow \mathbf{a}_{t-1 t-1}$ .
4. Optimise	Newton	$\mathbf{S}_t \leftarrow \mathbf{J}_t^{11} - \mathbf{J}_t^{12}(\mathbf{I}_{t-1 t-1} + \mathbf{J}_t^{22})^{-1} \mathbf{J}_t^{21} - \frac{d^2 \ell(\mathbf{y}_t   \mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}'_t}$ , $\mathbf{D}_t \leftarrow \mathbf{I}_{t-1 t-1} + \mathbf{J}_t^{22}$ , $\mathbf{G}_t^1 \leftarrow \mathbf{J}_t^1 + \frac{d\ell(\mathbf{y}_t   \mathbf{a}_t)}{d\mathbf{a}_t}$ , $\mathbf{G}_t^2 \leftarrow \mathbf{J}_t^2 - \mathbf{I}_{t-1 t-1}(\mathbf{a}_{t-1} - \mathbf{a}_{t-1 t-1})$ , $\mathbf{a}_t \leftarrow \mathbf{a}_t + \mathbf{S}_t^{-1} \mathbf{G}_t^1 - \mathbf{S}_t^{-1} \mathbf{J}_t^{12} \mathbf{D}_t^{-1} \mathbf{G}_t^2$ , $\mathbf{a}_{t-1} \leftarrow \mathbf{a}_{t-1} - \mathbf{D}_t^{-1} \mathbf{J}_t^{21} \mathbf{S}_t^{-1} \mathbf{G}_t^1 + (\mathbf{D}_t^{-1} + \mathbf{D}_t^{-1} \mathbf{J}_t^{21} \mathbf{S}_t^{-1} \mathbf{J}_t^{12} \mathbf{D}_t^{-1}) \mathbf{G}_t^2$ .
	Fisher	Like Newton's method, but with $\mathbf{S}_t$ adjusted to include $\mathbb{E}[d^2 \ell(\mathbf{y}_t   \mathbf{a}_t) / (d\mathbf{a}_t d\mathbf{a}'_t)   \mathbf{a}_t]$ .
5. Stop		Stop if some convergence criterion is satisfied or after a predetermined number of iterations.
6. Update		$\mathbf{a}_{t t} = \mathbf{a}_t$ and $\mathbf{a}_{t-1 t} = \mathbf{a}_{t-1}$ .
	Newton	$\mathbf{I}_{t t} = \mathbf{J}_t^{11} - \mathbf{J}_t^{12}(\mathbf{I}_{t-1 t-1} + \mathbf{J}_t^{22})^{-1} \mathbf{J}_t^{21} - \frac{d^2 \ell(\mathbf{y}_t   \mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}'_t} \Big _{\mathbf{a}_t = \mathbf{a}_{t t}, \mathbf{a}_{t-1} = \mathbf{a}_{t-1 t}}$
	Fisher	$\mathbf{I}_{t t} = \mathbf{J}_t^{11} - \mathbf{J}_t^{12}(\mathbf{I}_{t-1 t-1} + \mathbf{J}_t^{22})^{-1} \mathbf{J}_t^{21} - \mathbb{E} \left[ \frac{d^2 \ell(\mathbf{y}_t   \mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}'_t} \Big  \mathbf{a}_t \right] \Big _{\mathbf{a}_t = \mathbf{a}_{t t}, \mathbf{a}_{t-1} = \mathbf{a}_{t-1 t}}$
6. Proceed		Set $t = t + 1$ and return to step 2.

*Note:* The log-likelihood functions  $\ell(\mathbf{y}_t | \mathbf{a}_t)$  and  $\ell(\mathbf{a}_t | \mathbf{a}_{t-1})$  are known in closed form and can be read off from the data-generating process (1). Various derivatives of  $\ell(\mathbf{a}_t | \mathbf{a}_{t-1})$  are defined in equation (10). Two (intentionally vanilla) optimisation methods are listed under steps 4 and 6. Users may also implement more sophisticated and/or black-box optimisation methods based on maximisation (7).

obviating the need for manual computations; this will save researcher time but potentially increase the required computer time. The optimisation can be started using  $(\mathbf{a}_t, \mathbf{a}_{t-1}) \leftarrow (\mathbf{a}_{t|t-1}, \mathbf{a}_{t-1|t-1})$ , where  $\mathbf{a}_{t|t-1} := \arg \max_{\mathbf{a}} \ell(\mathbf{a} | \mathbf{a}_{t-1|t-1})$ , as indicated under steps 2 and 3 in Table 2. This prediction  $\mathbf{a}_{t|t-1}$  can often be computed in closed form.

To facilitate the proposed recursive method, the left-hand side of Bellman's equation (5) must also be approximated by a multivariate quadratic function. To this end, I compute the negative Hessian matrix (with respect to  $\mathbf{a}_t$ ) of the value function, i.e.  $V_t(\mathbf{a}_t) = \ell(\mathbf{y}_t | \mathbf{a}_t) + \max_{\mathbf{a}_{t-1}} \{\ell(\mathbf{a}_t | \mathbf{a}_{t-1}) + V_{t-1}(\mathbf{a}_{t-1})\}$ . The negative Hessian may be then be evaluated at the peak. Employing the second-order envelope theorem (Supplement C) yields

$$\mathbf{I}_{t|t} := \mathbf{J}_t^{11} - \mathbf{J}_t^{12}(\mathbf{I}_{t-1|t-1} + \mathbf{J}_t^{22})^{-1} \mathbf{J}_t^{21} - \frac{d^2 \ell(\mathbf{y}_t | \mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}'_t} \Big|_{\mathbf{a}_t = \mathbf{a}_{t|t}, \mathbf{a}_{t-1} = \mathbf{a}_{t-1|t}} \quad (11)$$

as shown in Table 2 under step 6. Fisher's version is obtained by taking a conditional expectation of the last term. For linear Gaussian state-space models, Newton and Fisher versions of update (11) are identical and equal to the information update of the Kalman filter (Supplement D). Update (11) can also be viewed as a 'realised' version of the recursion for the inverse of Cramér-Rao lower bounds (Tichavsky et al., 1998, eq. 21)—the difference being that equation (11) has no expectations. The predicted information  $\mathbf{I}_{t|t-1}$ , given in step 2 of Table 2, is similar in form and used for static-parameter estimation purposes in section 7.

The resulting Bellman filter in Table 2 has a computational complexity of  $O(m^3 t)$ , which is attributable to the need to invert  $m \times m$  matrices at every time step. This complexity matches that of (the information form of) the Kalman filter, thus offering scalability to at least moderately high dimensions  $m$ . I am unaware

of other approximate filters offering the same breadth of applicability and computational efficiency.<sup>2</sup>

### 3.2 Extension to the degenerate case

When some elements of  $\mathbf{a}_{t-1|t-1}$  are known to be pinpoint accurate, the corresponding diagonal values of the precision matrix  $\mathbf{I}_{t-1|t-1}$  in equation (8) are unbounded. Such infinite diagonal values make optimisation (7) easier rather than harder, as some elements of  $\mathbf{a}_{t-1}$  are constrained and need not be numerically optimised; rather, they can be fixed by hand. When the relevant restriction is implemented, the unbounded contributions in the quadratic term (8) can be dropped. Similarly, when the state-transition density  $\ell(\mathbf{a}_t|\mathbf{a}_{t-1})$  is degenerate, some elements of the current state are deterministic functions of the previous state. When these restrictions are implemented, the degenerate part of the transition density can be dropped. Indeed, this procedure will be used for the model in section 10, which involves degenerate state dynamics. Finally, when the observation density  $\ell(\mathbf{y}_t|\mathbf{a}_t)$  is degenerate, as when some elements of  $\mathbf{a}_t$  are fully revealed by the observation  $\mathbf{y}_t$ , optimisation (7) requires that some elements of  $\mathbf{a}_t$  take a specific functional form of  $\mathbf{y}_t$ . From an optimisation perspective, therefore, degeneracies correspond to equality constraints that can typically be implemented by hand, reducing the dimension of the numerical optimisation problem to be solved. This capacity to deal with (partially) deterministic state dynamics forms an advantage over e.g. particle-filtering methods, which may struggle in such situations.

## 4 Bellman filter for models with linear Gaussian state dynamics

This section applies the general idea developed in the previous section to models in which the state-transition equation remains linear and Gaussian. The advantage of this special case is that the ‘inner’ optimisation in Bellman’s equation (5), i.e. with respect to the lagged state  $\mathbf{a}_{t-1}$ , can now be performed in closed form. The ‘outer’ optimisation with respect to the current state  $\mathbf{a}_t$  remains numerical. Models with linear Gaussian state equations are written as in Koopman et al. (2015, 2016):

$$\mathbf{y}_t \sim p(\mathbf{y}_t | \boldsymbol{\alpha}_t), \quad \boldsymbol{\alpha}_{t+1} = \mathbf{c} + \mathbf{T} \boldsymbol{\alpha}_t + \boldsymbol{\eta}_{t+1}, \quad \boldsymbol{\eta}_t \sim \text{i.i.d. N}(\mathbf{0}, \mathbf{Q}), \quad \boldsymbol{\alpha}_1 \sim p(\boldsymbol{\alpha}_1), \quad (12)$$

where  $t = 1, \dots, n$ , and the state-transition equation contains the system vector  $\mathbf{c} \in \mathbb{R}^m$  and system matrix  $\mathbf{T} \in \mathbb{R}^{m \times m}$ . The state innovation  $\boldsymbol{\eta}_t$  is controlled by a positive semidefinite covariance matrix  $\mathbf{Q} \in \mathbb{R}^{m \times m}$ , which presents no loss of generality compared to authors who write the innovation as  $\mathbf{R}\boldsymbol{\eta}_t$  for some matrix  $\mathbf{R}$ .<sup>3</sup> The observation density  $p(\mathbf{y}_t | \boldsymbol{\alpha}_t)$  may still be non-Gaussian and involve nonlinearity.

<sup>2</sup>In related work, Koyama et al. (2010, p. 173) report a computational complexity of  $O(m^2t)$ , purportedly as  $O(m^2)$  is the ‘complexity of matrix manipulations’. This result comes with two important caveats. First, it relies on having a linear and Gaussian state equation; otherwise, their prediction step requires the (numerical) evaluation of an integral in  $m$  dimensions. Second, it overlooks the fact that the (dense) matrix inversion required by Newton’s method typically requires  $O(m^3)$  computational effort; not even the best linear solvers achieve  $O(m^2)$ .

<sup>3</sup>Indeed, my  $\mathbf{Q}$  could throughout be replaced by  $\mathbf{R}\mathbf{Q}\mathbf{R}'$ ; for a similar comment, see Durbin and Koopman (2000, p. 43).

## 4.1 Inner maximisation

Taking Bellman's equation (5), substituting the quadratic approximation (8) and the (similarly quadratic) logarithmic state-transition density from model (12) yields

$$V_t(\mathbf{a}_t) = \ell(\mathbf{y}_t|\mathbf{a}_t) + \max_{\mathbf{a}_{t-1} \in \mathbb{R}^m} \left\{ -\frac{1}{2}(\mathbf{a}_t - \mathbf{c} - \mathbf{T}\mathbf{a}_{t-1})' \mathbf{Q}^{-1} (\mathbf{a}_t - \mathbf{c} - \mathbf{T}\mathbf{a}_{t-1}) - \frac{1}{2}(\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1})' \mathbf{I}_{t-1|t-1} (\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1}) \right\} + \text{constants}, \quad \mathbf{a}_t \in \mathbb{R}^m. \quad (13)$$

While  $\mathbf{Q}^{-1}$  is assumed to exist in writing equation (13), the results derived below will remain valid when  $\mathbf{Q}$  is only positive semidefinite; this follows from standard limiting arguments (e.g. Chopin and Papaspiliopoulos, 2020, p. 78). Here I focus on the maximisation over the lagged state variable  $\mathbf{a}_{t-1}$ .

As the variable  $\mathbf{a}_{t-1}$  appears at most quadratically on the right-hand side of equation (13), its maximisation can be performed in closed form. Importantly for the development below, the solution, denoted  $\mathbf{a}_{t-1}^* \in \mathbb{R}^m$ , depends linearly on the variable  $\mathbf{a}_t \in \mathbb{R}^m$ , which is involved in the outer maximisation. Hence  $\mathbf{a}_{t-1}^*$  is a vector function  $\mathbf{a}_{t-1}^* : \mathbb{R}^m \rightarrow \mathbb{R}^m$ , whose expression following from the standard first-order condition can be usefully expressed (after some algebra, see Supplement E) as

$$\mathbf{a}_{t-1}^* = \mathbf{a}_{t-1|t-1} + \mathbf{I}_{t-1|t-1}^{-1} \mathbf{T}' \mathbf{I}_{t|t-1} (\mathbf{a}_t - \mathbf{a}_{t|t-1}), \quad (14)$$

which employs the definitions of the predicted state  $\mathbf{a}_{t|t-1}$  and the predicted precision matrix  $\mathbf{I}_{t|t-1}$  given under step 2 in Table 3. Expression (14) can be recognised the one-period version of RTS (1965) smoother, providing the best estimate of  $\mathbf{a}_{t-1}$  conditional on the best estimate of next state,  $\mathbf{a}_t$ , which at this point remains to be found; i.e. the optimal  $\mathbf{a}_{t-1}^*$  is a function of the (still to be optimised) state variable  $\mathbf{a}_t$ .

Regarding the predicted precision matrix  $\mathbf{I}_{t|t-1}$ , the first expression in step 2 of Table 3 relies on the positive definiteness of the matrix  $\mathbf{Q}$ . The second expression, which holds by the Woodbury matrix identity, remains valid even when  $\mathbf{Q}$  becomes singular; a similar argument is made in Chopin and Papaspiliopoulos (2020, p. 78). Hence the algorithm in Table 3 remains valid when  $\mathbf{Q}$  is singular. While the derivation here is different, the resulting prediction step 2 in Table 3 is in fact identical to that of the (information form of the) Kalman filter (e.g. Harvey, 1990, p. 106). Hence, while the usual derivation of the Kalman filter is based on taking expectations, the optimisation approach presented here yields the same result.

## 4.2 Outer maximisation

Substituting the vector function  $\mathbf{a}_{t-1}^* : \mathbb{R}^m \rightarrow \mathbb{R}^m$  of equation (14) back into Bellman's equation (13), we obtain (after some algebra, see Supplement F) the value function with a single argument,  $\mathbf{a}_t$ , as follows:

$$V_t(\mathbf{a}_t) = \ell(\mathbf{y}_t|\mathbf{a}_t) - \frac{1}{2}(\mathbf{a}_t - \mathbf{a}_{t|t-1})' \mathbf{I}_{t|t-1} (\mathbf{a}_t - \mathbf{a}_{t|t-1}) + \text{constants}, \quad \mathbf{a}_t \in \mathbb{R}^m, \quad (15)$$

where predicted quantities  $\mathbf{a}_{t|t-1} \in \mathbb{R}^m$  and  $\mathbf{I}_{t|t-1} \in \mathbb{R}^{m \times m}$  were derived above (see step 2 of Table 3). The (approximate) value function (15) involves two terms: (a) the log-likelihood contribution of  $\mathbf{y}_t$  evaluated at the state variable  $\mathbf{a}_t$  and (b) a quadratic term that penalises deviations of  $\mathbf{a}_t$  from  $\mathbf{a}_{t|t-1}$ . The filtered

Table 3: Bellman filter and smoother for model (12).

Step	Method	Computation
1. Initialise	Unconditional Estimation	Set $\mathbf{a}_{0 0} = (\mathbf{1}_{m \times m} - \mathbf{T})^{-1} \mathbf{c}$ and $\text{vec}(\mathbf{I}_{0 0}^{-1}) = (\mathbf{1}_{m^2 \times m^2} - \mathbf{T} \otimes \mathbf{T})^{-1} \text{vec}(\mathbf{Q})$ . Set $t = 1$ .
	Diffuse	Treat $\mathbf{a}_{0 0}$ as a static parameter to be estimated and set $\mathbf{I}_{0 0}$ equal to a large multiple of the identity matrix. Set $t = 1$ .
		Possible if $\arg \max_{\mathbf{a}} \ell(\mathbf{y}_1   \mathbf{a})$ exists. Set $\mathbf{I}_{0 0}$ equal to a small multiple of the identity matrix. Set $t = 1$ .
2. Predict		$\mathbf{a}_{t t-1} = \mathbf{c} + \mathbf{T} \mathbf{a}_{t-1 t-1}$ . $\mathbf{I}_{t t-1} = \mathbf{Q}^{-1} - \mathbf{Q}^{-1} \mathbf{T} (\mathbf{I}_{t-1 t-1} + \mathbf{T}' \mathbf{Q}^{-1} \mathbf{T})^{-1} \mathbf{T}' \mathbf{Q}^{-1} = (\mathbf{T} \mathbf{I}_{t-1 t-1}^{-1} \mathbf{T}' + \mathbf{Q})^{-1}$ .
3. Start		Set $\mathbf{a}_t \leftarrow \mathbf{a}_{t t-1}$ . Alternatively, set $\mathbf{a}_t \leftarrow \arg \max_{\mathbf{a}} \ell(\mathbf{y}_t   \mathbf{a})$ if this quantity exists.
4. Optimise	Newton	$\mathbf{a}_t \leftarrow \mathbf{a}_t + \left[ \mathbf{I}_{t t-1} - \frac{d^2 \ell(\mathbf{y}_t   \mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} \right]^{-1} \left[ \frac{d\ell(\mathbf{y}_t   \mathbf{a}_t)}{d\mathbf{a}_t} - \mathbf{I}_{t t-1} (\mathbf{a}_t - \mathbf{a}_{t t-1}) \right]$ .
	Fisher	Like Newton step, but replace $d^2 \ell(\mathbf{y}_t   \mathbf{a}_t) / (d\mathbf{a}_t d\mathbf{a}_t')$ by $\mathbb{E}[d^2 \ell(\mathbf{y}_t   \mathbf{a}_t) / (d\mathbf{a}_t d\mathbf{a}_t')   \mathbf{a}_t]$ .
	BHHH	Like Newton step, but replace $d^2 \ell(\mathbf{y}_t   \mathbf{a}_t) / (d\mathbf{a}_t d\mathbf{a}_t')$ by $-d\ell(\mathbf{y}_t   \mathbf{a}_t) / d\mathbf{a}_t \times d\ell(\mathbf{y}_t   \mathbf{a}_t) / d\mathbf{a}_t'$ .
5. Stop		Stop at if some convergence criterion is satisfied or after a predetermined number of iterations.
6. Update		$\mathbf{a}_{t t} = \mathbf{a}_t$ .
	Newton	$\mathbf{I}_{t t} = \mathbf{I}_{t t-1} - \frac{d^2 \ell(\mathbf{y}_t   \mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} \Big _{\mathbf{a}_t = \mathbf{a}_{t t}}$ if the realised information is positive semidefinite
	Fisher	Like Newton update, but replace $d^2 \ell(\mathbf{y}_t   \mathbf{a}_t) / (d\mathbf{a}_t d\mathbf{a}_t')$ by $\mathbb{E}[d^2 \ell(\mathbf{y}_t   \mathbf{a}_t) / (d\mathbf{a}_t d\mathbf{a}_t')   \mathbf{a}_t]$ .
7. Proceed	BHHH	Like Newton update, but replace $d^2 \ell(\mathbf{y}_t   \mathbf{a}_t) / (d\mathbf{a}_t d\mathbf{a}_t')$ by $-d\ell(\mathbf{y}_t   \mathbf{a}_t) / d\mathbf{a}_t \times d\ell(\mathbf{y}_t   \mathbf{a}_t) / d\mathbf{a}_t'$ .
		Set $t = t + 1$ and return to step 2.
8. Smooth		Run the Bellman filter and store $\mathbf{a}_{t t}$ , $\mathbf{P}_{t t} = \mathbf{I}_{t t}^{-1}$ and $\mathbf{P}_{t t-1} = \mathbf{I}_{t t-1}^{-1}$ for all $1 \leq t \leq n$ . Start with $t = n - 1$ and iterate the following recursions backwards until $t = 1$ is reached: $\mathbf{a}_{t n} = \mathbf{a}_{t t} + \mathbf{P}_{t t} \mathbf{T}' \mathbf{I}_{t+1 t} (\mathbf{a}_{t+1 n} - \mathbf{c} - \mathbf{T} \mathbf{a}_{t t})$ , and $\mathbf{P}_{t n} = \mathbf{P}_{t t} - \mathbf{P}_{t t} \mathbf{T}' \mathbf{I}_{t+1 t} (\mathbf{P}_{t+1 t} - \mathbf{P}_{t+1 n}) \mathbf{I}_{t+1 t} \mathbf{T} \mathbf{P}_{t t}$ .

*Note:* BHHH = Berndt-Hall-Hall-Hausman. The log-likelihood function  $\ell(\mathbf{y}_t | \mathbf{a}_t)$  is known in closed form and can be read off from the data-generating process (12). The corresponding score and the realised and expected information quantities are written as  $d\ell(\mathbf{y}_t | \mathbf{a}) / d\mathbf{a}$ ,  $-d^2 \ell(\mathbf{y}_t | \mathbf{a}) / (d\mathbf{a} d\mathbf{a}')$  and  $\mathbb{E}[-d^2 \ell(\mathbf{y}_t | \mathbf{a}) / (d\mathbf{a} d\mathbf{a}') | \mathbf{a}]$ , respectively, which are viewed as functions of  $\mathbf{a}$ , to be evaluated at some state estimate. Steps 4 and 6 list three (intentionally vanilla) optimisation methods, which may but need not be identical for both steps. Users may also implement more sophisticated optimisation methods based on the argmax (16). The expressions in the (optional) smoother step 8 are derived in section 6.

state at time  $t$  maximises the sum of both terms, i.e.

$$\mathbf{a}_{t|t} := \underset{\mathbf{a}_t \in \mathbb{R}^m}{\operatorname{argmax}} V_t(\mathbf{a}_t) = \underset{\mathbf{a}_t \in \mathbb{R}^m}{\operatorname{argmax}} \left\{ \ell(\mathbf{y}_t | \mathbf{a}_t) - \frac{1}{2} (\mathbf{a}_t - \mathbf{a}_{t|t-1})' \mathbf{I}_{t|t-1} (\mathbf{a}_t - \mathbf{a}_{t|t-1}) \right\}. \quad (16)$$

The optimisation can be performed in closed form when the observation density is Gaussian with mean  $\mathbf{d} + \mathbf{Z} \mathbf{a}_t$ , as in Corollary 1, in which case  $\ell(\mathbf{y}_t | \mathbf{a}_t)$  is multivariate quadratic in  $\mathbf{a}_t$ ; this yields the standard Kalman filter (see Supplement G for details). In general, the potentially complicated functional form of  $\ell(\mathbf{y}_t | \mathbf{a}_t)$  implies that optimisation (16) cannot be performed in closed form. Some plain-vanilla applications of optimisation methods are included in Table 3 under step 4. The presence of the score in this optimisation step is distinctive for the Bellman filter and guarantees its robustness if the observation density is heavy tailed. As before, the computational complexity of the resulting filter is  $O(m^3 t)$ .

A unique argmax (16) is guaranteed when the precision matrix  $\mathbf{I}_{t|t-1}$  is positive definite and the log-likelihood function  $\ell(\mathbf{y}_t | \mathbf{a}_t)$  is concave in the state variable  $\mathbf{a}_t \in \mathbb{R}^m$ . When the smallest eigenvalue of the precision matrix  $\mathbf{I}_{t|t-1}$  is sufficiently large, a unique argmax is still guaranteed to exist even when  $\ell(\mathbf{y}_t | \mathbf{a}_t)$  fails to be concave in  $\mathbf{a}_t$ . In the non-concave case, it is possible that  $\mathbf{I}_{t|t-1}$  is insufficiently ‘large’ to pin down the update. This may be solved by adding to  $\mathbf{I}_{t|t-1}$  some positive multiple of the identity matrix or

skip the optimisation altogether; in the simulation study in section 8, this situation never arose.

Before proceeding to the next time step, the value function (15) must be approximated by a multivariate quadratic function. Because constants are irrelevant and the argmax has already been found, what remains is to determine the negative matrix of second derivatives evaluated at the peak, denoted  $\mathbf{I}_{t|t}$ , as indicated in Table 3 under step 6. Intuitively, one expects  $\mathbf{I}_{t|t} \geq \mathbf{I}_{t|t-1}$ , where the weak inequality means that the left-hand side minus the right-hand side is positive semidefinite. The intuition derives from the fact that missing observations can be dealt with as in the Kalman filter by setting  $\mathbf{a}_{t|t} = \mathbf{a}_{t|t-1}$  and  $\mathbf{I}_{t|t} = \mathbf{I}_{t|t-1}$ . Any (existing) observation should be weakly more informative than a nonexistent one, implying  $\mathbf{I}_{t|t} \geq \mathbf{I}_{t|t-1}$ . The lower bound may be reached in the limit for extreme observations (i.e. outliers), which are uninformative. While Newton’s updating method under step 6 has the advantage of explicitly utilising the observation  $\mathbf{y}_t$ , enabling it to recognise that some observations carry little information, the inequality  $\mathbf{I}_{t|t} \geq \mathbf{I}_{t|t-1}$  is not guaranteed unless the realised information quantity is positive semidefinite. For Fisher’s updating method under step 6, the situation is reversed, failing to utilise the realisation  $\mathbf{y}_t$  while ensuring  $\mathbf{I}_{t|t} \geq \mathbf{I}_{t|t-1}$ . For some models it is possible to formulate a hybrid version, e.g. by taking a weighted average of Newton’s and Fisher’s updating methods, that achieves the best of both worlds (I use this hybrid method for some models in section 8).

### 4.3 Special cases of Bellman filter with linear Gaussian states

Special cases of the algorithm in Table 3 include the Kalman filter (Supplement G), the iterated extended Kalman filter (Supplement H), Fahrmeir’s (1992) approximate mode estimator (Supplement I), Koyama et al.’s (2010) Laplace Gaussian filter (Supplement J), and Toulis and Airoidi’s (2017) implicit stochastic gradient method for the estimation of states that are constant over time (Supplement K). The key difference with implicit stochastic gradient methods is that the Bellman filter, like the Kalman filter, generally remains perpetually responsive and does not converge to a ‘true’ parameter value.

## 5 Theory: Contractivity, error bounds and stability

This section investigates the theoretical properties of the Bellman filter derived in the previous section, i.e. under the assumption of linear and Gaussian state dynamics. Under appropriate conditions, this section will show that (a) at a fixed time step, the Bellman filtering step is contractive in quadratic mean to a small region around the true state, (b) over time, the mean squared filtering error remains uniformly bounded (i.e. approximation errors cannot accumulate), and (c) the effect of the initialisation of the filter vanishes asymptotically and exponentially fast, an important property known as invertibility (Straumann and Mikosch (2006) or stability (Koyama et al., 2010, Th. 4).

### 5.1 Contractivity at a fixed time step

Here the time step  $t \geq 1$  is considered fixed. Similarly, in the Bellman-filter update (16), predictions  $\mathbf{a}_{t|t-1} \in \mathbb{R}^m$  and  $\mathbf{I}_{t|t-1} \in \mathbb{R}^{m \times m}$  are fixed. Update (16) can generally be viewed as a stochastic version of Rockafellar’s (1976) proximal point algorithm, which similarly combines a target function to be optimised, in this case  $\ell(\mathbf{y}_t | \mathbf{a}_t)$ , with a quadratic penalty centred at a previous iterate, in this case  $\mathbf{a}_{t|t-1}$ . Indeed, optimisation (16) can be classed as a stochastic proximal point method (e.g. Ryu and Boyd, 2016, Bianchi, 2016, Patrascu and Necoara, 2018, Asi and Duchi, 2019). Its intuitive link with proximal optimisation

methods suggests that update (16) should remain both applicable and reasonably accurate outside the classic Kalman-filtering context. Theorem 1 below confirms this intuition.

**Notation:** For vectors  $\mathbf{x} \in \mathbb{R}^m$ , the Euclidean norm is denoted by  $\|\mathbf{x}\| := \sqrt{\mathbf{x}'\mathbf{x}}$ . For a positive definite weight matrix  $\mathbf{W} > \mathbf{0}$ , the weighted Euclidean vector norm is denoted  $\|\mathbf{x}\|_{\mathbf{W}} := \sqrt{\mathbf{x}'\mathbf{W}\mathbf{x}}$ , while for a matrix  $\mathbf{M} \in \mathbb{R}^{m \times m}$ , the induced matrix norm is denoted  $\|\mathbf{M}\|_{\mathbf{W}} := \max\{\|\mathbf{M}\mathbf{x}\|_{\mathbf{W}} : \|\mathbf{x}\|_{\mathbf{W}} = 1\}$  (see e.g. Jungers, 2009, Def. 2.8). The gradient and Hessian of  $\ell(\mathbf{y}|\mathbf{a})$  with respect to  $\mathbf{a}$  are written as  $\nabla\ell(\mathbf{y}|\mathbf{a})$  and  $\nabla^2\ell(\mathbf{y}|\mathbf{a})$ , respectively. The smallest and largest eigenvalues of a matrix  $\cdot$  are denoted  $\lambda_{\min}(\cdot)$  and  $\lambda_{\max}(\cdot)$ , respectively. The  $m \times m$  identity matrix is denoted by  $\mathbf{1}_{m \times m}$ .

**Assumption 1 (Concavity)** *With probability one in the random draw  $\mathbf{y}$ , the observation log density  $\ell(\mathbf{y}|\cdot)$  maps  $\mathbb{R}^m$  to  $\mathbb{R}$ , and is either (a) concave, or (b) strictly concave, or (c) strongly concave with parameter  $\epsilon > 0$ .*

**Assumption 2 (Differentiability)** *With probability one in the random draw  $\mathbf{y}$ , the observation log density  $\mathbf{a} \mapsto \ell(\mathbf{y}|\mathbf{a})$  is (a) once or (b) twice continuously differentiable on all of  $\mathbb{R}^m$ .*

**Assumption 3 (Bounded information)**  $\mathbb{E}[\|\nabla\ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)\|^2] \leq \sigma^2 < \infty$ , where  $\boldsymbol{\alpha}_t$  is the true latent state that generates  $\mathbf{y}_t \sim p(\mathbf{y}_t|\boldsymbol{\alpha}_t)$ .

**Theorem 1 (Contractivity of the mean squared error)** *Fix the time step  $t \geq 1$ . Let  $\mathbf{a}_{t|t-1} \in \mathbb{R}^m$  and  $\mathbf{I}_{t|t-1} \in \mathbb{R}^{m \times m}$  be given and fixed, where the latter is symmetric and positive definite with eigenvalues satisfying  $0 < \lambda_{\min}(\mathbf{I}_{t|t-1}) \leq \lambda_{\max}(\mathbf{I}_{t|t-1}) < \infty$ . Let update  $\mathbf{a}_{t|t}$  be defined by (16).*

1. **Boundedness of updates:** *Under Assumption 1a, with probability one, the update  $\mathbf{a}_{t|t}$  is well defined and satisfies*

$$\frac{1}{2} \|\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1}\|_{\mathbf{I}_{t|t-1}}^2 \leq \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) - \ell(\mathbf{y}_t|\mathbf{a}_{t|t-1}). \quad (17)$$

2. **Stability for a single time step:** *Let Assumption 2b hold. Let  $\lambda_{\min}(\mathbf{I}_{t|t-1}) > \max\{0, \lambda_{\max}(\nabla^2\ell(\mathbf{y}|\mathbf{a}))\}$  for all  $\mathbf{a} \in \mathbb{R}^m$  and with probability one in  $\mathbf{y}$ . Then, with probability one,*

$$\left\| \frac{d\mathbf{a}_{t|t}}{d\mathbf{a}'_{t|t-1}} \right\|_{\mathbf{I}_{t|t-1}} \leq 1 - \frac{\lambda_{\min}(-\nabla^2\ell(\mathbf{y}_t|\mathbf{a}_{t|t}))}{\lambda_{\max}(\mathbf{I}_{t|t-1}) + \lambda_{\max}(-\nabla^2\ell(\mathbf{y}_t|\mathbf{a}_{t|t}))}. \quad (18)$$

*The right-hand side does not exceed (is strictly less than) unity under the additional Assumption 1a (1b).*

3. **Contractivity of the quadratic error:** *Under Assumptions 1c, 2a and 3,*

$$\mathbb{E} \left( \|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1} + 2\epsilon \mathbf{1}_{m \times m}}^2 \right) \leq \mathbb{E} \left( \|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2 \right) + \frac{\sigma^2}{\lambda_{\min}(\mathbf{I}_{t|t-1})}. \quad (19)$$

The proof is presented in Supplement L. Compared with other results for approximate filters (e.g. Koyama et al., 2010), Theorem 1 is attractive because the assumptions are (a) more easily verifiable (relating to model inputs instead of outputs) and (b) less stringent. For example, Theorem 1 applies to the Kalman filter, while the theory developed in Koyama et al. (2010) does not.<sup>4</sup>

<sup>4</sup>Koyama et al. (2010) require logarithmic observation densities with five uniformly bounded derivatives, ruling out the Gaussian case in which the logarithmic density is quadratic, implying unbounded first derivatives on  $\mathbb{R}^m$ .

Part 1 of Theorem 1 indicates that the update is well-defined, while Part 2 demonstrates that the Bellman-filtered state  $\mathbf{a}_{t|t}$  is stable in the prediction  $\mathbf{a}_{t|t-1}$ . This stability property can be used to establish the stability of the Bellman filter (see section 5.3). Part 3 of Theorem 1 says that the quadratic filtering error is contractive in expectation towards a small region around the true state. Inequality (19) features a weighted norm on both sides, in which the predicted information matrix  $\mathbf{I}_{t|t-1}$  plays a key role. The weight matrix on the left-hand side of inequality (19) contains the additional term  $2\epsilon\mathbb{1}_{m\times m}$  such that the diagonal is ‘reinforced’: this drives the contraction. Intuitively, when the weight matrix is ‘bigger’ (i.e. has larger eigenvalues), the vector inside the norm must be ‘smaller’ in magnitude. Of course, an improvement is impossible when the prediction is perfect, such that the additive term  $\sigma^2/\lambda_{\min}(\mathbf{I}_{t|t-1})$  on the right-hand side of equation (19) is unavoidable. Hence updates are contractive in quadratic mean towards a ‘noise-dominated region’ (NDR) around the true state (e.g. Patrascu and Necoara, 2018, p. 3).

Theorem 1 also relates to Toulis et al. (2016, p. 1291), who present the seemingly stronger result that proximal updates are ‘contracting almost surely’ when the log-likelihood function is strongly concave; however, their result relies on a nonstandard definition of strong concavity that rules out important cases of interest, e.g. the Kalman filter (see Supplement M for a detailed comparison).

## 5.2 Error bounds over time

While Theorem 1 involved a fixed time step, it is equally important to investigate how filtered quantities behave over extended time periods. When the latent state is stationary, even a trivial filter may asymptotically achieve a bounded mean squared error (MSE), e.g. by setting the filter output equal to zero for all time steps. Hence a more pertinent question is whether the filter can asymptotically achieve a bounded MSE in the case of unit-root states. As this section shows, in the long run, the Bellman filter achieves a bounded MSE even if the true process is free to roam.

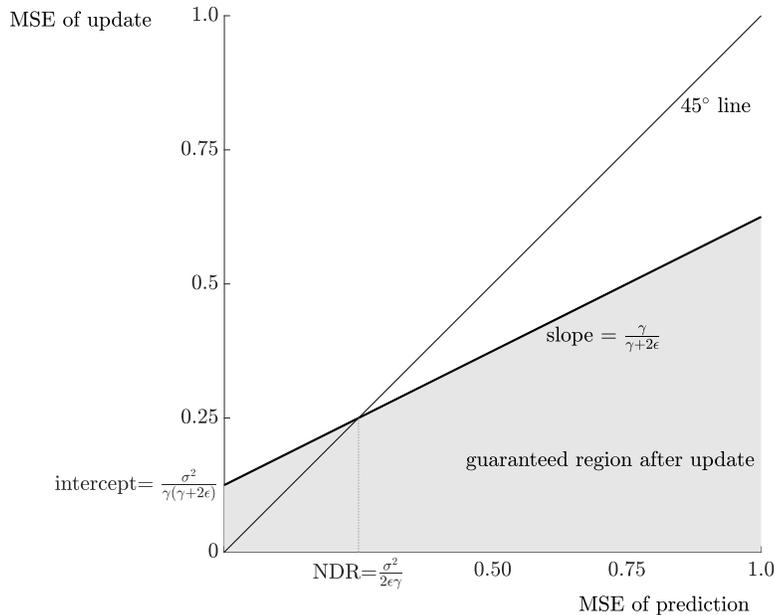
For simplicity I focus on the case in which  $\mathbf{I}_{t|t-1}$  is a constant multiple of the identity matrix; hence  $\mathbf{I}_{t|t-1} = \gamma\mathbb{1}_{m\times m}$ , where  $\gamma > 0$  can be interpreted as a smoothing parameter, and  $\lambda_{\min}(\mathbf{I}_{t|t-1}) = \lambda_{\max}(\mathbf{I}_{t|t-1}) = \gamma$ . The weighted MSE contraction (19) for a fixed time step then reduces to a standard MSE contraction:

$$\underbrace{\mathbb{E}\left(\|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|^2\right)}_{\text{MSE of update}} \leq \underbrace{\frac{\gamma}{\gamma + 2\epsilon}}_{<1} \left[ \underbrace{\mathbb{E}\left(\|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|^2\right)}_{\text{MSE of prediction}} + \underbrace{\frac{\sigma^2}{\gamma^2}}_{>0} \right]. \quad (20)$$

Inequality (20) features a multiplicative constant on its right-hand side that is strictly less than unity, which gives rise to the contraction. As illustrated in Figure 1, the inequality says that the MSE of the update is bounded above by a linear function of the MSE of the prediction. The slope of this line is  $\gamma/(\gamma + 2\epsilon) < 1$ , while the intercept is  $\sigma^2/(\gamma(\gamma + 2\epsilon)) > 0$ . The area below the line, shaded in grey, shows the contraction due to inequality (20). When the prediction error is large, the contractive property dominates and the update is expected to be beneficial: the grey area lies below the 45° line. When the prediction happens to be pinpoint accurate (i.e. the corresponding MSE is zero), the MSE of the update need not be zero, as can be seen in Figure 1 from the fact that the grey area stretches above the 45° line close to the origin. This is unavoidable with noisy data: when predictions are perfect, updates cannot be better. In the limit  $\epsilon \rightarrow 0$ , whereby the target function is concave but not strongly so, inequality (20) is closely related to Theorem 3.2 in Asi and Duchi (2019).

MSE contraction (20) is used below in Proposition 2 (see Supplement N for the proof) to demonstrate

Figure 1: Illustration of mean squared error (MSE) contraction due to inequality (20)



*Note:* NDR = noise-dominated region. The grey area corresponds to possible values of the MSE after updating, which is conditional on the MSE before updating. Purely for illustrative purposes, the parameters are  $\sigma = \epsilon = 1$  and  $\gamma = 2$ .

that the filtering MSE remains uniformly bounded over time. Proposition 2 applies to the Kalman filter, which can similarly track unit-root states in the long run, but holds more generally for strictly concave logarithmic observation densities.

**Proposition 2 (Uniformly bounded MSE)** *Assume  $\mathbf{\alpha}_t = \mathbf{\alpha}_{t-1} + \boldsymbol{\eta}_t$  with  $\boldsymbol{\eta}_t \sim$  i.i.d.  $(\mathbf{0}, \mathbf{Q})$ , which need not be Gaussian, and  $\sigma_\eta^2 = \text{Trace}(\mathbf{Q}) < \infty$ . Set  $\mathbf{a}_{t+1|t} = \mathbf{a}_{t|t}$  and take  $\mathbf{I}_{t+1|t} = \gamma \mathbf{1}_{m \times m}$  for some  $\gamma > 0$  and all  $t \geq 1$ . Let  $\mathbf{a}_{t|t}$  be given by update (16). Denote  $\text{MSE}_{t|t} := \mathbb{E} \|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|^2$  and  $\text{MSE}_{t|t-1} := \mathbb{E} \|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|^2$ . In the setting of part 3 of Theorem 1,*

$$\text{MSE}_{t|t} \leq \frac{\gamma}{\gamma + 2\epsilon} \left[ \text{MSE}_{t|t-1} + \frac{\sigma^2}{\gamma^2} \right], \quad \text{MSE}_{t+1|t} = \text{MSE}_{t|t} + \sigma_\eta^2, \quad t \geq 1. \quad (21)$$

*Irrespective of the initial value  $\text{MSE}_{1|0}$ , the long-run filtering error remains uniformly bounded:*

$$\limsup_{t \rightarrow \infty} \text{MSE}_{t|t} \leq \frac{\sigma^2}{2\gamma\epsilon} + \frac{\gamma\sigma_\eta^2}{2\epsilon}. \quad (22)$$

*Minimising the bound with respect to  $\gamma$  yields  $\gamma = \sigma/\sigma_\eta$ .*

### 5.3 Stability

As emphasised by Anderson and Moore (2012, p. 63), ‘a question of vital interest [...] is whether or not the filter is stable’. A filter can be considered stable if deviations in the initial conditions ‘tend to be reduced, rather than amplified, by conditioning on further observations’ (Koyama et al., 2010). To this end, it is sufficient that filtered paths with different initialisations—but based on identical data—converge exponentially fast over time, a concept known as ‘invertibility’ (e.g. Straumann and Mikosch, 2006). This section demonstrates the stability of a time-invariant version of the Bellman filter.

Stability analyses of the Kalman filter rely on the fact that, in the time-invariant version of the filter, the matrix  $d\mathbf{a}'_{t|t}/d\mathbf{a}_{t-1|t-1}$  is static, as  $\mathbf{a}_{t|t}$  is then a linear function of  $\mathbf{a}_{t-1|t-1}$  with a static coefficient matrix. Stability follows when the spectral radius of this coefficient matrix is strictly exceeded by one. Unfortunately, the stability analysis here is complicated by the fact that each derivative matrix  $d\mathbf{a}'_{t|t}/d\mathbf{a}_{t-1|t-1}$  is stochastic, depending on the observations as well as the filtered states. Moreover, an analysis based on the spectral radius is ruled out because it fails to be a norm. I follow the classic literature in investigating a time-invariant setting, which implies that the predicted information matrix  $\mathbf{I}_{t|t-1} = \mathbf{I} \in \mathbb{R}^{m \times m}$  is taken to be static over time. I deviate by basing the result not on the spectral radius but the (weighted) matrix norm  $\|\cdot\|_{\mathbf{I}}$ .

**Theorem 2 (Stability of the time-invariant Bellman filter.)** *Let the initialisation  $\mathbf{a}_{0|0} \in \mathbb{R}^m$  be given. For all  $t \geq 1$ , (a) set  $\mathbf{a}_{t|t-1} = \mathbf{c} + \mathbf{T}\mathbf{a}_{t-1|t-1}$ , where  $\mathbf{c} \in \mathbb{R}^m$  and  $\mathbf{T} \in \mathbb{R}^{m \times m}$  are given, and (b) let update  $\mathbf{a}_{t|t}$  be defined by maximisation (16), where  $\mathbf{I}_{t|t-1} = \mathbf{I} \in \mathbb{R}^{m \times m}$  is a time-invariant (i.e. static) positive-definite matrix with eigenvalues in the range  $(\nu_{\min}, \nu_{\max})$ . Assume that, with probability one, the observation log density  $\ell(\mathbf{y}|\mathbf{a})$  is twice continuously differentiable, while the negative Hessian matrix  $-\nabla^2 \ell(\mathbf{y}|\mathbf{a})$  has eigenvalues in the range  $(\mu_{\min}, \mu_{\max})$  uniformly for  $\mathbf{a} \in \mathbb{R}^m$ , where  $\max\{0, -\mu_{\min}\} < \nu_{\min}$ . Then, with probability one,*

$$\left\| \frac{d\mathbf{a}_{t|t}}{d\mathbf{a}'_{0|0}} \right\|_{\mathbf{I}} \leq \left( 1 - \min \left\{ \frac{\delta}{\nu_{\min}}, \frac{\delta}{\nu_{\max}} \right\} \right)^{t/2} \left( 1 - \frac{\mu_{\min}}{\nu_{\max} + \mu_{\max}} \right)^t, \quad (23)$$

where  $\delta := \lambda_{\min}(\mathbf{I} - \mathbf{T}'\mathbf{I}\mathbf{T}) \leq \nu_{\min}$ . As  $t \rightarrow \infty$ , exponential almost sure convergence to zero is guaranteed under the following sufficient condition:

$$\frac{1}{2} \log \left( 1 - \min \left\{ \frac{\delta}{\nu_{\min}}, \frac{\delta}{\nu_{\max}} \right\} \right) + \log \left( 1 - \frac{\mu_{\min}}{\nu_{\max} + \mu_{\max}} \right) < 0. \quad (24)$$

The proof is presented in Supplement O. Theorem 2 assumes that  $\mathbf{I}$  is positive definite while its smallest eigenvalue  $\nu_{\min} > 0$  is sufficiently large. For concave log densities (i.e.  $\mu_{\min} \geq 0$ ), it is required only that  $\nu_{\min} > 0$  such that  $\mathbf{I}$  is positive definite. For log densities that fail to be concave (i.e.  $\mu_{\min} < 0$ ), the stronger condition  $\nu_{\min} > \max\{0, -\mu_{\min}\}$  is imposed to ensure that optimisation problem (16) is well-defined and leads to unique solution  $\mathbf{a}_{t|t}$  for all  $t$ . The sufficient condition (24) for invertibility is automatically satisfied if the prediction and updating steps are both non-expansive (both  $\delta \geq 0$  and  $\mu_{\min} \geq 0$ ), while at least one is strictly contractive ( $\delta > 0$  and/or  $\mu_{\min} > 0$ ). For example, the observation log density could be strictly concave (i.e.  $\mu_{\min} > 0$ ) while  $\mathbf{T}$  is the identity matrix (in which case  $\delta = 0$ ); hence, unit root dynamics are permitted. Moreover, inequality (24) will always be satisfied if the observations point adequately to the underlying state. More specifically, if  $\mu_{\min}$  and  $\mu_{\max}$  approach infinity at the same rate (such that the measurement is exceedingly precise), then the second logarithm in condition (24) approaches negative infinity such that the condition is satisfied. For sufficiently informative observations, therefore, even explosive state dynamics may be accommodated.

## 6 Smoothing using Bellman's principle

Here the general method in section 2 is extended to present a unified method for both filtering and smoothing using Bellman's dynamic-programming principle. Readers purely interested in filtering can skip

this section without loss of continuity. While the approach below is general, I present the most explicit result in the case of a linear Gaussian state equation. This specialised setting allows me to show that the classic Rauch, Tung and Striebel (RTS, 1965) smoother expressions remain valid, albeit as approximations, for a general (i.e. non-Gaussian) observation density—an insight that may be useful in practice.

Below I introduce three value functions, based on (a) past data, (b) future data and (c) all data. All three are based on the partial log-likelihood function  $L_{t_1:t_2} : \Omega \times \mathbb{R}^m \times \dots \times \mathbb{R}^m \rightarrow \mathbb{R}$  involving states and observations from time  $t_1$  to  $t_2$  as follows:

$$L_{t_1:t_2}(\mathbf{a}_{t_1}, \dots, \mathbf{a}_{t_2}) := \sum_{i=t_1}^{t_2} \ell(\mathbf{y}_i | \mathbf{a}_i) + \sum_{i=t_1+1}^{t_2} \ell(\mathbf{a}_i | \mathbf{a}_{i-1}) + \mathbb{1}_{t_1=1} \ell(\mathbf{a}_1), \quad 1 \leq t_1 \leq t_2 \leq n, \quad (25)$$

where sums containing no terms are understood to be zero. Equation (25) generalises equation (2), which is a special case with  $t_1 = 1$  and  $t_2 = t$ . The new function  $L_{t_1:t_2}(\dots)$  depends on observations  $\mathbf{y}_{t_1}$  through  $\mathbf{y}_{t_2}$ , which are considered fixed, and involves  $t_2 - t_1$  state transitions from  $\mathbf{a}_{t_1}$  to  $\mathbf{a}_{t_2}$ . For definiteness, I assume that  $L_{t_1:t_2}(\cdot, \dots, \cdot)$  can be maximised with respect to each input argument; this assumption is too strong but sufficient for the development below.

**Assumption 4** *For all  $1 \leq t_1 \leq t_2 \leq n$ , the partial log-likelihood function  $L_{t_1:t_2}(\cdot, \dots, \cdot)$  defined in equation (25) has a unique maximum with respect to each state variable  $\mathbf{a}_t$ , i.e. for each  $t_1 \leq t \leq t_2$ .*

Assumption 4 allows us to define three value functions  $V_t(\cdot), W_t(\cdot), Z_t(\cdot) : \Omega \times \mathbb{R}^m \rightarrow \mathbb{R}$  as follows:

$$\text{using past data:} \quad V_t(\mathbf{a}_t) := \max_{\mathbf{a}_1, \dots, \mathbf{a}_{t-1}} L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t), \quad (26)$$

$$\text{using future data:} \quad W_t(\mathbf{a}_t) := \max_{\mathbf{a}_{t+1}, \dots, \mathbf{a}_n} L_{t:n}(\mathbf{a}_t, \dots, \mathbf{a}_n), \quad (27)$$

$$\text{using all data:} \quad Z_t(\mathbf{a}_t) := \max_{\mathbf{a}_1, \dots, \mathbf{a}_{t-1}, \mathbf{a}_{t+1}, \dots, \mathbf{a}_n} L_{1:n}(\mathbf{a}_1, \dots, \mathbf{a}_n), \quad (28)$$

where  $1 \leq t \leq n$ . Maximisations are written as  $\max_{\mathbf{a}}$  instead of  $\max_{\mathbf{a} \in \mathbb{R}^m}$ ; i.e. it is implicitly understood that each state variable takes values in the state space  $\mathbb{R}^m$ . The backward-looking value function  $V_t(\cdot)$  is identical to that in Definition 1. The forward-looking value function  $W_t(\cdot)$  is based on current and future data and specialises to that in Mayne (1966, eq. 18) for linear Gaussian state-space models. The convention that any maximisation involving no variables can be ignored gives the correct initial and terminal conditions for  $t = 1$  and  $t = n$ , respectively. Function  $Z_t(\cdot)$  is based on all data and implies a smoothed state estimate via  $\mathbf{a}_{t|n} := \operatorname{argmax}_{\mathbf{a}} Z_t(\mathbf{a})$ . The usefulness of the above definitions lies in the fact that the first two value functions satisfy forward and backward recursions, respectively, while jointly implying the third:

**Proposition 3 (Bellman's forward and backward recursions.)** *Let Assumption 4 hold. Then*

$$\text{forward recursion:} \quad V_t(\mathbf{a}_t) = \ell(\mathbf{y}_t | \mathbf{a}_t) + \max_{\mathbf{a}_{t-1}} \left\{ \ell(\mathbf{a}_t | \mathbf{a}_{t-1}) + V_{t-1}(\mathbf{a}_{t-1}) \right\}, \quad 1 < t \leq n, \quad (29)$$

$$\text{backward recursion:} \quad W_t(\mathbf{a}_t) = \ell(\mathbf{y}_t | \mathbf{a}_t) + \max_{\mathbf{a}_{t+1}} \left\{ \ell(\mathbf{a}_{t+1} | \mathbf{a}_t) + W_{t+1}(\mathbf{a}_{t+1}) \right\}, \quad 1 \leq t < n, \quad (30)$$

$$\text{relation between both:} \quad Z_t(\mathbf{a}_t) = V_t(\mathbf{a}_t) + \max_{\mathbf{a}_{t+1}} \left\{ \ell(\mathbf{a}_{t+1} | \mathbf{a}_t) + W_{t+1}(\mathbf{a}_{t+1}) \right\}, \quad 1 \leq t < n, \quad (31)$$

$$= W_t(\mathbf{a}_t) + \max_{\mathbf{a}_{t-1}} \left\{ \ell(\mathbf{a}_t | \mathbf{a}_{t-1}) + V_{t-1}(\mathbf{a}_{t-1}) \right\}, \quad 1 < t \leq n. \quad (32)$$

The proof, being a straightforward extension of that of Proposition 1, is omitted. Forward recursion (29) is identical that in Proposition 1, while backward recursion (30) can be derived using similar arguments; for linear Gaussian state-space models, the latter collapses to the backward recursion in Mayne (1966, eq. 27). Function  $Z_t(\cdot)$  can be constructed by combining the output of both recursions, where either the forward or backward recursion extends to time  $t$  as in equations (31) and (32), respectively. In both cases, a single-state transition log-density is added, followed by an optimisation involving a single state variable.

Interestingly, equations (31) and (32) do not (explicitly) contain the observation density. Instead, they contain only two value functions (one using past data, one using future data) that are linked through a single state-transition density. When both value functions are quadratic, and the state-transition equation is linear and Gaussian, such that  $\ell(\mathbf{a}_t|\mathbf{a}_{t-1})$  is also quadratic, then equations (31) and (32) contain only quadratic terms and should thus be analytically soluble. As illustrated below, this yields the classic RTS smoother expressions. However, the main innovation of this article is to consider quadratic value functions even when inexact. As the next proposition shows, if we are willing to accept that value functions may be reasonably approximated by quadratic functions, then the resulting expression is still given by the classic RTS smoother. This insight appears to be new, and considerably extends the domain of applicability of the RTS smoother, at least as an approximation. In practice, it means that the Bellman filter developed in section 4 can be executed and its output used in the standard RTS smoothing formulas to obtain approximate smoothed state estimates—which the simulation study in section 8 finds to be highly accurate.

**Proposition 4 (Bellman smoother with linear Gaussian state equation)** *Let Assumption 4 hold. Assume  $\boldsymbol{\alpha}_t = \mathbf{c} + \mathbf{T}\boldsymbol{\alpha}_{t-1} + \boldsymbol{\eta}_t$  with  $\boldsymbol{\eta}_t \sim i.i.d. N(\mathbf{0}, \mathbf{Q})$ . Suppose that both value functions on the right-hand side of equation (31) are approximated as quadratic functions; in particular let  $V_t(\cdot)$  have  $\text{argmax } \mathbf{a}_{t|t}$  and negative Hessian  $\mathbf{I}_{t|t} = \mathbf{P}_{t|t}^{-1} > \mathbf{0}$ . Under this approximation,  $Z_t(\cdot)$  on the left-hand side of equation (31) is also quadratic. Moreover, the  $\text{argmax } \mathbf{a}_{t|n}$  of  $Z_t(\cdot)$  can be expressed in terms of the  $\text{argmax } \mathbf{a}_{t+1|n}$  of  $Z_{t+1}(\cdot)$  as follows:*

$$\mathbf{a}_{t|n} = \mathbf{a}_{t|t} + \mathbf{P}_{t|t}\mathbf{T}'\mathbf{I}_{t+1|t}(\mathbf{a}_{t+1|n} - \mathbf{c} - \mathbf{T}\mathbf{a}_{t|t}), \quad (33)$$

$$\mathbf{P}_{t|n} = \mathbf{P}_{t|t} - \mathbf{P}_{t|t}\mathbf{T}'\mathbf{I}_{t+1|t}(\mathbf{P}_{t+1|t} - \mathbf{P}_{t+1|n})\mathbf{I}_{t+1|t}\mathbf{T}\mathbf{P}_{t|t}, \quad (34)$$

where  $\mathbf{I}_{t+1|t} := (\mathbf{T}\mathbf{P}_{t|t}\mathbf{T}' + \mathbf{Q})^{-1} > \mathbf{0}$  and  $\mathbf{I}_{t|n} = \mathbf{P}_{t|n}^{-1} > \mathbf{0}$  for  $t = 1, \dots, n$  is the negative Hessian of  $Z_t(\cdot)$ . Expressions (33) and (34) are identical to the classic RTS smoother expressions, but in a more general—i.e. possibly approximate—context.

The proof, presented in Supplement Q, employs only standard matrix algebra, including a simple lemma on multivariate quadratic functions in Supplement P. Exact solubility of equation (31) is clear given that all functions on its right-hand side are assumed to be quadratic; the crucial step is to relate the properties of  $Z_t(\cdot)$  to those of  $Z_{t+1}(\cdot)$  to obtain a backward recursion. The resulting RTS smoother (33) requires us to store the output of the filter for all time steps and subsequently to compute the smoothed state,  $\mathbf{a}_{t|n}$ , as a linear combination of the filtered state,  $\mathbf{a}_{t|t}$ , and the adjacent smoothed state,  $\mathbf{a}_{t+1|n}$ . The backward recursion can be initialised using the final filtered state,  $\mathbf{a}_{n|n}$ . The output of the backward matrix recursion (34), which provides a measure of uncertainty, is not required if one is merely interested in the smoothed state estimates (33).

## 7 Parameter estimation by likelihood approximation

This section presents a heuristic approach to the static-parameter estimation problem, as distinct from the filtering problem, in that we aim to estimate both the time-varying states and the static (hyper)parameter  $\psi$ . I deviate from the literature by decomposing the log-likelihood function of the data in terms of the ‘fit’ generated by the Bellman filter, penalised by a nonnegative term that resembles a ‘realised’ version of the Kullback-Leibler (KL, 1951) divergence between filtered and predicted states. Intuitively, this decomposition illustrates that we wish to maximise the congruence of the Bellman-filtered states and the data, while minimising the distance between the filtered and predicted states to prevent over-fitting.

The proposed pseudo log-likelihood decomposition has the advantage that all terms can be evaluated or approximated using the output of the Bellman filter; no sampling techniques or numerical integration methods are required. While no formal guarantees of convergence are provided, I analyse the statistical properties of the proposed static-parameter estimator in extensive simulation studies (see section 8) and find that it performs on par with simulation-based methods at a fraction of the computational cost. The development of an asymptotic theory remains unresolved.

To introduce the proposed decomposition, I focus on the log-likelihood contribution of a single observation,  $\ell(\mathbf{y}_t|\mathcal{F}_{t-1}) := \log p(\mathbf{y}_t|\mathcal{F}_{t-1})$ . The equalities below follow immediately from the definition of conditional densities and the assumption of the state-space model (1):

$$\ell(\mathbf{y}_t|\mathcal{F}_{t-1}) = \ell(\mathbf{y}_t, \boldsymbol{\alpha}_t|\mathcal{F}_{t-1}) - \ell(\boldsymbol{\alpha}_t|\mathbf{y}_t, \mathcal{F}_{t-1}) = \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) + \ell(\boldsymbol{\alpha}_t|\mathcal{F}_{t-1}) - \ell(\boldsymbol{\alpha}_t|\mathcal{F}_t). \quad (35)$$

While the above decomposition is valid for any  $\boldsymbol{\alpha}_t \in \mathbb{R}^m$ , the resulting expression is not a computable quantity, as the true latent state  $\boldsymbol{\alpha}_t$  remains unknown. It is practical to evaluate the expression at the Bellman-filtered state  $\mathbf{a}_{t|t}$  and swap the order of the last two terms, such that

$$\ell(\mathbf{y}_t|\mathcal{F}_{t-1}) = \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)\Big|_{\boldsymbol{\alpha}_t=\mathbf{a}_{t|t}} - \underbrace{\left\{ \ell(\boldsymbol{\alpha}_t|\mathcal{F}_t) - \ell(\boldsymbol{\alpha}_t|\mathcal{F}_{t-1}) \right\}}_{\text{‘realised’ KL divergence}}\Big|_{\boldsymbol{\alpha}_t=\mathbf{a}_{t|t}}. \quad (36)$$

The first term on the right-hand side,  $\ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)$  evaluated at  $\boldsymbol{\alpha}_t = \mathbf{a}_{t|t}$ , quantifies the congruence (or ‘fit’) between the Bellman-filtered state  $\mathbf{a}_{t|t}$  and the observation  $\mathbf{y}_t$ , which we wish to maximise. We simultaneously aim to minimise the term in curly brackets, i.e. the difference  $\ell(\boldsymbol{\alpha}_t|\mathcal{F}_t) - \ell(\boldsymbol{\alpha}_t|\mathcal{F}_{t-1})$  evaluated at  $\boldsymbol{\alpha}_t = \mathbf{a}_{t|t}$ . This difference can be viewed as a ‘realised’ version of the KL divergence between the filtered and predicted densities; intuitively, it indicates the level of ‘surprise’ associated with the filtered state  $\mathbf{a}_{t|t}$ . The standard KL divergence between filtered and predicted densities would have read  $\mathbb{E}[\log(\boldsymbol{\alpha}_t|\mathcal{F}_t) - \log(\boldsymbol{\alpha}_t|\mathcal{F}_{t-1})]$ , which involves an expectation operator that integrates out the state  $\boldsymbol{\alpha}_t$  using the true density  $p(\boldsymbol{\alpha}_t|\mathcal{F}_t)$ . Equation (36) contains no expectation but is simply evaluated at the filtered state  $\mathbf{a}_{t|t}$ ; hence, it can be viewed as a realised version. The trade-off in equation (36) between maximising the fit while minimising the surprise gives rise to a meaningful optimisation problem.

While decomposition (36) is exact, we do not generally have an exact expression for the terms in curly brackets. To ensure that the log-likelihood contribution (36) is computable, I now turn to approximating the realised KL divergence. In deriving the Bellman filter, I presumed that the researcher’s knowledge, as measured in log-likelihood space for each time step, could be approximated by a multivariate quadratic function. Extending this line of reasoning, I consider the following approximations of the two terms that

compose the realised KL divergence:

$$\ell(\boldsymbol{\alpha}_t|\mathcal{F}_t) \approx \frac{1}{2} \log \det\{\mathbf{I}_{t|t}/(2\pi)\} - \frac{1}{2}(\boldsymbol{\alpha}_t - \mathbf{a}_{t|t})' \mathbf{I}_{t|t} (\boldsymbol{\alpha}_t - \mathbf{a}_{t|t}), \quad (37)$$

$$\ell(\boldsymbol{\alpha}_t|\mathcal{F}_{t-1}) \approx \frac{1}{2} \log \det\{\mathbf{I}_{t|t-1}/(2\pi)\} - \frac{1}{2}(\boldsymbol{\alpha}_t - \mathbf{a}_{t|t-1})' \mathbf{I}_{t|t-1} (\boldsymbol{\alpha}_t - \mathbf{a}_{t|t-1}). \quad (38)$$

Here the state  $\boldsymbol{\alpha}_t$  is understood as a variable in  $\mathbb{R}^m$ , while  $\mathbf{a}_{t|t-1}$ ,  $\mathbf{a}_{t|t}$ ,  $\mathbf{I}_{t|t-1} \geq \mathbf{0}$  and  $\mathbf{I}_{t|t} \geq \mathbf{0}$  are known quantities determined by the Bellman filter in Table 2 or 3, depending on the context. If the model is linear and Gaussian, then the Bellman filter is exact (it is, in fact, the Kalman filter), as are equations (37)–(38). Based on approximations (37) and (38), the approximation of the realised KL divergence reads

$$\ell(\boldsymbol{\alpha}_t|\mathcal{F}_t) - \ell(\boldsymbol{\alpha}_t|\mathcal{F}_{t-1}) \Big|_{\boldsymbol{\alpha}_t=\mathbf{a}_{t|t}} \approx \frac{1}{2} \log \frac{\det(\mathbf{I}_{t|t})}{\det(\mathbf{I}_{t|t-1})} + \frac{1}{2}(\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1})' \mathbf{I}_{t|t-1} (\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1}), \quad (39)$$

where all constants involving  $\pi$  drop out. Nonnegativity of this quantity is guaranteed if  $\mathbf{I}_{t|t} \geq \mathbf{I}_{t|t-1}$ , which can be ensured in the implementation of the filter. Even when approximations (37)–(38) are somewhat inaccurate, it may be that the approximation of their difference in equation (39) is quite accurate. Intuitively, the realised KL divergence between two densities can be approximated to second order by considering the difference between both argmaxes and the sharpness of both peaks.

To define the proposed approximate maximum-likelihood estimator (MLE) for the static parameters, I take the usual definition  $\hat{\boldsymbol{\psi}} := \arg \max_{\boldsymbol{\psi}} \sum_t \ell(\mathbf{y}_t|\mathcal{F}_{t-1})$ . Then I substitute the (exact) decomposition (36) and the KL approximation (39), which gives

$$\hat{\boldsymbol{\psi}} := \arg \max_{\boldsymbol{\psi}} \sum_{t=t_0+1}^n \left\{ \underbrace{\ell(\mathbf{y}_t|\mathbf{a}_{t|t})}_{\text{'fit' of the filter}} - \underbrace{\left[ \frac{1}{2} \log \frac{\det(\mathbf{I}_{t|t})}{\det(\mathbf{I}_{t|t-1})} + \frac{1}{2}(\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1})' \mathbf{I}_{t|t-1} (\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1}) \right]}_{\geq 0, \text{ KL-type penalty}} \right\}, \quad (40)$$

where all terms on the right-hand side implicitly or explicitly depend on the (hyper)parameter  $\boldsymbol{\psi}$ . Time  $t_0 \geq 0$  is long enough to ensure the mode exists at time  $t_0$ . If model (12) is stationary and  $\boldsymbol{\alpha}_0$  is drawn from the unconditional distribution, as in the simulation studies in section 8, then  $t_0 = 0$ . The case  $t_0 > 0$  is analogous to that for the Kalman filter when the first  $t_0$  observations are used to construct a ‘proper’ prior (see Harvey, 1990, p. 123). The first term inside curly brackets, involving the observation density, is given by model (12). The remaining terms can be computed based on the output of the Bellman filter in Table 2 or 3. Expression (40) can be viewed as an alternative to the prediction-error decomposition for linear Gaussian state-space models (see e.g. Harvey, 1990, p. 126), the advantage being that estimator (40) remains applicable—albeit as an approximation—outside the classic context of linear Gaussian state-space models.

**Corollary 2** *Take the linear Gaussian state-space model specified in Corollary 1. Assume that the Kalman-filtered covariance matrices  $\{\mathbf{P}_{t|t}\}$  are positive definite. Estimator (40) then equals the MLE.*

Estimator (40) is only slightly more computationally demanding than static-parameter estimation using the Kalman filter. The sole source of additional computational complexity derives from the fact that the Bellman filter in Table 2 or 3 may perform several optimisation steps for each time step, while the Kalman filter performs only one. However, because each optimisation step is straightforward and few steps are typically required, the additional computational burden is negligible.

## 8 Simulation studies

### 8.1 Design

This section contains an extensive Monte Carlo study to investigate the performance of the Bellman filter for a range of data-generating processes (DGPs). I consider 10 DGPs with linear Gaussian state dynamics (12). (The empirical sections 9 and 10 consider high-dimensional and non-linear state dynamics, respectively.) The observation densities for this simulation study are listed in Supplement R, which also includes link functions, scores and other quantities used by the Bellman filter. To avoid selection bias, these DGPs have been taken from Koopman et al. (2016). While the numerically accelerated importance-sampling (NAIS) method in Koopman et al. (2015, 2016) has been shown to produce highly accurate results, the Bellman filter turns out to be equally (if not more) accurate at a fraction of the computational cost.

I add one DGP to the nine considered in Koopman et al. (2016): a local-level model with heavy-tailed observation noise. While a local-level model with additive Gaussian observation noise would be solved exactly by the Kalman filter, the latter does not adjust for heavy-tailed observation noise. Although the Kalman filter remains the best linear unbiased estimator of the state, the results below show that the (nonlinear) Bellman filter fares better.

The static (hyper)parameters for the first nine DGPs are taken from Koopman et al. (2016, Table 3). In particular, the state-transition equation (i.e.  $\alpha_t = c + T\alpha_{t-1} + \eta_t$  with  $\eta_t \sim N(0, \sigma_\eta^2)$ ) has parameters  $c = 0, T = \phi = 0.98$  and  $\sigma_\eta = 0.15$ , except for both dependence models, in which case  $c = 0.02, T = \phi = 0.98$  and  $\sigma_\eta = 0.10$ . In the observation densities (provided in Supplement R), the Student's  $t$  distributions have 10 degrees of freedom, i.e.  $\nu = 10$ , except for the local-level model, in which case  $\nu = 3$ . The remaining shape parameters are  $\kappa = 4$  for the negative binomial distribution,  $\kappa = 1.5$  for the Gamma distribution,  $\kappa = 1.2$  for the Weibull distribution and  $\sigma = 0.45$  for the local-level model.

For each of the 10 DGPs, I simulate 1,000 time series of length 5,000. I take the first 2,500 observations to represent the ‘in-sample’ period. For the purpose of static-parameter estimation, I use either (a) all 2,500 in-sample observations (long estimation window), (b) the last 1,000 in-sample observations (medium estimation window), or (c) the last 250 in-sample observations (short estimation window). Based on these parameter estimates, I run the Bellman filter and smoother in Table 3 on the entire dataset, including the out-of-sample period from  $t = 2,501$  through  $t = 5,000$ . For the Bellman filter, I also produce out-of-sample ‘smoothed’ state estimates  $a_{t|n}$  using parameters estimated from in-sample period, but including out-of-sample data for the purpose of smoothing.

I compute mean absolute errors (MAEs) and root mean squared errors (RMSEs) by comparing filtered and smoothed states against their true (simulated) counterparts.<sup>5</sup> For each DGP and each method, the reported average loss is based on  $2,500 \times 1,000 = 2.5$  million filtered states. I consider five methods:

1. **Infeasible mode estimator:** For filtering, I compute the mode using the true static parameters and a moving window of the most recent 250 observations; hence, 250 first-order conditions are solved for each time step (larger windows result in excessive computational times). The final state estimate  $a_{t|t}$  for each time  $t$  represents the filtered state. For smoothing, I use the mode estimator (3) based on the true parameters with  $t = n$  (i.e. based on the full sample).
2. **Bellman filter (BF):** The algorithm in Table 3 is initialised using the unconditional distribution.

---

<sup>5</sup>The Bellman filter, being based on the mode, is technically suboptimal for both loss functions.

Optimisation steps are performed until the estimated state is stable up to a tolerance of 0.0001 (on average,  $\sim 5$  iterations are needed). The logarithmic observation density is smooth and concave for the first seven DPGs, in which case optimisation (16) is strongly concave; quasi-Newton methods then quickly find the optimum (e.g. Nocedal and Wright, 2006). For simplicity, I pick Newton’s method which proved fast and stable. For the last three DPGs, the logarithmic observation density fails to be concave; in this case, I amend Newton’s method by replacing the Hessian of the logarithmic density by a weighted average of the Hessian and its expectation to ensure that the resulting expression is negative with probability one.<sup>6</sup> For these DPGs, the same weighting scheme ensures  $I_{t|t} \geq I_{t|t-1}$  as desired for the static-parameter estimator (40). Smoothed states are obtained as stated in Table 3.

3. **Particle filter (PF):** I follow Malik and Pitt’s (2011) implementation of the continuous sampling importance resampling (CSIR) particle filter, as it allows static parameters to be estimated using the same numerical optimisers employed for other methods. Experimentation suggests that using 1,000 particles is necessary to achieve a performance similar to that of the other methods. The seed that controls randomness is fixed beforehand, after which new random variates are drawn for each of the 1,000 time series; variations on this setup make no noticeable difference. The mean and the median of the particles at each time step are stored to compute RMSEs and MAEs, respectively.
4. **Numerically accelerated importance sampler (NAIS):** I follow Koopman et al. (2016), whose code is available online, deviating slightly by computing not only the weighted mean but also the weighted median of the (simulated) states. The resulting filtered states are used to compute RMSEs and MAEs, respectively.
5. **Kalman filter (KF):** I follow Ruiz (1994) and Harvey and Shephard (1996) in using quasi maximum-likelihood estimation (QMLE) to estimate the static parameters of both stochastic-volatility (SV) models. For both SV models, the observations are squared and taking the logarithm produces a linear state-space model, albeit with biased and non-Gaussian observation noise (for details, see Ruiz, 1994 or Harvey et al., 1994). For the local-level model with heavy-tailed observation noise, the Kalman filter is applied directly, i.e. without adjustments, and estimated by QMLE. For all three models, filtered and smoothed states are obtained, respectively, by the familiar Kalman filter and Rauch, Tung and Striebel smoother.

## 8.2 Results

This section compares (a) computational complexity, (b) quality of estimated (hyper)parameters, (c) quality of filtered and (d) smoothed state estimates, and (e) coverage (and length) of predicted, filtered and smoothed confidence intervals.

- a. **Computational complexity:** Table 4 shows average computation times (in seconds per sample) required for parameter estimation (based on the long estimation window) and filtering (based on all data) for three methods (BF, PF and NAIS). The BF is considerably faster than both simulation-based methods for the purposes of both parameter estimation and filtering. Compared to the NAIS

---

<sup>6</sup>For the dependence model with the Gaussian distribution, the weight placed on the expectation should weakly exceed  $1/2$ . For the Student’s  $t$  distribution, this generalises to  $1/2 \times (\nu + 4)/(\nu + 3)$ . For the local-level model with heavy-tailed noise, the weight given to the expectation should weakly exceed  $(1 + \nu/3)/(1 + 3\nu)$ .

Table 4: Average computing time (in seconds per sample) for parameter estimation and filtering

DGP Type	Distribution	Parameter estimation			Filtering		
		PF	NAIS	BF	NAIS	PF	BF
Count	Poisson	51	1.1	0.25	4.0	0.7	0.0024
Count	Negative binomial	146	3.1	0.64	5.2	1.0	0.0024
Intensity	Exponential	43	1.1	0.24	3.4	0.6	0.0022
Duration	Gamma	138	3.8	0.55	4.8	1.0	0.0026
Duration	Weibull	162	8.4	0.84	9.4	1.4	0.0060
Volatility	Gaussian	48	1.3	0.28	3.7	0.7	0.0023
Volatility	Student's $t$	95	2.7	0.70	5.2	1.0	0.0027
Dependence	Gaussian	69	2.4	0.57	5.5	0.8	0.0050
Dependence	Student's $t$	129	6.4	1.21	7.1	1.1	0.0060
Local level	Student's $t$	176	n/a	1.01	n/a	0.9	0.0029

*Note:* BF = Bellman filter. PF = particle filter. NAIS = numerically accelerated importance sampler. Computation times are measured on a computer running 64-bit Windows 8.1 Pro with an Intel(R) Core(TM) i7-4810MQ CPU @ 2.80GHz. Average parameter estimation times are based on the first 2,500 observations across 1,000 repetitions for each DGP. Average filtering times are based on filtering the entire sample of 5,000 observations across 1,000 repetitions for each DGP.

method, parameter estimation by the BF is faster by a factor 4 to 10, while filtering is faster by a factor between  $\sim 1,000$  and  $\sim 2,000$ . Compared to the PF, parameter estimation by the BF is faster by a factor between  $\sim 100$  and  $\sim 250$ , while filtering is faster by a factor between  $\sim 160$  and  $\sim 400$ .

- b. **(Hyper)parameter estimates:** Table 5 displays average (hyper)parameter estimates and root mean squared errors (RMSEs) versus the true parameters for three methods (BF, PF and NAIS) for the long estimation window. Parameter estimates for the short and medium windows are presented in Supplement S. The BF is about as accurate as both simulation-based methods for all three window sizes in terms of both average parameters and RMSEs relative to the true parameters. The average parameters are close to the true values and tend to be drawn even closer as the estimation window is increased, while the RMSEs decrease rapidly. These simulation results suggest that, for these models and sample sizes, any potential bias or loss of efficiency compared to the simulation-based methods under investigation is negligible.
- c. **Filtered state estimates:** Table 6 shows mean absolute errors (MAEs) of filtered states in the out-of-sample period, reported relative to the MAEs of the infeasible mode estimator, for four methods: BF, PF, NAIS and KF. The infeasible estimator uses true parameters and the same information set as the filtering methods. The main finding is that the BF, PF and NAIS perform near identically, while the KF, when applicable, lags substantially behind.<sup>7</sup> The out-of-sample performance of the BF based on the long estimation window falls within  $\sim 2\%$  of that of the infeasible state estimator across all DGPs. For this estimation window, the BF marginally outperforms the PF and NAIS for three DGPs (for the Poisson, negative binomial and exponential distributions). It performs on par with both these methods for four DGPs (with the Gamma/Weibull distributions and for the Gaussian volatility and Student's  $t$  dependence models), but is marginally outperformed for three DGPs (for the Student's  $t$  volatility, Gaussian dependence and local-level models), albeit by max  $\sim 0.3\%$ . Filtering results deteriorate by a few percentage points for the medium estimation window,

<sup>7</sup>This difference is not due to the choice of loss function; the relative performance of the KF deteriorates further when reporting RMSEs (see Supplement T).

Table 5: Average parameter estimates and RMSEs based on the long estimation window

DGP Type	Distribution	Truth		BF		PF		NAIS	
				Average	RMSE	Average	RMSE	Average	RMSE
Count	Poisson	$c$	0.000	-0.007	[0.008]	0.000	[0.003]	0.000	[0.003]
		$\phi$	0.980	0.977	[0.007]	0.978	[0.006]	0.978	[0.006]
		$\sigma_\eta$	0.150	0.153	[0.014]	0.152	[0.014]	0.149	[0.013]
Count	Negative Bin.	$c$	0.000	-0.004	[0.005]	0.000	[0.003]	0.000	[0.003]
		$\phi$	0.980	0.979	[0.006]	0.977	[0.007]	0.979	[0.006]
		$\sigma_\eta$	0.150	0.149	[0.015]	0.152	[0.016]	0.145	[0.015]
		$1/\kappa$	0.250	0.239	[0.036]	0.248	[0.031]	0.287	[0.049]
Intensity	Exponential	$c$	0.000	-0.007	[0.008]	0.000	[0.003]	0.000	[0.003]
		$\phi$	0.980	0.976	[0.008]	0.978	[0.007]	0.978	[0.007]
		$\sigma_\eta$	0.150	0.158	[0.017]	0.151	[0.014]	0.151	[0.014]
Duration	Gamma	$c$	0.000	0.007	[0.008]	0.000	[0.004]	0.000	[0.004]
		$\phi$	0.980	0.976	[0.007]	0.977	[0.006]	0.977	[0.006]
		$\sigma_\eta$	0.150	0.158	[0.015]	0.152	[0.013]	0.152	[0.013]
		$\kappa$	1.500	1.507	[0.043]	1.501	[0.043]	1.501	[0.043]
Duration	Weibull	$c$	0.000	0.009	[0.010]	0.000	[0.003]	0.000	[0.003]
		$\phi$	0.980	0.975	[0.008]	0.978	[0.006]	0.978	[0.006]
		$\sigma_\eta$	0.150	0.160	[0.018]	0.152	[0.013]	0.152	[0.013]
		$\kappa$	1.200	1.207	[0.023]	1.200	[0.021]	1.200	[0.021]
Volatility	Gaussian	$c$	0.000	0.007	[0.008]	0.000	[0.004]	0.000	[0.004]
		$\phi$	0.980	0.975	[0.010]	0.977	[0.008]	0.977	[0.008]
		$\sigma_\eta$	0.150	0.166	[0.026]	0.152	[0.018]	0.152	[0.018]
Volatility	Student's $t$	$c$	0.000	0.005	[0.006]	0.000	[0.004]	0.000	[0.004]
		$\phi$	0.980	0.975	[0.010]	0.977	[0.008]	0.977	[0.008]
		$\sigma_\eta$	0.150	0.162	[0.031]	0.153	[0.021]	0.153	[0.022]
		$1/\nu$	0.100	0.089	[0.030]	0.100	[0.010]	0.097	[0.023]
Dependence	Gaussian	$c$	0.020	0.021	[0.009]	0.024	[0.011]	0.024	[0.011]
		$\phi$	0.980	0.979	[0.008]	0.977	[0.010]	0.977	[0.010]
		$\sigma_\eta$	0.100	0.095	[0.020]	0.103	[0.024]	0.103	[0.024]
Dependence	Student's $t$	$c$	0.020	0.022	[0.010]	0.025	[0.013]	0.025	[0.014]
		$\phi$	0.980	0.977	[0.010]	0.975	[0.013]	0.975	[0.014]
		$\sigma_\eta$	0.100	0.098	[0.023]	0.106	[0.029]	0.107	[0.030]
		$1/\nu$	0.100	0.103	[0.012]	0.100	[0.006]	0.098	[0.025]
Level	Student's $t$	$c$	0.000	0.000	[0.004]	0.000	[0.003]		
		$\phi$	0.980	0.979	[0.005]	0.978	[0.005]		
		$\sigma_\eta$	0.150	0.139	[0.013]	0.151	[0.008]		
		$\sigma$	0.450	0.453	[0.025]	0.451	[0.027]		
		$1/\nu$	0.333	0.277	[0.066]	0.332	[0.024]		

*Note:* BF = Bellman filter. PF = Particle filter. NAIS = Numerically accelerated importance sampler. RMSE = root mean squared error. I simulated 1,000 time series each of length 5,000 for 10 data-generating processes with linear Gaussian state dynamics (12), i.e.  $\alpha_{t+1} = c + \phi\alpha_t + \eta_{t+1}$  with  $\eta_{t+1} \sim N(0, \sigma_\eta^2)$ . The observation densities are listed in Supplement R. The estimation of static parameters is based on the long estimation window, which consists of 2,500 observations. Parameter estimation is performed as follows: Bellman filter: based on estimator (40); Particle filter: as in Malik and Pitt (2011); Importance sampler: as in Koopman et al. (2015, 2016).

and by  $\sim 10\text{--}30\%$  for the short estimation window, in particular for both dependence models. Even for the short estimation window, the results for the BF, PF and NAIS are virtually identical with the KF lagging behind. The robustness of the BF means that it compares favourably with the KF for both the SV and local-level models: e.g. for the local-level model, the maximum absolute error in the out-of-sample period, averaged across 1,000 samples, is 1.80 for the KF; double that for the BF

Table 6: MAEs of filtered states in out-of-sample period

DGP		Infeasible estimator	Short estimation window (250 obs.)				Medium estimation window (1,000 obs.)				Long estimation window (2,500 obs.)			
			BF	PF	NAIS	KF	BF	PF	NAIS	KF	BF	PF	NAIS	KF
Type	Distribution	MAE	Relative MAE				Relative MAE				Relative MAE			
Count	Poisson	0.283	1.145	1.141	1.140		1.015	1.015	1.016		1.001	1.002	1.003	
Count	Neg. Bin.	0.300	1.159	1.154	1.155		1.018	1.019	1.020		1.005	1.006	1.007	
Intensity	Exponential	0.286	1.128	1.130	1.128		1.013	1.014	1.014		1.002	1.003	1.003	
Duration	Gamma	0.259	1.158	1.156	1.154		1.023	1.024	1.023		1.007	1.007	1.007	
Duration	Weibull	0.264	1.117	1.115	1.114		1.012	1.012	1.012		1.001	1.001	1.001	
Volatility	Gaussian	0.337	1.198	1.200	1.200	1.473	1.023	1.023	1.023	1.230	1.005	1.005	1.005	1.230
Volatility	Student's $t$	0.352	1.231	1.213	1.217	1.574	1.038	1.029	1.030	1.336	1.012	1.009	1.010	1.275
Dependence	Gaussian	0.288	1.291	1.296	1.290		1.056	1.056	1.055		1.018	1.016	1.016	
Dependence	Student's $t$	0.295	1.301	1.313	1.291		1.063	1.065	1.067		1.022	1.022	1.022	
Level	Student's $t$	0.159	1.059	1.045		1.196	1.014	1.004		1.128	1.003	1.000	1.122	

*Note:* MAE = mean absolute error. BF = Bellman filter. PF = particle filter. NAIS = numerically accelerated importance sampler. KF = Kalman filter. I simulated 1,000 time series each of length 5,000 for 10 data-generating processes of type (12); the observation densities are listed in Supplement R. The data is split in an ‘in-sample’ period (first 2,500 observations) and an ‘out-of-sample’ period (last 2,500 observations). The short, medium and long estimation windows consist of the 250, 1,000 or 2,500 observations, respectively, of the in-sample period. Filtered states based on simulation-based methods (importance sampler and particle filter) are computed by taking the median of the simulated states. In all cases, MAEs are computed by comparing the last 2,500 filtered states with their true (simulated) counterparts. MAEs are reported relative to the MAE of the infeasible mode estimator.

Table 7: MAEs of smoothed states in out-of-sample period

DGP		Infeasible estimator	Short estimation window (250 obs.)		Medium estimation window (1,000 obs.)		Long estimation window (2,500 obs.)	
			BF	KF	BF	KF	BF	KF
Type	Distribution	MAE	Relative MAE		Relative MAE		Relative MAE	
Count	Poisson	0.222	1.118		1.020		1.013	
Count	Neg. Bin.	0.236	1.139		1.018		1.009	
Intensity	Exponential	0.222	1.099		1.021		1.016	
Duration	Gamma	0.201	1.168		1.040		1.024	
Duration	Weibull	0.204	1.096		1.026		1.021	
Volatility	Gaussian	0.266	1.196	1.628	1.033	1.259	1.022	1.221
Volatility	Student's $t$	0.280	1.247	2.156	1.047	1.433	1.024	1.366
Dependence	Gaussian	0.240	1.359		1.056		1.018	
Dependence	Student's $t$	0.247	1.379		1.064		1.021	
Level	Student's $t$	0.126	1.035	1.154	1.017	1.131	1.015	1.129

*Note:* For the simulation setting, see the note to Table 6. For the SV models, the static parameters in the Kalman filter are estimated by QMLE as in Ruiz (1994), after which the RTS smoother is applied (Rauch et al., 1965). MAEs are reported relative to the MAE of the infeasible estimator (3).

(0.90). The BF is thus more robust in the face of heavy-tailed observation noise, while having only a single additional parameter to estimate (the degrees of freedom of the observation noise,  $\nu$ ).

- d. **Smoothed state estimates:** Table 7 shows the MAEs of smoothed states in the out-of-sample period obtained by the Bellman filter/smoothing combination in Table 3, where the static parameters are estimated based on three different in-sample estimation windows. The results are reported relative to those of the infeasible state estimator (3) with  $t = n$ , which similarly exploits all data and uses the true parameters. Where appropriate, results are also reported for the Kalman filter/smoothing. The performance of the Bellman filter/smoothing using the long estimation window lies within  $\sim 2\%$  of

Table 8: Coverage (in %) and average length (in square brackets) of Bellman-predicted, -filtered and -smoothed confidence intervals for different parameter-estimation windows

DGP Type	Distribution	Short estimation window (250 obs.)			Medium estimation window (1,000 obs.)			Long estimation window (2,500 obs.)		
		Predict	Filter	Smooth	Predict	Filter	Smooth	Predict	Filter	Smooth
Count	Poisson	90.2	90.6	92.5	94.7	94.8	94.7	95.2	95.3	94.9
		[1.52]	[1.41]	[1.17]	[1.51]	[1.41]	[1.11]	[1.51]	[1.41]	[1.11]
Count	Neg. Bin.	89.5	89.7	91.7	94.3	94.3	94.3	94.9	94.9	94.6
		[1.61]	[1.50]	[1.24]	[1.57]	[1.48]	[1.16]	[1.57]	[1.48]	[1.16]
Intensity	Exponential	90.8	91.1	93.4	95.4	95.4	95.5	95.8	95.8	95.5
		[1.56]	[1.46]	[1.20]	[1.57]	[1.47]	[1.16]	[1.57]	[1.47]	[1.15]
Duration	Gamma	90.8	90.9	92.1	95.2	95.2	94.9	95.7	95.7	95.3
		[1.43]	[1.31]	[1.06]	[1.44]	[1.32]	[1.04]	[1.44]	[1.33]	[1.03]
Duration	Weibull	92.4	92.6	94.3	95.6	95.6	95.5	96.0	95.9	95.5
		[1.50]	[1.37]	[1.12]	[1.48]	[1.36]	[1.07]	[1.48]	[1.36]	[1.06]
Volatility	Gaussian	88.1	88.4	90.8	95.3	95.3	95.5	96.1	96.0	95.8
		[1.81]	[1.73]	[1.47]	[1.84]	[1.76]	[1.42]	[1.84]	[1.77]	[1.41]
Volatility	Student's $t$	88.4	88.4	90.5	94.5	94.5	94.7	95.4	95.3	95.2
		[1.98]	[1.87]	[1.61]	[1.88]	[1.81]	[1.46]	[1.87]	[1.80]	[1.44]
Dependence	Gaussian	73.9	74.0	75.7	90.5	90.6	91.2	93.1	93.1	93.1
		[1.26]	[1.23]	[1.10]	[1.37]	[1.34]	[1.14]	[1.39]	[1.36]	[1.14]
Dependence	Student's $t$	71.9	71.9	73.5	90.4	90.4	91.2	93.0	93.1	93.4
		[1.28]	[1.25]	[1.13]	[1.42]	[1.40]	[1.20]	[1.43]	[1.41]	[1.19]
Level	Student's $t$	93.1	93.5	94.7	94.9	95.0	95.2	95.1	95.1	95.3
		[0.98]	[0.80]	[0.65]	[0.99]	[0.81]	[0.64]	[0.99]	[0.81]	[0.64]

Note: For the simulation setting, see the note to Table 6.

that of the infeasible state estimator across all DGPs. The performance compared with the filtering results in Table 6 is improved by  $\sim 20\%$ . This shows that smoothing has substantial benefits, which the Bellman filter/smoothing successfully exploits. The KF smoothing results are comparatively poor, especially for the short estimation window. Neither Malik and Pitt (2011) nor Koopman et al. (2016) present smoothing methods; hence, no PF or NAIS smoothing results are reported.

- e. **Coverage of confidence intervals:** Table 8 shows the coverage of approximate Bellman-predicted, -filtered and -smoothed confidence intervals with endpoints given by  $a_{t|t-1} \pm 2/\sqrt{I_{t|t-1}}$ ,  $a_{t|t} \pm 2/\sqrt{I_{t|t}}$  and  $a_{t|n} \pm 2/\sqrt{I_{t|n}}$ , respectively, as well as the average length of these intervals, where the estimation of static parameters is based on three possible window sizes. These confidence intervals are based on the quadratic approximation of the value function and are analogous to those in the Kalman filter. For brevity, both simulation-based approaches are excluded. The Bellman-predicted, -filtered and -smoothed confidence intervals based on the medium and long estimation windows tend to be fairly accurate, containing the true states  $\sim 93 - 96\%$  of the time for most DGPs and  $\sim 90 - 96\%$  for both dependence models. Confidence intervals based on the short estimation window tend to be overly optimistic, especially for the two dependence models. Finally, the length of confidence intervals based on the smoothed states is substantially reduced, while the coverage remains good for the medium and long estimation windows, further highlighting the benefits of smoothing.

## 9 Application I: High-dimensional state space

This section considers the modelling of high-dimensional cloud-intensity data from a regional climate model as in Katzfuss et al. (2020). In a simulation study with realistic parameter values, I demonstrate that the performance of the Bellman filter is unaffected as the dimension of the state increases from 10 to 150, while the performance of the standard (bootstrap) particle filter deteriorates sharply—even when using very many particles. When predicting real data, I show that the Bellman filter substantially outperforms the particle-ensemble Kalman filter in Katzfuss et al. (2020) and the exact approximation of the Rao-Blackwellised particle filter in Johansen et al. (2012).

### 9.1 Model

Following Katzfuss et al. (2020, p. 868), I consider a multivariate overdispersed Poisson density that generates an integer number of clouds recorded at adjacent locations over a period of time, in combination with a linear Gaussian state equation for the logarithmic cloud intensities. The model for  $t = 1, \dots, n$  reads

$$\mathbf{y}_t \sim \text{Poisson}(\exp(\boldsymbol{\beta}_t)), \quad \mathbf{y}_t \in \mathbb{N}^m, \boldsymbol{\beta}_t \in \mathbb{R}^m, \quad (41)$$

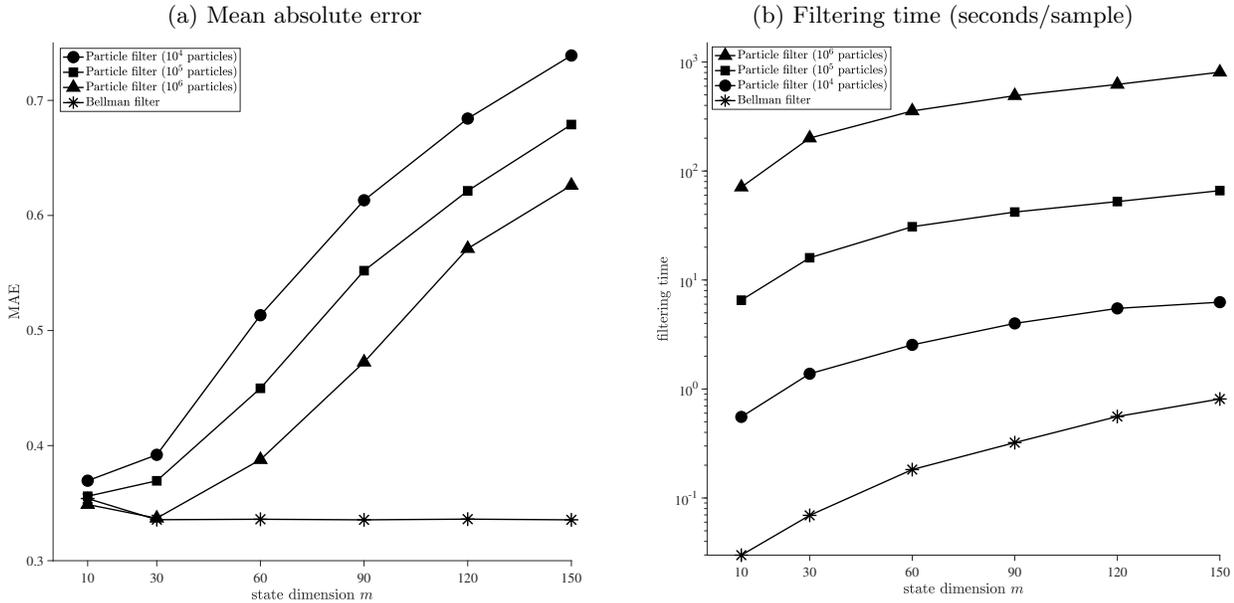
$$\boldsymbol{\beta}_t = \boldsymbol{\alpha}_t + \boldsymbol{\xi}_t, \quad \boldsymbol{\xi}_t \sim \text{i.i.d. } \mathcal{N}(\mathbf{0}_m, \sigma_\xi^2 \mathbb{1}_{m \times m}), \quad (42)$$

$$\boldsymbol{\alpha}_t = (\mathbb{1}_{m \times m} - \mathbf{T}) \mathbf{c} + \mathbf{T} \boldsymbol{\alpha}_{t-1} + \boldsymbol{\eta}_t, \quad \boldsymbol{\eta}_t \sim \text{i.i.d. } \mathcal{N}(\mathbf{0}_m, \mathbf{Q}), \quad (43)$$

where  $\boldsymbol{\alpha}_t \in \mathbb{R}^m$  is the latent state,  $\boldsymbol{\beta}_t \in \mathbb{R}^m$  is an overdispersed (i.e. noisy) realisation of  $\boldsymbol{\alpha}_t$  with overdispersion parameter  $\sigma_\xi \geq 0$ , and  $\mathbf{y}_t \in \mathbb{N}^m$  is a vector of  $m$  Poisson-generated counts with corresponding intensities  $\exp(\boldsymbol{\beta}_t)$ . The exponent of a vector in equation (41) is understood elementwise, i.e. observation  $y_{i,t}$  is drawn independently from a Poisson density with intensity  $\exp(\beta_{i,t})$  for each  $i = 1, \dots, m$ . When  $\sigma_\xi = 0$ , such that  $\boldsymbol{\alpha}_t = \boldsymbol{\beta}_t$  for all  $t$ , the model collapses to a standard state-space model with state vector  $\boldsymbol{\alpha}_t$  of length  $m$ . For  $\sigma_\xi > 0$ , the hierarchical structure (41)–(43) can be cast in the standard state-space format as I show below, where the dimension of the state is  $2m$ . Models with  $\sigma_\xi = 0$  and  $\sigma_\xi > 0$  are referred to as the ‘standard’ and ‘overdispersed’ versions of the model, respectively.

The system vectors and matrices in the state-transition equation are  $\mathbf{c} \in \mathbb{R}^m$  and  $\mathbf{T}, \mathbf{Q} \in \mathbb{R}^{m \times m}$ . Following Katzfuss et al. (2020), I assume that  $\mathbf{T}$  is tridiagonal with  $\gamma_1$  on the main diagonal,  $\gamma_2$  above the main diagonal, and  $\gamma_3$  below the main diagonal. Intuitively, these parameters govern the probability of cloud intensities staying in place or drifting left or right. As in Katzfuss et al. (2020), I assume new cloud formation to be more highly correlated at shorter distances. Specifically, the covariance matrix  $\mathbf{Q}$  is assumed to be a spatial Matèrn covariance matrix, with a smoothness of 1.5, spatial dependence parameter  $\lambda > 0$ , and overall scale governed by  $\tau > 0$ , i.e.  $(\mathbf{Q})_{ij} = \tau^2(1 + \sqrt{3}|i - j|/\lambda) \exp(-\sqrt{3}|i - j|/\lambda)$  for  $i, j = 1, \dots, m$ . While Katzfuss et al. (2020) set  $\mathbf{c} = \mathbf{0}_m$ , I consider the more general case  $\mathbf{c} \neq \mathbf{0}_m$ , where  $\mathbf{c}$  can be interpreted as the long-run average of  $\boldsymbol{\alpha}_t$  if the eigenvalues of  $\mathbf{T}$  lie inside the unit circle. For simplicity I set  $\mathbf{c} = c\mathbb{1}_m$ , where a single parameter  $c \in \mathbb{R}$  controls the overall level. Static parameters are collected in the vector  $\boldsymbol{\psi} = (c, \gamma_1, \gamma_2, \gamma_3, \tau, \lambda, \sigma_\xi)'$ .

Figure 2: MAE of filtered states and filtering times (in seconds per sample)



Note: MAE = mean absolute error. I simulated 100 instances of the model (41)–(43) with  $n = 80$  time steps and static parameters  $\psi = (0, 0.4, 0, 0.4, 0.8, 5, 0)'$  for various values of the state dimension  $m$ . Using the true static parameter for the purpose of filtering, I recorded the MAE of the filtered states  $\mathbf{a}_{t|t}$  relative to the true (simulated) states  $\mathbf{\alpha}_t$  and runtime in seconds per sample for the Bellman filter and particle filter, where the latter was implemented with  $10^4$ ,  $10^5$  and  $10^6$  particles.

## 9.2 State-space formulation and Bellman-filter implementation

For  $\sigma_\xi > 0$ , a standard state-space model can be obtained by writing the dynamics of  $\mathbf{\alpha}_t$  and  $\beta_t$  jointly as

$$\begin{bmatrix} \beta_t \\ \mathbf{\alpha}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{0}_m \\ (\mathbf{1}_{m \times m} - \mathbf{T})\mathbf{c} \end{bmatrix} + \begin{bmatrix} \mathbf{0}_{m \times m} & \mathbf{1}_{m \times m} \\ \mathbf{0}_{m \times m} & \mathbf{T} \end{bmatrix} \begin{bmatrix} \beta_{t-1} \\ \mathbf{\alpha}_t \end{bmatrix} + \begin{bmatrix} \xi_t \\ \eta_{t+1} \end{bmatrix}, \quad (44)$$

where  $\{\xi_t\}$  and  $\{\eta_t\}$  are series of i.i.d. disturbances with characteristics specified in equations (42)–(43). The state vector in the overdispersed model is  $(\beta'_t, \mathbf{\alpha}'_{t+1})' \in \mathbb{R}^{2m}$ , which is 120-dimensional when  $m = 60$  (as in Katzfuss et al., 2020). The Bellman filter in Table 3 is directly applicable after appropriate redefinitions; e.g.  $\mathbf{c}$  in Table 3 should be identified with the first vector on the right-hand side of equation (44).

The Bellman filter solves a high-dimensional optimisation problem at each time step. The logarithmic Poisson density is jointly concave in all elements of  $\beta_t$ . The Bellman-filtered state in equation (16) then is unique; it can typically be found using e.g. Newton steps. To avoid the need for repeated large-matrix inversions, however, I opted for the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm (e.g. Nocedal and Wright, 2006, §6.1), which proved both fast and stable. Indeed, at the estimated parameter values, executing the Bellman filter for the standard (overdispersed) model using data from Katzfuss et al. (2020), involving a 60-dimensional (120-dimensional) optimisation problem for each of 80 time steps, takes about  $\sim 0.25$  ( $\sim 0.60$ ) seconds. In both cases, convergence with a tolerance of  $10^{-5}$  at each time step is reached within  $\sim 12$  BFGS optimisation steps.

### 9.3 Simulation study with high-dimensional state space

This section investigates the performance of the Bellman filter in high-dimensional state spaces by performing a simulation study for the model (41)–(43) with varying spatial dimension  $m$ . I compare the Bellman filter’s performance against that of the standard (bootstrap) particle filter. For simplicity, the static parameter  $\psi$  is considered known and taken as  $\psi = (c, \gamma_1, \gamma_2, \gamma_3, \tau, \lambda, \sigma_\xi)' = (0, 0.4, 0, 0.4, 0.8, 5, 0)'$ , which is similar to the empirical parameter estimates obtained from real data. As in the real data, the relatively large value of  $\gamma_3 = 0.4$  reflects the fact that logarithmic cloud intensities tend to float from lower to higher location numbers, which may be due to a fixed wind direction during the observation period. The overdispersion parameter  $\sigma_\xi$  is set to zero, as my empirical study contains no evidence to suggest otherwise. For  $\sigma_\xi = 0$ , the state-augmentation procedure (44) is not required; hence, the dimension of the state space is simply  $m$ . I investigate cases where  $m$  equals 10, 30, 60, 90, 120 or 150, thus exploring different spatial dimensions beyond that of the real data set considered in Katzfuss et al. (2020), where  $m = 60$ . For each  $m$ , I simulate 100 datasets with 80 time steps, matching the time dimension of the real data.

The particle filter is subject to the curse of dimensionality and may struggle in higher dimensions (e.g. Surace et al., 2019). Hence, I experiment with  $10^4$ ,  $10^5$  and  $10^6$  particles; increasing this number further turns out to be computationally infeasible (see further discussion below). I compute the median of the particles as the filtered state. For both methods, mean absolute errors (MAEs) of filtered states are computed by taking the one-norm of the vector  $\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t \in \mathbb{R}^m$ , dividing this norm by  $m$ , and averaging the resulting quantity across 80 time steps and 100 simulated data sets.

Figure 2 (Panel A) shows that the MAE of the Bellman filter is almost entirely flat at  $\sim 0.34$ , independently of the dimension  $m$ . In fact, the MAE appears to improve slightly as the dimension  $m$  increases, possibly because the filter benefits from improved predictions: cloud observations even in distant locations may, due to wind conditions, be informative as to the possible future presence of clouds at other locations. In contrast, the MAE of the particle filter increases sharply with  $m$  and substantially exceeds that of the Bellman filter even at  $m = 60$  or  $m = 90$ . This heightened inaccuracy in higher dimensions materialises for any (fixed) number of particles. Even with  $10^6$  particles, the particle filter at  $m = 150$  produces an MAE of  $\sim 0.63$ , a factor  $\sim 1.8$  higher than that of the Bellman filter.

Figure 2 (Panel B) shows that using  $10^6$  particles in  $m = 150$  dimensions necessitates a filtering time of  $\sim 800$  seconds per simulated dataset, such that the total runtime for the particle filter across 100 simulations is  $100 \times 800$  seconds =  $\sim 22$  hours. The BFGS implementation of the Bellman filter required between 0.03 seconds (for  $m = 10$ ) and 0.80 seconds (for  $m = 150$ ), translating in the latter case to a total runtime across 100 simulations of only  $\sim 1.3$  minutes. Panel B also shows that the computational complexity of the particle filter scales with the number of particles employed: for  $10^6$  particles, the difference with the Bellman filter is around three orders of magnitude for any  $m$ . The relative accuracy and speed of the Bellman filter as demonstrated in this section can largely be attributed to its approach to optimisation, which is simpler than the sampling/integration approach used in the particle filter—especially in higher dimensions.

### 9.4 Real-data application with artificially missing data

For the real-data application, I take the cloud-motion data investigated by Katzfuss et al. (2020), which contains  $m = 60$  locations along a spatial transect (i.e. a line), where the number of visible clouds is

Table 9: Full-sample-with-missing-data parameter estimates for model (41)–(43)

	$c$	$\gamma_1$	$\gamma_2$	$\gamma_3$	$\tau$	$\lambda$	$\sigma_\xi$	MSE	CRPS
Standard model	-3.656 [0.242]	0.254 [0.053]	0.050 [0.040]	0.372 [0.056]	1.749 [0.100]	7.040 [0.471]		0.513	0.185
Standard model ( $c = 0$ )		0.260 [0.060]	0.127 [0.047]	0.482 [0.055]	1.771 [0.108]	8.295 [0.561]		0.547	0.192
Overdispersed model	-4.236 [0.072]	0.245 [0.025]	0.055 [0.033]	0.384 [0.027]	1.839 [0.053]	7.249 [0.053]	0.000 [0.018]	0.509	0.185
Overdispersed model ( $c = 0$ )		0.230 [0.055]	0.142 [0.045]	0.494 [0.047]	1.791 [0.102]	8.301 [0.346]	0.000 [0.035]	0.556	0.197

*Note:* MSE = mean squared error. CRPS = continuously ranked probability score. The standard model has  $\sigma_\xi = 0$ , while the overdispersed model has  $\sigma_\xi > 0$ . Numerical standard errors in square brackets are computed by taking the square root of diagonal elements of the inverse of the negative finite-difference Hessian matrix. Using the output of the Bellman filter at times and locations where observations were declared missing, I produce ‘nowcasts’ of missing data, the quality of which can be judged on the basis of MSE and CRPS values in the right-most columns.

recorded at each of  $n = 80$  time steps. Following their procedure, I artificially introduce ‘missing data’ by assuming that at each time step only 90% of the locations, i.e. 54 randomly selected locations, deliver a measurement that the researcher can use for parameter estimation and state filtering. The remaining  $80 \times 6 = 480$  observations are declared ‘missing’, but remain available for testing. For reproducibility, the same missing data are considered as in Katzfuss et al. (2020), whose code is available online. The aim is to ‘nowcast’ the (same) missing data by running the Bellman filter on the available data.

To implement the Bellman filter with missing data, I write the logarithm of the observation density at time  $t$  used in the Bellman-filter update (16) as

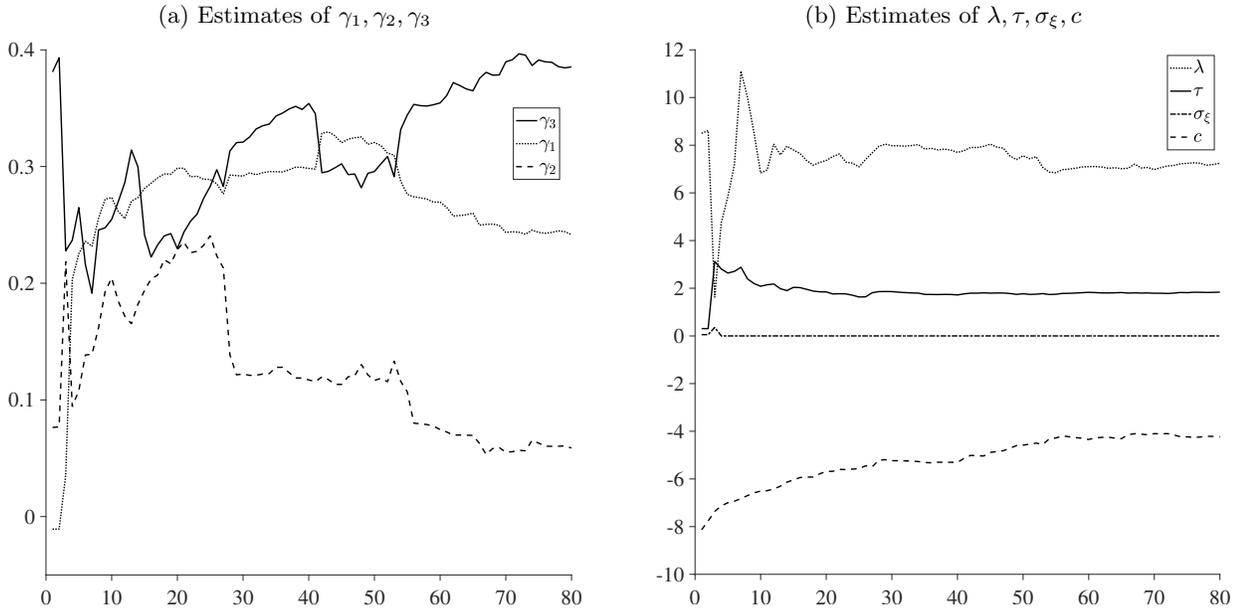
$$\log \text{Poisson}(\mathbf{y}_t | \exp(\boldsymbol{\beta}_t)) = \sum_{i \in \mathcal{O}_t} \log \text{Poisson}(y_{i,t} | \exp(\beta_{i,t})), \quad (45)$$

where  $\mathcal{O}_t$  is the set of available observations at time  $t$ ; i.e. log-likelihood contributions of missing data are excluded. The Bellman filter in Table 3 remains applicable as long as the score and (realised) information quantities are computed by taking derivatives of the logarithmic density on the right-hand side of equation (45). This implies that elements of the score vector corresponding to missing observations are set to zero. Nevertheless, the Bellman-filtered states at times and locations for which observations are declared missing remain non-trivial, because the filtered state—representing the solution to an optimisation problem—is affected by *all* available observations at a given time step. The Bellman filter in Table 3 is initialised with  $\mathbf{I}_{1|0}$  equal to a small multiple of the identity. The static parameter  $\boldsymbol{\psi}$  is estimated using the approximate maximum-likelihood estimator (40), employing equation (45) to exclude data declared missing.

## 9.5 Results: Full sample with missing data

Table 9 contains the resulting parameter estimates for various model specifications, where the parameter-estimation procedure used all data deemed available. Consistent with Katzfuss et al. (2020), in all specifications the relatively large estimate of  $\gamma_3$  picks up the drift of clouds along the spatial transect, indicating

Figure 3: Expanding-window parameter estimation results for model (41)–(43)



Note: Parameters estimated by an expanding window using cloud data from Katzfuss et al. (2020).

that clouds tend to float from lower to higher location numbers. While Katzfuss et al. (2020) investigated only the overdispersed model, our comparison of the overdispersed model and the standard model yields no evidence that the former is preferable to the latter: estimates of the overdispersion parameter  $\sigma_\xi$  are practically zero. On the other hand, the inclusion of an additional parameter  $c$  governing the overall level appears to be beneficial.

Running the Bellman filter on the entire sample with missing data produces filtered states at times and locations for which observations were declared missing. By taking the exponent, a filtered state translates to an intensity, which in turn equals the expected value of a draw from the relevant Poisson distribution. This allows us to produce both point and density ‘nowcasts’ of missing data conditional on the available data up to and including the relevant time step. Following Katzfuss et al. (2020), these point and density nowcasts can be compared with the actual observations using the mean squared error (MSE) and continuously ranked probability score (CRPS), respectively, which are reported in the right-most columns of Table 9. Depending on the model specification, the MSEs of the Bellman filter lie in the range  $\sim 0.51$ – $0.56$ , the CRPS in  $\sim 0.18$ – $0.20$ . These numbers are not (yet) directly comparable with those in Katzfuss et al. (2020), who use an expanding window for the purpose of parameter estimation. This is addressed in the next section.

## 9.6 Results: Expanding window with missing data

The highly parametrised model (41)–(43) allows us to estimate the static parameters in an expanding-window-with-missing-data setting, starting with a window of one time step. For the most general (i.e. overdispersed) version of model, Figure 3 shows the parameter estimates over time. At the end of the sample, the parameter estimates match the results in Table 9. For all time steps, the estimate of  $\sigma_\xi$  is practically zero. After some variation at the start of the sample, the estimates of  $\lambda, \tau$  and  $c$  converge relatively quickly. The estimates of  $\gamma_1, \gamma_2, \gamma_3$ , however, show considerable time variation even towards the

Table 10: Quality of nowcasts using an expanding window for parameter estimation and filtering

Model	Method	MSE	CRPS
Overdispersed ( $\sigma_\xi > 0$ )	Rao-Blackwellised particle filter ( $c = 0$ , Johansen et al., 2012)	1.26	0.33
	Particle ensemble Kalman filter ( $c = 0$ , Katzfuss et al., 2020)	0.75	0.25
	Bellman filter ( $c = 0$ )	0.554	0.194
	Bellman filter ( $c \neq 0$ )	0.519	0.188
Standard ( $\sigma_\xi = 0$ )	Bellman filter ( $c = 0$ )	0.556	0.196
	Bellman filter ( $c \neq 0$ )	0.525	0.190

*Note:* MSE = mean squared error. CRPS = continuously ranked probability score. The data (including the classification of training and test data) are available from Katzfuss et al. (2020). The first two rows are copied from Katzfuss et al. (2020), who consider only the overdispersed model with  $c = 0$ .

end of the sample, indicating that these parameters may not in fact be static. This may explain why the expanding-window results, discussed below, appear to be no worse than the full-sample results.

For the purpose of nowcasting missing data, Table 10 shows that both the standard ( $\sigma_\xi = 0$ ) and overdispersed ( $\sigma_\xi > 0$ ) versions of the model with  $c \neq 0$  achieve MSEs of  $\sim 0.52$ , with the particle ensemble Kalman filter and Rao-Blackwellised particle filter lagging behind by  $\sim 45\%$  and  $\sim 140\%$ , respectively. Irrespective of the exact specification, the Bellman filter achieves CRPS values of  $\sim 0.19$ , with the corresponding numbers for both particle-filtering methods inflated by  $\sim 30\%$  and  $\sim 75\%$ . This demonstrates that Bellman filter can outperform state-of-the-art particle filtering methods in high-dimensional settings, while the computational burden remains low.

## 10 Application II: Nonlinear and degenerate state dynamics

This section considers a recent state-space model in financial econometrics featuring multidimensional, nonlinear and degenerate state dynamics. A simulation study demonstrates that the Bellman filter outperforms the particle filter for the purposes of both parameter estimation and filtering, while an empirical application using real data yields similar results for both methods.

### 10.1 Model

Catania (2022, eq. 1) considers a stochastic-volatility model with a general leverage specification:

$$y_t = \mu + \exp(h_t/2) \varepsilon_t, \quad \varepsilon_t \sim \text{i.i.d. N}(0, 1), \quad (46)$$

$$h_t = c + \varphi h_{t-1} + \sigma_\eta \eta_t, \quad (47)$$

$$\eta_t = \sum_{j=0}^k \rho_j \varepsilon_{t-j} + \sigma_\xi \xi_t, \quad \xi_t \sim \text{i.i.d. N}(0, 1). \quad (48)$$

Here,  $y_t$  is a financial log return, with median (but not mean, as we shall see)  $\mu$ . The dynamics for the log-volatility process  $\{h_t\}$  feature the intercept  $c$ , persistence parameter  $|\varphi| < 1$  and variability  $\sigma_\eta > 0$ . The volatility shock  $\eta_t$  is a linear function of current and lagged return shocks, i.e.  $\varepsilon_t, \dots, \varepsilon_{t-k}$ , where  $k \geq 0$  represents the maximum lag length. Unlike in standard volatility models, the return shock  $\varepsilon_t$  and log-volatility  $h_t$  are generally dependent; both are related to  $\eta_t$  whenever  $\rho_0 \neq 0$ . When  $\rho_0 < 0$ , as is typical for financial returns, a negative return shock  $\varepsilon_t$  tends to coincide, contemporaneously, with a

positive volatility shock  $\eta_t$ . This is known as the ‘volatility-feedback effect’ (e.g. Carr and Wu, 2017) and implies that the distribution of  $y_t$  is negatively skewed, explaining why  $\mu$  is the median but not generally the mean. While Catania (2022) sets  $\mu = 0$ , the introduction of  $\mu$  enables a more accurate estimation of  $\rho_0$  by disentangling the location and scale. Parameters  $\rho_j \in (-1, 1)$  for  $j = 1, \dots, k$  quantify a generalised ‘leverage effect’: the impact of multiple lagged return shocks  $\varepsilon_{t-j}$  on the volatility shock  $\eta_t$ . Catania (2022) sets  $\sigma_\xi^2 = 1 - \sum_{j=0}^k \rho_j^2$  with  $\sum_{j=0}^k \rho_j^2 < 1$  to ensure that the unconditional variance of  $\eta_t$  is unity; this is required for the identification of  $\sigma_\eta$ .

## 10.2 State-space formulation

Model (46) through (48) can be written in the general state-space format (1) if the latent state is identified as  $\mathbf{a}_t = (h_t, h_{t-1}, \dots, h_{t-k})' \in \mathbb{R}^{k+1}$ , which contains the log volatility  $h_t$  as well as  $k$  lags. As shown in Supplement U, the probability density of  $y_t \in \mathbb{R}$  conditional on the (now multidimensional) state  $\mathbf{a}_t$  and the information set at time  $t - 1$  is Gaussian with mean  $\mu_{y,t}$  and standard deviation  $\sigma_{y,t}$  as follows:

$$p(y_t | \mathbf{a}_t, \mathcal{F}_{t-1}) = \frac{1}{\sigma_{y,t} \sqrt{2\pi}} \exp\left(-\frac{(y_t - \mu_{y,t})^2}{2\sigma_{y,t}^2}\right), \quad \sigma_{y,t} = \exp(h_t/2) \sqrt{1 - \frac{\rho_0^2}{1 - \sum_{j=1}^k \rho_j^2}}, \quad (49)$$

$$\mu_{y,t} = \mu + \frac{\rho_0}{1 - \sum_{j=1}^k \rho_j^2} \exp(h_t/2) \left[ \frac{h_t - c - \varphi h_{t-1}}{\sigma_\eta} - \sum_{j=1}^k \rho_j \frac{y_{t-j} - \mu}{\exp(h_{t-j}/2)} \right].$$

The mean  $\mu_{y,t}$  depends on the log volatility  $h_t$  as well as  $k$  of its lags (except when  $\rho_0 = 0$ ), such that  $y_t$  provides information about the entire state vector  $\mathbf{a}_t = (h_t, \dots, h_{t-k})'$ . This implies that, at each time step,  $k + 1$  logarithmic volatilities must be estimated; this insight will be important for the choice of estimation method. The density of the state vector  $\mathbf{a}_t$  conditional on the previous state and the information set  $\mathcal{F}_{t-1}$  is a degenerate Gaussian (for details, see Supplement U). The first element of  $\mathbf{a}_t$  (i.e.  $h_t$ ) has a proper distribution, while lagged versions of  $h_t$  are not random when the conditioning set includes the previous state  $\mathbf{a}_{t-1}$ :

$$p(\mathbf{a}_t | \mathbf{a}_{t-1}, \mathcal{F}_{t-1}) = \frac{1}{\sigma_{h,t} \sqrt{2\pi}} \exp\left(-\frac{(h_t - \mu_{h,t})^2}{2\sigma_{h,t}^2}\right) \times \prod_{j=1}^k \delta(a_{j+1,t} - a_{j,t-1}), \quad (50)$$

$$\mu_{h,t} = c + \varphi h_{t-1} + \sigma_\eta \sum_{j=1}^k \rho_j \frac{y_{t-j} - \mu}{\exp(h_{t-j}/2)}, \quad \sigma_{h,t} = \sigma_\eta \sqrt{1 - \sum_{j=1}^k \rho_j^2}.$$

Here,  $a_{j,t}$  denotes the  $j$ -th element of the state vector  $\mathbf{a}_t = (h_t, h_{t-1}, \dots, h_{t-k})'$ , and  $\delta(\cdot)$  denotes the Dirac delta function. The product of Dirac deltas ensures that the second element of  $\mathbf{a}_t$  equals the first element in  $\mathbf{a}_{t-1}$ , and so on. The resulting state dynamics are multidimensional, nonlinear and degenerate. This is problematic, as parameter estimation for multidimensional states (Kantas et al., 2015, p. 335) and/or degenerate state dynamics (Künsch, 2013, p. 1396) using particle-filtering methods remains a challenge that has not yet been fully resolved in the literature. For the same reasons, approximate filters such as that in Koyama et al. (2010) are ruled out.

Table 11: Average parameter estimates across 100 samples, standard deviations (in parentheses) and the average of numerical standard errors (in square brackets).

	Parameter estimates							MAE
	$\mu$	$c$	$\varphi$	$\sigma_\eta$	$\rho_0$	$\rho_1$	$\rho_2$	$h_{t t-1}$
True value $\rightarrow$	0.0015	-0.200	0.980	0.250	-0.700	-0.400	0.300	$h_t$
Bellman filter	0.0015 (0.0001) [0.0001]	-0.207 (0.038) [0.033]	0.979 (0.004) [0.003]	0.252 (0.024) [0.026]	-0.651 (0.089) [0.094]	-0.438 (0.115) [0.107]	0.294 (0.101) [0.102]	0.358
Particle filter	0.0016 (0.0002) [0.0001]	-0.262 (0.155) [0.004]	0.974 (0.016) [0.001]	0.279 (0.051) [0.004]	-0.739 (0.110) [0.004]	-0.109 (0.293) [0.005]	0.095 (0.203) [0.005]	0.382
True value $\rightarrow$	0.0015	-0.200	0.980	0.250	-0.400	-0.700	0.300	$h_t$
Bellman filter	0.0015 (0.0001) [0.0001]	-0.208 (0.033) [0.034]	0.979 (0.004) [0.004]	0.265 (0.034) [0.033]	-0.355 (0.083) [0.088]	-0.715 (0.062) [0.064]	0.306 (0.084) [0.089]	0.335
Particle filter	0.0015 (0.0003) [0.0001]	-0.242 (0.099) [0.005]	0.976 (0.010) [0.001]	0.250 (0.062) [0.006]	-0.471 (0.207) [0.007]	-0.441 (0.347) [0.007]	0.061 (0.258) [0.008]	0.358

*Note:* MAE = mean absolute error. For both sets of true parameter values, I simulate 100 samples of length 5,000 and compute parameter estimates based on the first 2,500 observations. For the Bellman filter, the proposed approximate estimator (40) is used. For the particle filter, I follow Catania (2022) in using Malik and Pitt’s (2011) continuous sampling importance resampling (CSIR) particle filter with 5,000 particles. For each sample I compute, in addition to parameter estimates, numerical standard errors by inverting the negative Hessian matrix evaluated at the peak and taking the square root of the diagonal. I exclude standard errors based on non-invertible Hessian matrices, which were encountered in  $\sim 40\%$  of samples based on the CSIR method. Using estimated parameters, I make out-of-sample predictions by running the filter on the entire data set, computing mean absolute errors (MAEs) by comparing out-of-sample predictions  $h_{t|t-1}$  with actual (simulated) values  $h_t$  for  $t > 2,500$ .

### 10.3 Parameter-estimation methods

Catania (2022) estimates the static parameters of the state-space model (49) and (50) using a univariate implementation of Malik and Pitt’s (2011) continuous sampling importance resampling (CSIR) method. The effect of this univariate approach on parameter estimation and model selection is a priori unclear. Moreover, this approach comes with three potential disadvantages. First, the univariate approach means that only the first element of the state vector  $\mathbf{a}_t = (h_t, h_{t-1}, \dots, h_{t-k})'$  is estimated at time  $t$ , while the other elements remain fixed at previously estimated values. However, the observation  $y_t$  contains information about the entire state vector  $\mathbf{a}_t$ , as can be seen from the observation density (49). While actual (i.e. true) lags of  $h_t$  are constant over time, the researcher’s estimates need not be. Even when focusing purely on the real-time estimation of  $h_t$ , the decision not to re-estimate the lags at each point in time may lead to an efficiency loss. Second, while the CSIR method guarantees a continuous approximation of the log-likelihood function, this approximation need not be smooth, potentially causing standard gradient-based optimisers to fail. I employ a grid search to identify promising areas of the parameter space, followed by a simplex-based optimisation algorithm that does not utilise gradients. Third, numerical standard errors derived from the inversion of negative Hessian matrices may be misleading when the objective function is nonsmooth. For a piecewise linear approximation as in the CSIR method, finite-difference Hessian matrices may be badly scaled when evaluated near kinks, or identically zero when evaluated on linear pieces. This may explain the exceedingly small standard errors reported in Catania (2022), as well as my finding that Hessian matrices based on the CSIR method frequently fail to be invertible.

In addition to the particle filter, I employ the general version of the Bellman filter (section 3.1) extended to account for degenerate state dynamics (section 3.2). The Bellman filter is implemented using closed-form expressions (given in Supplement V) for derivatives of the observation and state-transition log densities with respect to the entire state vector  $\mathbf{a}_t = (h_t, h_{t-1}, \dots, h_{t-k})$ ; hence, the entire  $(k+1)$ -dimensional state is estimated at each time  $t$ . I allow up to  $k_{\max} = 10$  lags, implying that the Bellman filter solves an optimisation problem with up to 11 dimensions at each time step. To estimate the static parameters, I identify promising starting values using a grid search, after which I implement estimator (40) using a gradient-based numerical optimiser. In the Bellman-filtering procedure, at each time step I execute Newton or Fisher optimisation steps when the search direction is well-defined; otherwise, the optimisation is skipped and the update is set equal to the prediction. This somewhat crude approach ensures that the filter runs smoothly even when using flawed parameter values, which may be encountered during the black-box estimation routine (40). At the optimal parameter values identified using this routine, the filter is convergent at every time step.

#### 10.4 Simulation results

To investigate the difference between the multivariate approach and the (one-dimensional) CSIR method, a simulation study is performed. Two sets of realistic parameter values are shown in Table 11. I generate 100 series of length 5,000, using the first half for parameter estimation. The results in Table 11 show that average parameter estimates of  $\rho_0, \rho_1$  and  $\rho_2$  obtained by the CSIR particle filter are inaccurate, while those based on the Bellman filter are relatively accurate. For example, the average estimate of  $\rho_2$  by the Bellman filter differs from the true value by no more than 0.01, compared to at least 0.20 for the particle filter. While Catania (2022) demonstrated that the CSIR method may produce accurate parameter estimates, this finding may partly be explained by the fact that the parameter-optimisation routine there was initialised using the true parameters, in which case the CSIR estimates typically remain close to the starting point. The results also show that the parameter estimates based on the particle filter vary greatly across samples, as can be seen from the large standard deviations in parentheses in Table 11, while parameter estimates based on the Bellman filter are relatively stable. Additionally, the average of numerically computed standard errors, in square brackets, indicates that standard errors are somewhat reliable for the Bellman filter, closely matching the actual variation across samples, but not for the CSIR method, where they are several orders of magnitude too small. This may be due to the nonsmooth approximation of the log-likelihood function in the CSIR method, and casts doubt on the validity of similarly small standard errors in Catania (2022). Finally, the right-most column shows that the improved parameter estimates lead to out-of-sample forecasting gains, which are consistent across samples (the Bellman filter produces better forecasts for each sample) and overwhelmingly statistically significant according to a standard Diebold-Mariano test (not shown).

#### 10.5 Empirical results

For the empirical application, I take log returns of the S&P500 from 3 Jan 1990 to 31 Dec 2019 (7,558 observations). Table 12 shows preferred models when using the Bayesian information criterion, which suggests setting  $k = 3$  lags for both parameter-estimation methods when up to 10 lags are allowed (full results are available in Supplement W). Parameter estimates for both methods are similar, perhaps due to the comparatively long dataset. Both methods indicate that volatility feedback and leverage play important

Table 12: Parameter estimates for preferred model specifications and numerical standard errors in square brackets

	$\mu$	$c$	$\varphi$	$\sigma_\eta$	$\rho_0$	$\rho_1$	$\rho_2$	$\rho_3$
Bellman filter	0.051 [0.008]	-0.001 [0.002]	0.982 [0.003]	0.258 [0.016]	-0.377 [0.049]	-0.583 [0.066]	-0.091 [0.099]	0.463 [0.060]
Particle filter	0.052 [0.004]	-0.006 [0.002]	0.983 [0.002]	0.239 [0.005]	-0.398 [0.009]	-0.571 [0.007]	-0.114 [0.007]	0.459 [0.005]

*Note:* For both parameter-estimation methods, the preferred model determined by the Bayesian information criterion (BIC) has three lags. Full parameter-estimation results with up to ten lags are available in Supplement W. The data are log returns of the S&P500 (multiplied by 100) from 3 Jan 1990 to 31 Dec 2019 (7,558 observations).

roles, with the positive estimate of  $\rho_3$  suggesting that the leverage effect is temporary: upward volatility shocks following negative returns may be partially reversed on day three. The small standard errors for the particle filter, similar to those reported in Catania (2022, table 2), may underestimate the true uncertainty surrounding the parameter estimates. Standard errors based on the Bellman filter, which are up to an order of magnitude higher for the parameters of interest, were in simulation studies found to be reasonably accurate.

## 11 Conclusion

The Bellman filter for state-space models as developed in this article generalises the Kalman filter and is equally computationally inexpensive in high-dimensional state spaces, but robust in the case of heavy-tailed observation noise and applicable to a wider range of (nonlinear and non-Gaussian) models. Under suitable conditions, the Bellman-filtered states are globally contractive to a small region around the true state at every time step, while filtering errors remain uniformly bounded over time. A second contribution is the development of a Bellman smoother that is mathematically equivalent to the classic Rauch, Tung and Striebel (1965) smoother, but applicable more generally—as an approximation—to state-space models with nonlinear and/or non-Gaussian observation equations. Third, the approximate static-parameter estimation procedure developed here is straightforward to implement and, again, computationally inexpensive; the resulting parameter estimates for various sample sizes appear to be no less accurate or efficient than those of (asymptotically exact) simulation-based methods.

In a simulation study involving a wide range of univariate models, the performance of the Bellman filter is near identical to those of state-of-the-art simulation-based methods in terms of parameter estimation and filtering, while additionally enabling smoothing. Filtering speeds are improved by factors up to  $\sim 160$  (compared to particle filters) and  $\sim 2,000$  (cf. importance samplers). Likewise, computation times for estimating the static parameters are reduced by factors up to  $\sim 10$  (cf. importance samplers) and  $\sim 400$  (cf. particle filters). In an application with a high-dimensional climate model, the tracking performance of the Bellman filter remains virtually unchanged as the dimension of the state space is increased from 10 to 150, while that of the particle filter deteriorates sharply—due to the curse of dimensionality—even when employing very many particles: e.g. with  $10^6$  particles in 150 spatial dimensions, the Bellman filter is both faster (by a factor  $\sim 1,000$ ) and more accurate (by a factor  $\sim 1.8$  in terms of mean absolute filtering error). In a second application with highly nonlinear and degenerate state dynamics, the Bellman filter outperforms the particle filter for the purposes of both parameter estimation and filtering.

## Acknowledgements

I thank Wisse Rutgers for research assistance, Serena Ng for helpful editorial guidance, and the anonymous AE and two referees for their valuable comments. Thanks are also due to Maksim Anisimov, Francisco Blasques, Leopoldo Catania, Dick van Dijk, Simon Donker van Heel, Jippe van Dunné, Dennis Fok, Maria Grith, Andrew Harvey, Christiaan Heij, Elwin Kardux, Matthias Katzfuss, Onno Kleen, Erik Kole, Siem Jan Koopman, Rutger Lit, Rasmus Lonn, André Lucas, Robin Lumsdaine, Jan Maciejowski, Andrea Naghi, Jochem Oorschot, Richard Paap, Andreas Pick, Krzysztof Postek, Rogier Quaedvlieg, Daniel Ralph, Bram van Os, Omiros Papaspiliopoulos, Marcel Scharth, Annika Schnücker, Ekaterina Smetanina, Panos Toulis, Stephen Thiele, Nando Vermeer, Sebastiaan Vermeulen, Michel van der Wel, Martina Zaharieva, Mikhail Zhelonkin and Chen Zhou. Finally, I thank participants of the 2021 North American summer meeting of the Econometric Society and the 27th international conference on Computing in Economics and Finance for stimulating discussions.

## References

- Anderson, B. D. and Moore, J. B. (2012) *Optimal Filtering*. Courier Corporation.
- Asi, H. and Duchi, J. C. (2019) Stochastic (approximate) proximal point methods: Convergence, optimality, and adaptivity. *SIAM Journal on Optimization*, **29**, 2257–2290.
- Baum, L. E. and Petrie, T. (1966) Statistical inference for probabilistic functions of finite state Markov chains. *The Annals of Mathematical Statistics*, **37**, 1554–1563.
- Bauwens, L. and Hautsch, N. (2006) Stochastic conditional intensity processes. *Journal of Financial Econometrics*, **4**, 450–493.
- Bauwens, L. and Veredas, D. (2004) The stochastic conditional duration model: A latent variable model for the analysis of financial durations. *Journal of Econometrics*, **119**, 381–412.
- Bellman, R. E. (1957) *Dynamic Programming*. PUP.
- Bertsekas, D. P. (2012) *Dynamic Programming and Optimal Control, Vol. II: Approximate Dynamic Programming*. Athena Scientific.
- Bianchi, P. (2016) Ergodic convergence of a stochastic proximal point algorithm. *SIAM Journal on Optimization*, **26**, 2235–2260.
- Carr, P. and Wu, L. (2017) Leverage effect, volatility feedback, and self-exciting market disruptions. *Journal of Financial & Quantitative Analysis*, **52**, 2119–2156.
- Catania, L. (2022) A stochastic volatility model with a general leverage specification. *Journal of Business & Economic Statistics*, **40**, 678–689.
- Chopin, N. and Papaspiliopoulos, O. (2020) *An Introduction to Sequential Monte Carlo*. Springer.
- Doucet, A., De Freitas, N. and Gordon, N. (2001) *Sequential Monte Carlo Methods in Practice*. Springer.
- Durbin, J. and Koopman, S. J. (1997) Monte Carlo maximum likelihood estimation for non-Gaussian state space models. *Biometrika*, **84**, 669–684.
- (2000) Time series analysis of non-Gaussian observations based on state space models from both classical and Bayesian perspectives. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **62**, 3–56.
- Fahrmeir, L. (1992) Posterior mode estimation by extended Kalman filtering for multivariate dynamic generalized linear models. *Journal of the American Statistical Association*, **87**, 501–509.
- Farmer, L. E. (2021) The discretization filter: A simple way to estimate nonlinear state space models. *Quantitative Economics*, **12**, 41–76.
- Frühwirth-Schnatter, S. and Wagner, H. (2006) Auxiliary mixture sampling for parameter-driven models of time series of counts with applications to state space modelling. *Biometrika*, **93**, 827–841.

- Fuh, C.-D. (2006) Efficient likelihood estimation in state space models. *The Annals of Statistics*, **34**, 2026–2068.
- Ghysels, E., Harvey, A. C. and Renault, E. (1996) Stochastic volatility. In *Handbook of Statistics, Vol. 14, Statistical Methods in Finance* (eds. G. Maddala and C. Rao), 119–191. Elsevier.
- Godsill, S. J., Doucet, A. and West, M. (2004) Monte Carlo smoothing for nonlinear time series. *Journal of the American Statistical Association*, **99**, 156–168.
- Hafner, C. M. and Manner, H. (2012) Dynamic stochastic copula models: Estimation, inference and applications. *Journal of Applied Econometrics*, **27**, 269–295.
- Hamilton, J. D. (1989) A new approach to the economic analysis of nonstationary time series and the business cycle. *Econometrica*, **57**, 357–384.
- Hansen, L. P. and Sargent, T. J. (2013) *Recursive Models of Dynamic Linear Economies*. PUP.
- Harvey, A. C. (1990) *Forecasting, Structural Time Series Models and the Kalman Filter*. CUP.
- Harvey, A. C., Ruiz, E. and Shephard, N. (1994) Multivariate stochastic variance models. *The Review of Economic Studies*, **61**, 247–264.
- Harvey, A. C. and Shephard, N. (1996) Estimation of an asymmetric stochastic volatility model for asset returns. *Journal of Business & Economic Statistics*, **14**, 429–434.
- Jacquier, E., Polson, N. G. and Rossi, P. E. (2002) Bayesian analysis of stochastic volatility models. *Journal of Business & Economic Statistics*, **20**, 69–87.
- Johansen, A. M., Whiteley, N. and Doucet, A. (2012) Exact approximation of Rao-Blackwellised particle filters. *IFAC Proceedings Volumes*, **45**, 488–493.
- Julier, S. J. and Uhlmann, J. K. (1997) New extension of the Kalman filter to nonlinear systems. In *Signal Processing, Sensor Fusion, and Target Recognition VI* (ed. I. Kadar), vol. 3068, 182–193. International Society for Optics and Photonics.
- Jungbacker, B. and Koopman, S. J. (2007) Monte Carlo estimation for nonlinear non-Gaussian state space models. *Biometrika*, **94**, 827–839.
- Jungers, R. (2009) *The Joint Spectral Radius: Theory and Applications*. Springer.
- Kalman, R. E. (1960) A new approach to linear filtering and prediction problems. *Journal of Basic Engineering*, **82**, 35–45.
- Kantas, N., Doucet, A., Singh, S. S., Maciejowski, J. and Chopin, N. (2015) On particle methods for parameter estimation in state-space models. *Statistical Science*, **30**, 328–351.
- Katzfuss, M., Stroud, J. R. and Wikle, C. K. (2020) Ensemble Kalman methods for high-dimensional hierarchical dynamic space-time models. *Journal of the American Statistical Association*, **115**, 866–885.
- Kitagawa, G. (1987) Non-Gaussian state-space modeling of nonstationary time series. *Journal of the American Statistical Association*, **82**, 1032–1041.
- Koopman, S. J., Lit, R. and Lucas, A. (2017) Intraday stochastic volatility in discrete price changes: The dynamic Skellam model. *Journal of the American Statistical Association*, **112**, 1490–1503.
- Koopman, S. J., Lucas, A. and Scharth, M. (2015) Numerically accelerated importance sampling for nonlinear non-Gaussian state-space models. *Journal of Business & Economic Statistics*, **33**, 114–127.
- (2016) Predicting time-varying parameters with parameter-driven and observation-driven models. *Review of Economics and Statistics*, **98**, 97–110.
- Koyama, S., Castellanos Pérez-Bolde, L., Shalizi, C. R. and Kass, R. E. (2010) Approximate methods for state-space models. *Journal of the American Statistical Association*, **105**, 170–180.
- Koyama, S. and Paninski, L. (2010) Efficient computation of the maximum a posteriori path and parameter estimation in integrate-and-fire and more general state-space models. *Journal of Computational Neuroscience*, **29**, 89–105.
- Kullback, S. and Leibler, R. A. (1951) On information and sufficiency. *The Annals of Mathematical Statistics*, **22**, 79–86.
- Künsch, H. R. (2001) State space and hidden Markov models. In *Complex Stochastic Systems* (eds. O. E. Barndorff-Nielsen and C. Kluppelberg), 109–174. Chapman & Hall/CRC.
- (2013) Particle filters. *Bernoulli*, **19**, 1391–1403.

- Liu, J. and West, M. (2001) Combined parameter and state estimation in simulation-based filtering. In *Sequential Monte Carlo Methods in Practice* (eds. A. Doucet, N. De Freitas and N. Gordon), 197–223. Springer.
- Liu, J. S. (2008) *Monte Carlo Strategies in Scientific Computing*. Springer.
- Liu, Q. and Ihler, A. (2013) Variational algorithms for marginal MAP. *The Journal of Machine Learning Research*, **14**, 3165–3200.
- Malik, S. and Pitt, M. K. (2011) Particle filters for continuous likelihood evaluation and maximisation. *Journal of Econometrics*, **165**, 190–209.
- Masreliez, C. (1975) Approximate non-Gaussian filtering with linear state and observation relations. *IEEE Transactions on Automatic Control*, **20**, 107–110.
- Mayne, D. Q. (1966) A solution of the smoothing problem for linear dynamic systems. *Automatica*, **4**, 73–92.
- Müller, U. K. and Petalas, P.-E. (2010) Efficient estimation of the parameter path in unstable time series models. *The Review of Economic Studies*, **77**, 1508–1539.
- Murphy, S. A. and Van der Vaart, A. W. (2000) On profile likelihood. *Journal of the American Statistical Association*, **95**, 449–465.
- Nocedal, J. and Wright, S. J. (2006) *Numerical Optimization*. Springer.
- Patrascu, A. and Necoara, I. (2018) Nonasymptotic convergence of stochastic proximal point methods for constrained convex optimization. *The Journal of Machine Learning Research*, **18**, 7204–7245.
- Rauch, H. E., Tung, F. and Striebel, C. T. (1965) Maximum likelihood estimates of linear dynamic systems. *AIAA Journal*, **3**, 1445–1450.
- Rockafellar, R. T. (1976) Monotone operators and the proximal point algorithm. *SIAM Journal on Control and Optimization*, **14**, 877–898.
- Ruiz, E. (1994) Quasi-maximum likelihood estimation of stochastic volatility models. *Journal of Econometrics*, **63**, 289–306.
- Ryu, E. K. and Boyd, S. (2016) Stochastic proximal iteration: A non-asymptotic improvement upon stochastic gradient descent. *Author website*.
- Singh, A. and Roberts, G. (1992) State space modelling of cross-classified time series of counts. *International Statistical Review*, **60**, 321–335.
- So, M. K. (2003) Posterior mode estimation for nonlinear and non-Gaussian state space models. *Statistica Sinica*, **13**, 255–274.
- Straumann, D. and Mikosch, T. (2006) Quasi-maximum-likelihood estimation in conditionally heteroscedastic time series: A stochastic recurrence equations approach. *The Annals of Statistics*, **34**, 2449–2495.
- Surace, S. C., Kutschireiter, A. and Pfister, J.-P. (2019) How to avoid the curse of dimensionality: Scalability of particle filters with and without importance weights. *SIAM Review*, **61**, 79–91.
- Taylor, S. J. (2008) *Modelling Financial Time Series*. World Scientific.
- Tichavsky, P., Muravchik, C. H. and Nehorai, A. (1998) Posterior Cramér-Rao bounds for discrete-time nonlinear filtering. *IEEE Transactions on Signal Processing*, **46**, 1386–1396.
- Toulis, P. and Airoidi, E. M. (2017) Asymptotic and finite-sample properties of estimators based on stochastic gradients. *Annals of Statistics*, **45**, 1694–1727.
- Toulis, P., Tran, D. and Airoidi, E. (2016) Towards stability and optimality in stochastic gradient descent. In *Artificial Intelligence and Statistics*, vol. 51, 1290–1298. PMLR.
- Viterbi, A. J. (1967) Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Transactions on Information Theory*, **13**, 260–269.
- West, M. (1981) Robust sequential approximate bayesian estimation. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, **43**, 157–166.
- Whittle, P. (1981) Risk-sensitive linear/quadratic/Gaussian control. *Advances in Applied Probability*, **13**, 764–777.
- (1996) *Optimal Control: Basics and Beyond*. Wiley.
- (2004) State structure, decision making and related issues. In *State space and unobserved component models: Theory and applications* (eds. A. Harvey, S. J. Koopman and N. Shephard), 26–39. CUP.

# Supplementary material for the article 'Bellman filtering and smoothing for state-space models'

Rutger-Jan Lange

Econometric Institute, Erasmus School of Economics, Rotterdam, The Netherlands

1 November 2023

## A Proof of Proposition 1

To understand how a recursive approach may be feasible, we start by noting that the joint log-likelihood function (2) satisfies a straightforward recursive relation for  $2 \leq t \leq n$  as follows:

$$L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t) = \ell(\mathbf{y}_t|\mathbf{a}_t) + \ell(\mathbf{a}_t|\mathbf{a}_{t-1}) + L_{1:t-1}(\mathbf{a}_1, \dots, \mathbf{a}_{t-1}). \quad (\text{A.1})$$

That is, in transitioning from time  $t-1$  to time  $t$ , two terms are added: one representing the state-transition density,  $\ell(\mathbf{a}_t|\mathbf{a}_{t-1})$ ; the other representing the observation density,  $\ell(\mathbf{y}_t|\mathbf{a}_t)$ . Next, standard dynamic-programming arguments imply

$$\begin{aligned} V_t(\mathbf{a}_t) &:= \max_{(\mathbf{a}_1, \dots, \mathbf{a}_{t-1}) \in \mathbb{R}^{m \times (t-1)}} L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t), && \text{by definition (4),} \\ &= \max_{\mathbf{a}_{1:t-1} \in \mathbb{R}^{m \times (t-1)}} \{ \ell(\mathbf{y}_t|\mathbf{a}_t) + \ell(\mathbf{a}_t|\mathbf{a}_{t-1}) + L_{1:t-1}(\mathbf{a}_1, \dots, \mathbf{a}_{t-1}) \}, && \text{by recursion (A.1),} \\ &= \max_{\mathbf{a}_{t-1} \in \mathbb{R}^m} \left\{ \ell(\mathbf{y}_t|\mathbf{a}_t) + \ell(\mathbf{a}_t|\mathbf{a}_{t-1}) + \max_{(\mathbf{a}_1, \dots, \mathbf{a}_{t-2}) \in \mathbb{R}^{m \times (t-2)}} L_{1:t-1}(\mathbf{a}_1, \dots, \mathbf{a}_{t-1}) \right\}, \\ &&& \text{by moving all but one maximisation inside curly brackets,} \\ &= \max_{\mathbf{a}_{t-1} \in \mathbb{R}^m} \{ \ell(\mathbf{y}_t|\mathbf{a}_t) + \ell(\mathbf{a}_t|\mathbf{a}_{t-1}) + V_{t-1}(\mathbf{a}_{t-1}) \}, && \text{again by definition (4),} \\ &= \ell(\mathbf{y}_t|\mathbf{a}_t) + \max_{\mathbf{a}_{t-1} \in \mathbb{R}^m} \{ \ell(\mathbf{a}_t|\mathbf{a}_{t-1}) + V_{t-1}(\mathbf{a}_{t-1}) \}. \end{aligned} \quad (\text{A.2})$$

Further, it is evident that

$$\mathbf{a}_{t|t} = \arg \max_{\mathbf{a}_t \in \mathbb{R}^m} V_t(\mathbf{a}_t) = \arg \max_{\mathbf{a}_t \in \mathbb{R}^m} \max_{(\mathbf{a}_1, \dots, \mathbf{a}_{t-1}) \in \mathbb{R}^{m \times (t-1)}} L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t). \quad (\text{A.3})$$

## B Block-matrix inversion

Consider the second diagonal block of the negative Hessian matrix in equation (9). Define this block as  $\mathbf{D}_t \in \mathbb{R}^{m \times m}$  and define its Schur complement  $\mathbf{S}_t \in \mathbb{R}^{m \times m}$  as follows:

$$\mathbf{D}_t := \mathbf{I}_{t-1|t-1} + \mathbf{J}_t^{22}, \quad \mathbf{S}_t := \mathbf{J}_t^{11} - \mathbf{J}_t^{12} \mathbf{D}_t^{-1} \mathbf{J}_t^{21} - \frac{d^2 \ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'}. \quad (\text{B.1})$$

As is standard (e.g. Bernstein, 2009, p. 108), the required block-matrix inverse can then be expressed as

$$\begin{bmatrix} \mathbf{J}_t^{11} - \frac{d^2 \ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} & \mathbf{J}_t^{12} \\ \mathbf{J}_t^{21} & \mathbf{I}_{t-1|t-1} + \mathbf{J}_t^{22} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{S}_t^{-1} & -\mathbf{S}_t^{-1} \mathbf{J}_t^{12} \mathbf{D}_t^{-1} \\ -\mathbf{D}_t^{-1} \mathbf{J}_t^{21} \mathbf{S}_t^{-1} & \mathbf{D}_t^{-1} + \mathbf{D}_t^{-1} \mathbf{J}_t^{21} \mathbf{S}_t^{-1} \mathbf{J}_t^{12} \mathbf{D}_t^{-1} \end{bmatrix}, \quad (\text{B.2})$$

as long as the required inverses exist.

## C Derivation of equation (11)

Here we compute the negative Hessian of the value function, i.e.

$$\begin{aligned} V_t(\mathbf{a}_t) &= \ell(\mathbf{y}_t|\mathbf{a}_t) + \max_{\mathbf{a}_{t-1} \in \mathbb{R}^m} \left\{ \ell(\mathbf{a}_t|\mathbf{a}_{t-1}) - \frac{1}{2}(\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1})' \mathbf{I}_{t-1|t-1} (\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1}) \right\}, \\ &= \ell(\mathbf{y}_t|\mathbf{a}_t) + \ell(\mathbf{a}_t|\mathbf{a}_{t-1}^*) - \frac{1}{2}(\mathbf{a}_{t-1}^* - \mathbf{a}_{t-1|t-1})' \mathbf{I}_{t-1|t-1} (\mathbf{a}_{t-1}^* - \mathbf{a}_{t-1|t-1}), \end{aligned} \quad (\text{C.1})$$

where the second line employs the definition

$$\mathbf{a}_{t-1}^* := \arg \max_{\mathbf{a}_{t-1} \in \mathbb{R}^m} \left\{ \ell(\mathbf{a}_t|\mathbf{a}_{t-1}) - \frac{1}{2}(\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1})' \mathbf{I}_{t-1|t-1} (\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1}) \right\}. \quad (\text{C.2})$$

We must keep in mind that  $\mathbf{a}_{t-1}^*$  depends on  $\mathbf{a}_t$ ; we could have written  $\mathbf{a}_{t-1}^*(\mathbf{a}_t)$ . Indeed, to compute the negative Hessian of  $V_t(\mathbf{a}_t)$ , we must account for the change in  $\mathbf{a}_{t-1}^*(\mathbf{a}_t)$  using the chain rule. The first-order condition satisfied by  $\mathbf{a}_{t-1}^*$ , i.e.

$$\mathbf{0} = \frac{d\ell(\mathbf{a}_t|\mathbf{a}_{t-1}^*)}{d\mathbf{a}_{t-1}^*} - \mathbf{I}_{t-1|t-1} (\mathbf{a}_{t-1}^* - \mathbf{a}_{t-1|t-1}), \quad (\text{C.3})$$

can be differentiated with respect to  $\mathbf{a}_t$  to obtain

$$\mathbf{0} = \left[ -\mathbf{J}_t^{21} - \mathbf{J}_t^{22} \frac{d\mathbf{a}_{t-1}^*}{d\mathbf{a}_t'} - \mathbf{I}_{t-1|t-1} \frac{d\mathbf{a}_{t-1}^*}{d\mathbf{a}_t'} \right]_{\mathbf{a}_{t-1} = \mathbf{a}_{t-1}^*}, \quad (\text{C.4})$$

where  $\mathbf{J}_t^{21}$  and  $\mathbf{J}_t^{22}$  are as in equation (10). Solving for the sensitivity of  $\mathbf{a}_{t-1}^*$  with respect to  $\mathbf{a}_t$ , we obtain

$$\frac{d\mathbf{a}_{t-1}^*}{d\mathbf{a}_t'} = [-(\mathbf{I}_{t-1|t-1} + \mathbf{J}_t^{22})^{-1} \mathbf{J}_t^{21}]_{\mathbf{a}_{t-1} = \mathbf{a}_{t-1}^*}. \quad (\text{C.5})$$

Next, the chain rule tells us that the Hessian with respect to  $\mathbf{a}_t$  can be computed as

$$\frac{d^2 \cdot}{d\mathbf{a}_t d\mathbf{a}_t'} = \left[ \frac{\mathbb{1}_{m \times m}}{d\mathbf{a}_{t-1}^*} \right]' \left[ \begin{array}{cc} \frac{\partial^2 \cdot}{\partial \mathbf{a}_t \partial \mathbf{a}_t'} & \frac{\partial^2 \cdot}{\partial \mathbf{a}_t \partial \mathbf{a}_{t-1}^*'} \\ \frac{\partial^2 \cdot}{\partial \mathbf{a}_{t-1}^* \partial \mathbf{a}_t'} & \frac{\partial^2 \cdot}{\partial \mathbf{a}_{t-1}^* \partial \mathbf{a}_{t-1}^*'} \end{array} \right] \left[ \frac{\mathbb{1}_{m \times m}}{d\mathbf{a}_{t-1}^*} \right], \quad (\text{C.6})$$

where instances of  $\partial$  and  $d$  denote ‘partial’ and ‘total’ derivatives, respectively, while  $\mathbb{1}_{m \times m}$  denotes an identity matrix of size  $m \times m$ . By the first-order envelope theorem, no first order derivative with respect to  $\mathbf{a}_{t-1}^*$  appears. The negative Hessian of  $V_t(\mathbf{a}_t)$  becomes

$$\begin{aligned} -\frac{d^2 V_t(\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} &= \left[ \frac{\mathbb{1}_{m \times m}}{d\mathbf{a}_{t-1}^*} \right]' \left[ \begin{array}{cc} \mathbf{J}_t^{11} - \frac{d^2 \ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} & \mathbf{J}_t^{12} \\ \mathbf{J}_t^{21} & \mathbf{I}_{t-1,t-1} + \mathbf{J}_t^{22} \end{array} \right] \left[ \frac{\mathbb{1}_{m \times m}}{d\mathbf{a}_{t-1}^*} \right]_{\mathbf{a}_{t-1} = \mathbf{a}_{t-1}^*}, \\ &= \mathbf{J}_t^{11} - \frac{d^2 \ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} - 2\mathbf{J}_t^{12} (\mathbf{I}_{t-1|t-1} + \mathbf{J}_t^{22})^{-1} \mathbf{J}_t^{21} + \frac{d\mathbf{a}_{t-1}^*}{d\mathbf{a}_t} (\mathbf{I}_{t-1,t-1} + \mathbf{J}_t^{22}) \frac{d\mathbf{a}_{t-1}^*}{d\mathbf{a}_t'} \Big|_{\mathbf{a}_{t-1} = \mathbf{a}_{t-1}^*}, \\ &= \mathbf{J}_t^{11} - \frac{d^2 \ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} - \mathbf{J}_t^{12} (\mathbf{I}_{t-1|t-1} + \mathbf{J}_t^{22})^{-1} \mathbf{J}_t^{12} \Big|_{\mathbf{a}_{t-1} = \mathbf{a}_{t-1}^*}. \end{aligned} \quad (\text{C.7})$$

Finally  $\mathbf{a}_{t-1}^*(\mathbf{a}_t) = \mathbf{a}_{t-1|t}$ , such that

$$-\frac{d^2 V_t(\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} \Big|_{\mathbf{a}_t} = \left[ \mathbf{J}_t^{11} - \frac{d^2 \ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} - \mathbf{J}_t^{12} (\mathbf{I}_{t-1|t-1} + \mathbf{J}_t^{22})^{-1} \mathbf{J}_t^{21} \right]_{\mathbf{a}_t = \mathbf{a}_t, \mathbf{a}_{t-1} = \mathbf{a}_{t-1|t}}, \quad (\text{C.8})$$

which confirms equation (11).

## D Kalman information update as a special case of (11)

For the linear Gaussian model in Corollary 1, we have  $\mathbf{J}_t^{11} = \mathbf{Q}^{-1}$ ,  $\mathbf{J}_t^{12} = \mathbf{Q}^{-1}\mathbf{T}$ ,  $\mathbf{J}_t^{21} = \mathbf{T}'\mathbf{Q}^{-1}$ ,  $\mathbf{J}_t^{22} = \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T}$  and  $d^2\ell(\mathbf{y}_t|\mathbf{a}_t)/(d\mathbf{a}_t d\mathbf{a}_t') = -\mathbf{Z}'\mathbf{H}^{-1}\mathbf{Z}$ . Substituting these equalities into the information update (11), we obtain

$$\begin{aligned} \mathbf{I}_{t|t} &= \mathbf{Q}^{-1} - \mathbf{Q}^{-1}\mathbf{T}(\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T})^{-1}\mathbf{T}'\mathbf{Q}^{-1} + \mathbf{Z}'\mathbf{H}^{-1}\mathbf{Z}, \\ &= \mathbf{I}_{t|t-1} + \mathbf{Z}'\mathbf{H}^{-1}\mathbf{Z}, \end{aligned} \quad (\text{D.1})$$

where  $\mathbf{I}_{t|t-1}$  is defined as

$$\mathbf{I}_{t|t-1} := \mathbf{Q}^{-1} - \mathbf{Q}^{-1}\mathbf{T}(\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T})^{-1}\mathbf{T}'\mathbf{Q}^{-1} = (\mathbf{T}\mathbf{I}_{t-1|t-1}^{-1}\mathbf{T}' + \mathbf{Q})^{-1}, \quad (\text{D.2})$$

and where the second equality follows by the Woodbury matrix equality (e.g. Henderson and Searle, 1981, eq. 1). Next, assuming the inverses  $\mathbf{P}_{t|t-1} := \mathbf{I}_{t|t-1}^{-1}$  and  $\mathbf{P}_{t|t} := \mathbf{I}_{t|t}^{-1}$  exist, using again Henderson and Searle (1981, eq. 1), we find

$$\mathbf{P}_{t|t} = \mathbf{I}_{t|t}^{-1} = (\mathbf{I}_{t|t-1} + \mathbf{Z}'\mathbf{H}^{-1}\mathbf{Z})^{-1} = \mathbf{P}_{t|t-1} - \mathbf{P}_{t|t-1}\mathbf{Z}'(\mathbf{Z}\mathbf{P}_{t|t-1}\mathbf{Z}' + \mathbf{H})^{-1}\mathbf{Z}\mathbf{P}_{t|t-1}, \quad (\text{D.3})$$

which is exactly the Kalman filter covariance matrix updating step (again, see Harvey, 1990, p. 106).

## E Derivation of equation (14)

The first-order condition for the maximisation over  $\mathbf{a}_{t-1}$  in equation (13) can be usefully manipulated as follows:

$$\begin{aligned} \mathbf{a}_{t-1}^* &= (\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T})^{-1}(\mathbf{I}_{t-1|t-1}\mathbf{a}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}(\mathbf{a}_t - \mathbf{c})), \\ &= \mathbf{a}_{t-1|t-1} + (\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T})^{-1}\mathbf{T}'\mathbf{Q}^{-1}(\mathbf{a}_t - \mathbf{c} - \mathbf{T}\mathbf{a}_{t-1|t-1}), \\ &= \mathbf{a}_{t-1|t-1} + \mathbf{I}_{t-1|t-1}^{-1}\mathbf{T}'(\mathbf{T}\mathbf{I}_{t-1|t-1}^{-1}\mathbf{T}' + \mathbf{Q})^{-1}(\mathbf{a}_t - \mathbf{c} - \mathbf{T}\mathbf{a}_{t-1|t-1}), \\ &= \mathbf{a}_{t-1|t-1} + \mathbf{I}_{t-1|t-1}^{-1}\mathbf{T}'\mathbf{I}_{t|t-1}(\mathbf{a}_t - \mathbf{a}_{t|t-1}), \end{aligned} \quad (\text{E.1})$$

which confirms equation (14) in the main text. This second line expresses  $\mathbf{a}_{t-1}^*$  as the sum of  $\mathbf{a}_{t-1|t-1}$  and a correction that is linear in the ‘innovation’  $\mathbf{a}_t - \mathbf{c} - \mathbf{T}\mathbf{a}_{t-1|t-1}$ . The third line uses matrix-inversion formulas by Henderson and Searle (1981, eqns. 9–11) to ensure that  $\mathbf{Q}^{-1}$  no longer appears, such that by a limiting argument the result remains valid even when  $\mathbf{Q}$  is singular. The last line employs the definitions of  $\mathbf{a}_{t|t-1}$  and  $\mathbf{I}_{t|t-1}$  in Table 3.

## F Derivation of equation (15)

Computing the first-order condition in equation (15), with respect to  $\mathbf{a}_{t-1}$ , we obtain

$$\mathbf{0} = \mathbf{T}'\mathbf{Q}^{-1}(\mathbf{a}_t - \mathbf{c} - \mathbf{T}\mathbf{a}_{t-1}) - \mathbf{I}_{t-1|t-1}(\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1}), \quad (\text{F.1})$$

the solution of which reads

$$\mathbf{a}_{t-1}^* = (\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T})^{-1}\{\mathbf{I}_{t-1|t-1}\mathbf{a}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}(\mathbf{a}_t - \mathbf{c})\}, \quad (\text{F.2})$$

which depends linearly on  $\mathbf{a}_t$ . In principle, equation (15) in the main text can be obtained by substituting equation (F.2) into equation (13) and performing algebraic manipulations. The desired result can be obtained more elegantly by ‘completing the square’ as follows. First, we replace  $\mathbf{a}_{t-1}$  with  $\mathbf{a}_{t-1}^*$  in equation (13), which then contains the following terms:

$$-\frac{1}{2}(\mathbf{a}_t - \mathbf{c} - \mathbf{T}\mathbf{a}_{t-1}^*)'\mathbf{Q}^{-1}(\mathbf{a}_t - \mathbf{c} - \mathbf{T}\mathbf{a}_{t-1}^*) - \frac{1}{2}(\mathbf{a}_{t-1}^* - \mathbf{a}_{t-1|t-1})'\mathbf{I}_{t-1|t-1}(\mathbf{a}_{t-1}^* - \mathbf{a}_{t-1|t-1}). \quad (\text{F.3})$$

Then we recall from equation (F.2) that  $\mathbf{a}_{t-1}^*$  is linear in  $\mathbf{a}_t$ , such that the collection of terms in equation (F.3) above is at most multivariate quadratic in  $\mathbf{a}_t$ . Hence, we should be able to rewrite equation (F.3) as a quadratic function (i.e., by completing the square) as follows:

$$-\frac{1}{2}(\mathbf{a}_t - \mathbf{a}_{t|t-1})'\mathbf{I}_{t|t-1}(\mathbf{a}_t - \mathbf{a}_{t|t-1}) + \text{constants}, \quad (\text{F.4})$$

for some vector  $\mathbf{a}_{t|t-1}$  to be found and some matrix  $\mathbf{I}_{t|t-1}$  to be determined.

To do this, we note that  $\mathbf{a}_{t|t-1}$  represents the argmax of equation (F.4), which can most readily be found by differentiating equation (F.3) with respect to  $\mathbf{a}_t$  and setting the result to zero. Using the envelope theorem, we need not account for the fact that  $\mathbf{a}_{t-1}^*$  depends on  $\mathbf{a}_t$  (the first derivative with respect to  $\mathbf{a}_{t-1}^*$  is zero because  $\mathbf{a}_{t-1}^*$  is optimal). Thus we set the derivative of equation (F.3) with respect to  $\mathbf{a}_t$  equal to zero, which gives  $\mathbf{0} = \mathbf{a}_t - \mathbf{c} - \mathbf{T}\mathbf{a}_{t-1}^*$ , or, by substituting  $\mathbf{a}_{t-1}^*$  from equation (F.2), we obtain

$$\begin{aligned} \mathbf{0} &= \mathbf{a}_t - \mathbf{c} - \mathbf{T}[\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T}]^{-1}\mathbf{I}_{t-1|t-1}\mathbf{a}_{t-1|t-1} \\ &\quad - \mathbf{T}[\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T}]^{-1}\mathbf{T}'\mathbf{Q}^{-1}(\mathbf{a}_t - \mathbf{c}). \end{aligned} \quad (\text{F.5})$$

The solution to this equation reads  $\mathbf{a}_{t|t-1} := \mathbf{T}\mathbf{a}_{t-1|t-1} + \mathbf{c}$ , which confirms the expression in Table 3.

Next, we compute the negative second derivative of equation (F.3) with respect to  $\mathbf{a}_t$ , which should give us  $\mathbf{I}_{t|t-1}$ . To account for the dependence of  $\mathbf{a}_{t-1}^*$  on  $\mathbf{a}_t$ , we use the chain rule. Specifically, in equation (F.2),  $\mathbf{a}_{t-1}^*$  is linear in  $\mathbf{a}_t$ , with the following Jacobian matrix:

$$\mathbf{J} := \frac{d\mathbf{a}_{t-1}^*}{d\mathbf{a}_t'} = [\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T}]^{-1}\mathbf{T}'\mathbf{Q}^{-1}. \quad (\text{F.6})$$

Next, the chain rule tells us that

$$\frac{d^2}{d\mathbf{a}_t d\mathbf{a}_t'} = \begin{bmatrix} \mathbf{1}_{m \times m} \\ \mathbf{J} \end{bmatrix}' \begin{bmatrix} \frac{\partial^2}{\partial \mathbf{a}_t \partial \mathbf{a}_t'} & \frac{\partial^2}{\partial \mathbf{a}_t \partial \mathbf{a}_{t-1}^{*'}} \\ \frac{\partial^2}{\partial \mathbf{a}_{t-1}^* \partial \mathbf{a}_t'} & \frac{\partial^2}{\partial \mathbf{a}_{t-1}^* \partial \mathbf{a}_{t-1}^{*'}} \end{bmatrix} \begin{bmatrix} \mathbf{1}_{m \times m} \\ \mathbf{J} \end{bmatrix}, \quad (\text{F.7})$$

where instances of  $\partial$  and  $d$  denote ‘partial’ and ‘total’ derivatives, respectively, while  $\mathbf{1}_{m \times m}$  denotes an identity matrix. As before, the envelope theorem ensures that no *first* derivative with respect to  $\mathbf{a}_t^*$  appears. When applying equation (F.7), we find that the negative second derivative of equation (F.3) becomes

$$\begin{aligned} & \begin{bmatrix} \mathbf{1}_{m \times m} \\ \mathbf{J} \end{bmatrix}' \begin{bmatrix} \mathbf{Q}^{-1} & -\mathbf{Q}^{-1}\mathbf{T} \\ -\mathbf{T}'\mathbf{Q}^{-1} & \mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T} \end{bmatrix} \begin{bmatrix} \mathbf{1}_{m \times m} \\ \mathbf{J} \end{bmatrix} \\ &= \mathbf{Q}^{-1} - \underbrace{\mathbf{Q}^{-1}\mathbf{T}\mathbf{J}} - \underbrace{\mathbf{J}'\mathbf{T}'\mathbf{Q}^{-1}} + \underbrace{\mathbf{J}'[\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T}]\mathbf{J}}, \\ &= \mathbf{Q}^{-1} - \mathbf{Q}^{-1}\mathbf{T}[\mathbf{I}_{t-1|t-1} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T}]^{-1}\mathbf{T}'\mathbf{Q}^{-1}. \end{aligned} \quad (\text{F.8})$$

In the last line, we have used the fact that all three terms with curly brackets equal  $\mathbf{Q}^{-1}\mathbf{T}[\mathbf{I}_{t|t} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T}]^{-1}\mathbf{T}'\mathbf{Q}^{-1}$ , such that two terms with curly brackets and opposite signs cancel, leaving only one term with a negative sign, which confirms the expression for  $\mathbf{I}_{t|t-1}$  in Table 3.

## G Kalman filter as a special case

Consider the linear Gaussian state-space model in Corollary 1. Suppose the inverse of the Kalman-filtered covariance matrix exists, i.e.  $\mathbf{P}_{t-1|t-1}^{-1} := \mathbf{I}_{t-1|t-1}$  exists. In Table 3, take the starting point  $\mathbf{a}_{t|t}^{(0)} = \mathbf{a}_{t|t-1}$ , and use Newton or Fisher optimisation steps. Given that the observation density is Gaussian, the log likelihood  $\ell(\mathbf{y}_t|\mathbf{a}_t)$  is multivariate quadratic in  $\mathbf{a}_t$ , such that the entire objective function (15) turns out to be multivariate quadratic in  $\mathbf{a}_t$ . The matrix of second derivatives is constant, such that Newton and Fisher optimisation steps are identical. Moreover, given the quadratic nature of the objective function, both methods find the location of the optimum in a single step. Indeed, the result is the classic Kalman filter, albeit written in the information form.

More explicitly, take  $\mathbf{y}_t = \mathbf{d} + \mathbf{Z}\boldsymbol{\alpha}_t + \boldsymbol{\varepsilon}_t$  with  $\boldsymbol{\varepsilon}_t \sim \text{i.i.d. } \mathbf{N}(\mathbf{0}, \mathbf{H})$ . Then

$$\ell(\mathbf{y}_t|\mathbf{a}_t) = -1/2(\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\mathbf{a}_t)'\mathbf{H}^{-1}(\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\mathbf{a}_t) + \text{constants}. \quad (\text{G.1})$$

The score and realised information are

$$\frac{d\ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t} = \mathbf{Z}'\mathbf{H}^{-1}(\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\mathbf{a}_t), \quad \frac{d^2\ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} = \mathbf{Z}'\mathbf{H}^{-1}\mathbf{Z}. \quad (\text{G.2})$$

As the realised information is constant, it equals the (expected) marginal information. Taking the starting point

$\mathbf{a}_{t|t}^{(0)} = \mathbf{a}_{t|t-1}$  for Newton's optimisation method, the estimate after a single Newton iteration reads

$$\mathbf{a}_{t|t}^{(1)} = \mathbf{a}_{t|t-1} + (\mathbf{I}_{t|t-1} + \mathbf{Z}'\mathbf{H}^{-1}\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{H}^{-1}(\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\mathbf{a}_{t|t-1}), \quad (\text{G.3})$$

which is exactly the Kalman filter level update written in information form. To see the equivalence with the covariance form of the Kalman filter, suppose that  $\mathbf{P}_{t|t-1} := \mathbf{I}_{t|t-1}^{-1}$  exists. Then, using a standard matrix-inversion formula (see e.g. Henderson and Searle, 1981, eqns. 9–10), the expression above is equivalent to

$$\mathbf{a}_{t|t}^{(1)} = \mathbf{a}_{t|t-1} + \mathbf{P}_{t|t-1}\mathbf{Z}'(\mathbf{Z}\mathbf{P}_{t|t-1}\mathbf{Z}' + \mathbf{H})^{-1}(\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\mathbf{a}_{t|t-1}), \quad (\text{G.4})$$

which is exactly the Kalman filter updating step (see e.g. Harvey, 1990, p. 106). For the information matrix update we have

$$\mathbf{I}_{t|t} = \mathbf{I}_{t|t-1} - \left. \frac{d^2 \ell(\mathbf{y}_t|\mathbf{a})}{d\mathbf{a} d\mathbf{a}'} \right|_{\mathbf{a}=\mathbf{a}_{t|t}} = \mathbf{I}_{t|t-1} + \mathbf{Z}'\mathbf{H}^{-1}\mathbf{Z}. \quad (\text{G.5})$$

If the inverses  $\mathbf{P}_{t|t-1} := \mathbf{I}_{t|t-1}^{-1}$  and  $\mathbf{P}_{t|t} := \mathbf{I}_{t|t}^{-1}$  exist, then, again using Henderson and Searle (1981, eq. 1), we find

$$\mathbf{P}_{t|t} = \mathbf{I}_{t|t}^{-1} = (\mathbf{I}_{t|t-1} + \mathbf{Z}'\mathbf{H}^{-1}\mathbf{Z})^{-1} = \mathbf{P}_{t|t-1} - \mathbf{P}_{t|t-1}\mathbf{Z}'(\mathbf{Z}\mathbf{P}_{t|t-1}\mathbf{Z}' + \mathbf{H})^{-1}\mathbf{Z}\mathbf{P}_{t|t-1}, \quad (\text{G.6})$$

which is exactly the Kalman filter covariance matrix updating step (again, see Harvey, 1990, p. 106).

## H Iterated extended Kalman filter as a special case

Consider the linear Gaussian state-space model in Corollary 1, except let  $\mathbf{y}_t = \mathbf{d} + \mathbf{Z}(\boldsymbol{\alpha}_t) + \boldsymbol{\varepsilon}_t$  for some nonlinear vector function  $\mathbf{Z}(\cdot)$  and  $\boldsymbol{\varepsilon}_t \sim \text{i.i.d. } \mathbf{N}(\mathbf{0}, \mathbf{H})$ . In Table 3, take the starting point  $\mathbf{a}_{t|t}^{(0)} = \mathbf{a}_{t|t-1}$  and perform Fisher optimisation steps, ignoring (i.e. setting to zero) all second-order derivatives of  $\mathbf{Z}(\cdot)$ . The iterated extended Kalman filter is then obtained as a special case.

More explicitly, take  $\mathbf{y}_t = \mathbf{d} + \mathbf{Z}(\boldsymbol{\alpha}_t) + \boldsymbol{\varepsilon}_t$  with  $\boldsymbol{\varepsilon}_t \sim \text{i.i.d. } \mathbf{N}(\mathbf{0}, \mathbf{H})$ . Here,  $\mathbf{Z}_t := \mathbf{Z}(\boldsymbol{\alpha}_t)$  is a column vector of the same size as  $\mathbf{y}_t$ , where each element of  $\mathbf{Z}_t$  depends on the elements of  $\boldsymbol{\alpha}_t$ . Then

$$\ell(\mathbf{y}_t|\mathbf{a}_t) = -1/2(\mathbf{y}_t - \mathbf{d} - \mathbf{Z}(\mathbf{a}_t))'\mathbf{H}^{-1}(\mathbf{y}_t - \mathbf{d} - \mathbf{Z}(\mathbf{a}_t)) + \text{constants}. \quad (\text{H.1})$$

The score and marginal information are similar to those in Appendix G, as long as  $\mathbf{Z}$  there is replaced by the Jacobian of the transformation from  $\boldsymbol{\alpha}_t$  to  $\mathbf{Z}_t$ , i.e.  $d\mathbf{Z}(\mathbf{a}_t)/d\mathbf{a}_t'$ . Hence

$$\frac{d\ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t} = \frac{d\mathbf{Z}'}{d\mathbf{a}_t} \mathbf{H}^{-1}(\mathbf{y}_t - \mathbf{d} - \mathbf{Z}(\mathbf{a}_t)), \quad (\text{H.2})$$

$$\frac{d^2\ell(\mathbf{y}_t|\mathbf{a}_t)}{d\mathbf{a}_t d\mathbf{a}_t'} = -\frac{d\mathbf{Z}'}{d\mathbf{a}_t} \mathbf{H}^{-1} \frac{d\mathbf{Z}}{d\mathbf{a}_t'} + \text{second-order derivatives}. \quad (\text{H.3})$$

The iterated extended Kalman filter (IEKF) is obtained from the Bellman filter by choosing Newton's method and by making one further simplifying approximation: namely that all second-order derivatives of elements of  $\mathbf{Z}_t$  with respect to the elements of  $\boldsymbol{\alpha}_t$  are zero. It is not obvious under what circumstances this approximation is justified, but here we are interested only in showing that the IEKF is a special case of the Bellman filter. Higher-order IEKFs may be obtained by retaining the second-order derivatives. If the observation noise  $\boldsymbol{\varepsilon}_t$  is heavy tailed, however, the Bellman filter in Table 3 suggests a 'robustified' version of the Kalman filter and its extensions, in which case the tail behaviour of  $p(\mathbf{y}_t|\mathbf{a}_t)$  is accounted for in the optimisation step by using the score  $d\ell(\mathbf{y}_t|\mathbf{a}_t)/d\mathbf{a}_t$ .

## I Fahrmeir's approximate mode estimator as a special case

When considering an observation density  $p(\mathbf{y}_t|\mathbf{a}_t)$  from the exponential family and taking just one optimisation step, we recover Fahrmeir's (1992) approximate mode estimator. Our analysis differs from Fahrmeir's in that (a) we show that online mode estimation can in theory be performed exactly by solving Bellman's equation, (b) we consider a general (rather than exponential) observation distribution, and (c) we allow more than one optimisation step.

## J Laplace Gaussian filter as a special case

When the state-transition density is linear and Gaussian, step 4 in the algorithm of Koyama et al. (2010) can be performed in closed form. The first-order Laplace Gaussian filter in step three of their algorithm is then equivalent to maximisation (16). Both algorithms differ when the state transition is nonlinear and/or non-Gaussian.

## K Implicit stochastic gradient method as a special case

In model (12), suppose that  $\mathbf{c} = \mathbf{0}$ ,  $\mathbf{Q} = \mathbf{0}$  and  $\mathbf{T} = \mathbb{1}_{m \times m}$ , where  $\mathbb{1}_{m \times m}$  is an  $m \times m$  identity matrix. The (constant) state  $\boldsymbol{\alpha}_t = \boldsymbol{\alpha}_1$  for all  $t = 1, 2, \dots$  now represents an unknown parameter to be estimated recursively over time. The prediction step of the Bellman filter simplifies to  $\mathbf{a}_{t|t-1} = \mathbf{a}_{t-1|t-1}$  and  $\mathbf{I}_{t|t-1} = \mathbf{I}_{t-1|t-1}$ , while update (16) equates to an implicit stochastic gradient method (e.g. Toulis and Airolidi, 2015, Toulis et al., 2016, Toulis and Airolidi, 2017, Toulis et al., 2021). In this case, the Bellman filter with BHHH updating steps becomes an implicit version of the (explicit) stochastic gradient methods in Amari et al. (2000, eq. 2.14) or Toulis and Airolidi (2017, eq. 11). While such methods are asymptotically convergent to the true parameter value, the Bellman filter typically remains perpetually responsive.

## L Proof of Theorem 1

1. The objective function  $V_t(\mathbf{a}) := \ell(\mathbf{y}_t|\mathbf{a}) - 1/2 \|\mathbf{a} - \mathbf{a}_{t|t-1}\|_{\mathbf{I}_{t|t-1}}^2$  is strongly concave with probability one because  $\ell(\mathbf{y}_t|\cdot)$  is concave with probability one (Assumption 1a), while  $-1/2 \|\mathbf{a} - \mathbf{a}_{t|t-1}\|_{\mathbf{I}_{t|t-1}}^2$  is strongly concave. Because the objective function is also real valued,  $\mathbf{a}_{t|t}$  is well defined. Moreover,  $V_t(\mathbf{a}_{t|t}) \geq V_t(\mathbf{a}_{t|t-1}) = \ell(\mathbf{y}_t|\mathbf{a}_{t|t-1})$ , i.e.

$$0 \leq V_t(\mathbf{a}_{t|t}) - V_t(\mathbf{a}_{t|t-1}) = \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) - \frac{1}{2} \|\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1}\|_{\mathbf{I}_{t|t-1}}^2 - \ell(\mathbf{y}_t|\mathbf{a}_{t|t-1}). \quad (\text{L.1})$$

Re-arranging gives

$$\frac{1}{2} \|\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1}\|_{\mathbf{I}_{t|t-1}}^2 \leq \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) - \ell(\mathbf{y}_t|\mathbf{a}_{t|t-1}). \quad (\text{L.2})$$

The right-hand side is bounded because the set  $\{\mathbf{a} \in \mathbb{R}^m : V_t(\mathbf{a}) \geq V_t(\mathbf{a}_{t|t-1})\}$  is bounded.

2. Assuming that  $\mathbf{a} \mapsto \ell(\mathbf{y}_t|\mathbf{a})$  is twice continuously differentiable (Assumption 2b), the following first- and second-order conditions must hold at the Bellman-filtered state  $\mathbf{a}_{t|t} \in \mathbb{R}^m$ :

$$\text{first-order condition:} \quad \nabla \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) - \mathbf{I}_{t|t-1}(\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1}) = \mathbf{0}_m, \quad (\text{L.3})$$

$$\text{second-order condition:} \quad \nabla^2 \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) - \mathbf{I}_{t|t-1} \leq \mathbf{0}_{m \times m}, \quad (\text{L.4})$$

where the weak inequality in the second line means the matrix on the left-hand side is negative semi-definite. Differentiating the first-order condition with respect to  $\mathbf{a}_{t|t-1}$ , we obtain

$$\nabla^2 \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) \frac{d\mathbf{a}_{t|t}}{d\mathbf{a}'_{t|t-1}} = \mathbf{I}_{t|t-1} \left[ \frac{d\mathbf{a}_{t|t}}{d\mathbf{a}'_{t|t-1}} - \mathbb{1}_{m \times m} \right], \quad (\text{L.5})$$

which can be re-written as

$$\frac{d\mathbf{a}_{t|t}}{d\mathbf{a}'_{t|t-1}} = [\mathbf{I}_{t|t-1} - \nabla^2 \ell(\mathbf{y}_t|\mathbf{a}_{t|t})]^{-1} \mathbf{I}_{t|t-1}, \quad (\text{L.6})$$

where the required inverse exists because  $\mathbf{I}_{t|t-1} - \nabla^2 \ell(\mathbf{y}_t|\mathbf{a}_{t|t})$  is positive definite by assumption.

Next, we use a result of Wang and Gong (1993, eq. 2), which says that  $\lambda_{\min}(\mathbf{A})\lambda_{\min}(\mathbf{B}) \leq \lambda_{\min}(\mathbf{AB})$  for two square, symmetric and positive semidefinite matrices  $\mathbf{A}$  and  $\mathbf{B}$ , where  $\lambda_{\min}(\cdot)$  denotes the smallest eigenvalue of a matrix. Denoting  $\mathbf{H}_t := -\nabla^2 \ell(\mathbf{y}_t|\mathbf{a}_{t|t})$  and applying this result to  $(\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1} \mathbf{I}_{t|t-1}$  yields

$$0 < \frac{\lambda_{\min}(\mathbf{I}_{t|t-1})}{\lambda_{\max}(\mathbf{I}_{t|t-1} + \mathbf{H}_t)} = \lambda_{\min}[(\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1}] \lambda_{\min}(\mathbf{I}_{t|t-1}) \leq \lambda_{\min}[(\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1} \mathbf{I}_{t|t-1}]. \quad (\text{L.7})$$

Hence, the eigenvalues of  $(\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1} \mathbf{I}_{t|t-1}$  are strictly positive. To show that the eigenvalues of  $(\mathbf{I}_{t|t-1} +$

$\mathbf{H}_t)^{-1}\mathbf{I}_{t|t-1}$  are bounded above by one, we note that

$$\begin{aligned}\lambda_{\max}[(\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1}\mathbf{I}_{t|t-1}] &= \lambda_{\max}[\mathbb{1}_{m \times m} - (\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1}\mathbf{H}_t], \\ &= 1 - \lambda_{\min}[(\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1}\mathbf{H}_t], \\ &\leq 1 - \lambda_{\min}[(\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1}]\lambda_{\min}(\mathbf{H}_t), \\ &= 1 - \frac{\lambda_{\min}(\mathbf{H}_t)}{\lambda_{\max}(\mathbf{I}_{t|t-1} + \mathbf{H}_t)} \leq 1 - \frac{\lambda_{\min}(\mathbf{H}_t)}{\lambda_{\max}(\mathbf{I}_{t|t-1}) + \lambda_{\max}(\mathbf{H}_t)},\end{aligned}\quad (\text{L.8})$$

which does not exceed (is strictly smaller than) than unity if  $\mathbf{H}_t \geq 0$  ( $\mathbf{H}_t > 0$ ). The conditions  $\mathbf{H}_t \geq 0$  or  $\mathbf{H}_t > 0$  are ensured, respectively, if the observation log density is concave (Assumption 1a) or strictly concave (Assumption 1b).

Next, we use the well known fact (e.g. Jungers, 2009, p. 39) that the induced matrix norm satisfies

$$\|\mathbf{M}\|_{\mathbf{W}} = \|\mathbf{W}^{1/2}\mathbf{M}\mathbf{W}^{-1/2}\| = \sqrt{\lambda_{\max}(\mathbf{W}^{1/2}\mathbf{M}\mathbf{W}^{-1}\mathbf{M}'\mathbf{W}^{1/2})} = \sqrt{\lambda_{\max}(\mathbf{M}\mathbf{W}^{-1}\mathbf{M}'\mathbf{W})},$$

where the last equality follows by cyclically rotating inside the  $\lambda_{\max}(\cdot)$  operator. Here  $\mathbf{M}, \mathbf{W} \in \mathbb{R}^{m \times m}$  and  $\mathbf{W} > \mathbf{0}$  is the positive definite weight matrix. Using this fact along with the symmetry of  $\mathbf{I}_{t|t-1}$  and  $\mathbf{H}_t$ , we then obtain

$$\begin{aligned}\left\| \frac{d\mathbf{a}_{t|t}}{d\mathbf{a}'_{t|t-1}} \right\|_{\mathbf{I}_{t|t-1}} &= \left\| (\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1}\mathbf{I}_{t|t-1} \right\|_{\mathbf{I}_{t|t-1}}, \\ &= \sqrt{\lambda_{\max} \left\{ (\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1}\mathbf{I}_{t|t-1}\mathbf{I}_{t|t-1}^{-1}\mathbf{I}_{t|t-1}(\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1}\mathbf{I}_{t|t-1} \right\}}, \\ &= \sqrt{\lambda_{\max} \left\{ [(\mathbf{I}_{t|t-1} + \mathbf{H}_t)^{-1}\mathbf{I}_{t|t-1}]^2 \right\}} \leq 1 - \frac{\lambda_{\min}(\mathbf{H}_t)}{\lambda_{\max}(\mathbf{I}_{t|t-1}) + \lambda_{\max}(\mathbf{H}_t)},\end{aligned}\quad (\text{L.9})$$

where we have used equation (L.8) along with the fact that the eigenvalues of the square of a matrix are equal to the squares of the eigenvalues of the original matrix. If additionally Assumption 1a (1b) holds, then we have  $\lambda_{\min}(\mathbf{H}_t) \geq 0$  ( $\lambda_{\min}(\mathbf{H}_t) > 0$ ), such that the right-hand side does not exceed (is strictly less than) unity.

3. Assuming that  $\mathbf{a} \mapsto \ell(\mathbf{y}_t|\mathbf{a})$  is strongly concave with parameter  $\epsilon > 0$  (Assumption 1c) and once continuously differentiable (Assumption 2a), standard arguments (e.g. Nesterov, 2003, eq. 2.1.17) give

$$\langle \mathbf{a}_t - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\mathbf{a}_t) - \nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) \rangle \leq -\epsilon \cdot \|\mathbf{a}_t - \boldsymbol{\alpha}_t\|^2, \quad \forall \mathbf{a}_t, \boldsymbol{\alpha}_t \in \mathbb{R}^m. \quad (\text{L.10})$$

Strong concavity means that equation (L.10) holds for all pairs  $\mathbf{a}_t, \boldsymbol{\alpha}_t \in \mathbb{R}^m$ , but we shall need it only when  $\boldsymbol{\alpha}_t$  is the true state. Assuming differentiability (Assumption 2a), the first-order condition  $\mathbf{I}_{t|t-1}(\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1}) = \nabla \ell(\mathbf{y}_t|\mathbf{a}_{t|t})$  is rewritten by pre-multiplying the equation by  $\mathbf{I}_{t|t-1}^{-1/2}$  and subtracting  $\mathbf{I}_{t|t-1}^{1/2}\boldsymbol{\alpha}_t - \mathbf{I}_{t|t-1}^{-1/2}\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)$  from both sides to obtain

$$\mathbf{I}_{t|t-1}^{1/2}(\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t) - \mathbf{I}_{t|t-1}^{-1/2} \{ \nabla \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) - \nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) \} = \mathbf{I}_{t|t-1}^{1/2}(\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t) + \mathbf{I}_{t|t-1}^{-1/2}\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t). \quad (\text{L.11})$$

Computing the quadratic norm on both sides and ignoring one term on the left, we obtain an inequality as follows:

$$\begin{aligned}\|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2 - 2 \langle \mathbf{a}_{t|t} - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) - \nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) \rangle \\ \leq \|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2 + 2 \langle \mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) \rangle + \|\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)\|_{\mathbf{I}_{t|t-1}^{-1}}^2.\end{aligned}$$

By strong concavity (L.10), we have

$$\begin{aligned}\|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2 + 2\epsilon \cdot \|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|^2 \\ \leq \|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2 + 2 \langle \mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) \rangle + \|\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)\|_{\mathbf{I}_{t|t-1}^{-1}}^2.\end{aligned}\quad (\text{L.12})$$

Taking expectations yields

$$\begin{aligned} \mathbb{E} \left( \|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2 \right) + 2\epsilon \mathbb{E} \left( \|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|^2 \right) \\ \leq \mathbb{E} \left( \|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2 \right) + \mathbb{E} \left( \|\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)\|_{\mathbf{I}_{t|t-1}^{-1}}^2 \right). \end{aligned} \quad (\text{L.13})$$

where we have used  $\mathbb{E} \langle \mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) \rangle = 0$ , which is obvious from the expectation of the score being zero, i.e.  $\mathbb{E}[\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)|\boldsymbol{\alpha}_t] = 0$ . Finally, the theorem is proved by noting that the left-hand side is  $\mathbb{E} \left( \|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1} + 2\epsilon \mathbf{1}_{m \times m}}^2 \right)$ , where  $\mathbf{1}_{m \times m}$  is an  $m \times m$  identity matrix, while Assumption 3 together with the assumed positive definiteness of  $\mathbf{I}_{t|t-1}$  implies that on the right-hand side we have

$$\mathbb{E} \left( \|\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)\|_{\mathbf{I}_{t|t-1}^{-1}}^2 \right) \leq \sigma^2 / \lambda_{\min}.$$

## M Comparison of Theorem 1 with Toulis et al. (2016)

This section casts light on the different definitions of strong concavity used in Theorem 1 and in Toulis et al. (2016). Here we show that Theorem 1 applies to e.g. the Kalman filter, while the seemingly stronger result in Toulis et al. (2016) does not.

By the combination of Assumptions 1c (strong concavity) and 2b (twice differentiability), part 3 of Theorem 1 assumes that the negative Hessian  $-\nabla^2 \ell(\mathbf{y}_t|\mathbf{a})$  is strictly positive definite with smallest eigenvalue  $\epsilon > 0$ . Standard arguments (e.g. Nesterov, 2003, eq. 2.1.17) imply that

$$\langle \mathbf{a}_t - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\mathbf{a}_t) - \nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) \rangle \leq -\epsilon \cdot \|\mathbf{a}_t - \boldsymbol{\alpha}_t\|^2, \quad \forall \mathbf{a}_t, \boldsymbol{\alpha}_t \in \mathbb{R}^m. \quad (\text{M.1})$$

Toulis et al. (2016) take a different view on strong concavity, defining a log-likelihood function to be strongly concave, for a typical observation  $\mathbf{y}_t \in \mathbb{R}^l$ , when

$$\text{strong concavity in Toulis et al. (2016): } \langle \mathbf{a}_t - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\mathbf{a}_t) \rangle \leq -\epsilon \cdot \|\mathbf{a}_t - \boldsymbol{\alpha}_t\|^2, \quad \forall \mathbf{a}_t, \boldsymbol{\alpha}_t \in \mathbb{R}^m, \quad (\text{M.2})$$

which differs from definition (M.1) in that the term  $\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)$  is no longer present. Inequality (M.2) appears in Remark 2 and equation 17 of the supplementary material to Toulis et al. (2016), where  $\mu_t > 0$  appears instead of our  $\epsilon$ , the random draw  $\xi_t$  appears instead of our  $\mathbf{y}_t$ ,  $\theta_t$  appears instead of our  $\mathbf{a}_t$ , the true value  $\theta_*$  appears instead of our  $\boldsymbol{\alpha}_t$ , their  $L$  is a negative log-likelihood function, and index  $n$  is used instead of our  $t$ . Toulis et al. (2016) permit the parameter of strong concavity to depend on the observation; for simplicity, we do not. The term  $\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t)$ , which appears in equation (M.1) but not equation (M.2), is the score function evaluated at the true parameter; hence, this term is zero on average. For many models of interest, however, realisations of the score are non-zero with probability one, such that definition (M.2) materially differs from (M.1).

While definition (M.1) of strong concavity was used in the proof of Theorem 1, definition (M.2) allows a stronger result due to Toulis et al. (2016) to be derived. First, the first-order condition corresponding to maximisation (16), i.e.  $\mathbf{I}_{t|t-1}(\mathbf{a}_{t|t} - \mathbf{a}_{t|t-1}) = \nabla \ell(\mathbf{y}_t|\mathbf{a}_{t|t})$ , is rewritten as

$$\mathbf{I}_{t|t-1}^{1/2}(\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t) - \mathbf{I}_{t|t-1}^{-1/2} \nabla \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) = \mathbf{I}_{t|t-1}^{1/2}(\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t). \quad (\text{M.3})$$

Computing the quadratic norm on both sides, we have

$$\|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2 - 2 \langle \mathbf{a}_{t|t} - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\mathbf{a}_{t|t}) \rangle + \|\nabla \ell(\mathbf{y}_t|\mathbf{a}_{t|t})\|_{\mathbf{I}_{t|t-1}^{-1}}^2 = \|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2. \quad (\text{M.4})$$

By strong concavity (M.2), it follows that

$$\|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2 + 2\epsilon \cdot \|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|^2 + \|\nabla \ell(\mathbf{y}_t|\mathbf{a}_{t|t})\|_{\mathbf{I}_{t|t-1}^{-1}}^2 \leq \|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2. \quad (\text{M.5})$$

Ignoring the third term on the left-hand side and combining terms, we find

$$\|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1} + 2\epsilon \mathbf{1}_{m \times m}}^2 \leq \|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|_{\mathbf{I}_{t|t-1}}^2, \quad (\text{M.6})$$

where  $\mathbf{1}_{m \times m}$  denotes an  $m \times m$  identity matrix. In Toulis et al. (2016, p. 1291) it holds that  $\mathbf{I}_{t|t-1} = \gamma^{-1} \mathbf{1}_{m \times m}$ ,

where  $\mathbf{1}_{m \times m}$  is an  $m \times m$  identity matrix and  $\gamma > 0$  is a learning parameter, in which case we obtain

$$\|\mathbf{a}_{t|t} - \boldsymbol{\alpha}_t\|^2 \leq \frac{1}{1 + 2\gamma\epsilon} \|\mathbf{a}_{t|t-1} - \boldsymbol{\alpha}_t\|^2, \quad (\text{M.7})$$

as in Toulis et al. (2016, p. 1291). This result is stronger than that in Theorem 1, because (M.7) holds for all realisations  $\mathbf{y}_t$ , without taking expectations. Inequality (M.7) implies that the update is ‘contracting almost surely’ (Toulis et al., 2016, p. 1291). Unfortunately, this desirable property is not observed in practice for e.g. the Kalman filter.

To explain why the Kalman filter fails to be almost surely contractive in the sense of Toulis et al. (2016), we observe that the Kalman filter satisfies our assumption (M.1) as used in Theorem 1, but *not* assumption (M.2) as used by Toulis et al. (2016). To demonstrate this, we take the linear Gaussian state-space model in Corollary 1, such that the observation density  $p(\mathbf{y}_t|\boldsymbol{\alpha}_t)$  is Gaussian with mean  $\mathbf{d} + \mathbf{Z}\boldsymbol{\alpha}_t$  and covariance matrix  $\mathbf{H}$ , which is assumed positive definite. The log-likelihood function and its gradient then read

$$\ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) = -\frac{1}{2}(\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\boldsymbol{\alpha}_t)' \mathbf{H}^{-1} (\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\boldsymbol{\alpha}_t) + \text{constants}, \quad (\text{M.8})$$

$$\nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) = \mathbf{Z}' \mathbf{H}^{-1} (\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\boldsymbol{\alpha}_t). \quad (\text{M.9})$$

The multivariate Gaussian is strongly concave according to our definition (M.1), because

$$\begin{aligned} \langle \mathbf{a}_t - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\mathbf{a}_t) - \nabla \ell(\mathbf{y}_t|\boldsymbol{\alpha}_t) \rangle &= \langle \mathbf{a}_t - \boldsymbol{\alpha}_t, \mathbf{Z}' \mathbf{H}^{-1} (\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\mathbf{a}_t) - \mathbf{Z}' \mathbf{H}^{-1} (\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\boldsymbol{\alpha}_t) \rangle, \\ &= -\langle \mathbf{a}_t - \boldsymbol{\alpha}_t, \mathbf{Z}' \mathbf{H}^{-1} \mathbf{Z} (\mathbf{a}_t - \boldsymbol{\alpha}_t) \rangle, \\ &= -\|\mathbf{a}_t - \boldsymbol{\alpha}_t\|_{\mathbf{Z}' \mathbf{H}^{-1} \mathbf{Z}}^2, \\ &\leq -\lambda_{\min}(\mathbf{Z}' \mathbf{H}^{-1} \mathbf{Z}) \cdot \|\mathbf{a}_t - \boldsymbol{\alpha}_t\|^2, \end{aligned} \quad (\text{M.10})$$

where  $\lambda_{\min}(\cdot)$  denotes the smallest eigenvalues of a matrix. Hence, condition (M.1) is satisfied with  $\epsilon = \lambda_{\min}(\mathbf{Z}' \mathbf{H}^{-1} \mathbf{Z}) > 0$ . Conversely, the multivariate Gaussian fails to be strongly concave when using the alternative definition (M.2) of Toulis et al. (2016), because

$$\langle \mathbf{a}_t - \boldsymbol{\alpha}_t, \nabla \ell(\mathbf{y}_t|\mathbf{a}_t) \rangle = \langle \mathbf{a}_t - \boldsymbol{\alpha}_t, \mathbf{Z}' \mathbf{H}^{-1} (\mathbf{y}_t - \mathbf{d} - \mathbf{Z}\mathbf{a}_t) \rangle \not\leq -\text{positive scalar} \cdot \|\mathbf{a}_t - \boldsymbol{\alpha}_t\|^2. \quad (\text{M.11})$$

Stepping back, it is not too surprising that the almost sure contractive property of Toulis et al. (2016) fails for the Kalman filter, because the Kalman filter can (and does) move in the wrong direction when confronted with atypical observations. The contribution of Theorem 1 is to demonstrate that, in a general context, such ‘bad’ behaviour does not dominate. Theorem 1 allows for the fact that updates may be less accurate than predictions, while still ensuring that the updates are contractive in quadratic mean towards a noise-dominated region around the true state, which is the situation that is relevant in practice.

## N Proof of Proposition 2

Repeated self-substitution of the recursions (21) yields:

$$\begin{aligned} \text{MSE}_{t|t} &\leq \left(\frac{\gamma}{\gamma + 2\epsilon}\right)^t \text{MSE}_{1|0} + \frac{\sigma^2}{\gamma^2} \sum_{i=1}^t \left(\frac{\gamma}{\gamma + 2\epsilon}\right)^i + \sigma_\eta^2 \sum_{i=1}^{t-1} \left(\frac{\gamma}{\gamma + 2\epsilon}\right)^i, \\ &= \left(\frac{\gamma}{\gamma + 2\epsilon}\right)^t \text{MSE}_{1|0} + \frac{\sigma^2}{\gamma^2} \left(\frac{\gamma}{\gamma + 2\epsilon}\right) \frac{1 - \left(\frac{\gamma}{\gamma + 2\epsilon}\right)^t}{1 - \frac{\gamma}{\gamma + 2\epsilon}} + \sigma_\eta^2 \left(\frac{\gamma}{\gamma + 2\epsilon}\right) \frac{1 - \left(\frac{\gamma}{\gamma + 2\epsilon}\right)^{t-1}}{1 - \frac{\gamma}{\gamma + 2\epsilon}}, \end{aligned}$$

where the second line employs  $\sum_{i=1}^t x^{i-1} = (1 - x^t)/(1 - x)$  for  $-1 < x < 1$ . Using  $\gamma, \epsilon > 0$  and taking the limit  $t \rightarrow \infty$  yields equation (22).

## O Proof of Theorem 2

By the chain rule, we have

$$\begin{aligned} \left\| \frac{d\mathbf{a}_{t|t}}{d\mathbf{a}'_{0|0}} \right\|_{\mathbf{I}} &= \left\| \frac{d\mathbf{a}_{t|t}}{d\mathbf{a}'_{t|t-1}} \frac{d\mathbf{a}_{t|t-1}}{d\mathbf{a}'_{t-1|t-1}} \times \dots \times \frac{d\mathbf{a}_{1|1}}{d\mathbf{a}'_{1|0}} \frac{d\mathbf{a}_{1|0}}{d\mathbf{a}'_{0|0}} \right\|_{\mathbf{I}} \leq \left\| \frac{d\mathbf{a}_{t|t}}{d\mathbf{a}'_{t|t-1}} \right\|_{\mathbf{I}} \|\mathbf{T}\|_{\mathbf{I}} \times \dots \times \left\| \frac{d\mathbf{a}_{1|1}}{d\mathbf{a}'_{1|0}} \right\|_{\mathbf{I}} \|\mathbf{T}\|_{\mathbf{I}}, \\ &\leq (\|\mathbf{T}\|_{\mathbf{I}})^t \prod_{\tau=1}^t \left( 1 - \frac{\lambda_{\min}(\mathbf{H}_{\tau})}{\lambda_{\max}(\mathbf{I}) + \lambda_{\max}(\mathbf{H}_{\tau})} \right) \leq (\|\mathbf{T}\|_{\mathbf{I}})^t \left( 1 - \frac{\mu_{\min}}{\nu_{\max} + \mu_{\max}} \right)^t. \end{aligned} \quad (\text{O.1})$$

The inequality in the first line holds by the sub-multiplicative property of the induced matrix norm in combination with the linear prediction step. The second line holds by equation (L.9), where  $\mathbf{H}_t := -\nabla^2 \ell(\mathbf{y}_t | \mathbf{a}_{t|t})$ . The last inequality holds because  $\lambda_{\max}(\mathbf{I}) = \nu_{\max}$  and  $0 \leq \mu_{\min} \leq \lambda_{\min}(\mathbf{H}_t) \leq \lambda_{\max}(\mathbf{H}_t) \leq \mu_{\max}$  by assumption.

To prove equation (23), we must still bound the term  $\|\mathbf{T}\|_{\mathbf{I}}$ . To this end, we define  $\delta := \lambda_{\min}(\mathbf{I} - \mathbf{T}'\mathbf{I}\mathbf{T})' \in \mathbb{R}$ , which could be positive or negative. Since  $\mathbf{I}$  is positive definite, we must have

$$\delta = \lambda_{\min}(\mathbf{I} - \mathbf{T}'\mathbf{I}\mathbf{T}) \leq \lambda_{\min}(\mathbf{I}) = \nu_{\min}, \quad (\text{O.2})$$

so  $\delta \leq \nu_{\min}$ . Next, we have the inequality

$$\mathbf{0} \leq \mathbf{I} - \delta \mathbf{1}_{m \times m} - \mathbf{T}'\mathbf{I}\mathbf{T}, \quad (\text{O.3})$$

as we will use below. As  $\mathbf{I}$  is positive definite with smallest and largest eigenvalues  $\nu_{\min}$  and  $\nu_{\max}$  respectively, we have

$$\frac{1}{\nu_{\max}} \mathbf{I} \leq \mathbf{1}_{m \times m} \leq \frac{1}{\nu_{\min}} \mathbf{I}.$$

When  $\delta > 0$ , multiplying this sequence of inequalities by  $-\delta$  yields

$$\frac{-\delta}{\nu_{\max}} \mathbf{I} \geq -\delta \mathbf{1}_{m \times m} \geq \frac{-\delta}{\nu_{\min}} \mathbf{I}, \quad \delta > 0.$$

When  $\delta < 0$ , we obtain instead

$$\frac{-\delta}{\nu_{\max}} \mathbf{I} \leq -\delta \mathbf{1}_{m \times m} \leq \frac{-\delta}{\nu_{\min}} \mathbf{I}, \quad \delta < 0.$$

Combining the last two results, we see that  $-\delta \mathbf{1}_{m \times m}$  is bounded above by  $-\delta/\nu_{\max} \mathbf{I}$  when  $\delta > 0$  and  $-\delta/\nu_{\min} \mathbf{I}$  when  $\delta < 0$ . This means that for all  $\delta \in \mathbb{R}$ , we can write

$$-\delta \mathbf{1}_{m \times m} \leq -\min \left\{ \frac{\delta}{\nu_{\min}}, \frac{\delta}{\nu_{\max}} \right\} \mathbf{I}, \quad \delta \in \mathbb{R}. \quad (\text{O.4})$$

Using inequality (O.4), inequality (O.3) can be further extended as

$$\mathbf{0} \leq \mathbf{I} - \delta \mathbf{1}_{m \times m} - \mathbf{T}'\mathbf{I}\mathbf{T} \leq \left( 1 - \min \left\{ \frac{\delta}{\nu_{\min}}, \frac{\delta}{\nu_{\max}} \right\} \right) \mathbf{I} - \mathbf{T}'\mathbf{I}\mathbf{T}. \quad (\text{O.5})$$

Equation (O.5) shows that  $z^2 \mathbf{I} - \mathbf{T}'\mathbf{I}\mathbf{T} \geq \mathbf{0}$  for a particular value of  $z$ . This is useful because from Jungers (2009, p. 39) we have

$$\|\mathbf{T}\|_{\mathbf{I}} = \inf \{ z \geq 0 : z^2 \mathbf{I} - \mathbf{T}'\mathbf{I}\mathbf{T} \geq \mathbf{0} \}, \quad (\text{O.6})$$

which says that  $\|\mathbf{T}\|_{\mathbf{I}}$  is the infimum of such values. Hence equations (O.5) and (O.6) together imply

$$\|\mathbf{T}\|_{\mathbf{I}} \leq \sqrt{1 - \min \left\{ \frac{\delta}{\nu_{\min}}, \frac{\delta}{\nu_{\max}} \right\}}. \quad (\text{O.7})$$

As a sanity check, we may verify that the right-hand side is nonnegative, as when  $\delta > 0$  we have  $\delta \leq \nu_{\min}$  by equation (O.2) above. Substituting equation (O.7) in equation (O.1) yields equation (23) in the main text.

To prove equation (24) in the main text, compute the derivative of the logarithm of the right-hand side of

equation (23) as follows:

$$\frac{d}{dt} \log \left[ \left(1 - \frac{\delta}{\nu_{\min}}\right)^{t/2} \left(1 - \frac{\mu_{\min}}{\nu_{\max} + \mu_{\max}}\right)^t \right] = \frac{1}{2} \log \left(1 - \frac{\delta}{\nu_{\min}}\right) + \log \left(1 - \frac{\mu_{\min}}{\nu_{\max} + \mu_{\max}}\right). \quad (\text{O.8})$$

When this quantity is strictly negative, exponential almost sure convergence to zero follows.

## P Lemma involving quadratic functions

**Lemma 1.** *Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^m$ . Let  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times m}$  be symmetric positive definite matrices. Define  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  as*

$$f(\mathbf{x}) := \max_{\mathbf{y}} \left\{ -\frac{1}{2} \mathbf{x}' \mathbf{A} \mathbf{x} - \frac{1}{2} \mathbf{y}' \mathbf{B} \mathbf{y} + \mathbf{x}' \mathbf{C} \mathbf{y} + \mathbf{a}' \mathbf{x} + \mathbf{b}' \mathbf{y} \right\}, \quad (\text{P.1})$$

$$= \max_{\mathbf{y}} \left\{ -\frac{1}{2} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}' \begin{bmatrix} \mathbf{A} & -\mathbf{C} \\ -\mathbf{C}' & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} + \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}' \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \right\}, \quad (\text{P.2})$$

for  $\mathbf{C}, \mathbf{a}, \mathbf{b}$  of appropriate size. Then  $f(\mathbf{x})$  is multivariate quadratic with negative Hessian matrix  $\mathbf{A} - \mathbf{C} \mathbf{B}^{-1} \mathbf{C}'$ . When this negative Hessian is positive definite, the argmax of  $f(\mathbf{x})$  over  $\mathbf{x}$  equals  $(\mathbf{A} - \mathbf{C} \mathbf{B}^{-1} \mathbf{C}')^{-1} (\mathbf{a} + \mathbf{C} \mathbf{B}^{-1} \mathbf{b})$ .

*Proof.* Take  $\mathbf{x}$  as fixed. The first-order condition for the maximisation over  $\mathbf{y}$  reads  $\mathbf{0} = -\mathbf{B} \mathbf{y} + \mathbf{b} + \mathbf{C}' \mathbf{x}$ , which leads to  $\mathbf{y} = \mathbf{B}^{-1} (\mathbf{b} + \mathbf{C}' \mathbf{x})$ . Substituting the optimised value of  $\mathbf{y}$  into the expression for  $f(\mathbf{x})$  gives

$$f(\mathbf{x}) = -\frac{1}{2} \mathbf{x}' \mathbf{A} \mathbf{x} - \frac{1}{2} (\mathbf{b} + \mathbf{C}' \mathbf{x})' \mathbf{B}^{-1} (\mathbf{b} + \mathbf{C}' \mathbf{x}) + \mathbf{x}' \mathbf{C} \mathbf{B}^{-1} (\mathbf{b} + \mathbf{C}' \mathbf{x}) + \mathbf{a}' \mathbf{x} + \mathbf{b}' \mathbf{B}^{-1} (\mathbf{b} + \mathbf{C}' \mathbf{x}).$$

Several terms cancel and remaining terms can be grouped as

$$f(\mathbf{x}) = -\frac{1}{2} \mathbf{x}' (\mathbf{A} - \mathbf{C} \mathbf{B}^{-1} \mathbf{C}') \mathbf{x} + (\mathbf{a} + \mathbf{C} \mathbf{B}^{-1} \mathbf{b})' \mathbf{x} + \text{constants},$$

where constants independent of  $\mathbf{x}$  are ignored. When  $\mathbf{A} - \mathbf{C} \mathbf{B}^{-1} \mathbf{C}'$  is positive definite, this quadratic function of  $\mathbf{x}$  is maximised at  $(\mathbf{A} - \mathbf{C} \mathbf{B}^{-1} \mathbf{C}')^{-1} (\mathbf{a} + \mathbf{C} \mathbf{B}^{-1} \mathbf{b})$ , completing the proof.  $\square$

## Q Proof of Proposition 4

To derive a relation between  $\mathbf{a}_{t|n}$  and  $\mathbf{a}_{t+1|n}$  in the context of approximately quadratic value functions, it is useful to define a new value function  $U_{t,t+1}(\cdot, \cdot) : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ , which takes two state variables as input. This value function is defined using the partial sum (25), and can be rewritten using the value functions  $V_t(\cdot)$  and  $W_{t+1}(\cdot)$  defined in equations (26) and (27), respectively, as follows:

$$U_{t,t+1}(\mathbf{a}_t, \mathbf{a}_{t+1}) := \max_{\mathbf{a}_1, \dots, \mathbf{a}_{t-1}, \mathbf{a}_{t+2}, \dots, \mathbf{a}_n} L_{1:n}(\mathbf{a}_1, \dots, \mathbf{a}_n), \quad (\text{Q.1})$$

$$= \max_{\mathbf{a}_1, \dots, \mathbf{a}_{t-1}, \mathbf{a}_{t+2}, \dots, \mathbf{a}_n} [L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t) + \ell(\mathbf{a}_{t+1} | \mathbf{a}_t) + L_{t+1:n}(\mathbf{a}_{t+1}, \dots, \mathbf{a}_n)], \quad (\text{Q.2})$$

$$= \left[ \max_{\mathbf{a}_1, \dots, \mathbf{a}_{t-1}} L_{1:t}(\mathbf{a}_1, \dots, \mathbf{a}_t) \right] + \ell(\mathbf{a}_{t+1} | \mathbf{a}_t) + \left[ \max_{\mathbf{a}_{t+2}, \dots, \mathbf{a}_n} L_{t+1:n}(\mathbf{a}_{t+1}, \dots, \mathbf{a}_n) \right], \quad (\text{Q.3})$$

$$= V_t(\mathbf{a}_t) + \ell(\mathbf{a}_{t+1} | \mathbf{a}_t) + W_{t+1}(\mathbf{a}_{t+1}), \quad (\text{Q.4})$$

$$= -\frac{1}{2} \|\mathbf{a}_t - \mathbf{a}_{t|t}\|_{\mathbf{I}_{t|t}}^2 - \frac{1}{2} \|\mathbf{a}_{t+1} - \mathbf{c} - \mathbf{T} \mathbf{a}_t\|_{\mathbf{Q}^{-1}}^2 - \frac{1}{2} \|\mathbf{a}_{t+1} - \hat{\mathbf{a}}_{t+1|t+1}\|_{\hat{\mathbf{I}}_{t+1|t+1}}^2. \quad (\text{Q.5})$$

In the last line, we take a linear Gaussian state equation as in Corollary 1, and use the assumption that  $V_t(\mathbf{a}_t)$  is multivariate quadratic with argmax  $\mathbf{a}_{t|t}$  and negative Hessian matrix  $\mathbf{I}_{t|t}$ , while  $W_{t+1}(\mathbf{a}_{t+1})$  is similarly multivariate quadratic with argmax  $\hat{\mathbf{a}}_{t+1|t+1}$  and negative Hessian matrix  $\hat{\mathbf{I}}_{t+1|t+1}$ . Here, hats denote ‘backward filtered’ quantities. It follows that  $U_{t,t+1}(\cdot, \cdot)$  is a multivariate quadratic function in two state variables,  $\mathbf{a}_t$  and  $\mathbf{a}_{t+1}$ .

From definition (Q.1), it is clear that  $Z_t(\cdot)$  and  $Z_{t+1}(\cdot)$  defined in equation (28) can be recovered from  $U_{t,t+1}(\cdot, \cdot)$  as follows:

$$Z_t(\mathbf{a}_t) = \max_{\mathbf{a}_{t+1}} U_{t,t+1}(\mathbf{a}_t, \mathbf{a}_{t+1}), \quad (\text{Q.6})$$

$$Z_{t+1}(\mathbf{a}_{t+1}) = \max_{\mathbf{a}_t} U_{t,t+1}(\mathbf{a}_t, \mathbf{a}_{t+1}). \quad (\text{Q.7})$$

Since  $\mathbf{a}_{t|n} := \arg \max_{\mathbf{a}} Z_t(\mathbf{a})$  while  $\mathbf{a}_{t+1|n} := \arg \max_{\mathbf{a}} Z_{t+1}(\mathbf{a})$ , it is clear that  $U_{t,t+1}(\cdot, \cdot)$  is maximised when  $\mathbf{a}_t = \mathbf{a}_{t|n}$  and  $\mathbf{a}_{t+1} = \mathbf{a}_{t+1|n}$ . We evaluate  $U_{t,t+1}(\cdot, \cdot)$  at  $\mathbf{a}_{t+1} = \mathbf{a}_{t+1|n}$ . Subsequently, the first-order condition with respect to  $\mathbf{a}_t$  reads

$$\mathbf{0} = \mathbf{I}_{t|t}(\mathbf{a}_t - \mathbf{a}_{t|t}) - \mathbf{T}'\mathbf{Q}^{-1}(\mathbf{a}_{t+1|n} - \mathbf{c} - \mathbf{T}\mathbf{a}_t).$$

Solving for  $\mathbf{a}_t$  yields  $\mathbf{a}_{t|n}$ , which can be usefully rewritten as

$$\mathbf{a}_{t|n} = (\mathbf{I}_{t|t} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T})^{-1} (\mathbf{I}_{t|t} \mathbf{a}_{t|t} + \mathbf{T}'\mathbf{Q}^{-1}(\mathbf{a}_{t+1|n} - \mathbf{c})), \quad (\text{Q.8})$$

$$= \mathbf{a}_{t|t} + (\mathbf{I}_{t|t} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T})^{-1} \mathbf{T}'\mathbf{Q}^{-1} (\mathbf{a}_{t+1|n} - \mathbf{c} - \mathbf{T}\mathbf{a}_{t|t}), \quad (\text{Q.9})$$

$$= \mathbf{a}_{t|t} + \mathbf{I}_{t|t}^{-1} \mathbf{T}' (\mathbf{T}\mathbf{I}_{t|t}^{-1}\mathbf{T}' + \mathbf{Q})^{-1} (\mathbf{a}_{t+1|n} - \mathbf{c} - \mathbf{T}\mathbf{a}_{t|t}), \quad (\text{Q.10})$$

$$= \mathbf{a}_{t|t} + \mathbf{I}_{t|t}^{-1} \mathbf{T}' \mathbf{I}_{t+1|t} (\mathbf{a}_{t+1|n} - \mathbf{a}_{t+1|t}). \quad (\text{Q.11})$$

This second line expresses  $\mathbf{a}_{t|n}$  as the sum of  $\mathbf{a}_{t|t}$  and a correction that is linear in  $\mathbf{a}_{t+1|n} - \mathbf{c} - \mathbf{T}\mathbf{a}_{t|t}$ . The third line uses matrix-inversion formulas by Henderson and Searle (1981, eqns. 9–11) to ensure that  $\mathbf{Q}^{-1}$  no longer appears, such that by a limiting argument the result remains valid even when  $\mathbf{Q}$  is singular. The last line employs the prediction step  $\mathbf{a}_{t+1|t} := \mathbf{c} + \mathbf{T}\mathbf{a}_{t|t}$  and  $\mathbf{I}_{t+1|t} := (\mathbf{T}\mathbf{I}_{t|t}^{-1}\mathbf{T}' + \mathbf{Q})^{-1}$ . Equation (Q.11) is the Rauch-Tung-Striebel smoother expression, given in the main article in equation (33).

To derive the backward recursion for the precision matrix, we note that  $U_{t,t+1}(\cdot, \cdot)$  in equation (Q.5) can be written using matrix notation as

$$\begin{aligned} U_{t,t+1}(\mathbf{a}_t, \mathbf{a}_{t+1}) = & -\frac{1}{2} \begin{bmatrix} \mathbf{a}_t \\ \mathbf{a}_{t+1} \end{bmatrix}' \begin{bmatrix} \mathbf{I}_{t|t} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T} & -\mathbf{T}'\mathbf{Q}^{-1} \\ -\mathbf{Q}^{-1}\mathbf{T} & \widehat{\mathbf{I}}_{t+1|t+1} + \mathbf{Q}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{a}_t \\ \mathbf{a}_{t+1} \end{bmatrix} \\ & + \begin{bmatrix} \mathbf{I}_{t|t}\mathbf{a}_{t|t} - \mathbf{T}'\mathbf{Q}^{-1}\mathbf{c} \\ \mathbf{Q}^{-1}\mathbf{c} + \widehat{\mathbf{I}}_{t+1|t+1}\mathbf{a}_{t+1|t+1:n} \end{bmatrix}' \begin{bmatrix} \mathbf{a}_t \\ \mathbf{a}_{t+1} \end{bmatrix} + \text{constants}, \end{aligned} \quad (\text{Q.12})$$

where any constants that do not depend on  $\mathbf{a}_t$  and  $\mathbf{a}_{t+1}$  are ignored. This representation together with Lemma 1 implies that  $Z_t(\cdot) := \max_{\mathbf{a}} U_{t,t+1}(\cdot, \mathbf{a})$  is multivariate quadratic functions with negative Hessian matrix given by the following Schur complement:

$$\mathbf{I}_{t|n} = \mathbf{I}_{t|t} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T} - \mathbf{T}'\mathbf{Q}^{-1}(\widehat{\mathbf{I}}_{t+1|t+1} + \mathbf{Q}^{-1})^{-1}\mathbf{Q}^{-1}\mathbf{T}, \quad (\text{Q.13})$$

$$= \mathbf{I}_{t|t} + \mathbf{T}'(\widehat{\mathbf{I}}_{t+1|t+1}^{-1} + \mathbf{Q})^{-1}\mathbf{T}, \quad (\text{Q.14})$$

where the second line employs the Woodbury matrix equality (e.g. Henderson and Searle, 1981, eq. 1). Similarly,  $Z_{t+1}(\cdot) := \max_{\mathbf{a}} U_{t,t+1}(\mathbf{a}, \cdot)$  is multivariate quadratic with a negative Hessian given by the other Schur complement as follows:

$$\mathbf{I}_{t+1|n} = \widehat{\mathbf{I}}_{t+1|t+1} + \mathbf{Q}^{-1} - \mathbf{Q}^{-1}\mathbf{T}(\mathbf{I}_{t|t} + \mathbf{T}'\mathbf{Q}^{-1}\mathbf{T})^{-1}\mathbf{T}'\mathbf{Q}^{-1}, \quad (\text{Q.15})$$

$$= \widehat{\mathbf{I}}_{t+1|t+1} + (\mathbf{T}\mathbf{I}_{t|t}^{-1}\mathbf{T}' + \mathbf{Q})^{-1}, \quad (\text{Q.16})$$

$$= \widehat{\mathbf{I}}_{t+1|t+1} + \mathbf{I}_{t+1|t}, \quad (\text{Q.17})$$

where the second line again follows by the Woodbury matrix identity, while the last line employs the definition  $\mathbf{I}_{t+1|t} := (\mathbf{T}\mathbf{I}_{t|t}^{-1}\mathbf{T}' + \mathbf{Q})^{-1}$ . To derive equation (34), we note that

$$\mathbf{I}_{t|n}^{-1} = [\mathbf{I}_{t|t} + \mathbf{T}'(\widehat{\mathbf{I}}_{t+1|t+1}^{-1} + \mathbf{Q})^{-1}\mathbf{T}]^{-1}, \quad (\text{Q.18})$$

$$= \mathbf{I}_{t|t}^{-1} - \mathbf{I}_{t|t}^{-1}\mathbf{T}'[\widehat{\mathbf{I}}_{t+1|t+1}^{-1} + \mathbf{T}\mathbf{I}_{t|t}^{-1}\mathbf{T}' + \mathbf{Q}]^{-1}\mathbf{T}\mathbf{I}_{t|t}^{-1}, \quad \text{by Woodbury,} \quad (\text{Q.19})$$

$$= \mathbf{I}_{t|t}^{-1} - \mathbf{I}_{t|t}^{-1}\mathbf{T}'[\widehat{\mathbf{I}}_{t+1|t+1}^{-1} + \mathbf{I}_{t+1|t}^{-1}]^{-1}\mathbf{T}\mathbf{I}_{t|t}^{-1}, \quad \text{by Woodbury,} \quad (\text{Q.20})$$

$$= \mathbf{I}_{t|t}^{-1} - \mathbf{I}_{t|t}^{-1}\mathbf{T}'[\mathbf{I}_{t+1|t} - \mathbf{I}_{t+1|t}(\widehat{\mathbf{I}}_{t+1|t+1}^{-1} + \mathbf{I}_{t+1|t})^{-1}\mathbf{I}_{t+1|t}]\mathbf{T}\mathbf{I}_{t|t}^{-1}, \quad \text{Woodbury again,} \quad (\text{Q.21})$$

$$= \mathbf{I}_{t|t}^{-1} - \mathbf{I}_{t|t}^{-1}\mathbf{T}'[\mathbf{I}_{t+1|t} - \mathbf{I}_{t+1|t}\mathbf{I}_{t+1|n}^{-1}\mathbf{I}_{t+1|t}]\mathbf{T}\mathbf{I}_{t|t}^{-1}, \quad \text{by equation (Q.17),} \quad (\text{Q.22})$$

$$= \mathbf{I}_{t|t}^{-1} - \mathbf{I}_{t|t}^{-1}\mathbf{T}'\mathbf{I}_{t+1|t}[\mathbf{I}_{t+1|t}^{-1} - \mathbf{I}_{t+1|n}^{-1}]\mathbf{I}_{t+1|t}\mathbf{T}\mathbf{I}_{t|t}^{-1}, \quad (\text{Q.23})$$

confirming equation (34) in the main text.

## R Simulation study: Observation densities

Table R.1: Overview of data-generating processes in simulation studies.

DGP		Link function	Density	Score	Realised information	Information
Type	Distribution		$p(\mathbf{y}_t \alpha_t)$	$\frac{d\ell(\mathbf{y}_t \alpha_t)}{d\alpha_t}$	$-\frac{d^2\ell(\mathbf{y}_t \alpha_t)}{d\alpha_t^2}$	$\mathbb{E}\left[-\frac{d^2\ell(\mathbf{y}_t \alpha_t)}{d\alpha_t^2}\middle \alpha_t\right]$
Count	Poisson	$\lambda_t = \exp(\alpha_t)$	$\frac{\lambda_t^{y_t} \exp(-\lambda_t)/y_t!}{\Gamma(\kappa + y_t) \left(\frac{\kappa}{\kappa + \lambda_t}\right)^\kappa \left(\frac{\lambda_t}{\kappa + \lambda_t}\right)^{y_t}}$	$y_t - \lambda_t$	$\lambda_t$	$\lambda_t$
Count	Negative bin.	$\lambda_t = \exp(\alpha_t)$	$\frac{\Gamma(\kappa)\Gamma(y_t + 1)}{\lambda_t \exp(-\lambda_t y_t)}$	$y_t - \frac{\lambda_t(\kappa + y_t)}{\kappa + \lambda_t}$	$\frac{\kappa \lambda_t(\kappa + y_t)}{(\kappa + \lambda_t)^2}$	$\frac{\kappa \lambda_t}{\kappa + \lambda_t}$
Intensity	Exponential	$\lambda_t = \exp(\alpha_t)$	$\frac{y_t^{\kappa-1} \exp(-y_t/\beta_t)}{\Gamma(\kappa)\beta_t^\kappa}$	$1 - \lambda_t y_t$	$\frac{y_t \lambda_t}{\beta_t}$	1
Duration	Gamma	$\beta_t = \exp(\alpha_t)$	$\frac{\kappa (y_t/\beta_t)^{\kappa-1}}{\beta_t \exp\{(y_t/\beta_t)^\kappa\}}$	$\frac{y_t}{\beta_t} - \kappa$	$\frac{y_t}{\beta_t}$	$\kappa$
Duration	Weibull	$\beta_t = \exp(\alpha_t)$	$\frac{\exp\{-y_t^2/(2\sigma_t^2)\}}{\{2\pi\sigma_t^2\}^{1/2}}$	$\kappa \left(\frac{y_t}{\beta_t}\right)^\kappa - \kappa$	$\kappa^2 \left(\frac{y_t}{\beta_t}\right)^\kappa$	$\kappa^2$
Volatility	Gaussian	$\sigma_t^2 = \exp(\alpha_t)$	$\frac{\Gamma(\frac{\nu+1}{2}) \left(1 + \frac{y_t^2}{(\nu-2)\sigma_t^2}\right)^{-\frac{\nu+1}{2}}}{\sqrt{(\nu-2)\pi}\Gamma(\nu/2)\sigma_t}$	$\frac{y_t^2}{2\sigma_t^2} - \frac{1}{\nu+1}$	$\frac{y_t^2}{2\sigma_t^2}$	$\frac{1}{2}$
Volatility	Student's $t$	$\sigma_t^2 = \exp(\alpha_t)$	$\frac{\exp\left\{-\frac{y_{1t}^2 + y_{2t}^2 - 2\rho_t y_{1t} y_{2t}}{2(1-\rho_t^2)}\right\}}{2\pi\sqrt{1-\rho_t^2}}$	$\frac{\omega_t y_t^2}{2\sigma_t^2} - \frac{1}{\nu+1}$	$\frac{\nu-2}{\nu+1} \frac{\omega_t^2 y_t^2}{2\sigma_t^2}$	$\frac{\nu}{2\nu+6}$
Dependence	Gaussian	$\rho_t = \frac{1 - \exp(-\alpha_t)}{1 + \exp(-\alpha_t)}$		$\omega_t := \frac{\nu+1}{\nu-2 + y_t^2/\sigma_t^2}$	$0 \not\leq \frac{1}{4} \frac{z_{1t}^2 + z_{2t}^2}{1-\rho_t^2} - \frac{1-\rho_t^2}{4}$	$\frac{1 + \rho_t^2}{4}$
				$z_{1t} := y_{1t} - \rho_t y_{2t}$ $z_{2t} := y_{2t} - \rho_t y_{1t}$		
Dependence	Student's $t$	$\rho_t = \frac{1 - \exp(-\alpha_t)}{1 + \exp(-\alpha_t)}$	$\frac{\nu \left(1 + \frac{y_{1t}^2 + y_{2t}^2 - 2\rho_t y_{1t} y_{2t}}{(\nu-2)(1-\rho_t^2)}\right)^{-\frac{\nu+2}{2}}}{2\pi(\nu-2)\sqrt{1-\rho_t^2}}$	$\frac{\rho_t}{2} + \frac{\omega_t}{2} \frac{z_{1t} z_{2t}}{1-\rho_t^2}$	$0 \not\leq \frac{\omega_t}{4} \frac{z_{1t}^2 + z_{2t}^2}{1-\rho_t^2} - \frac{1-\rho_t^2}{4} - \frac{1}{2} \frac{\omega_t^2}{\nu+2} \frac{z_{1t}^2 z_{2t}^2}{(1-\rho_t^2)^2}$	$\frac{2 + \nu(1 + \rho_t^2)}{4(\nu+4)}$
				$z_{1t} := y_{1t} - \rho_t y_{2t}$ $z_{2t} := y_{2t} - \rho_t y_{1t}$	$\omega_t := \frac{\nu+2}{\nu-2 + \frac{y_{1t}^2 + y_{2t}^2 - 2\rho_t y_{1t} y_{2t}}{1-\rho_t^2}}$	
Local level	Student's $t$	$\mu_t = \alpha_t$	$\frac{\Gamma(\frac{\nu+1}{2}) \left(1 + \frac{(y_t - \mu_t)^2}{(\nu-2)\sigma^2}\right)^{-\frac{\nu+1}{2}}}{\sqrt{(\nu-2)\pi}\Gamma(\frac{\nu}{2})\sigma}$	$\frac{1}{\sigma} \frac{(\nu+1)e_t}{\nu-2 + e_t^2}$ $e_t := \frac{y_t - \mu_t}{\sigma}$	$0 \not\leq \frac{\nu+1}{\sigma^2} \frac{\nu-2 - e_t^2}{(\nu-2 + e_t^2)^2}$	$\frac{\nu(\nu+1)}{\sigma^2(\nu-2)(\nu+3)}$

Note: The table contains ten data-generating processes (DGPs) and link functions, the first nine of which are adapted from Koopman et al. (2016). For each model, the DGP is given by the linear Gaussian state equation (12) in combination with the observation density and link functions indicated in the table. The table further displays scores, realised information quantities and expected information quantities. The realised information quantities are nonnegative except for the bottom three models.

## S Simulation study: Parameter-estimation results

Table S.1: Short-window parameter estimates

DGP Type	Distribution	Truth	BF		PF		NAIS		
			Average	RMSE	Average	RMSE	Average	RMSE	
Count	Poisson	$c$	0.000	-0.016	[0.088]	-0.003	[0.042]	-0.002	[0.040]
		$\phi$	0.980	0.932	[0.132]	0.941	[0.099]	0.945	[0.084]
		$\sigma_\eta$	0.150	0.182	[0.083]	0.170	[0.070]	0.168	[0.060]
Count	Negative Bin.	$c$	0.000	-0.019	[0.095]	-0.008	[0.080]	-0.001	[0.036]
		$\phi$	0.980	0.925	[0.147]	0.929	[0.153]	0.946	[0.099]
		$\sigma_\eta$	0.150	0.194	[0.123]	0.176	[0.098]	0.158	[0.055]
		$1/\kappa$	0.250	0.205	[0.138]	0.227	[0.122]	0.298	[0.141]
Intensity	Exponential	$c$	0.000	-0.006	[0.033]	0.000	[0.030]	0.000	[0.030]
		$\phi$	0.980	0.943	[0.070]	0.946	[0.079]	0.948	[0.064]
		$\sigma_\eta$	0.150	0.180	[0.070]	0.168	[0.063]	0.169	[0.059]
Duration	Gamma	$c$	0.000	0.002	[0.041]	-0.003	[0.036]	-0.003	[0.037]
		$\phi$	0.980	0.944	[0.072]	0.948	[0.072]	0.949	[0.062]
		$\sigma_\eta$	0.150	0.175	[0.062]	0.166	[0.054]	0.166	[0.054]
		$\kappa$	1.500	1.541	[0.160]	1.531	[0.156]	1.532	[0.155]
Duration	Weibull	$c$	0.000	0.005	[0.041]	-0.003	[0.034]	-0.003	[0.033]
		$\phi$	0.980	0.939	[0.079]	0.946	[0.069]	0.947	[0.064]
		$\sigma_\eta$	0.150	0.188	[0.075]	0.172	[0.064]	0.173	[0.060]
		$\kappa$	1.200	1.225	[0.080]	1.215	[0.075]	1.215	[0.075]
Volatility	Gaussian	$c$	0.000	0.000	[0.068]	-0.004	[0.063]	-0.003	[0.073]
		$\phi$	0.980	0.905	[0.200]	0.906	[0.218]	0.914	[0.184]
		$\sigma_\eta$	0.150	0.202	[0.119]	0.174	[0.112]	0.183	[0.099]
Volatility	Student's $t$	$c$	0.000	-0.010	[0.113]	-0.008	[0.106]	-0.005	[0.070]
		$\phi$	0.980	0.870	[0.261]	0.872	[0.311]	0.914	[0.162]
		$\sigma_\eta$	0.150	0.249	[0.198]	0.190	[0.151]	0.192	[0.116]
		$1/\nu$	0.100	0.063	[0.069]	0.088	[0.041]	0.082	[0.057]
Dependence	Gaussian	$c$	0.020	0.082	[0.103]	0.142	[0.292]	0.165	[0.350]
		$\phi$	0.980	0.916	[0.102]	0.859	[0.278]	0.834	[0.339]
		$\sigma_\eta$	0.100	0.124	[0.090]	0.155	[0.185]	0.144	[0.132]
Dependence	Student's $t$	$c$	0.020	0.148	[0.321]	0.263	[0.540]	0.189	[0.349]
		$\phi$	0.980	0.854	[0.303]	0.744	[0.501]	0.810	[0.344]
		$\sigma_\eta$	0.100	0.136	[0.128]	0.201	[0.225]	0.146	[0.139]
		$1/\nu$	0.100	0.100	[0.031]	0.096	[0.033]	0.091	[0.066]
Level	Student's $t$	$c$	0.000	0.000	[0.016]	0.000	[0.019]		
		$\phi$	0.980	0.965	[0.027]	0.959	[0.034]		
		$\sigma_\eta$	0.150	0.131	[0.028]	0.155	[0.027]		
		$\sigma$	0.450	0.433	[0.061]	0.484	[0.147]		
		$1/\nu$	0.333	0.237	[0.121]	0.324	[0.083]		

Note: BF = Bellman filter. PF = Particle filter. NAIS = Numerically accelerated importance sampler. RMSE = root mean squared error. For the simulation setting, see the note to Table 5 in the main text.

Table S.2: Medium-window parameter estimates

DGP Type	Distribution	Truth		BF		PF		NAIS	
				Average	RMSE	Average	RMSE	Average	RMSE
Count	Poisson	$c$	0.000	-0.007	[0.010]	0.000	[0.006]	0.000	[0.006]
		$\phi$	0.980	0.974	[0.013]	0.975	[0.011]	0.975	[0.011]
		$\sigma_\eta$	0.150	0.155	[0.023]	0.154	[0.022]	0.151	[0.021]
Count	Negative Bin.	$c$	0.000	-0.004	[0.008]	0.000	[0.007]	0.001	[0.006]
		$\phi$	0.980	0.976	[0.012]	0.974	[0.013]	0.976	[0.011]
		$\sigma_\eta$	0.150	0.152	[0.027]	0.155	[0.027]	0.147	[0.025]
		$1/\kappa$	0.250	0.236	[0.058]	0.245	[0.051]	0.288	[0.066]
Intensity	Exponential	$c$	0.000	-0.007	[0.010]	0.000	[0.007]	0.000	[0.007]
		$\phi$	0.980	0.972	[0.014]	0.974	[0.013]	0.974	[0.013]
		$\sigma_\eta$	0.150	0.162	[0.027]	0.154	[0.023]	0.154	[0.023]
Duration	Gamma	$c$	0.000	0.007	[0.010]	0.000	[0.007]	0.000	[0.007]
		$\phi$	0.980	0.973	[0.013]	0.974	[0.012]	0.974	[0.012]
		$\sigma_\eta$	0.150	0.159	[0.023]	0.154	[0.021]	0.153	[0.020]
		$\kappa$	1.500	1.510	[0.070]	1.503	[0.069]	1.503	[0.069]
Duration	Weibull	$c$	0.000	0.009	[0.012]	0.000	[0.007]	0.000	[0.007]
		$\phi$	0.980	0.971	[0.015]	0.974	[0.012]	0.974	[0.012]
		$\sigma_\eta$	0.150	0.163	[0.027]	0.154	[0.021]	0.154	[0.021]
		$\kappa$	1.200	1.209	[0.037]	1.201	[0.035]	1.202	[0.035]
Volatility	Gaussian	$c$	0.000	0.007	[0.010]	0.000	[0.007]	0.000	[0.007]
		$\phi$	0.980	0.970	[0.019]	0.973	[0.016]	0.973	[0.016]
		$\sigma_\eta$	0.150	0.169	[0.040]	0.156	[0.032]	0.156	[0.031]
Volatility	Student's $t$	$c$	0.000	0.004	[0.010]	0.000	[0.007]	0.000	[0.007]
		$\phi$	0.980	0.969	[0.023]	0.974	[0.015]	0.973	[0.015]
		$\sigma_\eta$	0.150	0.173	[0.059]	0.157	[0.037]	0.158	[0.038]
		$1/\nu$	0.100	0.083	[0.045]	0.098	[0.021]	0.094	[0.034]
Dependence	Gaussian	$c$	0.020	0.028	[0.024]	0.035	[0.055]	0.034	[0.039]
		$\phi$	0.980	0.972	[0.023]	0.965	[0.056]	0.966	[0.038]
		$\sigma_\eta$	0.100	0.101	[0.033]	0.113	[0.054]	0.113	[0.049]
Dependence	Student's $t$	$c$	0.020	0.034	[0.059]	0.042	[0.088]	0.039	[0.052]
		$\phi$	0.980	0.966	[0.063]	0.958	[0.082]	0.961	[0.053]
		$\sigma_\eta$	0.100	0.107	[0.044]	0.121	[0.072]	0.122	[0.074]
		$1/\nu$	0.100	0.102	[0.017]	0.099	[0.013]	0.095	[0.039]
Level	Student's $t$	$c$	0.000	0.000	[0.005]	0.000	[0.006]		
		$\phi$	0.980	0.979	[0.007]	0.975	[0.010]		
		$\sigma_\eta$	0.150	0.129	[0.023]	0.152	[0.012]		
		$\sigma$	0.450	0.431	[0.033]	0.455	[0.053]		
		$1/\nu$	0.333	0.246	[0.094]	0.330	[0.043]		

Note: BF = Bellman filter. PF = Particle filter. NAIS = Numerically accelerated importance sampler. RMSE = root mean squared error. For the simulation setting, see the note to Table 5 in the main text.

## T Simulation study: Root mean squared errors

Table T.1: Root mean squared errors (RMSEs) of filtered states in the out-of-sample period.

DGP		Infeasible estimator	Short estimation window (250 obs.)				Medium estimation window (1,000 obs.)				Long estimation window (2,500 obs.)			
			BF	PF	NAIS	KF	BF	PF	NAIS	KF	BF	PF	NAIS	KF
Type	Distribution	Absolute RMSE	Relative RMSE				Relative RMSE				Relative RMSE			
Count	Poisson	0.360	1.163	1.157	1.155	1.015	1.015	1.015	1.000	1.000	1.001			
Count	Neg. Bin.	0.379	1.177	1.171	1.173	1.019	1.019	1.020	1.005	1.005	1.006			
Intensity	Exponential	0.361	1.139	1.141	1.137	1.013	1.012	1.012	1.001	1.001	1.000			
Duration	Gamma	0.326	1.169	1.165	1.163	1.023	1.022	1.022	1.006	1.005	1.005			
Duration	Weibull	0.332	1.126	1.123	1.120	1.010	1.009	1.009	0.999	0.998	0.998			
Volatility	Gaussian	0.425	1.218	1.221	1.220	1.497	1.022	1.022	1.022	1.229	1.003	1.003	1.002	1.229
Volatility	Student's $t$	0.442	1.250	1.231	1.235	1.593	1.039	1.028	1.029	1.338	1.012	1.028	1.009	1.275
Dependence	Gaussian	0.362	1.307	1.313	1.321		1.057	1.056	1.054		1.017	1.014	1.014	
Dependence	Student's $t$	0.371	1.314	1.327	1.303		1.065	1.066	1.068		1.022	1.021	1.021	
Level	Student's $t$	0.204	1.058	1.045	n/a	1.233	1.007	1.000	n/a	1.156	0.998	0.996	n/a	1.148

Note: MAE = mean absolute error. BF = Bellman filter. PF = particle filter. NAIS = numerically accelerated importance sampler. KF = Kalman filter. See the note to Table 3 in the main text. The only difference is that here we report root mean squared errors (RMSEs), not mean absolute errors (MAEs).

## U Catania's (2022) model: State-space representation

Fix  $t > k+1$ . Conditional on the information set at time  $t-k-1$ , denoted  $\mathcal{F}_{t-k-1}$ , Catania's (2022) model (46)–(48) implies that the volatility shock  $\eta_t$  and the return shocks  $\varepsilon_t, \dots, \varepsilon_{t-k}$  are jointly normally distributed as

$$\begin{bmatrix} \eta_t \\ \varepsilon_t \\ \varepsilon_{t-1} \\ \vdots \\ \varepsilon_{t-k} \end{bmatrix} \Big|_{\mathcal{F}_{t-k-1}} \sim N \left( \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho_0 & \rho_1 & \dots & \rho_k \\ \rho_0 & 1 & 0 & \dots & 0 \\ \rho_1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_k & 0 & 0 & \dots & 1 \end{bmatrix} \right). \quad (\text{U.1})$$

Next, we compute the distribution of both current shocks, i.e.  $\eta_t$  and  $\varepsilon_t$ , conditional on the past shocks,  $\varepsilon_{t-1}, \dots, \varepsilon_{t-k}$ . From a well-known lemma regarding conditional Gaussian distributions (e.g. Harvey, 1990, p. 165), it follows that  $\eta_t, \varepsilon_t$  conditional on  $\varepsilon_{t-1}, \dots, \varepsilon_{t-k}$ , or, equivalently,  $\mathcal{F}_{t-1}$  and  $\mathbf{a}_{t-1}$ , are jointly normally distributed as

$$\begin{bmatrix} \eta_t \\ \varepsilon_t \end{bmatrix} \Big|_{\mathcal{F}_{t-1}, \mathbf{a}_{t-1}} \sim N \left( \begin{bmatrix} \sum_{j=1}^k \rho_j \varepsilon_{t-j} \\ 0 \end{bmatrix}, \begin{bmatrix} 1 - \sum_{j=1}^k \rho_j^2 & \rho_0 \\ \rho_0 & 1 \end{bmatrix} \right). \quad (\text{U.2})$$

The marginal distribution of  $\eta_t$  is again Gaussian, with a mean and variance that can be read off. Next, the state-transition equation implies that  $h_t = c + \varphi h_{t-1} + \sigma_\eta \eta_t$ , being a linear transformation of  $\eta_t$ , is distributed as

$$h_t \Big|_{\mathcal{F}_{t-1}, \mathbf{a}_{t-1}} \sim N(\mu_{h,t}, \sigma_{h,t}^2), \quad \text{where} \quad (\text{U.3})$$

$$\mu_{h,t} = c + \varphi h_{t-1} + \sigma_\eta \sum_{j=1}^k \rho_j \frac{y_{t-j} - \mu}{\exp(h_{t-j}/2)}, \quad \sigma_{h,t} = \sigma_\eta \sqrt{1 - \sum_{j=1}^k \rho_j^2}, \quad (\text{U.4})$$

where we have used  $\varepsilon_{t-j} = (y_{t-j} - \mu) \exp(-h_{t-j}/2)$  for  $j = 1, \dots, k$  in the expression for  $\mu_{h,t}$ . This confirms the non-degenerate part of the state-transition density (50). To derive the observation density, we note that the bivariate distribution (U.2) with another application of the conditional-Gaussian lemma (Harvey, 1990, p. 165) gives

$$\varepsilon_t \Big|_{\mathcal{F}_{t-1}, \mathbf{a}_{t-1}, \eta_t} \sim N(\mu_{\varepsilon,t}, \sigma_{\varepsilon,t}^2), \quad \text{where} \quad (\text{U.5})$$

$$\mu_{\varepsilon,t} = \frac{\rho_0}{1 - \sum_{j=1}^k \rho_j^2} \left( \eta_t - \sum_{j=1}^k \rho_j \varepsilon_{t-j} \right), \quad \sigma_{\varepsilon,t} = \sqrt{1 - \frac{\rho_0^2}{1 - \sum_{j=1}^k \rho_j^2}}. \quad (\text{U.6})$$

Noting that neither  $\mu_{\varepsilon,t}$  nor  $\sigma_{\varepsilon,t}$  depend on  $h_{t-k-1}$ , while  $\mathbf{a}_{t-1}$  and  $\eta_t$  together imply  $\mathbf{a}_t$ , the conditioning set  $(\mathcal{F}_{t-1}, \mathbf{a}_{t-1}, \eta_t)$  can be simplified to  $(\mathcal{F}_{t-1}, \mathbf{a}_t)$ . Further, by substituting  $\eta_t = (h_t - c - \varphi h_{t-1})/\sigma_\eta$  and  $\varepsilon_{t-j} = (y_{t-j} - \mu) \exp(-h_{t-j}/2)$  for  $j = 1, \dots, k$ , equations (U.5)–(U.6) become

$$\varepsilon_t | \mathcal{F}_{t-1}, \mathbf{a}_t, \sim \text{N}(\mu_{\varepsilon,t}, \sigma_{\varepsilon,t}^2), \quad \text{where} \quad (\text{U.7})$$

$$\mu_{\varepsilon,t} = \frac{\rho_0}{1 - \sum_{j=1}^k \rho_j^2} \left( \frac{h_t - c - \varphi h_{t-1}}{\sigma_\eta} - \sum_{j=1}^k \rho_j \frac{y_{t-j} - \mu}{\exp(h_{t-j}/2)} \right), \quad \sigma_{\varepsilon,t} = \sqrt{1 - \frac{\rho_0^2}{1 - \sum_{j=1}^k \rho_j^2}}. \quad (\text{U.8})$$

Finally, the distribution of the observation  $y_t = \mu + \exp(h_t/2)\varepsilon_t$  conditional on  $\mathcal{F}_{t-1}$  and  $\mathbf{a}_t$  is Gaussian with mean  $\mu_{y,t} = \mu + \exp(h_t/2)\mu_{\varepsilon,t}$  and variance  $\sigma_{y,t}^2 = \exp(h_t)\sigma_{\varepsilon,t}^2$ , where  $\mu_{\varepsilon,t}$  and  $\sigma_{\varepsilon,t}$  are given in equation (U.8). This confirms observation density (49).

## V Catania's (2022) model: Bellman-filter implementation

Bellman's equation (7) at time  $t$  involves the maximisation over two state variables, i.e.  $\mathbf{a}_t$  and  $\mathbf{a}_{t-1}$ , which in general contain independent components. For the specific case of Catania's (2022) model, as described in section 9, the state vector is  $\mathbf{a}_t = (h_t, h_{t-1}, \dots, h_{t-k})' \in \mathbb{R}^{k+1}$ , which contains the log-volatility  $h_t$  as well as  $k$  lags. This implies that the state variables  $\mathbf{a}_t$  and  $\mathbf{a}_{t-1}$  have  $k$  elements in common, namely  $h_{t-1}$  through  $h_{t-k}$ . Further,  $h_t$  appears only in  $\mathbf{a}_t$ , while  $h_{t-k-1}$  appears only in  $\mathbf{a}_{t-1}$ . Taking into account these restrictions, optimisation (7) specialised to Catania's (2022) model reads

$$\begin{bmatrix} \mathbf{a}_t | t \\ h_{t-k-1} | t \end{bmatrix} = \begin{bmatrix} h_t | t \\ h_{t-1} | t \\ \vdots \\ h_{t-k} | t \\ h_{t-k-1} | t \end{bmatrix} = \arg \max_{h_t, h_{t-1}, \dots, h_{t-k-1}} \left\{ \ell(y_t | \mathbf{a}_t, \mathcal{F}_{t-1}) + \ell(h_t | \mathbf{a}_{t-1}, \mathcal{F}_{t-1}) + V_{t-1}(\mathbf{a}_{t-1}) \right\}, \quad (\text{V.1})$$

where  $\ell(\cdot) := \log p(\cdot)$  and the observation and state-transition densities are given in equations (49) and (50), respectively. In equation (V.1), we have dropped the degenerate part of the state-transition density, which is permitted given that the optimisation variables are taken to be  $h_t, \dots, h_{t-k-1}$ , such that the restrictions on the components of  $\mathbf{a}_t$  and  $\mathbf{a}_{t-1}$  are automatically satisfied. Value function  $V_{t-1} : \mathbb{R}^{k+1} \rightarrow \mathbb{R}$  on the right-hand side is approximated by the quadratic form (8).

To simplify the analysis of optimisation (V.1), we introduce three notational conventions. First, the  $k+2$  optimisation variables in optimisation (V.1) are collected in a single vector:

$$\mathbf{x}_t := (h_t, h_{t-1}, \dots, h_{t-k-1})' = (h_t, \mathbf{a}'_{t-1})' = (\mathbf{a}'_t, h_{t-k-1})' \in \mathbb{R}^{k+2}. \quad (\text{V.2})$$

Second, we write the observation log density as  $f := \ell(y_t | \mathbf{a}_t, \mathcal{F}_{t-1})$ , such that by equation (49) we have

$$f(\mathbf{a}_t) := -\frac{1}{2} \log(2\pi) - \log(\sigma_{y,t}) - \frac{(y_t - \mu_{y,t})^2}{2\sigma_{y,t}^2}, \quad \sigma_{y,t} = \exp(h_t/2) \sqrt{1 - \frac{\rho_0^2}{1 - \sum_{j=1}^k \rho_j^2}}, \quad (\text{V.3})$$

$$\mu_{y,t} = \mu + \frac{\rho_0 \exp(h_t/2)}{1 - \sum_{j=1}^k \rho_j^2} \left[ \frac{h_t - c - \varphi h_{t-1}}{\sigma_\eta} - \sum_{j=1}^k \rho_j \frac{y_{t-j} - \mu}{\exp(h_{t-j}/2)} \right].$$

Third, for the state-transition log density we use the short-hand  $g := \ell(h_t | \mathbf{a}_{t-1}, \mathcal{F}_{t-1})$  and note from equation (50) that it does not depend on  $h_{t-k-1}$ , such that we may write  $g = g(\mathbf{a}_t)$  as follows:

$$g(\mathbf{a}_t) := -\frac{1}{2} \log(2\pi) - \log(\sigma_{h,t}) - \frac{(h_t - \mu_{h,t})^2}{2\sigma_{h,t}^2}, \quad (\text{V.4})$$

$$\mu_{h,t} = c + \varphi h_{t-1} + \sigma_\eta \sum_{j=1}^k \rho_j \frac{y_{t-j} - \mu}{\exp(h_{t-j}/2)}, \quad \sigma_{h,t} = \sigma_\eta \sqrt{1 - \sum_{j=1}^k \rho_j^2}.$$

Notation (V.2) through (V.4) allows us to write optimisation (V.1) as

$$\hat{\mathbf{x}}_{t|t} = \arg \max_{\mathbf{x}_t} \left\{ f(\mathbf{a}_t) + g(\mathbf{a}_t) - \frac{1}{2} (\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1})' \mathbf{I}_{t-1|t-1} (\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1}) \right\}. \quad (\text{V.5})$$

The Newton scoring algorithm for optimisation (V.5) reads

$$\mathbf{x}_t \leftarrow \mathbf{x}_t + \left[ \left( \begin{array}{c|c} -\frac{d^2 f}{d\mathbf{a}_t d\mathbf{a}'_t} - \frac{d^2 g}{d\mathbf{a}_t d\mathbf{a}'_t} & \mathbf{0}_{k+1} \\ \hline \mathbf{0}'_{k+1} & 0 \end{array} \right) + \left( \begin{array}{cc} 0 & \mathbf{0}'_{k+1} \\ \mathbf{0}_{k+1} & \mathbf{I}_{t-1|t-1} \end{array} \right) \right]^{-1} \left[ \left( \begin{array}{c} \frac{d(f+g)}{d\mathbf{a}_t} \\ 0 \end{array} \right) - \left( \begin{array}{c} 0 \\ \mathbf{I}_{t-1|t-1} (\mathbf{a}_{t-1} - \mathbf{a}_{t-1|t-1}) \end{array} \right) \right], \quad (\text{V.6})$$

where  $\mathbf{0}_{k+1}$  is a column vector consisting of  $k+1$  zeroes. Fisher scoring steps are obtained by replacing  $d^2 f / (d\mathbf{a}_t d\mathbf{a}'_t)$  by  $\mathbb{E}[d^2 f / (d\mathbf{a}_t d\mathbf{a}'_t) | \mathbf{a}_t, \mathcal{F}_{t-1}]$ . Iterating Newton step (V.6) or its Fisher equivalent requires (expectations of) first and second derivatives of  $f, g$ , as derived next.

**Derivatives of  $f$ :** By the chain rule, first and second derivatives of the function  $f$  defined in equation (V.3) with respect to  $\mathbf{a}_t = (h_t, \dots, h_{t-k})'$  read

$$\frac{df}{d\mathbf{a}_t} = \frac{df}{d\mu_{y,t}} \frac{d\mu_{y,t}}{d\mathbf{a}_t} + \frac{df}{d\sigma_{y,t}} \frac{d\sigma_{y,t}}{d\mathbf{a}_t}, \quad (\text{V.7})$$

$$\begin{aligned} \frac{d^2 f}{d\mathbf{a}_t d\mathbf{a}'_t} &= \frac{d^2 f}{(d\mu_{y,t})^2} \frac{d\mu_{y,t}}{d\mathbf{a}_t} \frac{d\mu_{y,t}}{d\mathbf{a}'_t} + \frac{d^2 f}{(d\sigma_{y,t})^2} \frac{d\sigma_{y,t}}{d\mathbf{a}_t} \frac{d\sigma_{y,t}}{d\mathbf{a}'_t} + \frac{d^2 f}{d\mu_{y,t} d\sigma_{y,t}} \frac{d\mu_{y,t}}{d\mathbf{a}_t} \frac{d\sigma_{y,t}}{d\mathbf{a}'_t} \\ &\quad + \frac{d^2 f}{d\mu_{y,t} d\sigma_{y,t}} \frac{d\sigma_{y,t}}{d\mathbf{a}_t} \frac{d\mu_{y,t}}{d\mathbf{a}'_t} + \frac{df}{d\mu_{y,t}} \frac{d^2 \mu_{y,t}}{d\mathbf{a}_t d\mathbf{a}'_t} + \frac{df}{d\sigma_{y,t}} \frac{d^2 \sigma_{y,t}}{d\mathbf{a}_t d\mathbf{a}'_t}. \end{aligned} \quad (\text{V.8})$$

$$\begin{aligned} \mathbb{E} \left[ \frac{d^2 f}{d\mathbf{a}_t d\mathbf{a}'_t} \middle| \mathbf{a}_t, \mathcal{F}_{t-1} \right] &= \frac{d^2 f}{(d\mu_{y,t})^2} \frac{d\mu_{y,t}}{d\mathbf{a}_t} \frac{d\mu_{y,t}}{d\mathbf{a}'_t} + \mathbb{E} \left[ \frac{d^2 f}{(d\sigma_{y,t})^2} \middle| \mathbf{a}_t, \mathcal{F}_{t-1} \right] \frac{d\sigma_{y,t}}{d\mathbf{a}_t} \frac{d\sigma_{y,t}}{d\mathbf{a}'_t} \\ &\quad + \mathbb{E} \left[ \frac{d^2 f}{d\mu_{y,t} d\sigma_{y,t}} \middle| \mathbf{a}_t, \mathcal{F}_{t-1} \right] \frac{d\mu_{y,t}}{d\mathbf{a}_t} \frac{d\sigma_{y,t}}{d\mathbf{a}'_t} + \mathbb{E} \left[ \frac{d^2 f}{d\mu_{y,t} d\sigma_{y,t}} \middle| \mathbf{a}_t, \mathcal{F}_{t-1} \right] \frac{d\sigma_{y,t}}{d\mathbf{a}_t} \frac{d\mu_{y,t}}{d\mathbf{a}'_t}. \end{aligned} \quad (\text{V.9})$$

Equation (V.9) contains two fewer terms than equation (V.8), because the expectation of the last two terms in equation (V.8) is zero. In equations (V.7) through (V.9), derivatives of  $f$  with respect  $\mu_{y,t}$  and  $\sigma_{y,t}$  are given by

$$\frac{df}{d\mu_{y,t}} = \frac{y_t - \mu_{y,t}}{\sigma_{y,t}^2}, \quad \frac{df}{d\sigma_{y,t}} = \frac{(y_t - \mu_{y,t})^2}{\sigma_{y,t}^3} - \frac{1}{\sigma_{y,t}}, \quad (\text{V.10})$$

$$\frac{d^2 f}{(d\mu_{y,t})^2} = \frac{-1}{\sigma_{y,t}^2}, \quad \frac{d^2 f}{d\mu_{y,t} d\sigma_{y,t}} = -2 \frac{y_t - \mu_{y,t}}{\sigma_{y,t}^3}, \quad \frac{d^2 f}{(d\sigma_{y,t})^2} = \frac{1}{\sigma_{y,t}^2} - \frac{3(y_t - \mu_{y,t})^2}{\sigma_{y,t}^4}, \quad (\text{V.11})$$

$$\mathbb{E} \left[ \frac{d^2 f}{d\mu_{y,t} d\sigma_{y,t}} \middle| \mathcal{F}_{t-1}, \mathbf{a}_t \right] = 0, \quad \mathbb{E} \left[ \frac{d^2 f}{(d\sigma_{y,t})^2} \middle| \mathcal{F}_{t-1}, \mathbf{a}_t \right] = \frac{-2}{\sigma_{y,t}^2}, \quad (\text{V.12})$$

where we also give expectations when relevant for Fisher scoring steps. In equations (V.7) and (V.8), first derivatives of  $\mu_{y,t}$  with respect to the elements of  $\mathbf{a}_t$  read

$$\frac{d\mu_{y,t}}{d\mathbf{a}_t} = \begin{bmatrix} (\mu_{y,t} - \mu)/2 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \frac{\rho_0 \exp(h_t/2)}{1 - \sum_{j=1}^k \rho_j^2} \begin{bmatrix} 1/\sigma_\eta \\ -\varphi/\sigma_\eta + \rho_1/2 \frac{y_{t-1} - \mu}{\exp(h_{t-1}/2)} \\ \rho_2/2 \frac{y_{t-2} - \mu}{\exp(h_{t-2}/2)} \\ \vdots \\ \rho_k/2 \frac{y_{t-k} - \mu}{\exp(h_{t-k}/2)} \end{bmatrix} =: \begin{bmatrix} (\mu_{y,t} - \mu)/2 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \mathbf{b}_t, \quad (\text{V.13})$$

where the second equality entails a definition of  $\mathbf{b}_t$ . For second derivatives of  $\mu_{y,t}$ , we have

$$\frac{d^2\mu_{y,t}}{d\mathbf{a}_t d\mathbf{a}_t'} = \text{diag} \begin{bmatrix} (\mu_{y,t} - \mu)/4 \\ 0 \\ \vdots \\ 0 \end{bmatrix} - \frac{1}{4} \frac{\rho_0 \exp(h_t/2)}{1 - \sum_{j=1}^k \rho_j^2} \text{diag} \begin{bmatrix} 0 \\ \rho_1 \frac{y_{t-1} - \mu}{\exp(h_{t-1}/2)} \\ \rho_2 \frac{y_{t-2} - \mu}{\exp(h_{t-2}/2)} \\ \vdots \\ \rho_k \frac{y_{t-k} - \mu}{\exp(h_{t-k}/2)} \end{bmatrix} + \begin{bmatrix} 1/2 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \mathbf{b}_t' + \mathbf{b}_t \begin{bmatrix} 1/2 & 0 & \dots & 0 \end{bmatrix}, \quad (\text{V.14})$$

where the diag operator creates a diagonal matrix from a given vector. The derivatives of  $\sigma_{y,t}$  read

$$\frac{d\sigma_{y,t}}{d\mathbf{a}_t} = \begin{bmatrix} \sigma_{y,t}/2 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \frac{d^2\sigma_{y,t}}{d\mathbf{a}_t d\mathbf{a}_t'} = \text{diag} \begin{bmatrix} \sigma_{y,t}/4 \\ 0 \\ \vdots \\ 0 \end{bmatrix}. \quad (\text{V.15})$$

All components of equations (V.7) and (V.8) have now been specified.

**Derivatives of  $g$ :** By the chain rule, first and second derivatives of the function  $g$  given in equation (V.4) with respect to  $\mathbf{a}_t = (h_t, \dots, h_{t-k})'$  are

$$\frac{dg}{d\mathbf{a}_t} = \frac{h_t - \mu_{h,t}}{\sigma_{h,t}^2} \begin{bmatrix} -1 \\ \varphi - \frac{\sigma_\eta}{2} \rho_1 \frac{y_{t-1} - \mu}{\exp(h_{t-1}/2)} \\ -\frac{\sigma_\eta}{2} \rho_2 \frac{y_{t-2} - \mu}{\exp(h_{t-2}/2)} \\ \vdots \\ -\frac{\sigma_\eta}{2} \rho_k \frac{y_{t-k} - \mu}{\exp(h_{t-k}/2)} \end{bmatrix} =: \frac{h_t - \mu_{h,t}}{\sigma_{h,t}^2} \mathbf{c}_t, \quad (\text{V.16})$$

$$\frac{d^2g}{d\mathbf{a}_t d\mathbf{a}_t'} = \frac{-1}{\sigma_{h,t}^2} \mathbf{c}_t \mathbf{c}_t' + \frac{h_t - \mu_{h,t}}{\sigma_{h,t}^2} \frac{\sigma_\eta}{4} \text{diag} \begin{bmatrix} 0 \\ \rho_1 \frac{y_{t-1} - \mu}{\exp(h_{t-1}/2)} \\ \rho_2 \frac{y_{t-2} - \mu}{\exp(h_{t-2}/2)} \\ \vdots \\ \rho_k \frac{y_{t-k} - \mu}{\exp(h_{t-k}/2)} \end{bmatrix}. \quad (\text{V.17})$$

Jointly, equations (V.7) through (V.17) specify all components of the Fisher scoring step (V.6).

Finally, the updated information matrix  $\mathbf{I}_{t|t}$  is determined by the Schur complement of the bottom-right element of the negative Hessian matrix used in Newton's scoring step, which is given by

$$\begin{pmatrix} -\frac{d^2f}{d\mathbf{a}_t d\mathbf{a}_t'} - \frac{d^2g}{d\mathbf{a}_t d\mathbf{a}_t'} & \mathbf{0}_{k+1} \\ \mathbf{0}'_{k+1} & 0 \end{pmatrix} + \begin{pmatrix} 0 & \mathbf{0}'_{k+1} \\ \mathbf{0}_{k+1} & \mathbf{I}_{t-1|t-1} \end{pmatrix},$$

Taking Schur complement of the bottom-right element and evaluating the result at the peak, i.e. at  $\mathbf{a}_{t|t}$ , gives the updated information matrix  $\mathbf{I}_{t|t}$ . The Fisher version of the updating steps is obtained by replacing  $d^2f/(d\mathbf{a}_t d\mathbf{a}_t')$  by  $\mathbb{E}[d^2f/(d\mathbf{a}_t d\mathbf{a}_t') | \mathbf{a}_t, \mathcal{F}_{t-1}]$ .

# W Full estimation results for the S&P500

Table W.1: Full estimation results for the Bellman filter (top panel) and particle filter (bottom panel).

$\mu$	$c$	$\varphi$	$\sigma_\eta$	$\rho_0$	$\rho_1$	$\rho_2$	$\rho_3$	$\rho_4$	$\rho_5$	$\rho_6$	$\rho_7$	$\rho_8$	$\rho_9$	$\rho_{10}$	LogL	BIC
.0696	.0004	.9839	.2006	-.7189											-9555.1	2.5344
.0519	-.0017	.9759	.2058	-.4830	-.4028										-9531.7	2.5294
.0518	-.0013	.9776	.2447	-.4020	-.5945	.2910									-9524.3	2.5286
.0513	-.0006	.9815	.2582	-.3770	-.5828	-.0913	.4633								-9503.2	<b>2.5242</b>
.0509	-.0003	.9826	.2456	-.3989	-.6108	-.0926	.3612	.1463							-9500.3	2.5246
.0509	-.0001	.9842	.2456	-.4016	-.6037	-.0962	.3665	-.0382	.2132						-9494.5	2.5243
.0503	.0002	.9852	.2412	-.4136	-.6107	-.0921	.3715	-.0424	.0808	.1616					-9490.9	2.5245
.0499	.0005	.9862	.2397	-.4193	-.6115	-.0936	.3750	-.0478	.0916	.0186	.1644				-9487.6	2.5248
.0508	.0002	.9867	.2376	-.4204	-.6163	-.0955	.3817	-.0511	.0968	.0159	.0540	.1242			-9482.0	2.5245
.0502	.0006	.9875	.2384	-.4223	-.6096	-.0897	.3791	-.0572	.0986	.0188	.0553	-.0462	.1901		-9477.4	2.5245
.0500	.0007	.9881	.2353	-.4309	-.6126	-.0912	.3828	-.0616	.1031	.0175	.0597	-.0471	.0804	.1277	-9474.5	2.5249
.0680	-.0042	.9850	.1926	-.7319											-9562.1	2.5362
.0517	-.0071	.9784	.1932	-.5071	-.4149										-9539.3	2.5314
.0511	-.0065	.9796	.2262	-.4278	-.5935	.2732									-9534.2	2.5312
.0519	-.0056	.9828	.2395	-.3979	-.5707	-.1141	.4593								-9516.9	<b>2.5278</b>
.0513	-.0065	.9826	.2420	-.3743	-.6300	-.0624	.4107	.0501							-9516.2	2.5288
.0502	-.0051	.9837	.2284	-.4059	-.6137	-.1062	.3489	.1464	.0044						-9515.1	2.5297
.0491	-.0041	.9853	.2267	-.4217	-.5909	-.1206	.3700	-.0808	.1629	.1019					-9509.1	2.5293
.0489	-.0038	.9860	.2301	-.4171	-.6001	-.1134	.3845	-.0756	.1106	-.0147	.1842				-9505.9	2.5296
.0495	-.0039	.9864	.2294	-.4165	-.5988	-.1126	.3838	-.0760	.1102	-.0146	.1846	.0001			-9505.9	2.5308
.0495	-.0039	.9863	.2294	-.4163	-.5991	-.1128	.3831	-.0761	.1104	-.0144	.1848	.0001	.0003		-9505.9	2.5320
.0471	-.0037	.9874	.2204	-.4107	-.6221	-.1563	.3621	.0545	.0495	.0157	.0727	-.0021	.0014	.1236	-9501.9	2.5321

*Note:* LogL = log likelihood. BIC = Bayesian information criterion. For each panel, the best BIC is indicated in bold. The data are  $100 \times$  the log returns of the S&P500 from 3 Jan 1990 to 31 Dec 2019 (7,558 observations). The Bellman filter is implemented as described in Appendix V and estimated using estimator (40). The particle filter is estimated as in Catania (2022), who uses the continuous sampling importance resampling (CSIR) method of Malik and Pitt (2011).

## References

- Amari, S.-i., Park, H. and Fukumizu, K. (2000) Adaptive method of realizing natural gradient learning for multilayer perceptrons. *Neural Computation*, **12**, 1399–1409.
- Bernstein, D. S. (2009) *Matrix Mathematics: Theory, Facts, and Formulas*. PUP.
- Catania, L. (2022) A stochastic volatility model with a general leverage specification. *Journal of Business & Economic Statistics*, **40**, 678–689.
- Fahrmeir, L. (1992) Posterior mode estimation by extended Kalman filtering for multivariate dynamic generalized linear models. *Journal of the American Statistical Association*, **87**, 501–509.
- Harvey, A. C. (1990) *Forecasting, Structural Time Series Models and the Kalman Filter*. CUP.
- Henderson, H. V. and Searle, S. R. (1981) On deriving the inverse of a sum of matrices. *SIAM Review*, **23**, 53–60.
- Jungers, R. (2009) *The Joint Spectral Radius: Theory and Applications*. Springer.
- Koopman, S. J., Lucas, A. and Scharth, M. (2016) Predicting time-varying parameters with parameter-driven and observation-driven models. *Review of Economics and Statistics*, **98**, 97–110.
- Koyama, S., Castellanos Pérez-Bolde, L., Shalizi, C. R. and Kass, R. E. (2010) Approximate methods for state-space models. *Journal of the American Statistical Association*, **105**, 170–180.
- Malik, S. and Pitt, M. K. (2011) Particle filters for continuous likelihood evaluation and maximisation. *Journal of Econometrics*, **165**, 190–209.
- Nesterov, Y. (2003) *Introductory Lectures on Convex Optimization: A Basic Course*. Springer.
- Toulis, P. and Airoldi, E. M. (2015) Scalable estimation strategies based on stochastic approximations: Classical results and new insights. *Statistics and Computing*, **25**, 781–795.
- (2017) Asymptotic and finite-sample properties of estimators based on stochastic gradients. *Annals of Statistics*, **45**, 1694–1727.
- Toulis, P., Horel, T. and Airoldi, E. M. (2021) The proximal Robbins–Monro method. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **83**, 188–212.
- Toulis, P., Tran, D. and Airoldi, E. (2016) Towards stability and optimality in stochastic gradient descent. In *Artificial Intelligence and Statistics*, vol. 51, 1290–1298. PMLR.
- Wang, B.-Y. and Gong, M.-P. (1993) Some eigenvalue inequalities for positive semidefinite matrix power products. *Linear Algebra and Its Applications*, **184**, 249–260.