

Manuscript version: Author's Accepted Manuscript

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/142420>

How to cite:

Please refer to published version for the most recent bibliographic citation information.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk.

Generalised Bayesian Filtering via Sequential Monte Carlo

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce a framework for inference in general state-space hidden Markov
2 models (HMMs) under likelihood misspecification. In particular, we leverage
3 the loss-theoretic perspective of Generalized Bayesian Inference (GBI) to define
4 generalised filtering recursions in HMMs, that can tackle the problem of inference
5 under model misspecification. In doing so, we arrive at principled procedures for
6 robust inference against observation contamination by utilising the β -divergence.
7 Operationalising the proposed framework is made possible via sequential Monte
8 Carlo methods (SMC), where most standard particle methods, and their associated
9 convergence results, are readily adapted to the new setting. We apply our approach
10 to object tracking and Gaussian process regression problems, and observe improved
11 performance over both standard filtering algorithms and other robust filters.

12 1 Introduction

13 Estimating the hidden states in dynamical systems is a long-standing problem in many fields of sci-
14 ence and engineering. This can be formulated as an inference problem of a general state-space hidden
15 Markov model (HMM) defined via two processes, *the hidden process* $(\mathbf{x}_t)_{t \geq 0}$, and *the observation pro-*
16 *cess* $(\mathbf{y}_t)_{t \geq 1}$. More precisely, we consider the general state-space hidden Markov models of the form

$$17 \quad \mathbf{x}_0 \sim \pi_0(\mathbf{x}_0), \quad (1) \quad \mathbf{x}_t | \mathbf{x}_{t-1} \sim f_t(\mathbf{x}_t | \mathbf{x}_{t-1}), \quad (2) \quad \mathbf{y}_t | \mathbf{x}_t \sim g_t(\mathbf{y}_t | \mathbf{x}_t), \quad (3)$$

18 where $\mathbf{x}_t \in \mathcal{X}$ for $t \geq 0$, $\mathbf{y}_t \in \mathcal{Y}$ for $t \geq 1$, f_t is a Markov kernel on \mathcal{X} and $g_t : \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{R}_+$ is the
19 likelihood function. We assume $\mathcal{X} \subseteq \mathbb{R}^{d_x}$ and $\mathcal{Y} \subseteq \mathbb{R}^{d_y}$ for convenience; however, the extension to
20 general Polish spaces follows directly. The key inference problem in this model class is estimating
21 the *filtering distributions*, i.e. the posterior distributions of the hidden states $(\mathbf{x}_t)_{t \geq 0}$ given the
22 observations $\mathbf{y}_{1:t}$ denoted as $(\pi_t(\mathbf{x}_t | \mathbf{y}_{1:t}))_{t \geq 1}$ — commonly known as *Bayesian filtering* [1, 2].

24 Under assumptions of linearity and Gaussianity, the inference problem for the hidden states of HMMs
25 can be solved analytically via the Kalman filter [3]. However, inference for general HMMs of the form
26 (1)–(3) with nonlinear, non-Gaussian transitions and likelihoods lacked a general, principled solution
27 until the arrival of the particle filtering schemes [4]. Particle filters (PFs) have become ubiquitous for
28 Bayesian filtering in the general setting. In short, the PFs retain a weighted collection of Monte Carlo
29 samples representing the filtering distribution $\pi_t(\mathbf{x}_t | \mathbf{y}_{1:t})$ and recursively approximate the sequence
30 of distributions $(\pi_t)_{t \geq 0}$ using a particle mutation-selection scheme [5].

31 While PFs (and other inference schemes for HMMs) implicitly assume that the assumed model
32 is well-specified, it is important to consider whether the proposed model class includes the true
33 data-generating mechanism (DGM). In particular, for general state-space HMMs, misspecification
34 can occur if the true dynamics of the hidden process significantly differ from the assumed model

35 f_t , or if the true observation model is markedly different from the assumed likelihood model g_t , e.g.
 36 corruption by heavy tailed noise. The latter case is of widespread interest within the field of *robust*
 37 *statistics* [6] and has recently attracted significant interest in the machine learning community [7]. It
 38 is the setting that this paper seeks to address.

39 When the true DGM cannot be modelled, one principled approach to address misspecification is
 40 Generalized Bayesian Inference (GBI) [8]. This approach views classical Bayesian inference as
 41 a loss minimisation procedure in the space of probability measures, a view first developed by [9].
 42 In particular, the standard Bayesian update can be derived from this view, where a loss function
 43 is constructed using the Kullback-Leibler (KL) divergence from the empirical distribution of the
 44 observations to the assumed likelihood [8]. The KL divergence is sensitive to outliers [10], hence
 45 the overall inference procedure is not robust to observations that are incompatible with the assumed
 46 model. A principled remedy is to replace the KL divergence with alternative discrepancy, such as the
 47 β -divergence, which makes the overall procedure more robust [11] while retaining interpretability.

48 Previous work on robust particle filters have been done for handling outliers, sensor failures and
 49 misspecification of the transition model [12, 13, 14, 15, 16, 17, 18, 19]. However, these approaches
 50 are either based on problem-specific heuristic outlier detection schemes, or make strong assumptions
 51 about the DGM in order to justify the use of heavy-tailed distributions [15]. This requires knowledge
 52 of the contamination mechanism that is implicitly embedded in the likelihood.

53 In this work we propose a principled approach to robust filtering that does not impose additional
 54 modelling assumptions. We adapt the GBI approach of [8] to the Bayesian filtering setting and develop
 55 sequential Monte Carlo (SMC) methods for inference. We illustrate the performance of this approach,
 56 using the β -divergence, to mitigate the effect of outliers. We show that this approach significantly
 57 improves the PF performance in settings with contaminated data, while retaining a general and
 58 principled approach to inference. We provide empirical results that demonstrate improvement over
 59 Kalman and particle filters for both linear and non-linear HMMs. We further provide comparisons
 60 with various robust schemes against heavy-tailed noise, including t-based likelihoods [15] or auxiliary
 61 particle filters (APFs) [12]. Finally, exploiting the state-space representations of Gaussian processes
 62 (GPs) [20], we demonstrate our framework on London air pollution data using robust GP regression
 63 which has linear time-complexity in the number of observations.

64 **Notation.** We denote the space of bounded, Borel measurable functions on X as $B(X)$. We denote the
 65 Dirac measure located at y as $\delta_y(dx)$ and note that $f(y) = \int f(x)\delta_y(dx)$ for $f \in B(X)$. We denote
 66 the Borel subsets of X as $\mathcal{B}(X)$ and the set of probability measures on $(X, \mathcal{B}(X))$ as $\mathcal{P}(X)$. For a
 67 probability measure $\mu \in \mathcal{P}(X)$ and $\varphi \in B(X)$, we write $\mu(\varphi) := \int \varphi(x)\mu(dx)$. Given a probability
 68 measure μ , we abuse the notation denoting its density with respect to the Lebesgue measure as $\mu(x)$.

69 2 Background

70 2.1 Generalized Bayesian Inference (GBI)

71 Bayesian inference implicitly assumes that the generative model is well-specified, in particular, the
 72 observations are generated from the assumed likelihood model. When this assumption is not expected
 73 to hold in real-world scenarios, one may wish to take into account the discrepancy between the true
 74 DGM and the assumed likelihood. GBI [8] is an approach to deal with such cases. Here, we present
 75 the main idea of GBI and refer the reader to the appendix for a more detailed description and to the
 76 original reference for a full-treatment.

77 For the simple Bayesian updating setup, consider a prior π_0 and the assumed likelihood function
 78 $g(y|x)$. The posterior $\pi(x|y) =: \pi(x)$ is given by Bayes rule $\pi(x) = \pi_0(x) \frac{g(y|x)}{Z}$, where $Z :=$
 79 $\int g(y|x)\pi_0(x)dx$. [9] and [8] showed that this update can be seen as a special case of a more general
 80 update rule, which can be described as a solution of an optimisation problem in the space of measures.
 81 This leads to a more general belief updating rule given by

$$\pi(x) = \pi_0(x) \frac{G(y|x)}{Z}, \quad (4)$$

82 with $G(y|x) := \exp(-\ell(x, y))$ where $\ell(x, y)$ is a loss function connecting the observations to the
 83 model parameters. Specifying $\ell(x, y)$ as the cross-entropy (from the KL-divergence) of the assumed
 84 likelihood relative to the empirical distribution of the data recovers the standard Bayes update.

85 As noted before, the standard Bayes update is not robust to outliers due to the properties of KL
 86 divergence [10]. Hence, substituting the cross-entropy with a more robust loss such as the β -cross-
 87 entropy [7], based on the β -divergence, can make the inference more robust. Specifically, in this
 88 setting the generalised Bayes update for the likelihood $g(\mathbf{y}|\mathbf{x})$ is written as $\pi(\mathbf{x}) = \pi_0(\mathbf{x}) \frac{G^\beta(\mathbf{y}|\mathbf{x})}{Z_\beta}$,
 89 where

$$G^\beta(\mathbf{y}|\mathbf{x}) = \exp\left(\frac{1}{\beta}g(\mathbf{y}|\mathbf{x})^\beta - \frac{1}{\beta+1} \int g(\mathbf{y}'|\mathbf{x})^{\beta+1} d\mathbf{y}'\right). \quad (5)$$

90 One can consider $G^\beta(\mathbf{y}|\mathbf{x})$ as a generalised likelihood, resulting from the use of a different loss
 91 function compared to the standard Bayes procedure. Here β is a hyperparameter that needs to be
 92 selected depending on the degree of misspecification. In general $\beta \in (0, 1)$ and $\lim_{\beta \rightarrow 0} G^\beta(\mathbf{y}|\mathbf{x}) =$
 93 $g(\mathbf{y}|\mathbf{x})$. Thus, intuitively, small β values are suitable for mild model misspecification and large β
 94 values are suitable when the assumed model is expected to significantly deviate from the true model.
 95 In the experimental section, we devote some attention to the selection of β and sensitivity analysis.

96 Generalised Bayesian updating is more robust against outliers if a suitable divergence is chosen
 97 [21, 22, 10]. We note that GBI is conceptually different from approximate Bayesian methods with
 98 alternative divergences such as [23, 24, 25, 26]. While these methods target approximate posteriors
 99 that minimise some discrepancy from the true posterior and are not necessarily robust, GBI methods
 100 change the inference target from the standard Bayesian posterior (obtained by setting $\ell(\mathbf{x}, \mathbf{y})$ to the
 101 KL divergence) to a different target distribution with more desirable properties such as robustness to
 102 outliers. Later, we demonstrate how the GBI approach can be used to construct robust PF procedures.

103 2.2 Sequential Monte Carlo for HMMs

104 Let $\mathbf{x}_{1:T}$ be a hidden process with $\mathbf{x}_t \in \mathsf{X}$ and $\mathbf{y}_{1:T}$ an observation process with $\mathbf{y}_t \in \mathsf{Y}$. Our goal is
 105 to conduct inference in HMMs of the form (1)–(3) where $\pi_0(\cdot)$ is a prior probability distribution on
 106 the initial state \mathbf{x}_0 , $f_t(\mathbf{x}|\mathbf{x}')$ is a Markov transition kernel on X and $g_t(\mathbf{y}_t|\mathbf{x}_t)$ is the likelihood for
 107 observation \mathbf{y}_t . The observation sequence $\mathbf{y}_{1:T}$ is assumed to be fixed but otherwise arbitrary.

108 The typical interest in probabilistic models is the estimation of expectations of general test functions
 109 with respect to the posterior distribution, in this case, of the hidden process $\pi_t(\mathbf{x}_t|\mathbf{y}_{1:t})$ and the
 110 associated joint distributions $\mathfrak{p}_t(\mathbf{x}_{0:t}|\mathbf{y}_{1:t})$. More precisely, given a bounded test function $\varphi \in B(\mathsf{X})$,
 111 we are interested in estimating integrals of the form

$$\pi_t(\varphi) = \int \varphi(\mathbf{x}_t) \pi_t(\mathbf{x}_t|\mathbf{y}_{1:t}). \quad (6)$$

112 Kalman filtering [3, 1] can be used to obtain closed form expressions for $(\pi_t, \mathfrak{p}_t)_{t \geq 0}$ if f_t and g_t are
 113 linear-Gaussian. However, for non-linear or non-Gaussian cases, the target distributions are almost
 114 always intractable, requiring an alternative approach, such as SMC methods [5, 27]. Known as Particle
 115 Filters (PFs) when employed in the HMM setting, SMC methods combine importance sampling and
 116 resampling algorithms tailored to approximate the solution of the filtering and smoothing problems.

117 In a typical iteration, a PF method proceeds as follows: given a collection of samples $\{\mathbf{x}_{t-1}^{(i)}\}_{i=1}^N$
 118 representing the posterior $\pi_{t-1}(\mathbf{x}_{t-1}|\mathbf{y}_{1:t-1})$, it first samples from a (possibly observation dependent)
 119 proposal $\bar{\mathbf{x}}_t^{(i)} \sim q_t(\mathbf{x}_t|\mathbf{x}_{1:t-1}^{(i)}, \mathbf{y}_{1:t})$. It then computes weights for each sample (particle) $\bar{\mathbf{x}}_{t-1}^{(i)}$
 120 in the collection for a given observation \mathbf{y}_t , evaluating its fitness with respect to the likelihood g_t as
 121 $\mathbf{w}_t^{(i)} \propto g_t(\mathbf{y}_t|\bar{\mathbf{x}}_t^{(i)}) \frac{f_t(\bar{\mathbf{x}}_t^{(i)}|\mathbf{x}_{t-1}^{(i)})}{q_t(\bar{\mathbf{x}}_t^{(i)}|\mathbf{x}_{1:t-1}^{(i)}, \mathbf{y}_t)}$, where $\sum_{i=1}^N \mathbf{w}_t^{(i)} = 1$. Finally, an optional resampling step¹
 122 is used to prevent degeneracy, leading to $\mathbf{x}_t^{(i)} \sim \sum_{i=1}^N \mathbf{w}_t^{(i)} \delta_{\bar{\mathbf{x}}_t^{(i)}}(d\mathbf{x}_t)$. One can then construct the
 123 empirical measure $\pi_t^N(d\mathbf{x}_t|\mathbf{y}_{1:t}) = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_t^{(i)}}(d\mathbf{x}_t)$, and the estimate of $\pi_t(\varphi)$ in (6) is given by

$$\pi_t^N(\varphi) = \frac{1}{N} \sum_{i=1}^N \varphi(\mathbf{x}_t^{(i)}). \quad (7)$$

¹In the simplest form, drawing N times with replacement from the weighted empirical measure to obtain an unweighted sample whose empirical distribution approximates the same target; see [28] for an overview of resampling schemes and their properties.

Algorithm 1 The generalised particle filter

Input: Observation sequence $\mathbf{y}_{1:T}$, number of samples N , proposal distributions $q_{1:T}(\cdot)$.

Initialize: Sample $\{\bar{\mathbf{x}}_0^{(i)}\}_{i=1}^N$ for the prior $\pi_0(\mathbf{x}_0)$.

for $t = 1$ **to** T **do**

Sample: $\bar{\mathbf{x}}_t^{(i)} \sim q_t(\mathbf{x}_t | \mathbf{x}_{1:t-1}^{(i)}, \mathbf{y}_t)$, for $i = 1, \dots, N$.

Weight: $w_t^{(i)} \propto \exp(-\ell(\bar{\mathbf{x}}_t^{(i)}, \mathbf{y}_t)) \frac{f_t(\bar{\mathbf{x}}_t^{(i)} | \mathbf{x}_{t-1}^{(i)})}{q_t(\bar{\mathbf{x}}_t^{(i)} | \mathbf{x}_{1:t-1}^{(i)}, \mathbf{y}_t)}$, for $i = 1, \dots, N$.

Resample: $\mathbf{x}_t^{(i)} \sim \sum_{i=1}^N w_t^{(i)} \delta_{\bar{\mathbf{x}}_t^{(i)}}(d\mathbf{x}_t)$, for $i = 1, \dots, N$.

end for

124 If the proposal is chosen as the transition density, i.e., $q_t(\mathbf{x}_t | \mathbf{x}_{1:t-1}^{(i)}, \mathbf{y}_t) = f_t(\mathbf{x}_t | \mathbf{x}_{t-1}^{(i)})$, we obtain
125 the bootstrap particle filter (BPF) [4]. This corresponds to the simple procedure of sampling $\bar{\mathbf{x}}_t^{(i)}$
126 from $f_t(\mathbf{x}_t | \mathbf{x}_{t-1}^{(i)})$, and setting its weight $w_t^{(i)} \propto g_t(\mathbf{y}_t | \bar{\mathbf{x}}_t^{(i)})$.

127 3 Generalised Bayesian filtering

128 3.1 A simple generalised particle filter

129 As explained in Section 2.1, given a standard probability model comprised of the prior $\pi_0(\mathbf{x})$ and a
130 likelihood $g(\mathbf{y} | \mathbf{x})$, the general Bayes update defines an alternative, generalised likelihood $G(\mathbf{y} | \mathbf{x})$.
131 The sequence of generalised likelihoods, denoted as $G_t(\mathbf{y}_t | \mathbf{x}_t)$ for $t \geq 1$, in an HMM yields a joint
132 generalised posterior density which factorises as

$$p_t(\mathbf{x}_{0:t} | \mathbf{y}_{1:t}) \propto \pi_0(\mathbf{x}_0) \prod_{k=1}^t f_k(\mathbf{x}_k | \mathbf{x}_{k-1}) G_k(\mathbf{y}_k | \mathbf{x}_k), \quad (8)$$

133 where $G_t(\mathbf{y}_t | \mathbf{x}_t) := \exp(-\ell_t(\mathbf{x}_t, \mathbf{y}_t))$. Inference can be done via SMC applied to this sequence of
134 twisted probabilities defining a Feynman-Kac flow in the terminology of [29].

135 Comparing the update rule in (4) to the standard Bayes update suggests a generalisation of the particle
136 filter. In particular, under the model in (1)–(3), one can perform generalised inference using $(f_t)_{t \geq 1}$
137 as usual, but replacing the likelihood with $(G_t)_{t \geq 1}$. Hence, a generalised sequential importance
138 resampling PF (given fully in Algorithm 1) keeps the sampling step intact, but applies a different
139 weight computation step $w_t^{(i)} \propto \exp(-\ell(\bar{\mathbf{x}}_t^{(i)}, \mathbf{y}_t)) \frac{f_t(\bar{\mathbf{x}}_t^{(i)} | \mathbf{x}_{t-1}^{(i)})}{q_t(\bar{\mathbf{x}}_t^{(i)} | \mathbf{x}_{1:t-1}^{(i)}, \mathbf{y}_t)}$. Indeed, most PFs (including the
140 APF, see Algorithm 3 in the appendix) and related algorithms can be adapted to the GBI context.

141 3.2 The β -BPF and the β -APF

142 The β -BPF is derived by selecting $\ell_t(\mathbf{x}_t, \mathbf{y}_t)$ as the β -divergence and applying the BPF procedure
143 with the associated generalised likelihood. In this case, the loss is

$$\ell_t^\beta(\mathbf{x}_t, \mathbf{y}_t) = \frac{1}{\beta + 1} \int g_t(\mathbf{y}'_t | \mathbf{x}_t)^{\beta+1} d\mathbf{y}'_t - \frac{1}{\beta} g_t(\mathbf{y}_t | \mathbf{x}_t)^\beta. \quad (9)$$

144 We can then construct the general β -likelihood as

$$G_t^\beta(\mathbf{y}_t | \mathbf{x}_t) \propto \exp(-\ell_t^\beta(\mathbf{x}_t, \mathbf{y}_t)). \quad (10)$$

145 In this instance, the use of the β -divergence provides the sampler with robust properties [11]. This
146 can informally be seen from the form of the loss function in (9), where small values of β temper
147 the likelihood extending its tails making the loss more forgiving to outliers. The β -BPF procedure
148 is given in Algorithm 2 in the appendix. The β -APF (Algorithm 3 in the appendix) is an Auxiliary
149 Particle Filter [12, 30] adapted to the GBI setting, and is derived similarly to the β -BPF.

150 Note that the integral term in (9) is independent of \mathbf{x}_t and can be absorbed, without evaluation, into
151 the normalising constant when \mathbf{x}_t is a location parameter for a symmetric $g_t(\cdot)$ and \mathbf{Y} is a linear
152 subspace of \mathbb{R}^{d_y} . More generally, if $g_t(\cdot)$ is a member of the exponential family, the integral can be

153 computed by identifying $g_t^\beta(\cdot)$ with the kernel of another member of the same family with canonical
 154 parameters scaled by β . The overhead of computing $G_t^\beta(\cdot)$ is negligible in this instance, which is
 155 not too restrictive in the context of misspecified models. For other likelihoods, unbiased estimators
 156 for $G_t^\beta(\cdot)$, e.g. Poisson estimator [31], can be used in a random weight particle filter framework
 157 [32], where the overhead of computing $G_t^\beta(\cdot)$ will depend on the variance of the estimator and the
 158 convergence results from this setting apply but as [32] demonstrate this cost need not be prohibitive.

159 3.3 Selecting β

160 It is often the case that the primary goal of inference, particularly in the presence of model misspeci-
 161 fication, is prediction. Hence, we propose choosing divergence parameters that lead to maximally
 162 predictive posterior belief distributions. In particular, for the β -BPF and β -APF, define $\mathcal{L}_\beta(\mathbf{y}_t, \hat{\mathbf{y}}_t)$ as
 163 a loss function of the observations \mathbf{y}_t and the predictions $\hat{\mathbf{y}}_t$. We propose to choose β as the solution
 164 to the following decision-theoretic optimisation problem:

$$\min_{\beta} \text{agg}_{t=1}^T (\mathbb{E}_{p(\hat{\mathbf{y}}_t | \mathbf{y}_{1:t-1})} \mathcal{L}_\beta(\mathbf{y}_t, \hat{\mathbf{y}}_t)), \quad (11)$$

165 where agg denotes an aggregating function. This approach requires some training data to allow the
 166 selection of β . In filtering contexts, this can be historical data from the same setting or other available
 167 proxies. For offline inference one could also employ the actual data within this framework. Since,
 168 this proposal relies on the quality of the observations, which in the case of outlier contamination is
 169 violated by definition. To remedy this, we propose choosing robust versions for agg and \mathcal{L} , e.g. the
 170 median and the (standardised) absolute error respectively.

171 4 Theoretical guarantees

172 Theoretical guarantees for SMC methods can be extended to the generalised Bayesian filtering
 173 setting. Since the generalised Bayesian filters can be seen as a standard SMC methods with modified
 174 likelihoods, the same analytical tools can be used in this setting. We provide guarantees for the β -BPF
 175 but emphasise that essentially the same results can be obtained much more broadly (including for the
 176 β -APF via the approach of [30]). We denote the generalised filters and generalised posteriors for the
 177 HMM in the β -divergence setting as π_t^β and \mathbf{p}_t^β respectively. Consequently, corresponding quantities
 178 constructed by the β -BPF are denoted as $\pi_t^{\beta, N}$ and $\mathbf{p}_t^{\beta, N}$.

179 Although the generalised likelihoods $G_t^\beta(\mathbf{y}_t | \mathbf{x}_t)$ are not normalised, they can be considered as
 180 potential functions [29]. Since $G_t^\beta(\mathbf{y}_t | \mathbf{x}_t) < \infty$ whenever $g_t(\mathbf{y}_t | \mathbf{x}_t) < \infty$ and β is fixed, we can
 181 adapt the standard convergence results into the generalised case.

182 **Assumption 1.** For a fixed arbitrary observation sequence $\mathbf{y}_{1:T} \in \mathcal{Y}^{\otimes T}$, the potential functions
 183 $(G_t^\beta)_{t \geq 1}$ are bounded and $G_t^\beta(\mathbf{y}_t | \mathbf{x}_t) > 0$, $\forall t \in \{1, \dots, T\}$ and $\mathbf{x}_t \in \mathcal{X}$.

184 This assumption holds for most used likelihood functions and their generalised extensions.

185 **Theorem 1.** For any $\varphi \in B(\mathcal{X})$ and $p \geq 1$, $\|\pi_t^{\beta, N}(\varphi) - \pi_t^\beta(\varphi)\|_p \leq \frac{c_{t,p,\beta} \|\varphi\|_\infty}{\sqrt{N}}$, where $c_{t,p,\beta} < \infty$
 186 is a constant independent of N .

187 The proof sketch and the constant $c_{t,p,\beta}$ are in the supplement. This L_p bound provides a theoretical
 188 guarantee on the convergence of particle approximations to generalised posteriors. The special case
 189 when $p = 2$ also provides the error bound for the mean-squared error. It is well known that Theorem 1
 190 with $p > 2$ leads to a law of large numbers via Markov's inequality and a Borel-Cantelli argument:

191 **Corollary 1.** Under the setting of Theorem 1, $\lim_{N \rightarrow \infty} \pi_t^{\beta, N}(\varphi) = \pi_t^\beta(\varphi)$ a.s., for $t \geq 1$.

192 Finally, a central limit theorem for estimates of expectations with respect to the smoothing distribu-
 193 tions can be obtained by considering the path space $\mathcal{X}^{\otimes t}$. Recall the joint posterior $\mathbf{p}_t^\beta(\mathbf{x}_{1:t} | \mathbf{y}_{1:t})$ and
 194 consider a test function $\varphi_t : \mathcal{X}^{\otimes t} \rightarrow \mathbb{R}$. We denote $\bar{\varphi}_t^\beta := \int \varphi_t^\beta(\mathbf{x}_{1:t}) \mathbf{p}_t^\beta(\mathbf{x}_{1:t} | \mathbf{y}_{1:t})$ and denote the
 195 β -BPF estimate of $\bar{\varphi}_t$ with $\bar{\varphi}_t^{\beta, N} := \int \varphi_t(\mathbf{x}_{1:t}) \mathbf{p}_t^{\beta, N}(\mathbf{x}_{1:t})$.

196 **Theorem 2.** Under the regularity conditions given in [33, Theorem 1], $\sqrt{N} (\bar{\varphi}_t^{\beta, N} - \bar{\varphi}_t^\beta) \xrightarrow{d}$
 197 $\mathcal{N}(0, \sigma_{t,\beta}^2(\varphi_t))$, as $N \rightarrow \infty$ where $\sigma_{t,\beta}^2(\varphi_t) < \infty$.

198 The expression for $\sigma_{t,\beta}^2(\varphi_t)$ can be found in the appendix. These results illustrate that the standard
 199 guarantees for generic particle filtering methods extend to our case.

200 5 Experiments

201 In this section, we focus on β -BPF illustrating its the properties and empirically verifying its robust-
 202 ness. We include three experiments in the main text and another in Appendix D. Furthermore, we
 203 specifically investigate the β -APF in Section 5.2 comparing its behaviour to the β -BPF. Throughout,
 204 we report the *normalised mean squared error (NMSE)* and the *90% empirical coverage* as goodness-
 205 of-fit measures. The NMSE scores indicate the mean fit for the inferred posterior distribution and
 206 the empirical coverage measures the quality of its uncertainty quantification. We note that any claim
 207 in performance difference is based on the Wilcoxon signed-rank test. Further results and in-depth
 208 details on the experimental setup are given in the supplementary material.

209 5.1 A Linear-Gaussian state-space model

210 The Wiener velocity model [34] is a standard model in the target tracking literature,
 211 where the velocity of a particle is modelled as a Wiener process. The discretised ver-
 212 sion of this model can be represented as a Linear-Gaussian State-Space model (LGSSM),

$$213 \quad \mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \boldsymbol{\nu}_{t-1}, \quad \boldsymbol{\nu}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}), \quad (12) \quad \mathbf{y}_t = \mathbf{H}\mathbf{x}_t + \boldsymbol{\epsilon}_t, \quad \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \Sigma), \quad (13)$$

214 where \mathbf{A}, \mathbf{Q} are state-transition pa-
 215 rameters dictated by the continuous-
 216 time model and \mathbf{H} is the observa-
 217 tion matrix (see Appendix). We simulate
 218 this model in two-dimensions with
 219 $\Sigma = \mathbf{I}$, contaminating the observa-
 220 tions with a large scale, zero-mean
 221 Gaussian, $\mathcal{N}(0, 100^2)$ with probabili-
 222 ty p_c . Our aim is to obtain the
 223 filtering density under the heavily-
 224 contaminated setting where optimal
 225 filters struggle to perform. We com-
 226 pare our scheme for a large range of
 227 β to the standard BPF with a Gaus-
 228 sian likelihood (BPF), as well as the
 229 (optimal) Kalman filter.

230 We shed light onto three questions on
 231 this simple setup: (a) Does the β -BPF
 232 produce accurate and well-calibrated
 233 posterior distributions in the presence
 234 of contaminated data? (b) Is it sen-
 235 sitive to the choice β ? (c) Does the
 236 method described in Section 3.3 for
 237 selecting β return a near optimal re-
 238 sult?

239 Figure 1 shows the results for $p_c =$
 240 0.1. We observe that (a) the β -BPF
 241 outperforms the Kalman filter and the standard BPF for $\beta \leq 0.2$
 242 while producing well-calibrated posteriors accounting for the uncertainty (for $\beta \in [0.01, 0.2]$
 243 the coverage approaches the 90% threshold), (b) we see drastic performance gains (with median NMSE
 244 scores around $10\times$ smaller than the BPF and $100\times$ smaller that the Kalman filter) for a large range
 of β values, (c) we also see that the β -choice heuristic² chooses a well-performing β (gold vertical

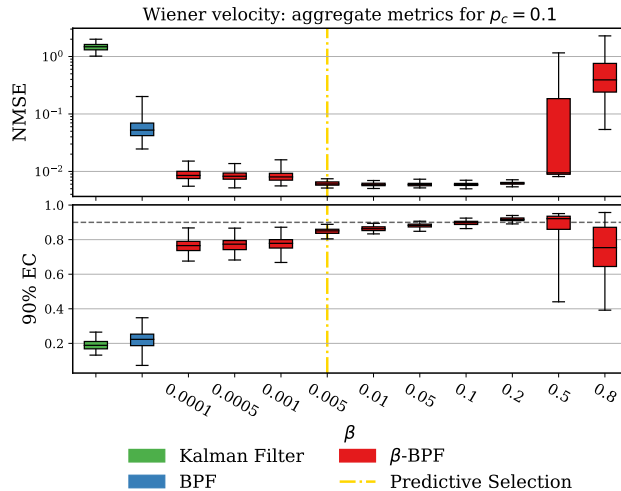


Figure 1: The mean metrics over state dimensions for the Wiener velocity example with $p_c = 0.1$. The top panel presents the NMSE results (lower is better) and the bottom panel presents the 90% empirical coverage results (higher is better), on 100 runs. The vertical dashed line in gold indicate the value of β chosen by the selection criterion in Section 3.3. The horizontal dashed line in black in the lower panel indicates the 90% mark for the coverage.

²We apply this choice criterion on an alternative dataset that is obtained from the same simulation but with 90% fewer observations.

245 lines in Figure 1). Note that, for most values of β , the β -BPF significantly outperforms both the
 246 Kalman filter and the standard BPF predictively. The full set of results for the predictive performance
 247 are presented in Table 4 in Appendix F.1.

248 5.2 Terrain Aided Navigation

249 Terrain Aided Navigation (TAN) is a challenging estimation problem, where the state evolution
 250 is defined as in (12) (in three dimensions), but with a highly non-linear observation model, $\mathbf{y}_t =$
 251 $h(\mathbf{x}_t) + \epsilon_t$, where $h(\cdot)$ is a non-linear function, typically including a non-analytic Digital Elevation
 252 Map (DEM). This problem simulates the trajectory of an aeroplane or a drone over a terrain map,
 253 where we observe its elevation over the terrain and its distance from its take-off hub from on-board
 254 sensors (see supplement for more details). We simulate transmission failure of the measurement
 255 system as impulsive noise on the observations, i.e., i.i.d. draws from a Student's t distribution with
 256 $\nu = 1$ degrees of freedom. In other words, we define $\epsilon_t \sim (1 - p_c)\mathcal{N}(0, 20^2) + p_c t_{\nu=1}(0, 20^2)$.

257 We apply both the β -BPF and the β -APF to this problem and compare them to the standard BPF
 258 with the Gaussian (BPF). We also compare against two other robust PF methods from the literature:
 259 Student's t (t-BPF) [15] and the APF [12]. We set the degrees of freedom for the t-BPF to the same
 260 value as the contamination $\nu = 1$.

261 From Figure 2, we observe
 262 that for low contamination, both
 263 the β -BPF and the β -APF out-
 264 perform the standard Gaussian
 265 BPF, the t-BPF and the APF.
 266 This shows that the use of t -
 267 distribution for the low contam-
 268 ination setting is inappropriate.
 269 This gap in the performance
 270 tightens, naturally, as p_c grows
 271 since t -distribution becomes a
 272 good model for the observations.
 273 Notably, the performance gaps
 274 between the standard PFs and
 275 their β -robustified counterparts
 276 are similar, indicating that the
 277 use of the β -divergence in a par-
 278 ticle filtering procedure does in-
 279 deed robustify the inference.

280 In Figure 3, we plot the filtering
 281 distributions for the sixth state
 282 dimension (vertical velocity) obtained from an illustrative run with $p_c = 0.1$. The top panel shows the
 283 filtering distributions from the (Gaussian) BPF (up) and the β -BPF (down). The locations of the most
 284 prominent outliers are marked with dashed vertical lines in black. Figure 3 displays the significant
 285 difference between the two approaches: while the uncertainty for the standard BPF collapses when
 286 it meets the outliers, e.g. around $t = 1700$, the β -BPF does not suffer from this problem. This
 287 performance difference is partly related to the stability of the weights. The lower panel in Figure 3
 288 demonstrates the effective sample size (ESS) with time for the two filters showing that the β -BPF
 289 consistently exhibits larger ESS values, avoiding particle degeneracy. The ESS values for the BPF,
 290 on the other hand, sharply decline when it meets outliers. A similar result is observed for the APF
 291 versus the β -APF in the figures in the Appendix F.2. Further results on predictive performance can be
 292 found in Appendix F.2.

293 5.3 London air quality Gaussian process regression

294 To measure air quality, London authorities use a network of sensors around the city recording pollutant
 295 measurements. Sensor measurements are susceptible to significant outliers due to environmental
 296 effects, manual calibration and sensor deterioration. In the experiment, we use Gaussian process (GP)
 297 regression to infer the underlying signal from a PM2.5 sensor.

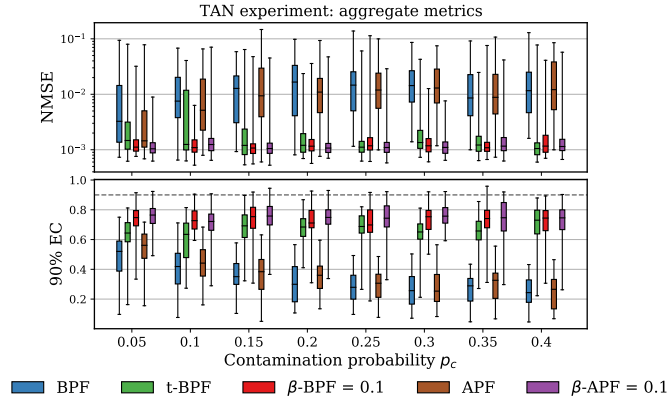


Figure 2: The mean metrics over state dimensions for the TAN example for different p_c . The top panel presents the NMSE results (lower is better) and the bottom panel presents the 90% empirical coverage results (higher is better), both evaluated on 50 runs. The horizontal dashed line in black in the lower panel indicate the 90% mark for the coverage.

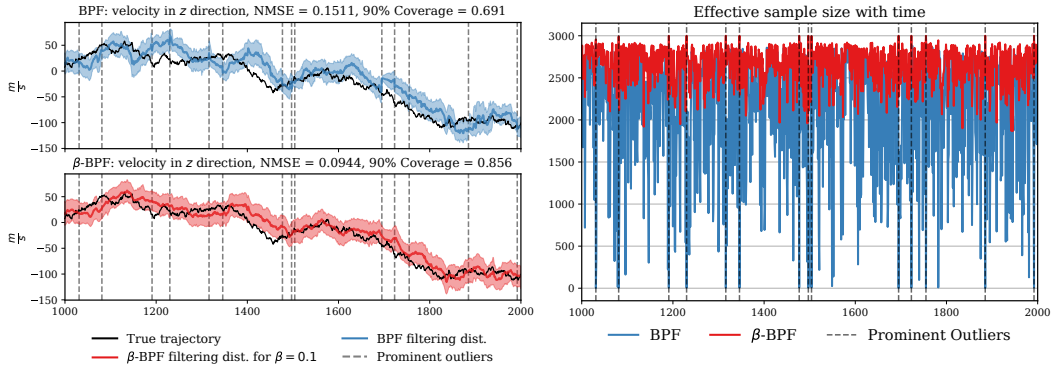


Figure 3: The left panel shows the inferred marginal filtering distributions for the velocity in the z direction for the BPF and β -BPF with $\beta = 0.1$. The right panel shows the effective sample size with time. The locations of the most prominent (largest deviation) outliers are shown as dashed vertical lines in black in both panels.

298 For 1-D time series data, GP inference [35] can be accelerated to linear time in the number of
 299 observations by formulating an equivalent stochastic differential equation whose solution precisely
 300 matches the GP under consideration³ [20]. The resulting model is a LGSSM of the form (12)–
 301 (13) where the smoothing distribution matches the GP marginals at discrete-times. One can then
 302 apply smoothing algorithms, such as Rauch Tung Striebel (RTS) [36] or Forward Filters Backward
 303 Smoothing (FFBS) [37], to obtain the GP posterior. These require a forward filtering step with the
 304 Kalman filter for RTS or a PF for FFBS. Here, we fit a Matérn 5/2 GP with known hyperparameters
 305 to a time series from one of the sensors. We plot the median of the signals from the wider sensor
 306 network to obtain a simple approximation of the ground truth.

307 To further investigate the GP solution of
 308 the β -BPF (FFBS), we show the fit for
 309 $\beta = 0.1$ and compare it with Kalman
 310 (RTS) smoothing. In Figure 24 we see
 311 that the latter is sensitive to outliers forcing
 312 the GP mean towards them while the
 313 β -BPF is robust and ignores them.

314 Table 1 compares results with a Gaus-
 315 sian likelihood for GP regression with
 316 Kalman (RTS) smoothing, the standard
 317 BPF (FFBS) and two runs for the β -BPF
 318 (FFBS) ($\beta = 0.1$ by predictive selection
 319 as Section 3.3 and $\beta = 0.2$ by overall
 320 best performance). For both choices of β , the β -BPF outperforms all other methods on both metrics .

Table 1: GP regression NMSE (higher is better) and 90% empirical coverage for the credible intervals of the posterior predictive distribution, on 100 runs. **Bold** indicates statistically significant best result from Wilcoxon signed-rank test. All presented results are statistically different from each other according to the test.

Filter (Smoother)	median (IQR)	
	NMSE	EC
Kalman (RTS)	0.144(0)	0.685(0)
BPF (FFBS)	0.116(0.015)	0.650(0.020)
($\beta = 0.1$)-BPF (FFBS)	0.061(0.003)	0.760(0.015)
($\beta = 0.2$)-BPF (FFBS)	0.059(0.002)	0.803(0.020)

321 6 Conclusions

322 We provided a generalised filtering framework based on GBI, which tackles likelihood misspecifi-
 323 cation in general state-space HMMs. Our approach leverages SMC methods, where we extended
 324 some analytical results to the generalised case. We presented the β -BPF, a simple instantiation of our
 325 approach based on the the β -divergence, developed an APF for this setting, and showed performance
 326 gains compared to other standard algorithms on a variety of problems and contamination settings.

327 This work opens up many exciting avenues for future research. Principle among which is online
 328 learning for model parameters (system identification) in the presence of misspecification. Our
 329 framework can directly incorporate most estimators found in the SMC literature and the computation
 330 of derivatives can be tackled with automatic differentiation tools.

³The SDE representation of a GP depends on the form of the covariance function. In this paper we use a GP with the Matérn 5/2 kernel, which admits a dual SDE representation.

331 7 Broader Impact

332 Robust inference in the context of misspecified models is a topic of broad interest. However, there are
333 a few robust generally-applicable methods which can be employed in the context of online inference
334 in time series settings. This paper provides a principled solution to this problem within a formal
335 framework backed by theoretical guarantees and opening up the benefits to multiple application
336 domains. The illustrative applications demonstrate the potential improvements in settings including
337 navigation and Gaussian process regression, which, if realised more widely, could have wide-reaching
338 impact. We hope that this inspires the community to build-on or apply our work to other challenging
339 real-world scenarios.

340 Of particular interest is the application of Robust SMC methods, like the β -BPF and the auxiliary
341 counterpart which were developed in this work, to impactful data-streaming applications in environ-
342 mental monitoring and forecasting. Indeed, our research in this area was motivated by a real-world
343 application in which existing techniques were inadequate (see *anonymized reference* for more details).
344 We have demonstrated the benefits such methods in proof-of-concept work and are incorporating the
345 resulting algorithms into a fully-developed platform, that has been in development for approximately
346 three years. We are partnering with local authorities to help in directly informing policy makers and
347 ultimately the general public.

348 More widely, this work provides an additional illustration that the GBI framework can provide
349 good solutions to challenging problems in the world of misspecified framework and hence provides
350 additional motivation to further investigate this extremely promising but rather new direction.

351 References

- 352 [1] Brian D O Anderson and John B Moore. *Optimal filtering*. Englewood Cliffs, N.J. Prentice
353 Hall, 1979.
- 354 [2] Simo Särkkä. *Bayesian Filtering and Smoothing*. Cambridge University Press, 2013.
- 355 [3] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. *Journal of*
356 *Fluids Engineering*, 82(1):35–45, 1960.
- 357 [4] Neil J Gordon, David J Salmond, and Adrian FM Smith. Novel approach to nonlinear/non-
358 Gaussian Bayesian state estimation. *IEE proceedings F (Radar and Signal Processing)*,
359 140(2):107–113, 1993.
- 360 [5] Arnaud Doucet, Simon Godsill, and Christophe Andrieu. On sequential Monte Carlo sampling
361 methods for Bayesian filtering. *Statistics and Computing*, 10(3):197–208, 2000.
- 362 [6] Peter J Huber. *Robust statistics*. Springer, 2011.
- 363 [7] Futoshi Futami, Issei Sato, and Masashi Sugiyama. Variational inference based on robust
364 divergences. In *International Conference on Artificial Intelligence and Statistics*, pages 813–
365 822, 2018.
- 366 [8] Pier Giovanni Bissiri, Chris C Holmes, and Stephen G Walker. A general framework for
367 updating belief distributions. *Journal of the Royal Statistical Society: Series B (Statistical*
368 *Methodology)*, 78(5):1103–1130, 2016.
- 369 [9] Arnold Zellner. Optimal information processing and Bayes’s theorem. *The American Statistician*,
370 42(4):278–280, 1988.
- 371 [10] Jeremias Knoblauch, Jack Jewson, and Theodoros Damoulas. Generalized Variational Inference:
372 Three arguments for deriving new Posteriors. *arXiv preprint arXiv:1904.02063*, 2019.
- 373 [11] Andrzej Cichocki and Shun-ichi Amari. Families of alpha-beta-and gamma-divergences:
374 Flexible and robust measures of similarities. *Entropy*, 12(6):1532–1568, 2010.
- 375 [12] Michael K Pitt and Neil Shephard. Filtering via simulation: Auxiliary particle filters. *Journal*
376 *of the American Statistical Association*, 94(446):590–599, 1999.

- 377 [13] Cristina S Maiz, Joaquin Miguez, and Petar M Djuric. Particle filtering in the presence of
378 outliers. In *2009 IEEE/SP 15th Workshop on Statistical Signal Processing*, pages 33–36. IEEE,
379 2009.
- 380 [14] Cristina S Maiz, Elisa M Molanes-Lopez, Joaquín Miguez, and Petar M Djuric. A particle
381 filtering scheme for processing time series corrupted by outliers. *IEEE Transactions on Signal*
382 *Processing*, 60(9):4611–4627, 2012.
- 383 [15] Dingjie Xu, Chen Shen, and Feng Shen. A robust particle filtering algorithm with non-Gaussian
384 measurement noise using student-t distribution. *IEEE Signal Processing Letters*, 21(1):30–34,
385 2013.
- 386 [16] Laurent E. Calvet, Veronika Czellar, and Elvezio Ronchetti. Robust filtering. *Journal of the*
387 *American Statistical Association*, 110(512):1591–1606, 2015.
- 388 [17] Francisco Curado Teixeira, João Quintas, Pramod Maurya, and António Pascoal. Robust
389 particle filter formulations with application to terrain-aided navigation. *International Journal of*
390 *Adaptive Control and Signal Processing*, 31(4):608–651, 2017.
- 391 [18] Xiao-Li Hu, Thomas B Schon, and Lennart Ljung. A robust particle filter for state estima-
392 tion—with convergence results. In *46th IEEE Conference on Decision and Control*, pages
393 312–317. IEEE, 2007.
- 394 [19] Ömer Deniz Akyildiz and Joaquín Míguez. Nudging the particle filter. *Statistics and Computing*,
395 30:305–330, 2020.
- 396 [20] Simo Särkkä, Arno Solin, and Jouni Hartikainen. Spatiotemporal learning via infinite-
397 dimensional Bayesian filtering and smoothing: A look at Gaussian process regression through
398 Kalman filtering. *IEEE Signal Processing Magazine*, 30(4):51–61, 2013.
- 399 [21] Abhik Ghosh and Ayanendranath Basu. Robust Bayes estimation using the density power
400 divergence. *Annals of the Institute of Statistical Mathematics*, 68(2):413–437, 2016.
- 401 [22] Jeremias Knoblauch, Jack E Jewson, and Theodoros Damoulas. Doubly robust Bayesian infer-
402 ence for non-stationary streaming data with β -divergences. In *Advances in Neural Information*
403 *Processing Systems*, pages 64–75, 2018.
- 404 [23] Tom Minka et al. Divergence measures and message passing. Technical report, Technical report,
405 Microsoft Research, 2005.
- 406 [24] Yingzhen Li and Richard E Turner. Rényi divergence variational inference. In *Advances in*
407 *Neural Information Processing Systems*, pages 1073–1081, 2016.
- 408 [25] Rajesh Ranganath, Dustin Tran, Jaan Altosaar, and David Blei. Operator variational inference.
409 In *Advances in Neural Information Processing Systems*, pages 496–504, 2016.
- 410 [26] Dilin Wang, Hao Liu, and Qiang Liu. Variational inference with tail-adaptive f-divergence. In
411 *Advances in Neural Information Processing Systems*, pages 5737–5747, 2018.
- 412 [27] Arnaud Doucet and Adam M Johansen. A tutorial on particle filtering and smoothing: Fifteen
413 years later. In D. Crisan and B. Rozovskiĭ, editors, *The Oxford Handbook of Nonlinear Filtering*,
414 pages 656–704. Oxford University Press, 2011.
- 415 [28] Mathieu Gerber, Nicolas Chopin, and Nick Whiteley. Negative association, ordering and
416 convergence of resampling methods. *Annals of Statistics*, 47(4):2236–2260, 2019.
- 417 [29] Pierre Del Moral. *Feynman-Kac formulae: Genealogical and interacting particle systems with*
418 *applications*. Springer, 2004.
- 419 [30] Adam M Johansen and Arnaud Doucet. A note on the auxiliary particle filter. *Statistics and*
420 *Probability Letters*, 78(12):1498–1504, September 2008.
- 421 [31] Alexandros Beskos, Omiros Papaspiliopoulos, Gareth O. Roberts, and Paul Fearnhead. Exact
422 and computationally efficient likelihood-based estimation for discretely observed diffusion
423 processes. *Journal of the Royal Statistical Society, Series B*, 68(3):333–382, 2006.

- 424 [32] Paul Fearnhead, Omiros Papaspiliopoulos, and Gareth O. Roberts. Particle filters for partially-
425 observed diffusion. *Journal of the Royal Statistical Society, Series B*, 70(4):755–777, 2008.
- 426 [33] Nicolas Chopin. Central limit theorem for sequential Monte Carlo methods and its application
427 to Bayesian inference. *The Annals of Statistics*, 32(6):2385–2411, 2004.
- 428 [34] Simo Särkkä and Arno Solin. *Applied Stochastic Differential Equations*, volume 10. Cambridge
429 University Press, 2019.
- 430 [35] Carl Edward Rasmussen and Christopher KI Williams. *Gaussian processes for machine learning*,
431 volume 1. MIT press Cambridge, 2006.
- 432 [36] Herbert E Rauch, F Tung, and Charlotte T Striebel. Maximum likelihood estimates of linear
433 dynamic systems. *American Institute of Aeronautics and Astronautics Journal*, 3(8):1445–1450,
434 1965.
- 435 [37] Mark Briers, Arnaud Doucet, and Simon Maskell. Smoothing algorithms for state space models.
436 *Annals of the Institute of Statistical Mathematics*, 62(1):61–89, 2010.
- 437 [38] John von Neumann and Oskar Morgenstern. *Theory of Games and Economic Behavior*. Prince-
438 ton University Press, 1947.
- 439 [39] Jack Jewson, Jim Q Smith, and Chris Holmes. Principles of Bayesian inference using general
440 divergence criteria. *Entropy*, 20(6):442, 2018.
- 441 [40] Pieralberto Guarniero, Adam M Johansen, and Anthony Lee. The iterated auxiliary particle
442 filter. *Journal of the American Statistical Association*, 112(520):1636–1647, 2017.
- 443 [41] Joaquín Míguez, Dan Crisan, and Petar M Djurić. On the convergence of two sequential Monte
444 Carlo methods for maximum a posteriori sequence estimation and stochastic global optimization.
445 *Statistics and Computing*, 23(1):91–107, 2013.