



**MONASH** University

**Australia**

Department of Econometrics and Business Statistics

<http://www.buseco.monash.edu.au/depts/ebs/pubs/wpapers/>

**Approximate Bayesian Computation in State Space  
Models**

**Gael M. Martin, Brendan P.M. McCabe, Worapree Maneesoonthorn  
and Christian P. Roberts**

**September 2014**

**Working Paper 20/14**

# Approximate Bayesian Computation in State Space Models\*

Gael M. Martin,<sup>†</sup> Brendan P.M. McCabe,<sup>‡</sup> Worapree Maneesoonthorn<sup>§</sup>  
and Christian P. Robert<sup>¶</sup>

September 29, 2014

## Abstract

A new approach to inference in state space models is proposed, based on approximate Bayesian computation (ABC). ABC avoids evaluation of the likelihood function by matching observed summary statistics with statistics computed from data simulated from the true process; exact inference being feasible only if the statistics are sufficient. With finite sample sufficiency unattainable in the state space setting, we seek asymptotic sufficiency via the maximum likelihood estimator (MLE) of the parameters of an auxiliary model. We prove that this auxiliary model-based approach achieves Bayesian consistency, and that - in a precise limiting sense - the proximity to (asymptotic) sufficiency yielded by the MLE is replicated by the score. In multiple parameter settings a separate treatment of scalar parameters, based on integrated likelihood techniques, is advocated as a way of avoiding the curse of dimensionality. Some attention is given to a structure in which the state variable is driven by a continuous time process, with exact inference typically infeasible in this case as a result of intractable transitions. The ABC method is demonstrated using the unscented Kalman filter as a fast and simple way of producing an approximation in this setting, with a stochastic volatility model for financial returns used for illustration.

*Keywords:* Likelihood-free methods, latent diffusion models, linear Gaussian state space models, asymptotic sufficiency, unscented Kalman filter, stochastic volatility.

*JEL Classification:* C11, C22, C58

---

\*This research has been supported by Australian Research Council (ARC) Future Fellowship FT0991045. We appreciate the input of various participants at *ABC in Rome*, June 2013, and *ABC in Sydney*, July 2014, on earlier drafts of the paper.

<sup>†</sup>Department of Econometrics and Business Statistics, Monash University, Melbourne, Australia. Corresponding author; email: gael.martin@monash.edu.

<sup>‡</sup>Management School, University of Liverpool, U.K.

<sup>§</sup>Melbourne Business School, University of Melbourne, Australia.

<sup>¶</sup>University of Paris Dauphine, Centre de Recherche en Économie et Statistique, and University of Warwick.

# 1 Introduction

Approximate Bayesian computation (ABC) (or likelihood-free inference) has become increasingly prevalent in areas of the natural sciences in which likelihood functions requiring integration over a large number of complex latent states are essentially intractable. (See Cornuet *et al.*, 2008, Beaumont, 2010, Marin *et al.*, 2011 and Sisson and Fan, 2011 for recent reviews.) The technique circumvents direct evaluation of the likelihood function by matching summary statistics calculated from the observed data with corresponding statistics computed from data simulated from the assumed data generating process. If such statistics are sufficient, the method yields an approximation to the exact posterior distribution of interest that is accurate, given an adequate number of simulations; otherwise, *partial* posterior inference, reflecting the information content of the set of summary statistics only, is the outcome.

The choice of statistics for use within the ABC method, in addition to techniques for determining the matching criterion, are clearly of paramount importance, with much recent research having been devoted to devising ways of ensuring that the information content of the chosen set is maximized, in some sense; e.g. Joyce and Marjoram (2008), Wegmann *et al.* (2009), Blum (2010a) and Fearnhead and Prangle (2012). Recent contributions here include those of Drovandi *et al.* (2011), Drovandi and Pettitt (2013), Gleim and Pigorsch (2013) and Creel and Kristensen (2014), in which the statistics are produced by estimating an approximating *auxiliary* model using both simulated and observed data. This approach mimics, in a Bayesian framework, the principle underlying the frequentist methods of indirect inference (II) (Gouriéroux *et al.* 1993, Smith, 1993, Heggland and Frigessi, 2004) and efficient method of moments (EMM) (Gallant and Tauchen, 1996), using, as it does, the approximating model to produce feasible, but sub-optimal, inference about an intractable true model. Whilst the price paid for the approximation in the frequentist setting is a possible reduction in efficiency, the price paid in the Bayesian case is posterior inference that is conditioned on statistics that are not sufficient for the parameters of the true model, and which amounts to only partial inference as a consequence.

Our paper continues in this spirit, but with particular focus given to the application of auxiliary model-based ABC methods in the state space model (SSM) framework. We begin by demonstrating that reduction to a set of sufficient statistics of fixed dimension relative to the sample size is *infeasible* in finite samples in SSMs. This key observation then motivates our decision to seek asymptotic sufficiency in the state space setting by using the MLE of the parameters of the auxiliary model as the (vector) summary statistic in the ABC matching criterion. We focus on two qualitatively different cases: 1) one in which the auxiliary model *coincides* with the true model, in which case asymptotic sufficiency for the true parameters is achievable via the proposed ABC technique; and 2) the more typical case in which the exact likelihood function is inaccessible, and the auxiliary model represents an approximation only. The first

case mimics that sometimes referenced in the II (or EMM) literature, in which the auxiliary model ‘nests’, or is equivalent to in some well-defined sense, the true model, and full asymptotic efficiency is achieved by the frequentist methods as a consequence. Investigation of this case allows us to document the maximum accuracy gains that are possible via the auxiliary model route, compared with ABC techniques based on alternative summaries, without the confounding effect of the error in the approximating model. The second case gives some insight into what can be achieved in a general non-linear state space setting when the investigator is forced to adopt an inexact approximating model in the implementation of auxiliary model-based ABC. We give emphasis here to non-linear models in which the state (and possibly the observed) is driven by a continuous time model, as this is the canonical example in which simulation from the true model is feasible (at least via an arbitrarily fine discretization), whilst the likelihood function is (typically) unavailable and exact posterior analysis thus not achievable.

We begin by considering the very concept of finite sample sufficiency in the state space context, and the usefulness of applying a typical ABC approach - based on *ad hoc* summary statistics - in this setting. Using the linear Gaussian model for illustration, we demonstrate the lack of reduction to a set of sufficient statistics of fixed dimension, this result providing motivation, as noted above, for the pursuit of asymptotic sufficiency via the auxiliary model method. We then proceed to demonstrate the Bayesian consistency of the auxiliary model approach, subject to the typical quasi-MLE form of conditions being satisfied. We also illustrate that to the order of accuracy that is relevant in establishing the theoretical properties of an ABC technique (i.e. allowing the tolerance used in the matching of the statistics to approach zero), a selection criterion based on the score of the auxiliary model - evaluated at the MLE computed from the observed data - yields equivalent results to a criterion based directly on the MLE itself. This equivalence is shown to hold in both the exactly and over-identified cases, and independently of any (positive definite) weighting matrix used to define the two alternative distance measures, and implies that the proximity to asymptotic sufficiency yielded by the use of the MLE in an ABC algorithm will be replicated by the use of the score. Given the enormous gain in speed achieved by avoiding optimization of the approximate likelihood at each replication of ABC, this is an critical result from a computational perspective. The application of the proposed method in multiple parameter settings is addressed, with separate treatment of scalar (or lower-dimensional blocks of) parameters, via marginal, or integrated likelihood principles advocated, as a possible way of avoiding the inaccuracy that plagues ABC techniques in high dimensions. (See Blum, 2010b and Nott *et al.*, 2014).

The results outlined in the previous paragraph are applicable to an auxiliary model-based ABC method applied in any context (subject to regularity) and, hence, are of interest in their own right. However, our particular interest, as already noted, is in applying the auxiliary model method - and thereby exploiting these properties - in the state space setting. For Case 1) we

choose to illustrate the approach using the linear Gaussian model, whereby the exact likelihood (and hence score) is accessible via the Kalman filter (KF), and asymptotic sufficiency thus achievable. For Case 2) we illustrate the approach via a particular choice of (approximating) auxiliary model for the continuous time SSM. Specifically, the approximating model is formed as a discretization of the true continuous time model, with the augmented unscented Kalman filter (AUKF) (Julier *et al.*, 1995, Julier and Uhlmann, 2004) used to evaluate the likelihood of that model. The general applicability, speed and simplicity of the AUKF calculations render the ABC scheme computationally feasible and relatively simple to implement. This particular approach to the definition of an auxiliary model also leads to a set of summary statistics of relatively small dimension. This is in contrast, for example, with an approach based on a highly parameterized (‘nesting’) approximating model (see, for example, Gleim and Pigorsch, 2013), in which the large number of auxiliary parameters - in principle sufficient for the parameters of the true latent diffusion model - is likely to yield a very inaccurate (non-parametric) estimate of the true posterior, due to the large dimension of the conditioning statistics. The equality between the number of parameters in our exact and approximating models also means that marginalization of the approximating model to produce a scalar matching criterion for each parameter of the true model is meaningful.

The paper proceeds as follows. In Section 2 we briefly summarize the basic principles of ABC as they would apply in a state space setting, including the role played by summary statistics and sufficiency. We demonstrate the lack of finite sample sufficiency reduction in an SSM, using the linear Gaussian model for illustration. In Section 3, we then proceed to demonstrate the properties of ABC based on the MLE, score and marginal score, respectively, of a generic approximating model followed, in Section 4, by an outline of a computationally feasible approximation - based on the AUKF - for use in the non-linear state space setting. Using (repeated samples of) artificially generated data, the accuracy with which the proposed technique reproduces the exact posterior distribution is assessed in Section 5. The ABC methods are based respectively on: i) the joint score; ii) the marginal score; iii) a (weighted) Euclidean metric based on statistics that are sufficient for an observed autoregressive model of order one; and iv) the dimension-reduction technique of Fearnhead and Prangle (2012), applied to this latter set of summary statistics. We conduct the assessment firstly within the context of the linear Gaussian model, with the issues of sufficiency and matching that are key to accurately reproducing the true posterior distribution able to be illustrated precisely in this setting. The overall superiority of both the joint and marginal score techniques over the ABC methods based on summary statistics is demonstrated numerically, as is the remarkable accuracy yielded by the marginal score technique in particular. This exercise thus forms a resounding proof-of-concept for the score-based ABC method, albeit in a case where the exact score is available.

We then proceed to assess performance in a particular non-linear latent diffusion model, in

which the degree of accuracy of the AUKF-based approximating model plays a role. A stochastic volatility for financial returns, in which the latent volatility is driven by a square root diffusion model, is adopted as the non-linear example, as the existence of known (non-central chi-squared) transition densities means that the exact likelihood function/posterior distribution is available, for the purpose of comparison. We apply the deterministic grid-based filtering method of Ng *et al.* (2013) - suitable for this particular setting - to produce the exact comparators for our ABC-based estimates of the relevant marginal posteriors, as well as the marginal posteriors associated with an Euler approximation to the true model. The score methods out-perform the summary statistic methods in the great majority of cases documented. Some gain in accuracy is still produced via the marginalization technique, although that gain is certainly less marked than in the linear Gaussian case, in which the exact score is accessible. Notably, all ABC-based approximations, which exploit simulation from the exact latent diffusion model, serve as more accurate estimates (overall) of the exact posteriors than do the AUKF and Euler approximations themselves. Section 6 concludes.

## 2 ABC in State Space Models

### 2.1 Outline of the basic approach

The aim of ABC is to produce draws from an approximation to the posterior distribution of a vector of unknowns,  $\boldsymbol{\theta}$ , given the  $T$ -dimensional vector of observed data  $\mathbf{y} = (y_1, y_2, \dots, y_T)'$ ,

$$p(\boldsymbol{\theta}|\mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta}),$$

in the case where both the prior,  $p(\boldsymbol{\theta})$ , and the likelihood,  $p(\mathbf{y}|\boldsymbol{\theta})$ , can be simulated. These draws are used, in turn, to approximate posterior quantities of interest, including marginal posterior moments, marginal posterior distributions and predictive distributions. The simplest (accept/reject) form of the algorithm (Tavaré *et al.* 1997, Pritchard, 1999) proceeds as follows:

Step 1 Simulate  $\boldsymbol{\theta}^i$ ,  $i = 1, 2, \dots, R$ , from  $p(\boldsymbol{\theta})$

Step 2 Simulate  $\mathbf{z}^i = (z_1^i, z_2^i, \dots, z_T^i)'$ ,  $i = 1, 2, \dots, R$ , from the likelihood,  $p(\cdot|\boldsymbol{\theta}^i)$

Step 3 Select  $\boldsymbol{\theta}^i$  such that:

$$d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^i)\} \leq \varepsilon, \tag{1}$$

where  $\boldsymbol{\eta}(\cdot)$  is a (vector) statistic,  $d\{\cdot\}$  is a distance criterion of some sort, and the tolerance level  $\varepsilon$  is arbitrarily small. In practice  $\varepsilon$  may be chosen such that, for a given value of  $R$ , a certain (small) proportion of draws of  $\boldsymbol{\theta}^i$  are selected.

The algorithm thus samples  $\boldsymbol{\theta}$  and  $\mathbf{z}$  from the joint posterior:

$$p_\varepsilon(\boldsymbol{\theta}, \mathbf{z}|\mathbf{y}) = \frac{p(\boldsymbol{\theta})p(\mathbf{z}|\boldsymbol{\theta})\mathbb{I}_{A_{\varepsilon,\mathbf{y}}}(\mathbf{z})}{\int p(\boldsymbol{\theta})p(\mathbf{z}|\boldsymbol{\theta})\mathbb{I}_{A_{\varepsilon,\mathbf{y}}}(\mathbf{z})},$$

where  $\mathbb{I}_B$  is the indicator function defined on the set  $B$  and  $A_{\varepsilon,\mathbf{y}} = \{\mathbf{z} \in \mathbf{S}; d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z})\} \leq \varepsilon\}$ . Clearly, when  $\boldsymbol{\eta}(\cdot)$  is sufficient and  $\varepsilon$  arbitrarily small,

$$p_\varepsilon(\boldsymbol{\theta}|\mathbf{y}) = \int p_\varepsilon(\boldsymbol{\theta}, \mathbf{z}|\mathbf{y})d\mathbf{z}$$

approximates the true posterior,  $p(\boldsymbol{\theta}|\mathbf{y})$ , and draws from  $p_\varepsilon(\boldsymbol{\theta}, \mathbf{z}|\mathbf{y})$  can be used to estimate features of the true posterior. In practice however, the complexity of the models to which ABC is applied implies, almost by definition, that sufficiency is unattainable. Hence, in the limit, as  $\varepsilon \rightarrow 0$ , the draws can be used only to approximate features of  $p(\boldsymbol{\theta}|\boldsymbol{\eta}(\mathbf{y}))$ .

Adaptations of the basic rejection scheme have involved post-sampling corrections of the draws using kernel methods (Beaumont *et al.*, 2002, Blum, 2010a, Blum and François, 2010), or the insertion of Markov chain Monte Carlo (MCMC) and/or sequential Monte Carlo (SMC) steps (Marjoram *et al.*, 2003, Sisson *et al.*, 2007, Beaumont *et al.*, 2009, Toni *et al.*, 2009 and Wegmann *et al.*, 2009), to improve the accuracy with which  $p(\boldsymbol{\theta}|\boldsymbol{\eta}(\mathbf{y}))$  is estimated, for any given number of draws. Focus is also given to choosing  $\boldsymbol{\eta}(\cdot)$  and/or  $d\{\cdot\}$  so as to render  $p(\boldsymbol{\theta}|\boldsymbol{\eta}(\mathbf{y}))$  a closer match to  $p(\boldsymbol{\theta}|\mathbf{y})$ , in some sense; see Joyce and Marjoram (2008), Wegmann *et al.*, Blum (2010b) and Fearnhead and Prangle (2012). In the latter vein, Drovandi *et al.* (2011) argue, in the context of a specific biological model, that the use of  $\boldsymbol{\eta}(\cdot)$  comprised of the MLEs of the parameters of a well-chosen approximating model, may yield posterior inference that is conditioned on a large portion of the information in the data and, hence, be close to exact inference based on  $p(\boldsymbol{\theta}|\mathbf{y})$ . (See also Drovandi and Pettitt, 2013, Gleim and Pigorsch, 2013, and Creel and Kristensen, 2014). It is the spirit of this approach that informs the current paper, but with our attention given to rendering the approach feasible in a *general* state space framework that encompasses a large number of the models that are of interest to practitioners, including continuous time models.

Our focus then is on the application of ABC in the context of a general SSM with measurement and transition distributions,

$$p(y_t|x_t, \boldsymbol{\phi}) \tag{2}$$

$$p(x_t|x_{t-1}, \boldsymbol{\phi}) \tag{3}$$

respectively, where  $\boldsymbol{\phi}$  is a  $p$ -dimensional vector of static parameters, elements of which may characterize either the measurement or state relation, or both. For expositional simplicity, and without loss of generality, we consider the case where both  $y_t$  and  $x_t$  are scalars. In financial applications it is common that both the observed and latent processes are driven by continuous

time processes, with the transition distribution in (3) being unknown (or, at least, computationally challenging) as a consequence. Bayesian inference would then typically proceed by invoking (Euler) discretizations for both the measurement and state processes and applying MCMC- or SMC-based techniques, with such methods being tailor-made to suit the features of the particular (discretized) model at hand; see Giordini *et al.* (2011) for a recent review. In some models expressed initially in discrete time, it may also be the case that the conditional distribution in (2) is unavailable in closed form, such as when empirically relevant distributions for financial returns are adopted (e.g. Peters *et al.*, 2012). In such cases, SMC- or MCMC-based inferential methods are typically infeasible. In contrast, in all of these cases the proposed ABC method *is* feasible, as long as simulation from the true model (at least via an arbitrarily fine discretization, in the continuous time case) is possible.

The full set of unknowns thus constitutes the augmented vector  $\boldsymbol{\theta} = (\boldsymbol{\phi}', \mathbf{x}'_c)'$  where, in the case where  $x_t$  evolves in continuous time,  $\mathbf{x}_c$  represents the infinite-dimensional vector comprising the continuum of unobserved states over the sample period. However, to fix ideas, we define  $\boldsymbol{\theta} = (\boldsymbol{\phi}', \mathbf{x}')'$ , where  $\mathbf{x} = (x_1, x_2, \dots, x_T)'$  is the  $T$ -dimensional vector comprising the time  $t$  states for the  $T$  observation periods in the sample.<sup>1</sup> Implementation of the algorithm thus involves simulating from  $p(\boldsymbol{\theta})$  by simulating  $\boldsymbol{\phi}$  from the prior  $p(\boldsymbol{\phi})$ , followed by simulation of  $x_t$  via the process for the state, conditional on the draw of  $\boldsymbol{\phi}$ , and subsequent simulation of artificial data  $z_t$  conditional on the draws of  $\boldsymbol{\phi}$  and the state variable.

Crucially, the focus in this paper is on inference about  $\boldsymbol{\phi}$  only; hence, only draws of  $\boldsymbol{\phi}$  are retained (via the selection criterion) and those draws used to produce an estimate of the marginal posterior,  $p(\boldsymbol{\phi}|\mathbf{y})$ , and with sufficiency to be viewed as relating to  $\boldsymbol{\phi}$  only. Hence, from this point onwards, when we reference a vector of summary statistics,  $\boldsymbol{\eta}(\mathbf{y})$ , it is the information content of that vector with respect to  $\boldsymbol{\phi}$  that is of importance, and the proximity of  $p(\boldsymbol{\phi}|\boldsymbol{\eta}(\mathbf{y}))$  to the marginal posterior of  $\boldsymbol{\phi}$  that is under question. We comment briefly on state inference in Section 6.

Before outlining the proposed methodology for the model in (2) and (3) in Section 3 we highlight the key observation that motivates our approach, namely that reduction to sufficiency in finite samples is not possible in state space settings. We use a linear Gaussian state space model to illustrate this result, as closed-form expressions are available in this case; however, as highlighted at the end of the section, the result is, in principle, applicable to any SSM.

## 2.2 Lack of finite sample sufficiency reduction

Sufficient statistics are useful for inference about a (vector) parameter since, being in possession of the sufficient set means that the data itself may be discarded for inference purposes. When

---

<sup>1</sup>For example, in a continuous time stochastic volatility model such values may be interpreted as end-of-day volatilities.

the cardinality of the sufficient set is small relative to the sample size a significant reduction in complexity is achieved and in the case of ABC, conditioning on the sufficient statistics leads to no loss of information, and the method produces a simulation-based estimate of the true posterior. The difficulty that arises is that only distributions that are members of the exponential family (EF) possess sufficient statistics that achieve a reduction to a fixed dimension relative to the sample size. In the context of the general SSM described by (2) and (3) the effective use of sufficient statistics is problematic. For any  $t$  it is unlikely that the marginal distribution of  $y_t$  will be a member of the EF, due to the vast array of non-linearities that are possible, in either the measurement or state equations, or both. Moreover, even if  $y_t$  were a member of the EF for each  $t$ , to achieve a sufficiency reduction it is required that the *joint* distribution of  $\mathbf{y} = \{y_t; t = 1, 2, \dots, T\}$  also be in the EF. For example, even if  $y_t$  were Gaussian, it does not necessarily follow that the joint distribution of  $\mathbf{y}$  will achieve a sufficiency reduction. The most familiar example of this is when  $\mathbf{y}$  follows a Gaussian moving average (MA) process and consequently only the whole sample is sufficient.

Even the simplest SSMs generate MA-like dependence in the data. Consider the linear Gaussian SSM, expressed in regression form as

$$y_t = x_t + e_t \tag{4}$$

$$x_t = \delta + \rho x_{t-1} + v_t, \tag{5}$$

where the disturbances are respectively independent  $N(0, \sigma_e^2 = 1)$  and  $N(0, \sigma_v^2)$  variables. In this case, the joint distribution of the vector of  $y_t$ 's (which are marginally normal and members of the EF) is  $\mathbf{y} \sim \mathbf{N}[\boldsymbol{\mu}, \sigma_x^2 (r\mathbf{I} + \mathbf{V})]$ , where  $r = \sigma_e^2/\sigma_x^2$  is the inverse of the signal-to-noise (SN) ratio and  $\mathbf{V}$  is the familiar Toeplitz matrix associated with an autoregressive (AR) model of order 1. To construct the sufficient statistics we need to evaluate  $(r\mathbf{I} + \mathbf{V})^{-1}$ , which appears in the quadratic form of the multivariate normal density, with the structure of  $(r\mathbf{I} + \mathbf{V})^{-1}$  determining the way in which sample information about the parameters is accumulated and, hence, the sufficiency reduction that is achievable. (See, for example, Anderson, 1958, Chp 6.) Representing  $(r\mathbf{I} + \mathbf{V})^{-1}$  as

$$(r\mathbf{I} + \mathbf{V})^{-1} = \mathbf{V}^{-1} - r\mathbf{V}^{-2} + r^2\mathbf{V}^{-3} - \dots, \tag{6}$$

it is straightforward to show that as the order of the approximation is increased by retaining more terms, the extent of the accumulation across successive observations is reduced, with the full sample of observations on  $y_t$  ultimately being needed to attain sufficiency. Given that the magnitude of  $r$  determines how many terms in (6) are required for the approximation to be accurate, we see that the SN ratio determines how well the set of summary statistics *that would be sufficient* for an observed AR(1) process (with  $r = 0$ ), namely,

$$s_1 = \sum_{t=2}^{T-1} y_t, \quad s_2 = \sum_{t=2}^{T-1} y_t^2, \quad s_3 = \sum_{t=2}^T y_t y_{t-1}, \quad s_4 = y_1 + y_T, \quad s_5 = y_1^2 + y_T, \quad (7)$$

approximates the information content of the true set of sufficient statistics, that is, the full sample. If the SN ratio is large (i.e.  $r$  is small) then using the set in (7) as summary statistics may produce a reasonable approximation to sufficiency. However, as the SN ratio declines (and higher powers of  $r$  cannot be ignored as a consequence) then this set deviates further and further from sufficiency.

This same qualitative problem would also characterize any SSM nested in (2) and (3), with the only difference being that, in any particular case there would not necessarily be an analytical link between the SN ratio and the lack of sufficiency associated with any finite set of statistics calculated from the observations. The quest for an accurate ABC technique in a state space setting - based on an arbitrary set of statistics - is thus not well-founded and this, in turn, motivates the search for asymptotic sufficiency via the MLE.

### 3 Auxiliary model-based ABC

#### 3.1 Theoretical properties

The asymptotic Gaussianity of the MLE for the parameters of (2) and (3) (under regularity) implies that the MLE satisfies the factorization theorem and is thereby asymptotically sufficient for the parameters of that model. (See Cox and Hinkley, 1974, Chp. 9 for elucidation of this matter.) Denoting the log-likelihood function by  $L(\mathbf{y}; \boldsymbol{\phi})$ , maximizing  $L(\mathbf{y}; \boldsymbol{\phi})$  with respect to  $\boldsymbol{\phi}$  yields  $\hat{\boldsymbol{\phi}}$ , which could, in principle, be used to define  $\boldsymbol{\eta}(\cdot)$  in an ABC algorithm. For large enough  $T$  the algorithm would produce draws from the exact posterior. Indeed, in arguments that mirror those adopted by Gallant and Tauchen (1996) and Gouriéroux *et al.* (1993) for the EMM and II estimators respectively, Gleim and Pigorsch (2013) demonstrate that if  $\boldsymbol{\eta}(\cdot)$  is chosen to be the MLE of an auxiliary model that ‘nests’ the true model in some well-defined way, asymptotic sufficiency will still be achieved; see also Gouriéroux and Monfort (1995) on this point.

Of course, if the SSM in question is such that the exact likelihood is accessible, the model is likely to be tractable enough to preclude the need for treatment via ABC. Further, as we allude to in the Introduction, the quest for asymptotic sufficiency via a (possibly large) nesting auxiliary model conflicts with the quest for an accurate non-parametric estimate of the posterior using the ABC draws; a point that, to our knowledge, has not been noted in the literature. Hence, in practice, the appropriate goal in the ABC context is to define a *parsimonious*, analytically tractable (and computationally efficient) model that *approximates* the (generally intractable) data generating process in (2) and (3) as well as possible, and use that model as the basis

for constructing a summary statistic within an ABC algorithm. If the approximating model is ‘accurate enough’ as a representation of the true model, such an approach will yield, via the ABC algorithm, an estimate of the posterior distribution that is conditioned on a statistic that is ‘close to’ being sufficient, at least for a large enough sample.

Despite the loss of full (asymptotic) sufficiency associated with the use of an approximating model to generate the matching statistics in an ABC algorithm, we show here that Bayesian consistency will still be achieved, subject to certain regularity conditions. Define the choice criterion in (1) as

$$d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^i)\} = \sqrt{\left[\widehat{\boldsymbol{\beta}}(\mathbf{y}) - \widehat{\boldsymbol{\beta}}(\mathbf{z}^i)\right]' \boldsymbol{\Omega} \left[\widehat{\boldsymbol{\beta}}(\mathbf{y}) - \widehat{\boldsymbol{\beta}}(\mathbf{z}^i)\right]} \leq \varepsilon, \quad (8)$$

where  $\widehat{\boldsymbol{\beta}}(\cdot)$  is the MLE of the parameter vector  $\boldsymbol{\beta}$  of the auxiliary model with log-likelihood function,  $L_a(\mathbf{y}; \boldsymbol{\beta})$ , and  $\boldsymbol{\Omega}$  is some positive definite matrix. The quadratic function under the square root essentially mimics the criterion used in the II technique, in which case  $\boldsymbol{\Omega}$  would assume the sandwich form of variance-covariance estimator - appropriate for when the auxiliary model does not coincide with the true model - and optimization with respect to the parameters of the latter is the goal.<sup>2</sup> In Bayesian analyses, in which (8) is used to produce ABC draws,  $\boldsymbol{\Omega}$  may also be defined as the sandwich estimator (Drovandi and Pettit, 2013, and Gleim and Pigorsch, 2013), or simply as the inverse of the (estimated) variance-covariance matrix of  $\widehat{\boldsymbol{\beta}}$ , evaluated at  $\widehat{\boldsymbol{\beta}}(\mathbf{y})$  (Drovandi *et al.*, 2011). However, in common with the frequentist proof of consistency of the II estimator, Bayesian consistency - whereby the posterior for  $\boldsymbol{\phi}$  is degenerate as  $T \rightarrow \infty$  at the true  $\boldsymbol{\phi} = \boldsymbol{\phi}_0$  - is invariant to the choice of the (positive definite)  $\boldsymbol{\Omega}$ . The demonstration follows directly from the same arguments used to prove consistency of the II estimator (see, for e.g. Gouriéroux and Monfort, 1996, Appendix 4A.1), with the following regularity conditions required:

**(A1)**  $\lim_{T \rightarrow \infty} T^{-1} L_a(\mathbf{z}(\boldsymbol{\phi}); \boldsymbol{\beta}) = L_\infty(\boldsymbol{\phi}, \boldsymbol{\beta})$ , uniformly in  $\boldsymbol{\beta}$ , where  $L_\infty(\boldsymbol{\phi}, \boldsymbol{\beta})$  is a deterministic limit function.

**(A2)**  $L_\infty(\boldsymbol{\phi}, \boldsymbol{\beta})$  has a unique maximum with respect to  $\boldsymbol{\beta} : \mathbf{b}(\boldsymbol{\phi}) = \arg \max_{\boldsymbol{\beta}} L_\infty(\boldsymbol{\phi}, \boldsymbol{\beta})$ .

**(A3)** The equation  $\boldsymbol{\beta} = \mathbf{b}(\boldsymbol{\phi})$  admits a unique solution in  $\boldsymbol{\phi}$ , for all  $\boldsymbol{\phi}$ .

Under these conditions it follows that

$$\widehat{\boldsymbol{\beta}}(\mathbf{y}) = \arg \max_{\boldsymbol{\beta}} T^{-1} L_a(\mathbf{y}(\boldsymbol{\phi}_0); \boldsymbol{\beta}) \xrightarrow{a.s.} \arg \max_{\boldsymbol{\beta}} L_\infty(\boldsymbol{\phi}_0, \boldsymbol{\beta}) = \mathbf{b}(\boldsymbol{\phi}_0) \quad (9)$$

and

$$\widehat{\boldsymbol{\beta}}(\mathbf{z}) = \arg \max_{\boldsymbol{\beta}} T^{-1} L_a(\mathbf{z}(\boldsymbol{\phi}); \boldsymbol{\beta}) \xrightarrow{a.s.} \arg \max_{\boldsymbol{\beta}} L_\infty(\boldsymbol{\phi}, \boldsymbol{\beta}) = \mathbf{b}(\boldsymbol{\phi}). \quad (10)$$

---

<sup>2</sup>In practice the implementation of II involves the use of a simulated sample in the computation of  $\widehat{\boldsymbol{\beta}}(\mathbf{z}^i)$  that is a multiple of the size of the empirical sample.

Hence, in the ABC context, in which the generic parameter  $\phi$  represents a draw  $\phi^i$  from the prior,  $p(\phi)$ , we see that as  $T \rightarrow \infty$  the choice criterion in (8) approaches

$$d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^i)\} = \sqrt{[\mathbf{b}(\phi_0) - \mathbf{b}(\phi^i)]' \boldsymbol{\Omega} [\mathbf{b}(\phi_0) - \mathbf{b}(\phi^i)]} \leq \varepsilon.$$

As  $\varepsilon \rightarrow 0$ , being the relevant addition limiting condition required in the ABC setting, we see that irrespective of the form of  $\boldsymbol{\Omega}$ , the only values of  $\phi^i$  that will be selected and, hence, be used to construct an estimate of the posterior distribution, are values such that  $\mathbf{b}(\phi_0) = \mathbf{b}(\phi^i)$ . Given the assumption of the uniqueness of the solution of  $\mathbf{b}(\cdot)$  for  $\phi$  (or  $\phi_0$ ),  $\mathbf{b}(\phi_0) = \mathbf{b}(\phi^i)$  if and only if  $\phi_0 = \phi^i$ . Hence, ABC produces draws that produce a degenerate distribution at the true parameter  $\phi_0$ , as required by the Bayesian consistency property. Once again, this is despite the fact that asymptotic sufficiency will not be achieved in the typical case in which the approximating model is in error, a result that is analogous to the frequentist finding of consistency for the II estimator, without full (Cramer Rao) efficiency obtaining.

We conclude this section by citing related work in hidden Markov models (e.g. Yildirim *et al.*, 2013; Dean *et al.*, 2014), in which ABC principles that avoid summarization have been advocated. Specifically, the difference between the observed data and the simulated pseudo-data is operated time step by time step, as in  $\prod_{t=1}^T \mathbb{I}_{d\{z_t^i, y_t\} \leq \varepsilon}$ . This form of ABC approximation also allows for the derivation of consistency properties (in the number of observations) of the ABC estimates. In particular, using such a distance in the algorithm allows for the approximation to converge to the genuine posterior when the tolerance  $\varepsilon$  goes to zero.<sup>3</sup> One problem with this approach, however, is that the acceptance rate decreases quickly with  $T$ , unless  $\varepsilon$  is increasing with  $T$ . Jasra *et al.* (2014) provide some solutions here, but within an observation-driven model context only. Finally, looking at the application of this form of ABC approximation in a particle MCMC (PMCMC) setting, Jasra *et al.* (2013) and Martin *et al.* (2014) establish convergence (to the exact posterior), in connection with the alive particle filter (Le Gland and Oudjane, 2006).

## 3.2 Score-based implementation

With large computational gains,  $\boldsymbol{\eta}(\cdot)$  in (1) can be defined using the score of the auxiliary model. That is, the score vector associated with the approximating model, when evaluated at the simulated data, and with  $\widehat{\boldsymbol{\beta}}(\mathbf{y})$  substituted for  $\boldsymbol{\beta}$ , will be closer to zero the ‘closer’ is the simulated data to the true. Hence, the choice criterion in (1) for an ABC algorithm can be based on  $\boldsymbol{\eta}(\cdot) = \mathbf{S}(\cdot; \cdot)|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}(\mathbf{y})}$ , where

$$\mathbf{S}(\mathbf{z}^i; \boldsymbol{\beta}) = T^{-1} \frac{\partial L_a(\mathbf{z}^i; \boldsymbol{\beta})}{\partial \boldsymbol{\beta}}, \quad (11)$$

---

<sup>3</sup>This is also the setting in which Fearnhead and Prangle (2012) show that noisy ABC (Wilkinson, 2013) is well-calibrated, i.e. converges to the relevant posterior distribution (determined by the choice of summary statistics in this case). See also, Dean and Singh (2011) and Dean *et al.* (2014).

yielding

$$d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^i)\} = \sqrt{\left[\mathbf{S}(\mathbf{z}^i; \widehat{\boldsymbol{\beta}}(\mathbf{y}))\right]' \boldsymbol{\Sigma} \left[\mathbf{S}(\mathbf{z}^i; \widehat{\boldsymbol{\beta}}(\mathbf{y}))\right]} \leq \varepsilon, \quad (12)$$

where  $\boldsymbol{\Sigma}$  denotes an arbitrary positive definite weighting matrix. Implementation of ABC via (12) is faster (by many orders of magnitude) than the approach based upon  $\boldsymbol{\eta}(\cdot) = \widehat{\boldsymbol{\beta}}(\cdot)$ , due to the fact that maximization of the approximating model is required only once, in order to produce  $\widehat{\boldsymbol{\beta}}(\cdot)$  from the observed data  $\mathbf{y}$ . All other calculations involve simply the *evaluation* of  $\mathbf{S}(\cdot; \cdot)$  at the simulated data, with a numerical differentiation technique invoked to specify  $\mathbf{S}(\cdot; \cdot)$ .

Once again in line with the proof of the consistency of the relevant frequentist (EMM) estimator, the Bayesian consistency result in Section 3.1 could be re-written in terms  $\boldsymbol{\eta}(\cdot) = \mathbf{S}(\cdot; \cdot)|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}(\mathbf{y})}$ , upon the addition of a differentiability condition regarding  $L_a(\mathbf{z}^i; \boldsymbol{\beta})$  and the assumption that  $\boldsymbol{\beta} = b(\boldsymbol{\phi})$  is the unique solution to the limiting first-order condition,  $\partial L_\infty(\boldsymbol{\phi}, \boldsymbol{\beta})/\partial \boldsymbol{\beta} = \lim_T T^{-1} \partial L_a(\mathbf{z}(\boldsymbol{\phi}); \boldsymbol{\beta})/\partial \boldsymbol{\beta}$  and the convergence is uniform in  $\boldsymbol{\beta}$ . In brief, given that  $\widehat{\boldsymbol{\beta}}(\mathbf{y}) \xrightarrow{a.s.} \mathbf{b}(\boldsymbol{\phi}_0)$ , as  $T \rightarrow \infty$  the choice criterion in (12) approaches

$$d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^i)\} = \sqrt{\left[\partial L_\infty(\boldsymbol{\phi}^i; \mathbf{b}(\boldsymbol{\phi}_0))/\partial \boldsymbol{\beta}\right]' \boldsymbol{\Sigma} \left[\partial L_\infty(\boldsymbol{\phi}^i; \mathbf{b}(\boldsymbol{\phi}_0))/\partial \boldsymbol{\beta}\right]} \leq \varepsilon.$$

As  $\varepsilon \rightarrow 0$ , irrespective of the form of  $\boldsymbol{\Sigma}$ , the only values of  $\boldsymbol{\phi}^i$  that will be selected via ABC are values such that  $\mathbf{b}(\boldsymbol{\phi}_0) = \mathbf{b}(\boldsymbol{\phi}^i)$ , which, given Assumption (A3), implies  $\boldsymbol{\phi}_0 = \boldsymbol{\phi}^i$ .

Hence, Bayesian consistency is maintained through the use of the score. However, a remaining pertinent question concerns the impact on sufficiency (or, more precisely, on *the proximity to asymptotic sufficiency*) of the use of the score instead of the MLE. In practical terms this question can be re-phrased as: does the selection criterion based on  $\mathbf{S}(\cdot; \cdot)$  yield identical draws of  $\boldsymbol{\phi}$  to those yielded by the selection criterion based on  $\widehat{\boldsymbol{\beta}}$ ? If the answer is yes then, unambiguously, for large enough  $T$  and for  $\varepsilon \rightarrow 0$ , the score- and MLE-based ABC criteria will yield equivalent estimates of the exact posterior, with the accuracy of those (equivalent) estimates dependent, of course, on the nature of the auxiliary model itself.

For any auxiliary model (satisfying identification and regularity conditions) with unknown parameter vector  $\boldsymbol{\beta}$ , we expand the score function in (11), evaluated at  $\widehat{\boldsymbol{\beta}}(\mathbf{y})$ , around the point  $\widehat{\boldsymbol{\beta}}(\mathbf{z}^i)$  (with scaling via  $T^{-1}$  having been introduced at the outset in the definition of the score in (11)),

$$\mathbf{S}(\mathbf{z}^i; \widehat{\boldsymbol{\beta}}(\mathbf{y})) = \mathbf{S}(\mathbf{z}^i; \widehat{\boldsymbol{\beta}}(\mathbf{z}^i)) + \mathbf{D} \left[ \widehat{\boldsymbol{\beta}}(\mathbf{y}) - \widehat{\boldsymbol{\beta}}(\mathbf{z}^i) \right] = \mathbf{D} \left[ \widehat{\boldsymbol{\beta}}(\mathbf{y}) - \widehat{\boldsymbol{\beta}}(\mathbf{z}^i) \right],$$

where

$$\mathbf{D} = T^{-1} \frac{\partial^2 L_a(\mathbf{z}^i; \widetilde{\boldsymbol{\beta}}(\mathbf{z}))}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} \quad (13)$$

and  $\widetilde{\boldsymbol{\beta}}(\mathbf{z}^i)$  denotes an (unknown) intermediate value between  $\widehat{\boldsymbol{\beta}}(\mathbf{y})$  and  $\widehat{\boldsymbol{\beta}}(\mathbf{z}^i)$ . Hence, the (scaled)

criterion in (12) becomes

$$\begin{aligned} & \sqrt{\left[\mathbf{S}(\mathbf{z}^i; \hat{\boldsymbol{\beta}}(\mathbf{y}))\right]' \boldsymbol{\Sigma} \left[\mathbf{S}(\mathbf{z}^i; \hat{\boldsymbol{\beta}}(\mathbf{y}))\right]} \\ &= \sqrt{\left[\hat{\boldsymbol{\beta}}(\mathbf{y}) - \hat{\boldsymbol{\beta}}(\mathbf{z}^i)\right]' \mathbf{D}' \boldsymbol{\Sigma} \mathbf{D} \left[\hat{\boldsymbol{\beta}}(\mathbf{y}) - \hat{\boldsymbol{\beta}}(\mathbf{z}^i)\right]'} \leq \varepsilon. \end{aligned} \quad (14)$$

Subject to standard conditions regarding the second derivatives of the auxiliary model, the matrix  $\mathbf{D}$  in (13) will be of full rank and as  $T \rightarrow \infty$ ,  $\mathbf{D}' \boldsymbol{\Sigma} \mathbf{D} \rightarrow$  some positive definite matrix (given the positive definiteness of  $\boldsymbol{\Sigma}$ ) that is some function of  $\boldsymbol{\phi}^i$ . Hence, whilst for any  $\varepsilon > 0$ , the presence of  $\mathbf{D}$  affects selection, as it is a function of the drawn value  $\boldsymbol{\phi}^i$  (through  $\mathbf{z}^i$ ), as  $\varepsilon \rightarrow 0$ ,  $\boldsymbol{\phi}^i$  will be selected via (14) if and only if  $\hat{\boldsymbol{\beta}}(\mathbf{y})$  and  $\hat{\boldsymbol{\beta}}(\mathbf{z}^i)$  are equal. Similarly, irrespective of the form of the (positive definite) weighting matrix in (8), the MLE criterion will produce these same selections. This result pertains no matter what the dimension of  $\boldsymbol{\beta}$  relative to  $\boldsymbol{\phi}$ , i.e. no matter whether the true parameters are exactly or over-identified by the parameters of the auxiliary model. This result thus goes beyond the comparable result regarding the II/EMM estimators (see, for e.g. Gouriéroux and Monfort, 1996), in that the equivalence is independent of the form of weighting matrix used *and* the form of identification that prevails.

Of course in practice, ABC is implemented with  $\varepsilon > 0$ , at which point the two ABC criteria will produce different draws. However, for the models entertained in this paper, preliminary investigation has assured us that the difference between the ABC estimates of the posteriors yielded by the alternative criteria is negligible for small enough  $\varepsilon$ . Hence, we proceed to operate solely with the score-based approach as the computationally feasible method of extracting approximate asymptotic sufficiency in the state space setting.

The actual selection and optimization of the tolerance level,  $\varepsilon$ , has been the subject of intense scrutiny in the recent years (see, for example, Marin *et al.*, 2011 for a detailed survey). What appears to be the most fruitful path to the calibration of the tolerance is to firmly set it within the realm of non-parametric statistics (Blum and François, 2010) as this provides proper convergence rates for the tolerance (rates that differ between standard and noisy ABC; see Fearnhead and Prangle, 2012) and shows that the optimal value stays away from zero for a given sample size. In addition, the practical constraints imposed by finite computing time and the necessity to produce an ABC sample of reasonable length lead us to follow the recommendations found in Biau *et al.* (2012), namely to analyze the ABC approximation as a k-nearest neighbour technique and to exploit this perspective to derive a practical value for the tolerance.

### 3.3 Dimension reduction via marginal likelihood techniques

An ABC algorithm induces two forms of approximation error. Firstly, and most fundamentally, the use of a vector of summary statistics  $\boldsymbol{\eta}(\mathbf{y})$  to define the selection criterion in (1) means that a simulation-based estimate of the posterior of interest is the outcome of the exercise. Only

if  $\boldsymbol{\eta}(\mathbf{y})$  is sufficient for  $\boldsymbol{\phi}$  is  $p(\boldsymbol{\phi}|\boldsymbol{\eta}(\mathbf{y}))$  equivalent to the exact posterior  $p(\boldsymbol{\phi}|\mathbf{y})$ ; otherwise the exact posterior is necessarily estimated with error because of the analytical difference between the exact density and the *partial* posterior density  $p(\boldsymbol{\phi}|\boldsymbol{\eta}(\mathbf{y}))$ . Secondly, the partial posterior density itself,  $p(\boldsymbol{\phi}|\boldsymbol{\eta}(\mathbf{y}))$ , will be estimated with simulation error. Critically, as highlighted by Blum (2010b), the accuracy of the simulation-based estimate of  $p(\boldsymbol{\phi}|\boldsymbol{\eta}(\mathbf{y}))$  will be less, all other things given, the larger the dimension of  $\boldsymbol{\eta}(\mathbf{y})$ . This ‘curse of dimensionality’ obtains even when the parameter  $\boldsymbol{\phi}$  is a scalar, and relates solely to the dimension of  $\boldsymbol{\eta}(\mathbf{y})$ . As elaborated on further by Nott *et al.* (2014), this problem is exacerbated as the dimension of  $\boldsymbol{\phi}$  itself increases, firstly because an increase in the dimension of  $\boldsymbol{\phi}$  brings with it a concurrent need for an increase in the dimension of  $\boldsymbol{\eta}(\mathbf{y})$  and, secondly, because the need to estimate a multi-dimensional density (for  $\boldsymbol{\phi}$ ) brings with it its own problems related to dimension.

As a potential solution to the inaccuracy induced by the dimensionality of the problem, Nott *et al.* (2014) suggest allocating (via certain criteria) a subset of the full set of summary statistics to each element of  $\boldsymbol{\phi}$ ,  $\phi_j$ ,  $j = 1, 2, \dots, p$ , using kernel density techniques to estimate each marginal density,  $p(\phi_j|\boldsymbol{\eta}_j(\mathbf{y}))$ , and then using standard techniques to retrieve a more accurate estimate of the joint posterior,  $p(\boldsymbol{\phi}|\boldsymbol{\eta}(\mathbf{y}))$ , if required. However, the remaining problem associated with the (possibly still high) dimension of each  $\boldsymbol{\eta}_j(\mathbf{y})$ , in addition to the very problem of defining an appropriate set  $\boldsymbol{\eta}_j(\mathbf{y})$  for each  $\phi_j$ , remains unresolved. See Blum *et al.* (2013) for further elaboration on the dimensionality issue in ABC and a review of current approaches for dealing with the problem.

The principle advocated in this paper is to exploit the information content in the MLE of the parameters of an auxiliary model,  $\boldsymbol{\beta}$ , to yield ‘approximate’ asymptotic sufficiency for  $\boldsymbol{\phi}$ . Within this framework, the dimension of  $\boldsymbol{\beta}$  determines the dimension of  $\boldsymbol{\eta}(\mathbf{y})$  and the curse of dimensionality thus prevails for high-dimensional  $\boldsymbol{\beta}$ . However, in this case a solution is available, at least when the dimensions of  $\boldsymbol{\beta}$  and  $\boldsymbol{\phi}$  are equivalent and there is a one-to-one match between the elements of the two parameter vectors. This is clearly so for the two cases tackled in this paper. In the linear Gaussian model investigated as Case 1) the auxiliary model coincides exactly with the true model, in which case  $\boldsymbol{\beta} = \boldsymbol{\phi}$ . In the stochastic volatility SSM investigated as Case 2), we produce an auxiliary model by discretizing the latent diffusion (and evaluating the resultant likelihood via the AUKF); hence,  $\dim(\boldsymbol{\beta}) = \dim(\boldsymbol{\phi})$  and there is a natural one-to-one mapping between the parameters of the (true) continuous time and (approximating) discretized models. In both of these examples then, marginalizing the auxiliary likelihood function with respect to all parameters other than  $\beta_j$  and then producing the score of this function with respect to  $\beta_j$  (as evaluated at the marginal MLE from the observed data,  $\widehat{\beta}_j(\mathbf{y})$ ), yields, by construction, an obvious *scalar* statistic for use in selecting draws of  $\phi_j$  and, hence, a method for estimating  $p(\phi_j|\mathbf{y})$ . If the marginal posteriors only are of interest, then all  $p$  marginals can be estimated in this way, with  $p$  applications of  $(p - 1)$ -dimensional integration required at

each step within ABC to produce the relevant score statistics. Importantly, we do not claim here that the ‘proximity’ to sufficiency (for  $\phi$ ) of the vector statistic  $\eta(\mathbf{y})$ , translates into an equivalent relationship between the score of the marginalized (auxiliary) likelihood function and the corresponding scalar parameter, nor that the associated product density is coherent with a joint probability distribution. If the joint posterior (of the full vector  $\phi$ ) is of particular interest, the sort of techniques advocated by Nott et al. (2014), amongst others, can be used to yield joint inference from the estimated marginals.

In Section 5 we explore the benefits of marginalization, in addition to the increase in accuracy yielded by using a score-based ABC method (either joint or marginal) rather than an ABC algorithm based on a more *ad hoc* choice of summary statistics. However, prior to that we provide details in the following section of the form of auxiliary model advocated for the non-linear SSM case.

## 4 The AUKF approximation for the general non-linear SSM

When the SSM defined in (2) and (3) is analytically intractable, an approximating model is needed to drive the score-based ABC technique. In the canonical non-linear example being emphasized in the paper, in which either  $y_t$  or  $x_t$  (or both) is driven by a continuous time process, this approximation begins with the specification of a discretized version of (2) and (3), expressed generically using a regression formulation as:

$$y_t = h_t(x_t, e_t, \phi) \tag{15}$$

$$x_t = k_t(x_{t-1}, v_t, \phi), \tag{16}$$

for  $t = 1, 2, \dots, T$ , where the  $\{e_t\}$  and  $\{v_t\}$  are assumed to be sequences of *i.i.d.* random variables. This formulation is general enough to include, for example, independent random jump components in either the measurement or state equations (subsumed under  $e_t$  and  $v_t$  respectively), but does exclude cases where the nature of the model is such that a regression formulation in discrete time is not feasible. We note here that in producing (15) and (16) either the observation or the state variable, or both, may need to be transformed; however, for notational simplicity we continue to use the same symbols,  $x_t$  and  $y_t$ , as are used to denote the variables in the true model in (2) and (3). The nature of the discretization affects the functional form of  $h_t(\cdot)$  and  $k_t(\cdot)$ . In Section 5.2 we illustrate the method using a model in which the true model for  $y_t$  is already expressed in discrete time (i.e. there is no discretization error via the measurement process) and in which a (first-order) Euler process is initially used to approximate a square root diffusion model for the state.

The log-likelihood function associated with the approximate model in (15) and (16) is defined by

$$L_a(\mathbf{y}; \boldsymbol{\phi}) = \ln p(y_1) + \sum_{t=1}^{T-1} \ln p(y_{t+1} | \mathbf{y}_{1:t}). \quad (17)$$

where  $\mathbf{y}_{1:t} = (y_1, y_2, \dots, y_t)'$ . Only if the approximate model is linear and Gaussian or the state variable is discrete on a finite support, are the components used to define (17) available in closed form. Given the nature of the problem we are tackling here, namely one in which (15) and (16) are produced as discrete approximations to a continuous time model (or, indeed, one in which (15) and (16) represent an initial discrete-time formulation that is non-linear and/or non-Gaussian), it is reasonable to assume that (17) cannot be evaluated exactly. Whilst several methods (see, e.g. Simon, 2006) are available to approximate  $L_a(\mathbf{y}; \boldsymbol{\phi})$  including, indeed, simulation-based methods such as particle filtering and SMC, the fact that this computation is to be *embedded* within the ABC algorithm makes it essential that the technique is both fast and numerically stable. The AUKF satisfies these criteria and, hence, is our method of choice for this illustration.

In brief, the unscented Kalman filter (UKF) is based on the theory of unscented transformations, which is a method for calculating the moments of a non-linear transformation of a random variable. It involves selecting a set of points on the support of the random variable, called *sigma points*, according to a predetermined and deterministic criterion. These sigma points yield, in turn, a ‘cloud’ of transformed points through the non-linear function, which are then used to produce approximations to the moments of the state variable  $x_t$  used in implementing the KF, via simple weighted sums of the transformed sigma points. (See Haykin, 2001, Chapter 7, for an excellent introduction.) The *augmented* version of the UKF - the AUKF - generalizes the filter to the case where the state and measurement errors are non-additive, applying the principles of unscented transformations to the *augmented* state vector  $s_t = (x_t, \nu_t, e_t)'$ . The computational burden of the filter is thus minimal - comprising the calculation of updated sigma points for the time varying state  $x_t$  at each  $t$ , the computation of the relevant means and variances (for both  $x_t$  and  $y_t$ ) using simple weighted sums, and the use of the usual KF up-dating equations. Details of the specification of the sigma points, plus the steps involved in estimating (17) are provided in the Appendix.

In implementing the score-based ABC method in a non-linear continuous time setting, there are two aspects of the approximation used therein to consider: 1) the accuracy of the discretization; and 2) the accuracy with which the likelihood function of the approximate (discretized) model is evaluated via the AUKF and, hence, the accuracy of the resultant estimate of the MLE. In addressing the first aspect, one has access to the existing literature in which various discretized versions of continuous time models are derived, to different orders of accuracy. With regard to the accuracy of the AUKF evaluation of the likelihood function, we make reference to relevant results (see, e.g. Haykin, 2001) that document the higher-order accuracy (relative,

say, to the extended KF) of the AUKF-based estimates of the mean and variance of the filtered and predictive (for both state and observed) distributions. We also advocate here using any transformation of the measurement (or state) equation that renders the Gaussian approximations invoked by the AUKF more likely to be accurate. However, beyond that, the accuracy of the Gaussian assumption embedded in the AUKF-based likelihood estimate is case-specific.<sup>4</sup>

## 5 Numerical assessment of alternative ABC methods

We now undertake a numerical exercise in which the accuracy of the score-based methods of ABC (both joint and marginal) is compared with that of ABC methods based on a set of summary statistics that are chosen without reference to a model. We conduct this exercise firstly within the context of the linear Gaussian model, in which case the exact score is accessible via the KF. The results here thus provide evidence on two points of interest: 1) whether or not accuracy can be increased by accessing the asymptotic sufficiency of the (exact) MLE in an ABC treatment of a state space model (compared to the use of other summary statistics); and 2) whether or not the curse of dimensionality can be obviated via marginalization, or integration, of the exact likelihood.

In Section 5.2 we then assess accuracy in the typical setting in which an approximating model is used to generate the score. The square root volatility model is used as the example, with the approximate score produced by evaluating the likelihood function of a discretized version of that model using the AUKF. Whilst recognizing that the results relate to one particular model and approximation thereof, they do, nevertheless, serve to illustrate that score-based methods can dominate the summary statistic-based techniques (overall), even when the approximating model used is not particularly accurate.

### 5.1 Case 1: Linear Gaussian model

#### 5.1.1 Data generation and computational details

In this section we simulate a sample of size  $T = 400$  from the linear Gaussian (LG) model in (4) and (5), based on the parameter settings:  $\rho = 0.7$ ,  $\delta = 0.1$  and  $\sigma_v^2 = 1$ , with the three-dimensional parameter  $\phi = (\rho, \delta, \sigma_v)'$  to be estimated. The value of the measurement error variance,  $\sigma_e^2$ , is set in order to fix the SN ratio. We compare the performance of the (joint and marginal) score-based techniques with that of more conventional ABC methods based on summary statistics that may be deemed to be a sensible choice in this setting. Given the relationship between the LG model

---

<sup>4</sup>One potential benefit of the AUKF-based approach is that higher-order discretizations, which produce non-linearity in the relevant error terms could be adopted (e.g. Milstein, 1998), as the relevant sigma points that underpin the AUKF method need only be substituted into this additional non-linear function. We do not explore this option in the current paper, basing the numerical demonstration in Section 5.2 on a first-order Euler discretization only.

and an observable AR(1) process, it seems sensible to propose a set of summary statistics that are sufficient for the latter, as given in (7). Two forms of distances are used. Firstly, we apply the conventional Euclidean distance, with each summary statistic also weighted by the inverse of the variance of the values of the statistic across the ABC draws. That is, we define

$$d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^i)\} = \left[ \sum_{j=1}^5 (s_j^i - s_j^{obs})^2 / \text{var}(s_j) \right]^{1/2} \quad (18)$$

for ABC iteration  $i = 1, 2, \dots, R$ , where  $\text{var}(s_j)$  is the variance (across  $i$ ) of the  $s_j^i$ , and  $s_j^{obs}$  is the observed value of the  $j$  statistic. Secondly, we use a distance measure proposed in Fearnhead and Prangle (2012) which, as made explicit in Blum *et al.* (2013), is a form of dimension reduction method. We explain this briefly as follows. Given the vector of observations  $\mathbf{y}$ , the set of summary statistics in (7) are used to produce an estimate of  $E(\phi_j|\mathbf{y})$ ,  $j = 1, 2, 3$ , which, in turn, is used as the summary statistic in a subsequent ABC algorithm. The steps of the procedure (as modified for this context) described for selection of the scalar parameter  $\phi_j$ ,  $j = 1, 2, 3$ , are as follows:

1. Simulate  $\phi_j^i$ ,  $i = 1, 2, \dots, R$ , from  $p(\phi_j)$  and, subsequently, simulate  $\mathbf{x}^i = (x_1^i, x_2^i, \dots, x_T^i)'$  from (24) using the exact transitions, and pseudo data,  $\mathbf{z}^i$  using the conditional Gaussian form of  $p(\mathbf{z}|\mathbf{x})$ .
2. For  $\mathbf{z}^i$ ,  $i = 1, 2, \dots, R$ , calculate

$$\mathbf{s}^i = [s_1^i, s_2^i, s_3^i, s_4^i, s_5^i]' \quad (19)$$

3. Define  $\boldsymbol{\phi}_j = (\phi_j^1, \phi_j^2, \dots, \phi_j^R)'$ ,  $\mathbf{X} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ \mathbf{s}^1 & \mathbf{s}^2 & \dots & \mathbf{s}^R \end{bmatrix}'$  and

$$\boldsymbol{\phi}_j = E[\boldsymbol{\phi}_j|\mathbf{Z}] + \mathbf{e} = \mathbf{X} \begin{bmatrix} \alpha \\ \boldsymbol{\beta} \end{bmatrix} + \mathbf{e},$$

where  $\mathbf{Z} = [\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^R]$  and  $\boldsymbol{\beta}$  is of dimension  $(5 \times 1)$ .

4. Use OLS to estimate  $E[\boldsymbol{\phi}_j|\mathbf{Z}]$  as  $\widehat{E}[\boldsymbol{\phi}_j|\mathbf{Z}] = \widehat{\alpha} + [\mathbf{s}^1 \quad \mathbf{s}^2 \quad \dots \quad \mathbf{s}^R]' \widehat{\boldsymbol{\beta}}$
5. Define:

$$\boldsymbol{\eta}(\mathbf{z}^i) = \widehat{E}(\boldsymbol{\phi}_j|\mathbf{z}^i) = \widehat{\alpha} + \mathbf{s}^{i'} \widehat{\boldsymbol{\beta}} \text{ and } \boldsymbol{\eta}(\mathbf{y}) = \widehat{E}(\boldsymbol{\phi}_j|\mathbf{y}) = \widehat{\alpha} + \mathbf{s}^{obs'} \widehat{\boldsymbol{\beta}},$$

where  $\mathbf{s}^{obs}$  denotes the vector of summary statistics in (19) calculated from the vector of observed returns, and use:

$$d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^i)\} = \left| \widehat{E}(\boldsymbol{\phi}_j|\mathbf{y}) - \widehat{E}(\boldsymbol{\phi}_j|\mathbf{z}^i) \right| = \left| \mathbf{s}^{i'} \widehat{\boldsymbol{\beta}} - \mathbf{s}^{obs'} \widehat{\boldsymbol{\beta}} \right| \quad (20)$$

as the selection criterion for  $\phi_j$  at each iteration  $i$ .

The joint score-based method uses the distance measure in (12), but with the score in this case computed from the *exact* model, evaluated using the KF. The weighting matrix  $\Sigma$  is set equal to the Hessian-based estimate of the covariance matrix of the (joint) MLE estimator of  $\beta = \phi$ , evaluated at the MLE computed from the observed data,  $\hat{\phi}(\mathbf{y})$ . The marginal score-based method for estimating the marginal posterior for the  $j$ th element of  $\phi$ ,  $\phi_j$ ,  $j = 1, 2, 3$ , is based on the distance

$$d\{\boldsymbol{\eta}(\mathbf{y}), \boldsymbol{\eta}(\mathbf{z}^i)\} = |S(\mathbf{y}; \phi_j) - S(\mathbf{z}^i; \phi_j)|, \quad (21)$$

where

$$S(\mathbf{z}^i; \phi_j) = T^{-1} \frac{\partial L(\mathbf{z}^i; \phi_j)}{\partial \phi_j} \quad (22)$$

and  $L(\mathbf{z}^i; \phi_j)$  is produced by integrating (numerically) the exact likelihood function (evaluated via the KF) with respect to all parameters other than  $\phi_j$ , and taking the logarithm.

### 5.1.2 Numerical results for the LG model

We produce results that compare the performance of the four different methods: the joint score-based ABC (‘ABC-joint score’ in all figures); the marginal score-based ABC (‘ABC-marg score’); the summary statistic-based ABC using the Euclidean metric in (18) (‘ABC-sum stats’); and the approach of Fearnhead and Prangle (2012) based on the metric in (20) (‘ABC-FP’). Marginal density estimates are produced initially for a single run of ABC, based on 50,000 replications of the accept/reject algorithm detailed in Section 2.1, and with  $\varepsilon$  defined as the 5th percentile of the 50,000 draws. The true data is generated from a process in which the SN ratio is high (i.e.  $[\sigma_v^2/(1 - \rho^2)]/\sigma_e^2 = 20$ ), with the true marginal posteriors computed by normalizing the likelihood function evaluated using the KF (and multiplied by a uniform prior), then marginalizing using deterministic integration over a very fine grid for  $\phi$ . The score-based methods also use the exact likelihood function to compute the score. All three sets of marginal posteriors, for  $\rho$ ,  $\delta$  and  $\sigma_v$ , are produced in Figure 1, Panels A, B and C respectively. We then summarize the results for 100 replications of ABC, using box plots, in Figures 2 to 4. We produce estimates of the 5th, 25th, 50th, 75th and 95th percentiles of each true (KF-based) posterior density, with the exact percentile represented by the horizontal dotted line. The short-hand notation used to denote the form of ABC method corresponds to that used in Figure 1.

Figure 1, Panels A to C, highlight graphically - for a single ABC run - the performance of the four ABC methods in reproducing the true marginal posteriors. The two most notable features of all three plots are: i) the remarkable accuracy of the marginal score approach; and ii) the very poor performance of the summary statistic approach based on the Euclidean metric. These features are replicated across the 100 runs of ABC, as evidenced by Figures 2 to 4. For all three parameters, the marginal score approach produces percentile estimates that are extremely accurate in terms of both mean location and spread. For the parameter  $\rho$  the joint score approach

Figure 1: Marginal posterior densities for the three parameters in the linear Gaussian state space model in (4) and (5). As per the key, the graphs reproduced are the exact posterior, in addition to the four ABC-based estimates, as detailed in the text.

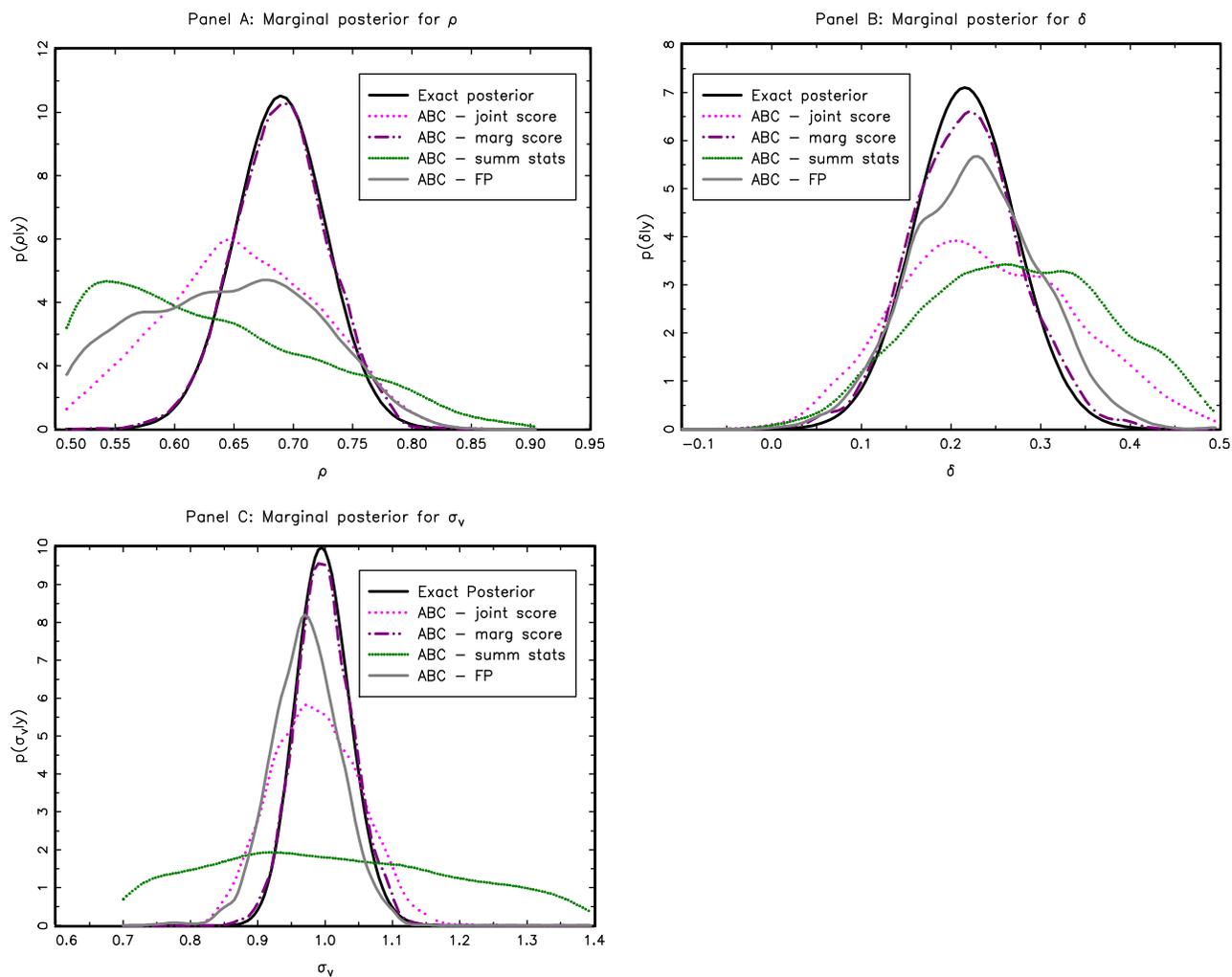
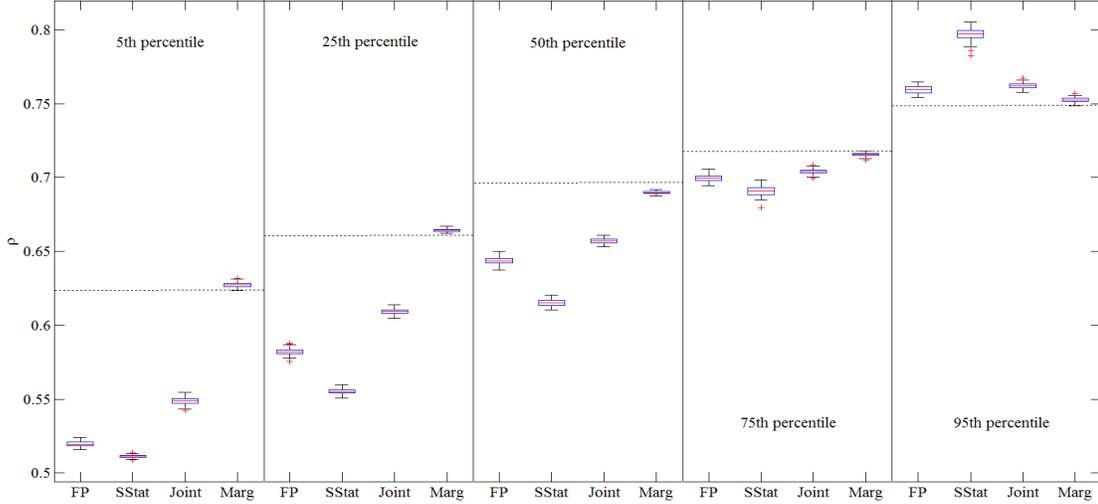


Figure 2: Box plots for 100 replications of ABC, with each replication based on 50,000 draws. Estimation of the marginal posterior density of the state intercept parameter  $\rho$  in the linear Gaussian model in (4) and (5).  $T = 400$ . High SN ratio. The percentiles for the exact marginal posterior are represented by the horizontal dotted lines; and the four ABC methods referenced correspond to those described in the text.



is the next most accurate technique, followed by the FP method. For the parameter  $\sigma_v$  there is little to choose between the latter two methods, both of which produce quite accurate estimates of the true percentiles, although both being still less accurate than the marginal score technique. For  $\delta$  the FP method tends to outperform the joint score approach, but with neither method competing with the accuracy yielded by the marginal score. For all three parameters, the summary statistic approach based on the Euclidean metric produces the poorest results overall, despite the high SN ratio.

Further results (not documented here, for reasons of space) demonstrate that decreasing the SN ratio has no qualitative effective on the performance of the score-based results, including the marked accuracy of the marginal score method. This robustness of the score methods to the SN ratio is to be anticipated, give that the likelihood function for the full state space model has been used to generate the matching statistics. The impact of the change in the SN ratio on the summary statistic-based methods is not uniform, with there certainly being no clear tendency for the results the worsen. This suggests that the accuracy with which  $p(\phi_j|\boldsymbol{\eta}(\mathbf{y}))$  itself is estimated (and, hence, the issue of dimensionality), rather than the relationship between  $p(\phi_j|\boldsymbol{\eta}(\mathbf{y}))$  and  $p(\phi_j|\mathbf{y})$  (and, hence the ‘closeness’ of  $\boldsymbol{\eta}(\mathbf{y})$  to sufficiency), remains a dominant influence on final accuracy, no matter what the SN value.

Figure 3: Box plots for 100 replications of ABC, with each replication based on 50,000 draws. Estimation of the marginal posterior density of the state intercept parameter  $\delta$  in the linear Gaussian model in (4) and (5).  $T = 400$ . High SN ratio. The percentiles for the exact marginal posterior are represented by the horizontal dotted lines; and the four ABC methods referenced correspond to those described in the text.

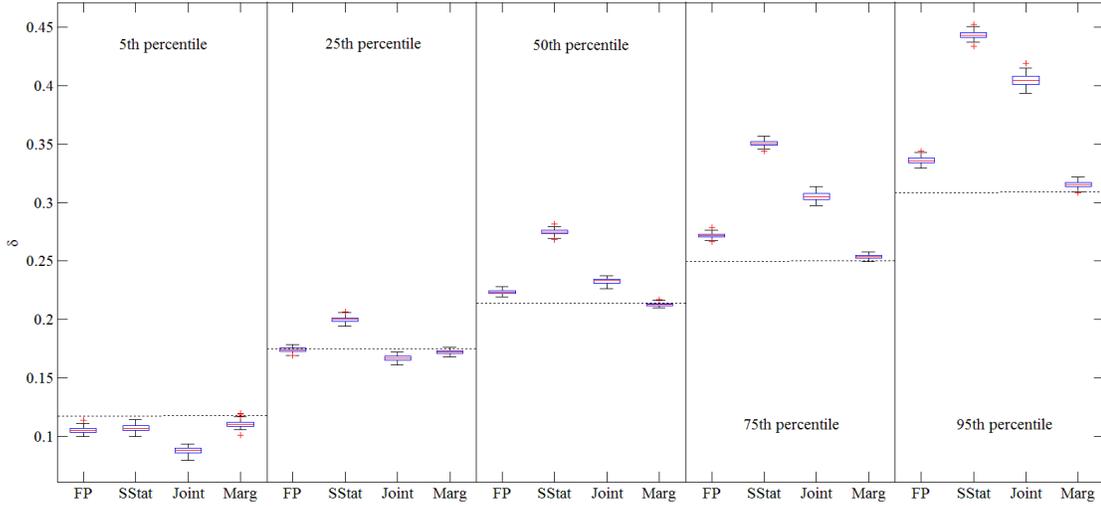
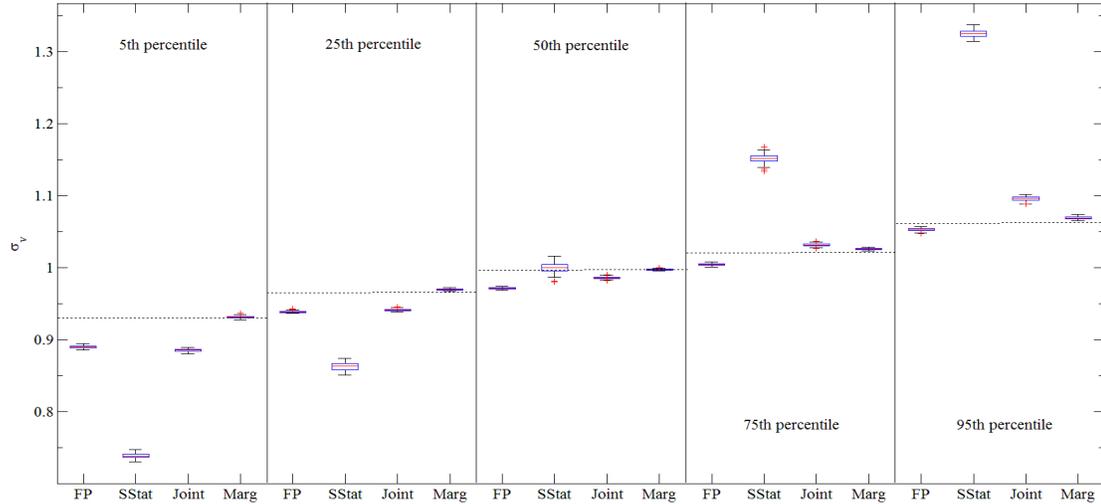


Figure 4: Box plots for 100 replications of ABC, with each replication based on 50,000 draws. Estimation of the marginal posterior density of the state intercept parameter  $\sigma_v$  in the linear Gaussian model in (4) and (5).  $T = 400$ . High SN ratio. The percentiles for the exact marginal posterior are represented by the horizontal dotted lines; and the four ABC methods referenced correspond to those described in the text.



## 5.2 Case 2: Heston stochastic volatility (SV) model

### 5.2.1 Data generation and computational details

Given the critical role played by volatility in asset pricing, portfolio management and the calculation of risk measures, a large segment of the empirical finance literature has been devoted to the construction and analysis of volatility models. Three decades of empirical studies have demonstrated that the *constant* volatility feature of a geometric Brownian motion process for an asset price is inconsistent with both the observed time-variation in return volatility and the non-Gaussian characteristics of empirical distributions of returns; see Bollerslev, Chou and Kroner (1992) for a review. Empirical regularities documented in the option pricing literature, most notably implied volatility ‘smiles’, are also viewed as evidence that asset prices deviate from the geometric Brownian motion assumption underlying the Black and Scholes (1973) option price; see, for example, Bakshi *et al.* (1997) and Lim *et al.* (2005), or Garcia *et al.* (2010) for a recent review.

In response to these now well-established empirical findings, many alternative time-varying volatility models have been proposed, with continuous time stochastic volatility (SV) models - often augmented by random jump processes - being particularly prominent of late. This focus on the latter form of models is due, in part, to the availability of (semi-) closed-form option prices, with variants of the SV model of Heston (1993) becoming the workhorse of the empirical option pricing literature (e.g. Eraker, 2004, Forbes *et al.* 2007, Broadie *et al.* 2007, Johannes *et al.* 2009), and MCMC and particle filtering techniques the typical numerical methods of choice. It is of interest, therefore, to assess the performance of the proposed ABC method in the context of this form of model.

We adopt here the simplest version of Heston ‘square root’ SV model, given by:

$$r_t = \sqrt{V_t}\epsilon_t \tag{23}$$

$$dV_t = (\delta - \alpha V_t) dt + \sigma_v \sqrt{V_t} dW_t, \tag{24}$$

where  $r_t$  denotes the (demeaned) logarithmic return on an asset price over period (day say)  $t$ ,  $\epsilon_t \sim i.i.d.N(0, 1)$ ,  $W_t$  is a standard Wiener process, and the restriction  $2\delta \geq \sigma_v^2$  ensures the positivity of the stochastic variance  $V_t$  ( $= x_t$  in our previous generic notation). For  $\alpha > 0$ ,  $V_t$  is mean reverting and as  $t \rightarrow \infty$  the variance approaches a steady state gamma distribution, with  $E[V_t] = \delta/\alpha$  and  $var(V_t) = \sigma_v^2 \delta / 2\alpha^2$ . The transition density for  $V_t$ , conditional on  $V_{t-1}$ , is

$$p(V_t|V_{t-1}) = c \exp(-u - v) \left(\frac{v}{u}\right)^{q/2} I_q(2(uv)^{1/2}), \tag{25}$$

where  $c = 2\alpha/\sigma_v^2(1 - \exp(-\alpha))$ ,  $u = cV_{t-1} \exp(-\alpha)$ ,  $v = cV_t$ ,  $q = \frac{2\delta}{\sigma_v^2} - 1$ , and  $I_q(\cdot)$  is the modified Bessel function of the first kind of order  $q$ . The conditional distribution function is non-Central chi-square,  $\chi^2(2cV_t; 2q + 2, 2u)$ , with  $2q + 2$  degrees of freedom and non-centrality parameter

$2u$ . The discrete-time model for  $r_t$  in (23) can be viewed as a discretized version of a diffusion process for returns (or returns can be viewed as being *inherently* discretely observed) whilst we retain the diffusion model for the latent variance. That is, we eschew the discretization of the variance process that would typically be used, with simulation of  $V_t$  in Step 2 of the ABC algorithm occurring exactly, through the treatment of the  $\chi^2(2cV_t; 2q + 2, 2u)$  random variable as a composition of central  $\chi^2(2q + 2 + 2j)$  and  $j \sim \text{Poisson}(u)$  variables.<sup>5</sup>

For the purpose of this illustration we set the parameters in (24) to values that produce simulated values of both  $r_t$  and  $V_t$  that match the characteristics of (respectively) daily returns and daily values of realized volatility (constructed from 5 minute returns) for the S&P500 stock index over the 2003-2004 period, namely

$$\rho = 1 - \alpha = 0.92; \quad \delta = 0.0024; \quad \sigma_v = 0.062. \quad (26)$$

This relatively calm period in the stock market is deliberately chosen as a reference point, as the inclusion of price and volatility jumps, and/or a non-Gaussian conditional distribution in the model would be an empirical necessity for any more volatile period, such as that witnessed during the recent 2008/2009 financial crisis.

In order to implement the score-based ABC method, using the AUKF algorithm to evaluate an auxiliary likelihood, we invoke the following discretization, based on (an exact) transformation of the measurement equation and an Euler approximation of the state equation,

$$\begin{aligned} \ln(r_t^2) &= \ln(V_t) + \ln(\varepsilon_t^2) \\ &\Rightarrow \\ y_t &= \ln(V_t) + e_t \end{aligned} \quad (27)$$

$$V_t = \delta + \rho V_{t-1} + \sigma_v \sqrt{V_{t-1}} v_t, \quad (28)$$

where  $v_t$  is treated as a truncated Gaussian variable with lower bound,

$$v_t > \frac{-(\delta + \rho V_{t-1})}{\sigma_v \sqrt{V_{t-1}}}. \quad (29)$$

Directing the reader to the Appendix for the detailed outline of the AUKF approach, we note that sigma points that span the support of  $e_t$  are defined by calculating  $E(e_t)$  and  $\text{var}(e_t)$  using deterministic integration and the closed form of  $p(e_t)$ , with the specification of  $a_e = b_e = \sqrt{3}$  adopted for convenience. Those for  $V_t$ ,  $t = 0, 1, \dots, T$ , are defined as:

$$V_t^1 = E(V_t^1 | \cdot); \quad V_t^2 = E(V_t^1 | \cdot) + \sqrt{3} \sqrt{\text{var}(V_t^1 | \cdot)}; \quad V_t^3 = 0.00001,$$

---

<sup>5</sup>As noted earlier, even in the typical case in which the transition distribution of the diffusion is unknown, an arbitrarily fine discretization, limited only by computer power, effectively enables draws from the exact diffusion to be produced.

where  $E(V_t^1|.)$  and  $var(V_t^1|.)$  respectively denote the mean and variance of the relevant distribution of  $V_t$  (marginal, filtered or predictive, depending on the particular step in the AUKF algorithm). The sigma points for (29) are then defined using the mean and variance of the truncated normal distribution, with the value of  $V_{t-1}$  in (29) represented using the relevant sigma point for  $V_{t-1}$ , and (with reference to the Appendix)  $a_v = b_v = \sqrt{3}$  specified.

In order to evaluate the accuracy of the estimate of the posterior produced using the ABC method, we produce the exact joint posterior distribution for  $\phi = (\rho, \delta, \sigma_v)'$  via the deterministic non-linear filtering method of Ng *et al.* (2013). In brief, this method represents the recursive filtering and prediction distributions used to define the likelihood function as the numerical solutions of integrals defined over the support of  $e_t$  in (27), with deterministic integration used to evaluate the relevant integrals, and the *exact* transitions in (25) used in the specification of the filtering and up-dating steps.<sup>6</sup> Whilst lacking the general applicability of the ABC-based method proposed here, this deterministic filtering method is ideal for the particular model used in this illustration, and can be viewed as producing a very accurate estimate of the exact density, without any of the simulation error that would be associated with an MCMC-based comparator, for instance. We refer the reader to Ng *et al.* for more details of the technique; see also Kitagawa (1987).<sup>7</sup> The likelihood function, evaluated via this method, is then multiplied by a uniform prior that imposes the restrictions:  $0 < \rho < 1$ ,  $\delta > 0$ ,  $\sigma_v^2 > 0$  and  $2\delta \geq \sigma_v^2$ . The three marginal posteriors are then produced via deterministic numerical integration (over the parameter space), with a very fine grid on  $\phi$  being used to ensure accuracy.

We compare the auxiliary model-based ABC technique with ABC approaches based on summary statistics, as discussed in the linear Gaussian context in Section 5.1. For want of a better choice we use the vector of statistics given attention therein, as well as the two alternative distance measures described there. We also compute marginal posterior densities estimated using the AUKF-based approximation (of the likelihood) itself. That is, using the AUKF to evaluate the likelihood function associated with the discretized model in (27) and (28) and normalizing (using a uniform prior) produces an approximation of the posterior which can, in principle, be invoked as an approximation in its own right, independently of its subsequent use as a score generator with an ABC algorithm. Finally, we compute the marginal posteriors (again, based on a uniform prior) using the likelihood function of the Euler approximation in (27) and (28)

---

<sup>6</sup>A numerically efficient and stable algorithm for evaluating the transitions densities for the exact model - in which a non-central chi-squared density is represented as an infinite mixture of Poisson and central chi-squared densities - is used. This enables the numerical problems typically associated with the direct computation of Bessel functions to be avoided.

<sup>7</sup>We note that the application of this filter in Ng *et al.* is to a non-parametric representation of  $e_t$ . In the current setting, in which  $e_t$  is specified parametrically, the known form of the distribution of  $e_t$  is used directly in the evaluation of the relevant integrals. We refer the reader to Section 2.2. of that paper for a full description of the algorithm. Preliminary experimentation with the number of grid points used in the deterministic integration was undertaken in order to ensure that the resulting estimate of the likelihood function/posterior stabilized, with 100 grid points underlying the final results documented here.

evaluated using the Ng *et al.* (2013) filtering method, with Gaussian transitions used in the filtering and up-dating steps. When normalized, this density can be viewed as the quantity that a typical MCMC scheme (as based on the equivalent prior) would be targeting, given that the tractability of Gaussian approximations to the transitions would typically be exploited in structuring an MCMC algorithm.

As a concluding note on computational matters, we re-iterate that the time taken to evaluate the AUKF-based approximate likelihood function at any point in the parameter space is roughly comparable to that required for KF evaluation, thus rendering it a feasible method to be inserted within the ABC algorithm. In contrast, evaluation of the Euler-based likelihood via the Ng *et al.* (2013) technique, whilst producing, in the main, (as will be seen in the following section) a more accurate estimate of the exact posterior than the AUKF method, is many orders of magnitude slower and, hence, simply infeasible as a score generator within ABC.

### 5.2.2 Numerical results for the Heston SV model

In order to abstract initially from the impact of dimensionality on the ABC methods, we first report results for each single parameter of the Heston model, keeping the remaining two parameters fixed at their true values. Three ABC-based estimates of the relevant exact (univariate) posterior, invoking a uniform prior, are produced in this instance. Three matching statistics are used, respectively: 1) the (uni-dimensional) auxiliary score based on the approximating model (ABC-score); 2) the summary statistics in (7), matched via the Euclidean distance measure in (18) (ABC-summ stats); and 3) the summary statistics in (7), matched via the FP distance measure in (20) (ABC-FP). We produce representative posterior (estimates) in each case, to give some visual idea of the accuracy (or otherwise) that is achievable via the ABC methods. We then summarize accuracy by reporting the average (over the 100 runs) of the root mean squared error (RMSE) of each ABC-based estimate of the exact posterior for a given parameter, computed as:

$$RMSE = \sqrt{\frac{1}{G} \sum_{g=1}^G (\hat{p}_g - p_g)^2}, \quad (30)$$

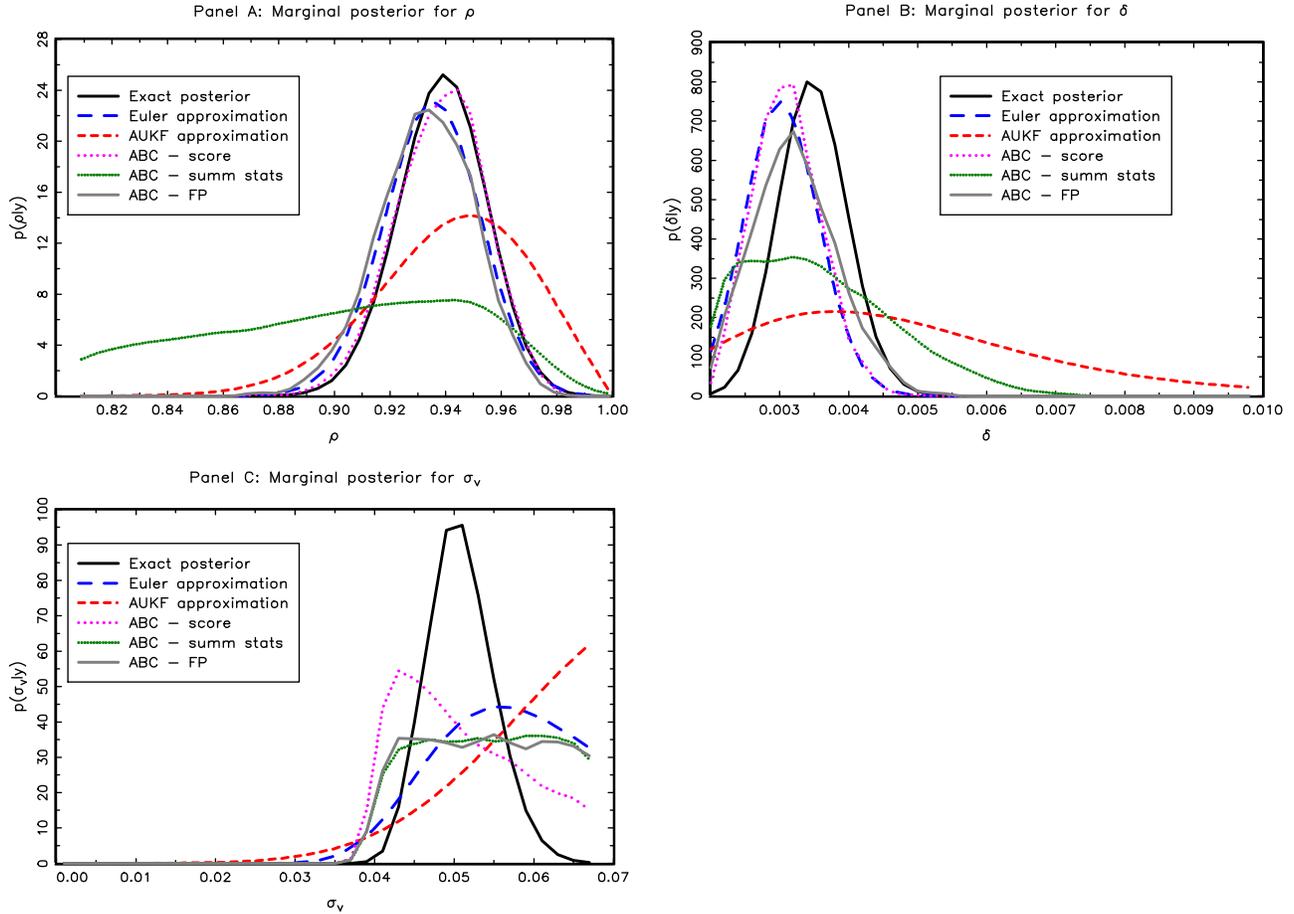
where  $\hat{p}_g$  is the ordinate of the ABC density estimate and  $p_g$  the ordinate of the exact posterior density, at the  $g$ th grid-point used to produce the plots.<sup>8</sup> All single parameter results are documented in Panel A of Table 1. We also tabulate there, as benchmarks of a sort, the RMSEs associated with the (one-off) AUKF- and Euler-based approximations of each univariate density.

Figure 5, Panel A reproduces the exact posterior of (the single unknown parameter)  $\rho$ , the posteriors associated with the AUKF- and Euler-based approximations, and the three ABC-based estimates. As is clear, the AUKF-based approximation is reasonably inaccurate, in terms of replicating the location and shape of the exact posterior - an observation that is interesting in

---

<sup>8</sup>The associated box plots are also available on request, but were not included here due to space considerations.

Figure 5: Posterior densities for each single unknown parameter of the model in (27) and (28), with the other two parameters set to their true values. As per the key, the graphs reproduced are the exact, Euler- and AUKF-based approximations, in addition to the three ABC-based estimates. Both the exact and Euler posteriors are evaluated using the grid-based non-linear filter of Ng *et al.* (2013).



its own right, given the potential for such a simple and computationally efficient approximation method to be used to evaluate likelihood functions (and posterior distributions) in non-linear state space models such as the one under consideration. However, once the approximation is embedded within an ABC scheme, in the manner described in Section 4, the situation is altogether different, with the (pink) dotted line (denoted by ‘ABC-score’ in the key) providing a remarkably accurate estimate of the exact posterior, using only 50,000 replications of the simplest rejection-based ABC algorithm, and fifteen minutes of computing time on a desktop computer. It is worth noting that the ABC-based estimate is also more accurate in this case than the Euler approximation, where we highlight, once again, that production of the latter still requires the application of the much more computationally burdensome non-linear filtering method. Most notably, the ABC method based on the summary statistics, combined using a Euclidean distance measure, performs very badly, although the dimensional reduction technique of Fearnhead and

Prangle (2012), applied to this same set of summary statistics, yields a reasonable estimate of the exact posterior in this instance.

Comparable graphs are produced for the single parameters  $\delta$  and  $\sigma_v$  in Panels B and C respectively of Figure 5, with the remaining pairs of parameters ( $\rho$  and  $\sigma_v$ , and  $\rho$  and  $\delta$  respectively) held fixed at their true values. In the case of  $\delta$ , the score-based method arguably provides the best representation of the shape of the exact posterior, despite being slightly inaccurate in terms of location. The Fearnhead and Prangle (2012) method also provides a reasonable estimate, whilst the summary statistic approach using the Euclidean distance, once again performs very poorly. For the parameter  $\sigma_v$ , *only* the score based method yields a density with a well-defined shape, with the two summary statistic-based techniques essentially producing uniform densities that reflect little more than the restricted support imposed on the parameter draws. Interestingly, the Euler approximation itself provides a quite poor representation of the exact marginal, a result which has not, as far as we know, been remarked upon in the literature, given that a typical MCMC scheme (as noted above) would in fact be targeting the Euler density itself, as the best representation of the true model - the exact posterior remaining unaccessed due to the difficulty of devising an effective MCMC scheme which uses the exact transitions. For neither  $\delta$  nor  $\sigma_v$  is the AUKF approximation *itself* particularly accurate, despite the fact that it respects the non-linearity in the true state space model.

The RMSE results recorded in Panel A of Table 1 confirm the qualitative nature of the single-run graphical results. For  $\delta$  and  $\sigma_v$ , all three ABC-based estimates are seen to produce lower RMSE values (sometimes an order of magnitude lower) than the AUKF approximation, indicating that *any* of the ABC procedures would yield gains over the use of the unscented filtering method itself. For  $\rho$ , the AUKF approximation is better than that of the summary statistic-based ABC estimate, but in part because the latter is so poor. For  $\rho$  and  $\sigma_v$ , the score-based ABC method is the most accurate, and is most notably *very* precise for the case of the persistence parameter  $\rho$ . For the parameter  $\delta$  the FP method is the most accurate according to this measure although, as indicated by the nature of the graphs in Panel B of Figure 5, this result does tend to understate the ability of the score-based method to capture the basic shape of the exact posterior. For  $\rho$  and  $\delta$ , the (Euclidean) summary-statistic method is an order of magnitude more inaccurate than the other two ABC methods, whilst also exhibiting no ability to identify the shape of the true posterior in the case of  $\sigma_v$ . Once again as is consistent with the graphs in Figure 5, the Euler approximation for  $\rho$  is reasonably accurate, but not as accurate as the score-based ABC estimate. For  $\delta$  and  $\sigma_v$ , the Euler approximation, whilst being more accurate than the AUKF approximation, is dominated by all three ABC estimates.

Table 1: RMSE of an estimated marginal and the exact marginal: average RMSE value over multiple runs of ABC using 50,000 replications. ‘Score’ refers to the ABC method based on the score of the AUKF model; ‘SS’ refers to the ABC method based on a Euclidean distance for the summary statistics in (7); ‘FP’ refers to the Fearnhead and Praugle ABC method, based on the summary statistics in (7). For the single parameter case, the (single) score method is documented in the row denoted by ‘ABC-Marginal Score’, whilst in the multi-parameter case, there are results for both the joint and marginal score methods. The RMSE of the AUKF and Euler approximations (computed once only, using the observed data) are recorded as benchmarks, in the top two rows of each panel. For the single and dual parameter cases, 100 runs of ABC were used to produce the results, whilst for the three parameter case, 50 runs were used. The smallest RMSE figure in each column is highlighted in bold.

| Approximate Density | Panel A       |               |               | Panel B       |               | Panel C       |               | Panel D        |               |               |
|---------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|----------------|---------------|---------------|
|                     | One unknown   |               |               | Two unknowns  |               | Two unknowns  |               | Three unknowns |               |               |
|                     | $\rho$        | $\delta$      | $\sigma_v$    | $\rho$        | $\sigma_v$    | $\rho$        | $\delta$      | $\rho$         | $\delta$      | $\sigma_v$    |
| AUKF                | 0.0263        | 0.0535        | 0.0935        | 0.0529        | 0.0370        | 0.0185        | 0.0798        | 0.0201         | 0.0862        | 0.0287        |
| Euler               | 0.0096        | 0.0434        | 0.0664        | 0.0072        | 0.0308        | 0.0175        | 0.0818        | 0.0120         | 0.0459        | 0.0242        |
| ABC-Joint Score     | -             | -             | -             | 0.0054        | <b>0.0217</b> | 0.0101        | 0.0441        | <b>0.0063</b>  | <b>0.0124</b> | <b>0.0166</b> |
| ABC-Marginal Score  | <b>0.0028</b> | 0.0301        | <b>0.0392</b> | <b>0.0045</b> | 0.0219        | <b>0.0048</b> | 0.0480        | 0.0085         | 0.0381        | 0.0167        |
| ABC-SS              | 0.0353        | 0.0316        | 0.0427        | 0.0310        | 0.0234        | 0.0119        | <b>0.0312</b> | 0.0109         | 0.0389        | 0.0170        |
| ABC-FP              | 0.0178        | <b>0.0215</b> | 0.0431        | 0.0145        | 0.0233        | 0.0093        | 0.0358        | 0.0124         | 0.0407        | 0.0168        |

In Panels B and C respectively of Table 1, we record all RMSE results for the case when two, then all three parameters are unknown, with a view to gauging the relative performance of the ABC methods when multiple matches are required. In these multiple parameter cases, a preliminary run of ABC, based on uniform priors defined over the domains defined above, has been used to determine the high mass region of the joint posterior. This information has been used to further truncate the priors in a subsequent ABC run, and final marginal posterior estimates then produced. The results recorded in Panels B and C highlight that when two parameters are unknown (either  $\rho$  and  $\sigma_v$  or  $\rho$  and  $\delta$ ), the score-based ABC method produces the most accurate density estimates in three of the four cases. Marginalization produces an improvement in accuracy for  $\rho$ . For the other two parameters marginalization does not yield an increase in accuracy; however, the differences between the joint and marginal score estimates are minimal. Only in one case (as pertains to  $\delta$ ) does an ABC method based on summary statistics outperform the score-based methods. In all four cases, the AUKF estimate is inferior to all other comparators, and the Euler approximation also inferior to all ABC-based estimates in three of the four cases. As is seen in Panel D, when all three parameters are to be estimated, the score-based ABC estimates remain the most accurate, with the joint score method superior overall and yielding notably improvements in accuracy over both the AUKF and Euler approximations.

## 6 Conclusions and Discussion

This paper has explored the application of approximate Bayesian computation in the state space setting. Certain fundamental results have been established, namely the lack of reduction to finite sample sufficiency and the Bayesian consistency of the auxiliary model-based method. The (limiting) equivalence of ABC estimates produced by the use of both the maximum likelihood and score-based summary statistics has also been demonstrated. The idea of tackling the dimensionality issue that plagues the application of ABC in high dimensional problems via an integrated likelihood approach has been proposed. The approach has been shown to work extremely well in the case in which the auxiliary model is exact, and to yield some benefits otherwise. However, a much more comprehensive analysis of different non-linear settings (and auxiliary models) would be required for a definitive conclusion to be drawn about the trade-off between the gain to be had from marginalization and the loss that may stem from integrating over an *inaccurate* auxiliary model.

Indeed, the most important challenge that remains, as is common to the related frequentist techniques of indirect inference and efficient methods of moments, is the specification of a computationally efficient and accurate approximating model. Given the additional need for parsimony, in order to minimize the number of statistics used in the matching exercise, the principle of aiming for a large nesting model, with a view to attaining full asymptotic sufficiency, is not an

attractive one. We have illustrated the use of one simple approximation approach based on the unscented Kalman filter. The relative success of this approach in the particular example considered, certainly in comparison with methods based on other more *ad hoc* choices of summary statistics, augers well for the success of score-based methods in the non-linear setting. Further exploration of approximation methods in other non-linear state space models is the subject of on-going research. (See also Creel and Kristensen, 2014, for some contributions on this front.)

Finally, we note that despite the focus of this paper being on inference about the static parameters in the state space model, there is nothing to preclude marginal inference on the states being conducted, at a second stage. Specifically, conditional on the (accepted) draws used to estimate  $p(\phi|\mathbf{y})$ , existing filtering and smoothing methods (including the recent methods that exploit ABC at the filtering/smoothing level; see, for example, Jasra *et al.*, 2010, Calvet and Czellar, 2014, Martin *et al.*, 2014) could be used to yield draws of the states, and (marginal) smoothed posteriors for the states produced via the usual averaging arguments. With the asymptotic properties of both approaches established (under relevant conditions), of particular interest would be a comparison of both the finite sample accuracy and computational burden of the ABC-PMCMC method developed Martin *et al.* (2014), with that of the method proposed herein, in which  $p(\phi|\mathbf{y})$  is targeted more directly via the score-based approach.

## References

- [1] Anderson, T.W. 1958. *The Statistical Analysis of Time Series*, John Wiley and Sons.
- [2] Bakshi, G., Cao, C. and Chen, Z. 1997. Empirical Performance of Alternative Option Pricing Models, *Journal of Finance*, 52, 2003-2049.
- [3] Black, F. and Scholes, M. 1973. The Pricing of Options and Corporate Liabilities, *Journal of Political Economy*, 81, 637-659.
- [4] Bollerslev, T., Chou, R.Y. and Kroner, K.F. 1992. ARCH Modelling in Finance: A Review of the Theory and Empirical Evidence, *Journal of Econometrics*, 52, 5-59.
- [5] Beaumont, M.A. 2010. Approximate Bayesian Computation in Evolution and Ecology, *Annual Review of Ecology, Evolution, and Systematic*, 41, 379-406.
- [6] Beaumont, M.A., Cornuet, J-M., Marin, J-M. and Robert, C.P. 2009. Adaptive Approximate Bayesian Computation, *Biometrika* 96, 983-990.
- [7] Beaumont, M.A., Zhang, W. and Balding, D.J. 2002. Approximate Bayesian Computation in Population Genetics, *Genetics* 162, 2025-2035.

- [8] Biau, G., Cérou, F. and Guyader, A. 2014. New Insights into Approximate Bayesian Computation. *Annales de l'IHP (Probability and Statistics)*, In press.
- [9] Blum, M.G.B. 2010a. Choosing the Summary Statistics and the Acceptance Rate in Approximate Bayesian Computation, *Proceedings of Compstat 2010*.
- [10] Blum, M.G.B. 2010b. Approximate Bayesian Computation: a Nonparametric Perspective, *Journal of the American Statistical Association* 105, 1178-1187.
- [11] Blum, M.G.B. and François, O. 2010. Non-linear Regression Models for Approximate Bayesian Computation, *Statistics and Computing* 20, 63-73.
- [12] Blum, M.G.B., Nunes, M.A., Prangle, D. and Sisson, S.A. 2013. A Comparative Review of Dimension Reduction Methods in Approximate Bayesian Computation, *Statistical Science*, 28, 189-208.
- [13] Broadie, M., Chernov, M. and Johannes, M. 2007. Model Specification and Risk Premia: Evidence from Futures Options, *The Journal of Finance*, LXII: 1453-1490.
- [14] Calvet, C. and Czellar, V. 2014. Accurate Methods for Approximate Bayesian Computation Filtering. *Journal of Econometrics* (to appear).
- [15] Cornuet, J-M., Santos, F., Beaumont, M.A., Robert, C.P., Marin, J-M., Balding, D.J., Guillemand, T. and Estoup, A. 2008. Inferring Population History with DIY ABC: a User-friendly Approach to Approximate Bayesian Computation, *Bioinformatics* 24, 2713-2719.
- [16] Cox, D.R. and Hinkley, D.V. 1974. *Theoretical Statistics*, Chapman and Hall, London.
- [17] Creel, M. and Kristensen, D. 2014. ABC of SV: Limited Information Likelihood Inference in Stochastic Volatility Jump-Diffusion Models, *Draft paper*.
- [18] Dean, T. and Singh, S. 2011. Asymptotic Behaviour of Approximate Bayesian Estimators, *Technical Report, University of Cambridge*.
- [19] Dean, T., Singh, S., Jasra, A. and Peters, G. 2014. Parameter Inference for Hidden Markov Models with Intractable Likelihoods, *Scand. J. Statist.* (to appear).
- [20] Drovandi, C.C., Pettitt, A.N. and Faddy, M.J. 2011. Approximate Bayesian Computation Using Indirect Inference, *JRSS(C)*, 60 1 - 21.
- [21] Drovandi, C. and Pettitt, A. 2013. Bayesian Indirect Inference. <http://eprints.qut.edu.au/63767/>.

- [22] Eraker, B. 2004. Do Stock Prices and Volatility Jump? Reconciling Evidence from Spot and Option Prices, *The Journal of Finance*, LIX: 1367-1403.
- [23] Fearnhead, P, Prangle, D. 2012. Constructing Summary Statistics for Approximate Bayesian Computation: Semi-automatic Approximate Bayesian Computation. *Journal of the Royal Statistical Society, Series B*. 74: 419–474.
- [24] Forbes C.S., Martin, G.M. and Wright J. 2007. Inference for a Class of Stochastic Volatility Models Using Option and Spot Prices: Application of a Bivariate Kalman Filter, *Econometric Reviews, Special Issue on Bayesian Dynamic Econometrics*, 26: 387-418.
- [25] Gallant, A.R. and Tauchen, G. 1996. Which Moments to Match, *Econometric Theory* 12, 657-681.
- [26] Garcia, R., Ghysels, E. and Renault, E. 2010. The Econometrics of Option Pricing, *Handbook of Financial Econometrics*, Vol 1. (eds Ait-Sahalia, Y. and Hansen, L.), North Holland.
- [27] Giordini, P, Pitt, M, Kohn R. 2011. Bayesian Inference for Time series State Space Models. *The Oxford Handbook of Bayesian Econometrics* (Eds. Geweke, Koop, van Dijk), New York.
- [28] Gleim, A, Pigorsch, C. 2013. Approximate Bayesian Computation with Indirect Summary Statistics. Draft paper: <http://ect-pigorsch.mee.uni-bonn.de/data/research/papers/>.
- [29] Gouriéroux, C. and Monfort, A. 1995. *Statistics and Econometric Models*. CUP.
- [30] Gouriéroux, C. and Monfort, A. 1996. *Simulation-based Econometric Methods*, OUP.
- [31] Gouriéroux, C., Monfort, A. and Renault, E. 1993. Indirect Inference, *Journal of Applied Econometrics*, 85, S85-S118.
- [32] Haykin, S. 2001. *Kalman Filtering and Neural Networks*, John Wiley and Sons.
- [33] Heggland, K. and Frigessi, A. 2004. Estimating Functions in Indirect Inference, *JRSS(B)* 66, 447-462.
- [34] Heston, S.L. 1993. A Closed-form Solution for Options with Stochastic Volatility with Applications to Bond and Currency Options, *The Review of Financial Studies*, 6, 327-343.
- [35] Jasra, A, Kantas, N. and Ehrlich, E. 2014. Approximate Inference for Observation Driven Time Series Models with Intractable Likelihoods. *TOMACS* (to appear).
- [36] Jasra, A, Lee, A, Yau, C. and Zhang, X. 2013. The Alive Particle Filter. *arXiv preprint*.

- [37] Jasra, A, Singh, S, Martin, J., McCoy, E. 2010. Filtering via Approximate Bayesian Computation. *Statistics and Computing* 22, 1223-1237.
- [38] Johannes, M., Polson, N.G. and Stroud, J.R. 2009. Optimal Filtering of Jump-Diffusions: Extracting Latent States from Asset Prices, *Review of Financial Studies*, 22: 2759-2799.
- [39] Joyce, P. and Marjoram, P. 2008. Approximately Sufficient Statistics and Bayesian Computation. *Statistical applications in genetics and molecular biology*, 7, 1-16.
- [40] Julier, S.J. and Uhlmann, J.K. 2004. Unscented Filtering and Nonlinear Estimation. *Proceedings of the IEEE* 92, 401-422.
- [41] Julier, S.J., Uhlmann, J.K. and Durrant-Whyte, H.F. 1995. A New Approach for Filtering Nonlinear Systems. *Proceedings of the American Control Conference*, 1628-1632.
- [42] Julier, S.J., Uhlmann, J.K. and Durrant-Whyte, H.F. 2000. A New Method for the Nonlinear Transformation of Means and Covariances in Filters and Estimators, *IEEE Transactions on Automatic Control* 45, 477-481.
- [43] Le Gland, F. and Oudjane, N. 2006. A Sequential Particle Algorithm that Keeps the Particle System Alive. In *Stochastic Hybrid Systems : Theory and Safety Critical Applications* (H. Blom and J. Lygeros, Eds), Lecture Notes in Control and Information Sciences 337, 351–389, Springer: Berlin.
- [44] Lim, G.C., Martin, G.M. and Martin, V.L. 2005. Parametric Pricing of Higher Order Moments in S&P500 Options, *Journal of Applied Econometrics*, 20, 377-404.
- [45] Kitagawa, G. 1987. Non-Gaussian State Space Modeling of Nonstationary Time Series. *Journal of the American Statistical Association* 76, 1032-1064.
- [46] Marin, J-M, Pudlo, P, Robert C, Ryder, R. 2011. Approximate Bayesian Computation Methods. *Statistics and Computing* 21, 289–291.
- [47] Marjoram, P., Molitor, J., Plagnol, V. and Tavaré, S. 2003. Markov Chain Monte Carlo Without Likelihoods, *Proceedings of the National Academie of Science USA* 100, 15324-15328.
- [48] Martin, J. S., Jasra, A., Singh, S. S., Whiteley, N., Del Moral, P. and McCoy, E. 2014. Approximate Bayesian Computation for Smoothing, *Stoch. Anal. Appl.* (to appear).
- [49] Milstein, G. 1978. A Method of Second Order Accuracy Integration of Stochastic Differential Equations, *Theory of Probability and Its Applications* 23, 396-401.

- [50] Ng, J., Forbes, C.S., Martin, G.M. and McCabe, B.P.M. 2013. Non-parametric Estimation of Forecast Distributions in Non-Gaussian, Non-linear State Space Models, *International Journal of Forecasting* 29, 411–430
- [51] Nott D, Fan, Y, Marshall, L, Sisson, S. 2014. Approximate Bayesian Computation and Bayes Linear Analysis: Towards High-dimensional ABC, *Journal of Computational and Graphical Statistics*, 23, 65-86.
- [52] Peters G, Sisson, S, Fan, Y. 2012. Likelihood-free Bayesian Inference for Alpha-stable Models. *Computational Statistics and Data Analysis* 56, 3743-3756.
- [53] Ponomareva, K. and Date, P. 2010. Some Results on Theory and Applications of Higher-order Sigma Point Filter. *Draft Paper*.
- [54] Pritchard, J.K., Seilstad, M.T., Perez-Lezaun, A. and Feldman, M.W. 1999. Population Growth of Human Y Chromosomes: A Study of Y Chromosome Microsatellites, *Molecular Biology and Evolution* 16 1791-1798.
- [55] Simon, D. 2006. *Optimal State Estimation*. John Wiley and Sons. New Jersey.
- [56] Sisson S. and Fan, Y. 2011. Likelihood-free Markov Chain Monte Carlo. In *Handbook of Markov Chain Monte Carlo* (Eds. Brooks, Gelman, Jones, Meng). Chapman and Hall/CRC Press.
- [57] Sisson, S., Fan, Y. and Tanaka, M. 2007. Sequential Monte Carlo without Likelihoods, *Proceedings of the National Academie of Science USA* 104, 1760-1765.
- [58] Smith, A.A. 1993. Estimating Non-linear Time Series Models Using Vector Autoregressions: Two Approaches, *Journal of Applied Econometrics* 8, 63-84.
- [59] Tavaré, S., Balding, D.J., Griffiths, R.C. and Donnelly, P. 1997. Inferring Coalescence Times from DNA Sequence Data, *Genetics* 145, 505-518.
- [60] Toni, T., Welch, D., Strelkowa, N., Ipsen, A. and Stumpf, M.P.H. 2009. Approximate Bayesian Computation Scheme for Parameter Inference and Model Selection in Dynamical Systems, *JRSS (Interface)* 6, 187-202.
- [61] Wegmann, D., Leuenberger, C. and Excoffier, L. 2009. Efficient Approximate Bayesian Computation Coupled with Markov chain Monte Carlo with Likelihood, *Genetics* 182, 1207-1218.
- [62] Wilkinson, R.D. 2013. Approximate Bayesian Computation (ABC) Gives Exact Results under the Assumption of Model Error. *Statist. Appl. Genetics Mole. Biol.*

[63] Yildirim, S, Singh, S.S, Dean, T. and Jasra, A. 2013. Parameter Estimation in Hidden Markov Models with Intractable Likelihoods Using Sequential Monte Carlo. *arXiv preprint*.

### Appendix: Implementation details for the AUKF

Given the assumed invariance (over time) of both  $\nu_t$  and  $e_t$  in (15) and (16), the sigma points are determined as:

$$e^1 = E(e_t); e^2 = E(e_t) + a_e \sqrt{\text{var}(e_t)}; e^3 = E(e_t) - b_e \sqrt{\text{var}(e_t)}$$

and

$$v^1 = E(v_t); v^2 = E(v_t) + a_v \sqrt{\text{var}(v_t)}; v^3 = E(v_t) - b_v \sqrt{\text{var}(v_t)}$$

respectively, and propagated at each  $t$  through the relevant non-linear transformations,  $h_t(\cdot)$  and  $k_t(\cdot)$ . The values  $a_e$ ,  $b_e$ ,  $a_v$  and  $b_v$  are chosen according to the assumed distribution of  $e_t$  and  $v_t$ , with a Gaussian assumption for both variables yielding values of  $a_e = b_e = a_v = b_v = \sqrt{3}$  as being ‘optimal’. Different choices of these values are used to reflect higher-order distributional information and thereby improve the accuracy with which the mean and variance of the non-linear transformations are estimated; see Julier *et al.* (2000) and Ponomareva and Date (2010) for more details. Restricted supports are also managed via appropriate truncation of the sigma points. The same principles are applied to produce the mean and variance of the time varying state  $x_t$ , except that the sigma points need to be recalculated at each time  $t$  to reflect the up-dated mean and variance of  $x_t$  as each new value of  $y_t$  is realized.

In summary, the steps of the AUKF applied to evaluate the likelihood function of (15) and (16) are as follows:

1. Use the (assumed) marginal mean and variance of  $x_t$ , along with the invariant mean and variance of  $v_t$  and  $e_t$  respectively, to create the  $(3 \times 7)$  matrix of augmented sigma points for  $t = 0$ ,  $X_{a0}$ , as follows. Define:

$$E(X_{a0}) = \begin{bmatrix} E(x_t) \\ E(v_t) \\ E(e_t) \end{bmatrix}, P_{a0} = \begin{bmatrix} \text{var}(x_t) & 0 & 0 \\ 0 & \text{var}(v_t) & 0 \\ 0 & 0 & \text{var}(e_t) \end{bmatrix}, \quad (31)$$

and  $\sqrt{P_{a0j}}$  as the  $j$ th column of the Cholesky decomposition (say) of  $P_{a0}$ . Given the diagonal form of  $P_{a0}$  (in this case), we have

$$\sqrt{P_{a01}} = \begin{bmatrix} \sqrt{\text{var}(x_t)} \\ 0 \\ 0 \end{bmatrix}; \sqrt{P_{a02}} = \begin{bmatrix} 0 \\ \sqrt{\text{var}(v_t)} \\ 0 \end{bmatrix}; \sqrt{P_{a03}} = \begin{bmatrix} 0 \\ 0 \\ \sqrt{\text{var}(e_t)} \end{bmatrix}.$$

The seven columns of  $X_{a0}$  are then generated by

$$E(X_{a0}); E(X_{a0}) + a_j\sqrt{P_{a0j}}; \text{ for } j = 1, 2, 3; E(X_{a0}) - b_j\sqrt{P_{a0j}}; \text{ for } j = 1, 2, 3,$$

where  $a_1 = a_x$ ,  $a_2 = a_v$  and  $a_3 = a_e$ , and the corresponding notation is used for  $b_j$ ,  $j = 1, 2, 3$ .

2. Propagate the  $t = 0$  sigma points through the transition equation as  $X_{x1} = k_1(X_{a0}, \phi)$  and estimate the predictive mean and variance of  $x_1$  as:

$$E(x_1|y_0) = \sum_{i=1}^7 w_i X_{x1}^i \quad (32)$$

$$var(x_1|y_0) = \sum_{i=1}^7 w_i (X_{x1}^i - E(x_1|y_0))^2, \quad (33)$$

where  $X_{x1}^i$  denotes the  $i$ th element of the  $(1 \times 7)$  vector  $X_{x1}$  and  $w_i$  the associated weight, determined as an appropriate function of the  $a_j$  and  $b_j$ ; see Ponomareva and Date (2010).

3. Produce a new matrix of sigma points,  $X_{a1}$ , for  $t = 1$  generated by

$$E(X_{a1}); E(X_{a1}) + a_j\sqrt{P_{a1j}}; \text{ for } j = 1, 2, 3; E(X_{a1}) - b_j\sqrt{P_{a1j}}; \text{ for } j = 1, 2, 3,$$

using the updated formulae for the mean and variance of  $x_t$  from (32) and (33) respectively, in the calculation of  $E(X_{a1})$  and  $P_{a1}$ .

4. Propagate the  $t = 1$  sigma points through the measurement equation as  $X_{y1} = h_1(X_{a1}, \phi)$  and estimate the predictive mean and variance of  $y_1$  as:

$$E(y_1|y_0) = \sum_{i=1}^7 w_i X_{y1}^i \quad (34)$$

$$var(y_1|y_0) = \sum_{i=1}^7 w_i (X_{y1}^i - E(y_1|y_0))^2, \quad (35)$$

where  $X_{y1}^i$  denotes the  $i$ th element of the  $(1 \times 7)$  vector  $X_{y1}$  and  $w_i$  is as defined in Step 3.

5. Estimate the first component of the likelihood function,  $p(y_1|y_0)$ , as a Gaussian distribution with mean and variance as given in (34) and (35) respectively.
6. Given observation  $y_1$  produce the up-dated filtered mean and variance of  $x_t$  via the usual KF up-dating equations:

$$\begin{aligned} E(x_1|y_1) &= E(x_1|y_0) + M_1(y_1 - E(y_1|y_0)) \\ var(x_1|y_1) &= var(x_1|y_0) - M_1^2 var(y_1|y_0), \end{aligned}$$

where:

$$M_1 = \frac{\sum_{i=1}^7 w_i (X_{x_1}^i - E(x_1|y_0))(X_{y_1}^i - E(y_1|y_0))}{var(y_1|y_0)}$$

and the  $X_{x_1}^i$ ,  $i = 1, 2, \dots, 7$  are as computed in Step 3.

7. Continue as for Steps 2 to 6, with the obvious up-dating of the time periods and the associated indexing of the random variables and sigma points, and with the likelihood function in evaluated as the product of the components produced in each implementation of Step 5, and the log-likelihood in (17) produced accordingly.