

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

Bayesian Estimation of Agent-Based Models

This is the author's manuscript

Original Citation:

Availability:

This version is available <http://hdl.handle.net/2318/1619141> since 2017-02-23T14:32:33Z

Terms of use:

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

Bayesian Estimation of Agent-Based Models

Jakob Grazzini^{a,b,c}, Matteo Richiardi^{d,e}, and Mike Tsionas^f

^aCatholic University of Milan, Department of Economics and Finance, Italy.

^bCatholic University of Milan, Complexity Lab in Economics, Italy.

^cCESifo, Munich, Germany

^dInstitute for New Economic Thinking and Nuffield College, Oxford, UK.

^eUniversity of Torino, Department of Economics and Statistics, and Collegio Carlo Alberto, Italy

^fLancaster University Management School, UK

February 7, 2017

Abstract

We consider Bayesian inference techniques for Agent-Based (AB) models, as an alternative to simulated minimum distance (SMD). Three computationally heavy steps are involved: (i) simulating the model, (ii) estimating the likelihood and (iii) sampling from the posterior distribution of the parameters. Computational complexity of AB models implies that efficient techniques have to be used with respect to points (ii) and (iii), possibly involving approximations. We first discuss non-parametric (kernel density) estimation of the likelihood, coupled with Markov chain Monte Carlo sampling schemes. We then turn to parametric approximations of the likelihood, which can be derived by observing the distribution of the simulation outcomes around the statistical equilibria, or by assuming a specific form for the distribution of external deviations in the data. Finally, we introduce Approximate Bayesian Computation techniques for likelihood-free estimation. These allow embedding SMD methods in a Bayesian framework, and are particularly suited when robust estimation is needed. These techniques are first tested in a simple price discovery model with one parameter, and then employed to estimate the behavioural macroeconomic model of De Grauwe (2012), with nine unknown parameters.

JEL codes: C11, C15.

1 Introduction

Agent-based (AB) models are structural dynamical models characterized by three features: (i) there are a multitude of objects that interact with each other and with the environment, (ii) these objects are autonomous, that is there is no central, or ‘top-down’ control over their behaviour and more generally on the dynamics of the system (e.g. a Walrasian auctioneer), and (iii) aggregation is performed numerically (Richiardi, 2012). AB models are increasingly used in disciplines as diverse as geography, anthropology, sociology, biology, political science, epidemiology (Macal, 2016). In macroeconomics, they have been proposed as an alternative to dynamic stochastic general equilibrium (DSGE) models, where the tenet of rational expectations is replaced by an explicit modelling of learning and selection, thus allowing for more heterogeneity in agents’ characteristics and behaviour and a more detailed account of the (physical and institutional) environment, leading to more complex interaction patterns (Richiardi, 2016).¹

There is a small but growing literature on estimation of AB models. This is surveyed in Grazzini and Richiardi (2015), who also show how to apply simulated minimum distance (SMD) techniques to estimate the parameters, following a frequentist approach. The method of simulated moments (MSM) and indirect inference (II), among other techniques, fall in this general class. Typically, certain summary statistics have to be selected in advance in order to implement the minimum-distance estimator, requiring additional sensitivity analysis to understand their properties. Good properties of the summary statistics ensure good properties of the estimator, e.g. consistency. Recently, Kukacka and Barunik (2016) have proposed non-parametric simulated maximum likelihood, and apply it to the estimation of the model described in Brock and Hommes (1998). Other contributions have introduced sophisticated validation methods aimed at measuring the distance between real and simulated data, where the simulated data can either come from the original AB model (Fabretti, 2013; Marks, 2013; Lamperti, 2015; Recchioni et al., 2015; Barde, 2016; Guerini and Moneta, 2016), or from a surrogate model (Salle and Yildizoglu, 2014; Sani et al., 2016). Minimisation of such distance metrics within a SMD approach leads to alternative estimators with respect to those commonly employed in the literature, such as MSM or II, although the properties of these estimators, including consistency, and their associated

¹For a plea for the AB approach, and a critique of the DSGE approach to macroeconomics, see Fagiolo and Roventini (2016).

uncertainty have yet to be fully assessed.²

While the literature on calibration and SMD estimation of AB models is moving fast, the Bayesian approach common for instance in the DSGE (dynamic stochastic general equilibrium) literature has so far received less attention. In this note, we show how Bayesian methods can be used to perform statistical inference in AB models. The advantages of Bayesian methods with respect to the frequentist and the calibration approaches are twofold: (i) they do not require to pre-select moments (MSM) or an auxiliary model (II), or other metrics (calibration) to evaluate the distance between the real and the simulated time series, and (ii) they allow to incorporate prior information, leading to a proper statistical treatment of the uncertainty of our knowledge, and how it is updated given the available observations. Moreover, with respect to methods that require the use of summary statistics (as MSM), the Bayesian approach fully exploits the informational content of the data, hence achieving, at least asymptotically, greater efficiency.³

On the other hand, the Bayesian approach can mask identification problems (Canova, 2008; Canova and Sala, 2009) by artificially adding curvature to the posterior with appropriately selected (and often poorly justified) priors, in presence of a flat likelihood. This is a malpractice that afflicts DSGE models (Fagiolo and Roventini, 2012, 2016), but that should not let us jump to the conclusion to “throw the baby out with the bathwater”. A more serious potential disadvantage of Bayesian methods, with respect to calibration techniques, is the computational burden of estimating the likelihood function by simulation. To save on these computational costs, in addition to using efficient sampling schemes in models with large parameters’ space, several approximations can be introduced, whose appropriateness should be evaluated on a case-by-case basis. These approximations might also involve giving up (i), and resorting again to make inference based on the informational content of (generally insufficient) summary statistics, an appropriate choice of which can also result in more robust estimation.

²As such, these techniques might be classified as advanced calibration. Kydland and Prescott (1996, p. 74) distinguish calibration from estimation by suggesting that calibration is concerned with data tracking (finding the values of the parameters that make the model behave as close as possible to the real data), while estimation refers to the inferential effort to learn about the underlying values of the parameters: “[D]ata are used to calibrate the model economy so that it mimics the world as close as possible along a limited, but clearly specified, number of dimensions. Note that calibration is not an attempt at assessing the size of something: it is not estimation. *Estimation* is the determination of the approximate quantity of something”. Other authors, however, have a less clear-cut understanding of the two terms. For instance, Hansen and Heckman (1996, p. 91) say that “[T]he distinction drawn between calibrating and estimating the parameters of a model is artificial at best. Moreover, the justification for what is called calibration is vague and confusing. In a profession that is already too segmented, the construction of such artificial distinctions is counterproductive.” This applies even more to a Bayesian approach to estimation, where no “true” values of the parameters are assumed to exist.

³The Bayesian estimator minimises the posterior expected loss with quadratic loss functions (mean squared error, MSE), where the expectation is taken over the posterior distributions of the parameters.

Bayesian methods are commonly employed for estimating DSGE models.⁴ However, two features of DSGE models make Bayesian estimation simpler: (i) they produce analytical expressions for the behaviour of the agents around the steady state, and (ii) they involve only a limited number of different agents, hence equations (e.g. textbook-version NK models have just three equations). Having analytical expressions for the steady state behaviour allows (log-) linearisation⁵ and the application of a simple Kalman filter to derive the likelihood and perform exact inference (on the approximated model). A limited number of equations implies that even if linearisation is not imposed, and a more complicated filter is used, simulation of the model is relatively fast. Still, the computational time required to repeatedly solve the model is an issue in DSGE modelling, as recognised by Fernández-Villaverde et al. (2016): “[C]losed-form solutions are the exception and typically not available for models used in serious empirical applications. [...] This ultimately leads to a trade-off: given a fixed amount of computational resources, the more time is spent on solving a model conditional on a particular θ , e.g., through the use of a sophisticated projection technique, the less often an estimation objective function can be evaluated. For this reason, much of the empirical work relies on first-order perturbation approximations of DSGE models, which can be obtained very quickly. The estimation of models solved with numerically sophisticated projection methods is relatively rare, because it requires a lot of computational resources.”

AB models almost never lead to closed form equations to describe the behaviour of the agents in the stationary state (if any);⁶ moreover, heterogeneity is such that each individual agent has to be simulated. Consequently, applications of Bayesian techniques requiring a high number of simulations of the model have so far been scant.⁷

Another approach, which we follow in this paper, is to save on the time to compute the likelihood, by approximating it.

The paper is structured as follows. Section 2 formalises AB models; section 3 describes

⁴See the recent review by Fernández-Villaverde et al. (2016).

⁵Linearisation however is not neutral: it eliminates asymmetries, threshold effects and many other interesting phenomena (Rubio-Ramirez and Fernández-Villaverde, 2005).

⁶An exception is the class of models considered in Alfarano et al. (2005, 2006, 2007).

⁷Ward et al. (2016) discuss the use of the the ensemble Kalman filter (EnKF) and present an application to a model of individual mobility. EnKF (Evensen, 2009) is a particular class of particle filters (see section 3.5), and is widely used in weather forecasting, where the models are of extremely high order and non-linear, the initial states are highly uncertain, and a large number of measurements are available. EnKF makes the assumption that all probability distributions involved are Gaussian, though non-Gaussian processes can be modelled via Gaussian mixtures (Stordal et al., 2011).

The terms of the trade-off described by Fernández-Villaverde et al. (2016) between misspecification due to over simplistic assumptions and poor identification due to insufficient sampling of the parameters space can be altered by reducing the time required to simulate the model and obtain a measure of fit with the observations.

the basic Bayesian techniques that can be used to make inference in AB models; section 4 offers a comparison of the different techniques in the simple model of price learning used in Grazzini and Richiardi (2015), with one parameter; section 5 applies the most general of the estimation strategies considered here to the macro model of De Grauwe (2012), with nine estimated parameters; section 6 offers our concluding remarks.

2 AB models

An Agent-Based (AB) model is a Markov chain where the state of the system at time t is given by the collection of all micro-states at time t $\mathbf{X}_t \equiv \{\mathbf{x}_{it}\}$, $i = 1 \dots N$, $t = 1 \dots T$, where the vector \mathbf{x}_{it} represents the micro-state of agent i in time period t . The evolution of each agent can be written as

$$\mathbf{x}_{i,t+1} = \mathbf{f}_i(\mathbf{X}_t, \boldsymbol{\Xi}_t, \boldsymbol{\theta}) \quad (1)$$

where \mathbf{f}_i is a function taking values in \mathbb{R}^k , $\boldsymbol{\Xi}_t \equiv \{\xi_{it}\}$ is a vector of stochastic elements and $\boldsymbol{\theta} \in \Theta$ is a parameter vector, with Θ being a compact subset of \mathbb{R}^Q . We assume for simplicity that the model is ergodic, that is, the effects of the random draws $\boldsymbol{\Xi}_t$ fade away with time.⁸

Equation (1) allows us to identify the main differences with respect to DSGE models: the functions \mathbf{f} are typically complicated, possibly heterogeneous and involving discontinuities, if-else statements, etc.; even when they are simple, there can be many of them (one for each agent); no equilibrium can be defined in terms of consistency of individual choices.⁹

A set of K aggregate statistics $\mathbf{y}_t \equiv \{y_{kt}\}$, $k = 1 \dots K$ can then be defined over \mathbf{X}_t :

$$\mathbf{y}_t = \mathbf{m}(\mathbf{X}_t). \quad (2)$$

where \mathbf{m} is a function transforming the collection of all micro-states \mathbf{X}_t in the aggregate statistics \mathbf{y}_t . Eqs. (1)-(2) together give

$$\mathbf{y}_{t+1} = \mathbf{g}(\mathbf{X}_t, \boldsymbol{\Xi}_t, \boldsymbol{\theta}) \quad (3)$$

One strand of the literature has looked at the use of Bayesian optimisation techniques to build emulators of AB models: the emulators then replace the AB models in a SMD framework. An application to demographic research is provided by Bijak et al. (2013).

⁸The theory of estimation of non-ergodic AB models is still under-explored, and it is the object of our current research. Tests for non-ergodicity of the time series produced by an AB model are developed in Grazzini (2012).

⁹In rational expectation models, equilibrium is defined as a consistency condition in the behavioural equations: agents (whether representative or not) must act consistently with their expectations, and the actions of all the agents must be mutually consistent. The system is therefore always in equilibrium, even during a phase of adjustment after a shock. By converse, equilibria in AB models are defined only at the aggregate level and only in statistical terms, when macro-observables (eq. 2) become stationary (Grazzini and Richiardi, 2015).

where \mathbf{X}_t is predetermined. If the model is stationary, a long-run statistical equilibrium –called ‘absorbing’ in Grazzini and Richiardi (2015)– is reached after T^* periods, where the state of the system is independent of the initial conditions \mathbf{X}_0 and the seed s which governs the random disturbances Ξ :¹⁰

$$\mathbf{y}^* = E[\mathbf{y}_t | t > T^*] = \mathbf{g}^*(\boldsymbol{\theta}). \quad (4)$$

Typically, we observe data

$$\mathbf{y}_{t+1}^R = \mathbf{g}^R(\mathbf{X}_t^R, \Xi_t, \boldsymbol{\theta}^R, \mathbf{u}_t) \quad (5)$$

and, in equilibrium,

$$\mathbf{y}_{t+1}^R = \mathbf{g}^{*R}(\boldsymbol{\theta}^R, \mathbf{u}_t) \quad (6)$$

where \mathbf{u}_t is a vector of disturbances (accounting for measurement errors, specification errors, etc.), and the superscript R refers to the observed data. Note that, in general, the micro states \mathbf{X}_t^R are not observable, and only aggregate data \mathbf{y}_t^R might be available.

SMD techniques work by comparing theoretical constructs computed over \mathbf{y}_t , $\boldsymbol{\mu}(\mathbf{y}_t(\boldsymbol{\theta}))$, which depend on the structural parameters $\boldsymbol{\theta}$, with their observed counterparts $\boldsymbol{\mu}^R$, computed on \mathbf{y}^R : a value of $\boldsymbol{\theta}$ is selected in order to minimise the distance between the theoretical and observed quantities. Because no closed form expression for the theoretical quantities can be found, they are estimated by simulation, in the artificial data produced by the model. The method of simulated moments (MSM), where the model is summarised by longitudinal moments¹¹, and indirect inference (II), where the model is summarised by the estimated coefficients of an auxiliary model, both fall into this class of estimators, as does simulated maximum likelihood (SML), where the model is summarised by the probability of reproducing the raw observed data. With respect to SML, standard Bayesian methods use the likelihood to derive a posterior distribution of the parameters, as described in the next section.

¹⁰Non-ergodicity implies the existence of multiple equilibria. In a Bayesian framework, if a model is non-ergodic the simulated data should come from different replications of the model with different random seeds, rather than from one (longer) simulation run. By converse, if the model is ergodic fixing the random seed allows a more precise estimation, as the changes in model behaviour when the parameters are changed are smoother.

Note that \mathbf{y} can be any transformation of the data. As an extreme example, in chaotic systems they may be properties of the orbits / attractors.

¹¹that is, moments computed over time

3 Methods

The fundamental equation for Bayesian methods is a simple application of the Bayes theorem:

$$p(\boldsymbol{\theta}|\mathbf{Y}^R) \propto \mathcal{L}(\boldsymbol{\theta}; \mathbf{Y}^R) p(\boldsymbol{\theta}) \quad (7)$$

where $p(\boldsymbol{\theta})$ is the prior distribution of the parameters, $\mathcal{L}(\boldsymbol{\theta}; \mathbf{Y}^R) \equiv p(\mathbf{Y}^R|\boldsymbol{\theta})$ is the likelihood of observing the data $\mathbf{Y}^R \equiv \{\mathbf{y}_t^R\}$, $t = 1 \dots T$ given the value of the parameters, and $p(\boldsymbol{\theta}|\mathbf{Y}^R)$ is the posterior distribution, that is the updated distribution once the information coming from the observed data is properly considered. The prior distribution typically comes from other studies or subjective evaluations. A uniform distribution in the allowed range of the parameters is often used as a way to introduce ‘uninformative’ priors, though not such a thing as an uninformative prior actually exists (Bernardo, 1997).¹² What matters, the prior is a distribution, which through application of Bayes theorem produces another distribution as an output (by converse, maximising the likelihood would only produce a point estimate).

Sampling the posterior distribution $p(\boldsymbol{\theta}|\mathbf{Y}^R)$ involves two computationally intensive steps: (i), for given values of $\boldsymbol{\theta}$, obtaining the likelihood \mathcal{L} , (ii) iterating over different values of $\boldsymbol{\theta}$. The resulting posterior is then normalized to have unit integral—for instance via a Simpson’s rule, a common method for numerical integration.

Sections 3.1 to 3.3 below deal with problem (i), while section 3.4 is devoted to likelihood-free methods and section 3.5 discusses problem (ii).

3.1 Non-parametric density estimation

Computation of the likelihood, for any given value of $\boldsymbol{\theta}$, is conceptually straightforward. Assuming we are in a statistical equilibrium, the probability of observing the whole (unordered) series of data $\mathbf{Y}^R \equiv \{\mathbf{y}_t^R\}$ is simply

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{Y}^R) \propto \prod_{t=1}^T f(\mathbf{y}_t^R|\boldsymbol{\theta}). \quad (8)$$

where we assume to be ‘blind’ with respect to all the other data points $\{\mathbf{y}_{-t}^R\}$ when we evaluate \mathbf{y}_t^R . Any autocorrelation in $\{\mathbf{y}_t\}$ would then lead to an increase in the variance of the estimated

¹²Sometimes a maximum entropy distribution is used as a ‘minimally informative’ prior; this however requires (a) some level of information on some moments of the prior to specify the constraints and (b) the choice of a reference measure in continuous settings.

distribution of \mathbf{y}_t , and in turn to an increase in the variance of $\mathcal{L}(\boldsymbol{\theta}; \mathbf{Y}^R)$.¹³

All we need is an estimate for the distribution $\tilde{f}(\boldsymbol{\theta})$; we then evaluate the estimated distribution $\tilde{f}(\boldsymbol{\theta})$ at each observed \mathbf{y}_t^R , and compute $\prod \tilde{f}(\mathbf{y}_t^R | \boldsymbol{\theta})$.

Estimation of the density, for any value of $\boldsymbol{\theta}$, is done by simulation: in a statistical equilibrium, the outcome fluctuates around a stationary level $\mathbf{y}^*(\boldsymbol{\theta}) = E[\mathbf{y}_t(\boldsymbol{\theta}) | t > \bar{T}]$. If we collect the artificial data produced by the model in such a statistical equilibrium, we can construct a probability distribution around \mathbf{y}^* , and therefore evaluate the density at each observed data point \mathbf{y}_t^R . If the outcomes $\mathbf{y}(\boldsymbol{\theta})$ were discrete, we would only have to count the frequency of occurrence of each observed value \mathbf{y}_t^R . With continuous $\mathbf{y}(\boldsymbol{\theta})$, the likelihood has to be estimated either non-parametrically or parametrically, under appropriate distributional assumptions. A traditional *non-parametric* method is kernel density estimation (KDE), which basically produces histogram-smoothing: the artificial data are grouped in bins (the histogram), and then a weighted moving average of the frequency of each bin is computed.¹⁴ The approximation bias introduced by KDE can be reduced by using a large number of very small bins, but then the variance in the estimate of the density grows.

3.2 Parametric estimation of the likelihood

The main problem with KDE is its computational cost (see section 4.2 below). A (faster) alternative is to assume a *parametric* distribution for the density around $\mathbf{y}^*(\boldsymbol{\theta})$:

$$\mathbf{y}_{t+1} = \mathbf{g}^*(\boldsymbol{\theta}) + \boldsymbol{\epsilon}_t. \quad (10)$$

Imposing additional information about the distribution of the simulation output can help generate better estimates from a limited number of simulation runs. On the other hand, those estimates may increase the bias if the assumed distribution does not conform to the true density

¹³As an example, think of an AR(1) process with mean 0: $y_{t+1} = \rho y_t + \varepsilon_t$, where ρ is the autocorrelation coefficient and $\varepsilon_t \sim \mathcal{N}(0, \sigma_\varepsilon^2)$. The unconditional distribution of y_t is $y_t \sim \mathcal{N}\left(0, \frac{\sigma_\varepsilon^2}{1-\rho^2}\right)$: an increase in the autocorrelation coefficient ρ increases the variance of the distribution.

¹⁴More formally, kernel density estimation (KDE), given a simulated time series $\mathbf{y}(\boldsymbol{\theta}) \equiv \{\mathbf{y}_s\}$, $s = 1 \dots S$, approximates the density $f(\mathbf{y}_t^R, \boldsymbol{\theta})$, for each observed data point \mathbf{y}_t^R , with:

$$\tilde{f}(\mathbf{y}_t^R | \boldsymbol{\theta}) = \frac{1}{Sh} \sum_{s=1}^S \mathcal{K} \left\{ \frac{\sqrt{\sum_{k=1}^K \sum_{s=1}^S (y_{ks}^S - y_{kt}^R)^2}}{h} \right\} \quad (9)$$

where \mathcal{K} is a kernel function that places greater weight on values y_{ks} that are closer to y_{kt}^R , is symmetric around zero and integrates to one, and h is the bandwidth. Note that in eq.(9) we assumed, for simplicity, that the multi-variate density is symmetric, and h is the same for each dimension k (intuitively, we are assuming that each dimension k has the same marginal variance). Algorithms for KDE are available in most statistical packages.

of the model output. Such an assumption can of course be tested in the artificial data, in the relevant range of θ (i.e. where the estimated coefficients lie). Use of parametric methods leads to *synthetic likelihood* or *pseudo-likelihood* estimation (Wood, 2010; Hartig et al., 2011).

The parametric and non-parametric methods discussed above use the output variability that is predicted by the model, and use this information for inference. However, it is possible that the model shows much less variability than the data. This can be due to fundamental *specification errors* (the model is only a poor approximation of the real process, so that the mean model predictions do not fit to the data), *simplification errors* (the model is a good approximation of the real process, but there are additional stochastic processes that have acted on the data and are not included in the model) or *measurement errors* (there is uncertainty about the data). While the first type of errors calls for a re-specification of the model, the latter types could in principle be dealt with by including additional processes that explain this variability within the model. “However, particularly, when those processes are not of interest for the scientific question asked, it is simpler and more parsimonious to express this unexplained variability outside the stochastic simulation. One way to do this is adding an *external error model* with a tractable likelihood on top of the results of the stochastic simulation” (Hartig et al., 2011):

$$\mathbf{y}_{t+1}^R = \mathbf{g}^*(\boldsymbol{\theta}^R) + \mathbf{u}_t. \quad (11)$$

The validity of such a strategy depends on the quality of the assumption about the distribution of these external errors, given the model and the data. This assumption cannot be tested *per se*, as the variability in the real data comes both from the *explained* components (i.e. the model) and from the *unexplained* ones (the external errors), the external errors being defined as a residual.

The two parametric strategies (modelling the variability of model outcome and modelling external errors that might affect the real data, in the stationary state) are often explored separately. For instance, most studies that use the augmentation by external errors approach then treat the model outcome as deterministic, in the stationary state, whilst most studies that employ a synthetic likelihood do not consider external errors.¹⁵

Under the assumption that the distribution for $\boldsymbol{\epsilon}$ or \mathbf{u} is Gaussian, the likelihood is also Gaussian, with a mean which depends on $\boldsymbol{\theta}$. We are then able to write down an expression for

¹⁵It is in principle possible to combine the two approaches together, and explicit distributional assumptions for both the model outcomes and the external errors, though there are no practical advantages.

the likelihood –hence the posterior distribution, under a convenient specification of the priors (see the Appendix). However, the likelihood needs in any case to be simulated, as $\mathbf{g}^*(\boldsymbol{\theta})$ has no closed form expression. Therefore, the assumption of normality is easy to be replaced, if other distributions appear to fit the data better than the normal, or there are theoretical reasons to believe that the external errors are non-Gaussian.

3.3 Unconditional vs. conditional estimation

An important difference between estimation of AB models and estimation of models with exact aggregation results (as in DSGE models) is that in the latter case the evolution of the aggregate state of the system, \mathbf{y}_{t+1} , depends only on the current macro-state \mathbf{y}_t and the parameters $\boldsymbol{\theta}$, with the set of micro states \mathbf{X}_t not providing any further information. In terms of eq. (3), we have

$$\mathbf{y}_{t+1} = \mathbf{g}(\mathbf{y}_t, \boldsymbol{\xi}_t, \boldsymbol{\theta}) \quad (3')$$

This can be restated in more general terms by saying that the aggregate representation of the system, which is a projection of the original Markov chain where the state of the system at time t is given by the collection of all the micro-states, is still a Markov chain –a condition called *lumpability* (Kemeny and Snell, 1976).¹⁶ In such a case a better estimation of the likelihood can be obtained.

Let $f(\mathbf{y}_{t+1}|\mathbf{y}_t, \boldsymbol{\theta})$ denote the distribution of \mathbf{y} at time $t + 1$ given the observables at t ; the likelihood function can be expressed as

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{Y}^R) = f(\mathbf{y}_0^R|\boldsymbol{\theta}) \prod_{t=0}^{T-1} f(\mathbf{y}_{t+1}^R|\mathbf{y}_t^R, \boldsymbol{\theta}). \quad (8')$$

Assuming $f(\mathbf{y}_0^R|\boldsymbol{\theta}) \propto \text{const.}$ or that the initial state of the system is known with certainty, it is possible to initialise at any time step t the simulations with the observed values \mathbf{y}_t^R , run D different one-period ahead simulations (rather than 1 ‘long’ simulation lasting D periods) and estimate the conditional densities using $\tilde{f}(\mathbf{y}_{t+1}^R|\mathbf{y}_t^R, \boldsymbol{\theta})$ instead of $\tilde{f}(\mathbf{y}_t^R|\boldsymbol{\theta})$ in (8). The scheme is more efficient because the simulated time series are constrained to remain closer to the observed ones, hence the likelihood is estimated with more precision.¹⁷ It also allows to relax the weakly

¹⁶There are necessary and sufficient conditions for lumpability, and they have to do with symmetries in the micro-state space (Banish et al., 2012).

¹⁷To continue with our AR(1) example, the conditional distribution of y_{t+1} given y_t and ρ is $y_{t+1} \sim \mathcal{N}(\rho y_t, \sigma_\varepsilon^2)$.

stationarity assumption, as conditioning on \mathbf{y}_t^R controls for persistence in the process.¹⁸

Note that eq. (8') holds generally. However, when the system is non lumpable the precise way all the individual states \mathbf{X}_t are re-initialised, in order to simulate the system, does matter. Moreover, there might be multiple combinations of micro states \mathbf{X}_t which give the same macro state \mathbf{y}_t^R , but different evolution \mathbf{y}_{t+1} .¹⁹ In general, if eq. 3' does not hold (and in AB models it typically does not hold, otherwise it would be more convenient to use a representative agent formulation), and the macro-state \mathbf{y}_{t+1} depends on the *distribution* of the unobserved micro-states \mathbf{X}_t , the unconditional approach is the only one feasible.

3.4 Likelihood-free methods

As we have seen, obtaining a non-parametric estimate of the likelihood can be computationally heavy. Turning to parametric estimates, under the assumption of a fixed distributional form of the variable of interest around a long-term stationary state predicted by the model —where the variability is produced either by model uncertainty or external errors— can sometimes be too restrictive. Originating from population genetics (Tavaré et al., 1997; Fu and Li, 1997), where the task of estimating the likelihood of the observed changes in DNA is impervious, a new set of methods have appeared in the last fifteen years to produce approximations of the posterior distributions without relying on the likelihood. These methods are labelled ‘likelihood-free’ methods, and the best known class is approximate Bayesian computation (ABC).²⁰

In standard Bayesian methods, it is the likelihood function that provides the fit of the model with the data —describing how plausible a particular parameter set θ is. The likelihood is however often computationally impractical to evaluate. The basic idea of ABC is to replace the evaluation of the likelihood with a 0-1 indicator, describing whether the model outcome is close enough to the observed data. To allow such an assessment, the model outcome and the data must first be summarised. Then, a distance between the simulated and the real data is computed. The model is assumed to be close enough to the data if the distance falls within the admitted tolerance. As such, there are three key ingredients in ABC: (i) the selection of *summary statistics*, (ii) the definition of a *distance measure*, (iii) the definition of a *tolerance threshold*. The choice of a distance measure is usually the least controversial point (the Euclidean distance

¹⁸If persistence is of order higher than 1, conditioning on periods previous to t is required.

¹⁹The latter would not be a problem if only we could be sure to draw the micro states \mathbf{X}_t randomly from the set of micro states which aggregate to \mathbf{y}_t^R , but in practice we can never be sure that this is the case.

²⁰See Marin et al. (2011); Turner and Zandt (2012). An application to an AB model of cancer development can be found in Sottoriva and Tavaré (2010).

or weighted Euclidean distance, where the weights are given by the inverse of the standard deviation of each summary statistics, is generally used). The choice of a tolerance threshold, as we shall see, determines the trade-off between *sampling error* and *approximation error*, given computing time. The choice of summary statistics is the most challenging, and we will discuss it in greater details later in the section. Note that, by performing inference based on summary statistics, ABC resembles SMD. Indeed, ABC can be seen as SMD embedded in a Bayesian framework, which allows incorporation of the priors and quantification of uncertainty by means of posterior distributions of the parameters. As in SMD, moments selection is left to the intuition of the researcher, and then validated by sensitivity analysis on the model behaviour; however, in complicated models finding the right moments could be a serious challenge. On the other hand, the need to analyse the behaviour of the summary statistics prior to estimation prompts a better understanding of the model behaviour, and can lead to more robust estimation.²¹

The basic ABC algorithm works as follows:

1. a candidate vector $\boldsymbol{\theta}^c$ is drawn from a prior distribution;
2. a simulation is run with parameters vector $\boldsymbol{\theta}^c$, obtaining simulated data from the model density $p(\mathbf{y}|\boldsymbol{\theta}^c)$;
3. the candidate vector is either retained or dismissed depending whether the distance between the summary statistics computed on the artificial data $\boldsymbol{\mu}(\mathbf{y}(\boldsymbol{\theta}))$ and summary statistics computed on the real data $\boldsymbol{\mu}(\mathbf{y}^R)$ is within or outside the admitted tolerance h : $d(\boldsymbol{\mu}, \boldsymbol{\mu}^R) \leq h$.

This is repeated N times; the retained values of the parameters define an empirical approximated posterior distribution. KDE can then be applied to smooth out the resulting histogram, and obtain an estimate of the theoretical approximated posterior. Approximation error refers to the fact that the posterior is approximated; sampling error refers to the fact that we learn about the approximated posterior from a limited set of data. It is easy to see where the approximation error comes from. While the true posterior distribution is $p(\boldsymbol{\theta}|\mathbf{y} = \mathbf{y}^R)$, in ABC we get $p(\boldsymbol{\theta}|\boldsymbol{\mu}(\mathbf{y})) \approx \boldsymbol{\mu}(\mathbf{y}^R)$.

If we set the tolerance threshold $h \rightarrow 0$, and our statistics were *sufficient summary statis-*

²¹This has also been noted in DSGE modelling. Ruge-Murcia (2007) compares SMD to Bayesian techniques and concludes that moment-based estimation methods compare very favourably to the more widely used likelihood-based methods, being more robust to misspecification and less affected by stochastic singularity.

*tics*²², we would get back to standard Bayesian inference, and sample from the exact posterior distribution. However —and this is the whole point— because of the complexity of the underlying model the likelihood of observing the real data is tiny everywhere, so that acceptances are impossible, or at least very rare. When h is too small, the distribution of accepted values of θ is closer to the true posterior, and the approximation error is smaller; however, the number of acceptances is usually too small to obtain a precise estimate of the (approximated) posterior distribution: the sampling error increases. On the other hand, when h is too large, the precision of the estimate improves because we have more accepted values (the sampling error goes down), but the approximation error gets bigger. In other words, we obtain a better estimate of a worse object.

An alternative to choosing h in advance is to specify the number of acceptances n required (e.g. $n = 500$): then, h is chosen (after the distance for every draw is computed) in order to achieve that number of acceptances. Finally, note that the trade-off between sampling error and approximation error is for a given number of draws (hence, a given computing time). Drawing more candidates allows to reduce the approximation error (by decreasing h) without increasing the sampling error. Stated more formally, ABC converges to the true posterior as $h \rightarrow 0, N \rightarrow \infty$. In particular, it can be shown that the bias of the ABC estimate is asymptotically proportional to h^2 , as $h \rightarrow 0$. On the other hand, the computational cost for generating n accepted ABC samples is proportional to nh^{-q} , where q is the dimension of the observations (Barber et al., 2015). The same authors show that the mean squared error (MSE) of the estimates for constant expected computational costs is given by

$$\text{MSE}(h) = ah^{-2} + bh^4 \tag{12}$$

for some constants a and b . Minimisation of the MSE gives the optimal level of tolerance, given a fixed amount of computational resources.

The choice of summary statistics is at the same time the weak point of ABC and a great source of flexibility. For instance, by choosing as summary statistics longitudinal moments, or the coefficients of an appropriate auxiliary model, it allows to embed the method of simulated moments and indirect inference in a Bayesian setting, incorporating prior information. Also,

²²A summary statistics is said to be sufficient if “no other statistic that can be calculated from the same sample provides any additional information as to the value of the parameter” (Fisher, 1922). Sufficient statistics satisfy $p(\theta|\mu(\mathbf{y}^R)) = p(\theta|\mathbf{y}^R)$.

an appropriate choice of the summary statistics allows to make *conditional forecasts* about the evolution of the real world. Suppose an extreme case where we only condition on the state of the system at time t : we wish to project the likely evolution of a system given \mathbf{y}_t . We can then simply set our summary statistics $\boldsymbol{\mu}(\mathbf{y}) = \mathbf{y}_t^R$: the ABC algorithm will retain any simulated trajectory that passes for \mathbf{y}_t^R , producing not only a (quite poor, in this case) approximation of the posterior, but also conditional projections about future states.²³

Any condition can in principle be used as a summary statistics: of course, the lower the informational content of the condition, the poorer the approximation. However, there is also a drawback in increasing the informational content of the summary statistics, and it comes again from the trade-off between sampling error and approximation error. As Beaumont et al. (2002, p. 2026) put it, “A crucial limitation of the [...] method is that only a small number of summary statistics can usually be handled. Otherwise, either acceptance rates become prohibitively low or the tolerance [...] must be increased, which can distort the approximation.” This is because the asymptotic rate of converge of ABC to the true posterior distribution, as $h \rightarrow 0, N \rightarrow \infty$, worsens with $\dim(\boldsymbol{\mu})$. The problem of choosing appropriate *low dimensional* summary statistics that are informative about $\boldsymbol{\theta}$ is an open issue in ABC. “The insidious issue is that it is rarely possible to verify either sufficiency or insufficiency. Furthermore, if they are insufficient, it is usually not possible to determine how badly they have distorted results. Said another way, you know you are probably making errors, but you don’t know how large they are” (Holmes, 2015).

The topic is an active area of research. Recent years have seen the development of techniques that provide guidance in the selection of the summary statistics (see e.g. Fearnhead and Prangle, 2012). Also, post-processing of the results can improve the quality of the approximation by correcting the distribution of $\boldsymbol{\theta}$ by the difference between the observed and simulated summary statistics (Beaumont et al., 2002). Efficiency can be improved by assigning a *continuation probability* to each simulation: the idea is to stop prematurely simulations that are likely to end up in a rejection, and has originated the *lazy ABC* approach (Prangle, 2014).

3.5 Sampling from the posterior distribution

Application of the Bayes theorem, once the likelihood is known, allows to get a density for the posterior distribution, at *one* given value of $\boldsymbol{\theta}$. However, to recover the whole shape of

²³The case when it is possible to kill two birds (inference and conditional statements) with one stone is of course quite a lucky one. More in general, when the condition in the conditional statement is too poor to allow for good inference, we should keep the two problems separate: *first*, get a good approximation of the posterior (by

the posterior distribution, many values have to be sampled. In simple models, exploration of the parameters space can be accomplished by ‘brute force’ grid exploration: the parameters space is sampled at regular (small) intervals. For instance, if there are two parameters that can potentially vary continuously between 0 and 1, and we set the value of the step to .1, we have 11 values to consider for each parameter, and their combination gives 121 points to sample: by discretising the parameters, we have reduced the size of the parameters space from \mathbb{R}^2 to 121 points. Multi-level grid search, where the grid is explored at smaller intervals in ranges of the parameters’ space on the bases of the results of previous, looser, grid explorations, can be devised to improve efficiency. However, the curse of dimensionality –the fact that when the dimensionality increases, the volume of the parameters’ space increases so fast that sampling becomes sparse– precludes adopting such an approach except when the number of the parameters is small. Grid exploration involves evaluating the density of the posterior distribution at many points where it is practically zero, while more likely values of θ , where a finer search might be valuable, are sampled with the same probability. Efficient sampling involves devising algorithms where the sampling probability, rather than being constant, is proportional to the posterior density.

There are four main classes of *efficient sampling schemes*, to obtain samples from a function of θ , the *target distribution* (the posterior, in our case), which is unknown analytically but can be evaluated point-wise for each θ : (i) rejection sampling, (ii) importance sampling, (iii) Sequential Monte Carlo, and (iv) Markov chain Monte Carlo.²⁴ Here we provide only an intuition of how they work, drawing extensively from the excellent survey by Hartig et al. (2011).²⁵

Rejection sampling (RS).* The simplest possibility of generating a distribution that approximates $\mathcal{L}(\theta)$ is to sample random parameters θ and accept those proportionally to their (point-wise approximated) value of $\mathcal{L}(\theta)$. This approach can be slightly improved by importance sampling or stratified sampling methods such as the Latin hypercube design, but rejection approaches encounter computational limitations when the dimensionality of the parameter space becomes larger than typically 10-15 parameters.

Importance sampling (IS). The intuition behind importance sampling is to study the

selecting appropriate summary statistics); *then*, sample from the estimated posterior and select the trajectories that fulfil the condition.

²⁴The literature on efficient sampling is huge —see for instance Kroese et al. (2011)— and many variants have been proposed to the basic algorithms, including the use of quasi-Monte Carlo methods.

²⁵The entries marked with a ‘*’ are excerpt from Hartig et al. (2011), where we have replaced ϕ in their notation with θ , in order to maintain consistency. The entries marked with a ‘**’ are based on Hartig et al. (2011), appropriately integrated.

distribution $\mathcal{L}(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\mathbf{y})$ while sampling from another, simpler distribution $q(\boldsymbol{\theta})$ (called *importance distribution*). This technique was born as a variance reduction technique, aimed at increasing the likelihood to sample from an ‘important’ but small region by sampling from a different distribution that overweights the important region (hence the name). Having over-sampled the important region, we have to adjust our estimate somehow to account for having sampled from this other distribution. This is done by re-weighting the sampled values by the adjustment factor $p(\boldsymbol{\theta})/(q(\boldsymbol{\theta}))$. Importance sampling and rejection sampling are similar in as much both distort a sample from one distribution in order to sample from another. They also share the limitation that they do not work well in high dimensions.

Sequential Monte Carlo methods (SMC).** Particle filters or sequential Monte Carlo methods (SMCs) work by filtering proposed values for $\boldsymbol{\theta}$ to arrive at a sample of values drawn from the desired distribution. In SMC each step of the algorithm contains N parameter combinations $\boldsymbol{\theta}_i$ (particles), that are assigned weights ω_i proportional to their likelihood or posterior value $\mathcal{L}(\boldsymbol{\theta}_i)$ (see Arulampalam et al., 2002). When starting with a random sample of parameters, many particles may be assigned close to zero weights, meaning that they carry little information for the inference (degeneracy). To avoid this, a resampling step is usually added where a new set of particles is created based on the current weight distribution. The traditional motivation for a particle filter is to include new data in each filter step, but the filter may also be used to work on a fixed dataset or to subsequently add independent subsets of the data.

The advantage of SMC –and of Markov chain Monte Carlo, see below– is that the time needed to obtain acceptable convergence is typically much shorter than for rejection sampling, because the sampling effort is concentrated in the areas of high likelihood or posterior density.

Markov chain Monte Carlo (MCMC).** MCMC sampling also try to concentrate the sampling effort in the areas of high likelihood or posterior density, based on previous samples. MCMC algorithms construct a Markov chain of parameter values $(\boldsymbol{\theta}_1 \dots \boldsymbol{\theta}_n)$, where the next parameter combination $\boldsymbol{\theta}_{i+1}$ is chosen by proposing a random move conditional on the last parameter combination $\boldsymbol{\theta}_i$, and accepted conditional on the ratio of $\mathcal{L}(\boldsymbol{\theta}_{i+1})/\mathcal{L}(\boldsymbol{\theta}_i)$. Given that certain conditions are met (see, e.g. Andrieu et al., 2003), the Markov chain of parameter values will eventually converge to the target distribution $\mathcal{L}(\boldsymbol{\theta})$.

There are a number of MCMC samplers, the most popular of which is the Metropolis-Hastings algorithm. In its simplest form, the *random-walk Metropolis-Hastings*, in each period a candidate $\boldsymbol{\theta}^c \sim \mathcal{N}(\boldsymbol{\theta}^{(s)}, \mathbf{V})$ is drawn, given the current value $\boldsymbol{\theta}^{(s)}$. The candidate is accepted

with probability

$$\alpha = \min \left\{ 1, \frac{p(\boldsymbol{\theta}^c | \mathbf{y}_R)}{p(\boldsymbol{\theta}^{(s)} | \mathbf{y}_R)} \right\} \quad (13)$$

in which case we set $\boldsymbol{\theta}^{(s+1)} = \boldsymbol{\theta}^c$; else, we set $\boldsymbol{\theta}^{(s+1)} = \boldsymbol{\theta}^{(s)}$ and we repeat the previous candidate.

With multiple parameters, a slightly more sophisticated version of the basic Metropolis-Hastings algorithm can be used, which resembles the Gibbs sampler and is most useful when choosing a particular covariance matrix \mathbf{V} is not easy. The Gibbs sampler relies on drawing a sample $\{\boldsymbol{\theta}^{(s)}, s = 1 \dots S\}$ by drawing repeatedly from the conditional distributions

$$\theta_j | \mathbf{Y}^R, \boldsymbol{\theta}_{-j}, j = 1 \dots k$$

When these posterior conditional distributions are not known families, we use the following.

Draw a candidate θ_j^c from the stochastic process

$$\theta_j^c = \theta_j^s + \eta_j \varepsilon \quad (14)$$

where ε is random noise, and η_j is a parameter governing the acceptance rate of the MCMC.

The candidate is accepted with probability

$$\alpha = \min \left\{ 1, \frac{p(\theta_j^c | \mathbf{Y}^R, \boldsymbol{\theta}_{-j}^{(s)})}{p(\theta_j^{(s)} | \mathbf{Y}^R, \boldsymbol{\theta}_{-j}^{(s)})} \right\} \quad (15)$$

in which case we set $\theta_j^{(s+1)} = \theta_j^c$, otherwise we repeat the previous accepted value, $\theta_j^{(s+1)} = \theta_j^{(s)}$. Each element of the parameter vector, $\theta_j \in \boldsymbol{\theta}$ is drawn sequentially, and the posterior probability associated to each new candidate is computed holding the other parameters at their last accepted values. The parameters η_j are adjusted so that approximately 25% of the candidates are accepted.

A final note concerns efficient sampling in an ABC setting. The standard scheme for ABC is, as we have seen, rejection sampling. Candidates are drawn from the prior distribution, and only those that ‘perform well’ are retained. This is not very efficient, especially if the prior distribution differs significantly from the posterior. However, it is possible to employ ABC with more efficient sampling schemes (see Sisson et al., 2016). For instance, rather than sampling from the prior one could sample from an importance distribution $q(\boldsymbol{\theta})$. Candidate are then

accepted if $d(\boldsymbol{\mu}(\boldsymbol{\theta}), \boldsymbol{\mu}^R) \leq h$, with a weight $p(\boldsymbol{\theta})/q(\boldsymbol{\theta})$. SMC methods can then be employed to adaptively refine both the threshold and the importance distribution. MCMC methods can also be employed, where new candidates depend on the current value of $\boldsymbol{\theta}$ and are accepted with a modified Metropolis-Hastings rule:

$$\alpha = \min \left\{ 1, \frac{p(\boldsymbol{\theta}^c) \mathbb{1}[d(\mathbf{S}(\boldsymbol{\theta}^c), \mathbf{S}^R) \leq h]}{p(\boldsymbol{\theta}^{(s)}) \mathbb{1}[d(\mathbf{S}(\boldsymbol{\theta}^{(s)}), \mathbf{S}^R) \leq h]} \right\} \quad (16)$$

where $p(\boldsymbol{\theta})$ is the prior density, and $\mathbb{1}[d(\mathbf{S}(\boldsymbol{\theta}^c), \mathbf{S}^R) \leq h]$ is an indicator whether the simulation outcome falls within the tolerance radius.

4 Comparison of different techniques on a simple model with 1 parameter

4.1 A price learning mechanism

As a testbed for the application of some of the Bayesian techniques described above, we first consider the stock market model proposed in Cliff and Bruten (1997) and used by Grazzini and Richiardi (2015). Participants in a stock market can trade by placing ask or bid limit orders in an order book. Traders do not know the demand and supply schedule but they can observe the order book; they cannot lend or borrow. The limit price of agent i in period t is

$$p_{i,t} = v_i(1 + \mu_{i,t}) \quad (17)$$

where v_i is the (constant) subjective value of the traded asset for agent i and $\mu_{i,t}$ is a profit margin (positive for sellers and negative for buyers). In each period traders look at the book and update their target price $\tau_{i,t}$: traders increase their target price if the last trade occurred at a high price, and lower it otherwise. This in turns determines a change in the limit price:

$$\Delta_{i,t} = \beta_i(\tau_{i,t} - p_{i,t}) \quad (18)$$

where β_i is a behavioural parameter determining how traders react to a difference between the target price and the current price $p_{i,t}$. The new limit price is:

$$p_{i,t+1} = p_{i,t} + \Delta_{i,t} \quad (19)$$

and the profit margin is coherently updated as follows:

$$\mu_{i,t+1} = \frac{p_{i,t} + \Delta_{i,t}}{v_i} - 1 \quad (20)$$

For simplicity, it is assumed that $\beta_i = \beta, \forall i = 1, \dots, N$. β is therefore the only parameter to be estimated in the model.

We perform Monte Carlo experiments where the pseudo-true data are generated, as in Grazzini and Richiardi (2015), with $\beta = 0.55$ and $N = 11$ buyers and sellers.²⁶ We adopt a uniform prior in $[0,1]$. We first estimate the model by KDE, with ‘brute force’ sampling of the parameters space. Second, we replace grid exploration with MCMC sampling. We then turn to parametric estimation of the likelihood, by considering augmentation of the model with Gaussian errors. Finally, we experiment with a simple rejection sampling ABC algorithm to embed the MSM strategy of Grazzini and Richiardi (2015) in a Bayesian framework.

4.2 Kernel density estimation

The model is simulated, for every value of the parameter β , for 150,000 periods, once the long-run stationary state has been reached.²⁷ We employ a Gaussian kernel, with optimal bandwidth $h = 1.06\hat{\sigma}S^{-0.2}$, where $\hat{\sigma}$ is the estimated standard deviation in the simulated sample of size S (Silverman, 1986, p. 45).²⁸ Figure 1 depicts the (normalized) posterior obtained for one series of pseudo-true data, computed on 100 equispaced values of β in the interval $[.01,1]$, for different numbers of simulated trading days.

The posterior shrinks as the number of observations increases, though it converges to a value slightly below 0.55. This is the effect of the bias introduced by KDE, and goes in an *a priori* unknown direction.²⁹ Experiments with different random seeds confirm the presence of a (downward) bias (figure 2).

²⁶Given the smooth behaviour of the model with respect to the parameter —see Grazzini and Richiardi (2015), in particular their figure 1 (p. 160)— estimation with other pseudo-true values of the parameters leads to similar performances.

²⁷To determine when the long-run stationary state is reached, we apply the Runs test described in Grazzini (2012). It turns out that discarding the first 500 periods of the simulation is enough to reach the stationary state, for any value of the parameter above $\beta = .01$. Given the pseudo-true value used in the exercise, slightly restricting the range of our prior from $[0,1]$ to $[.01,1]$ makes no harm. As an alternative one could estimate the time required to reach stationarity as a function of β and adopt a variable transient, or (ii) initialise the model at the equilibrium, hence getting rid of the transient altogether.

²⁸The optimal bandwidth, also known as *Silverman’s rule of thumb*, is the value of h that minimises the mean integrated squared error (Silverman, 1998, p. 45). Further improvements can be obtained by adopting a variable bandwidth, depending upon the location of either the estimate (balloon estimator) or the samples (pointwise estimator). This is known as adaptive or variable bandwidth kernel density estimation.

²⁹By converse, the small sample bias associated with MSM, arising from non-linearities of the moment func-

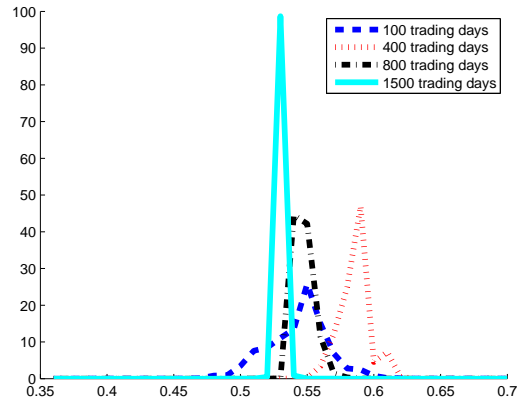


Figure 1: Posterior distribution of the learning parameter β with kernel density estimation and grid exploration of the parameter space (100 equispaced values in $[.01,1]$). Different numbers of trading days are assumed to be observed. The pseudo-true value is $\beta = 0.55$.

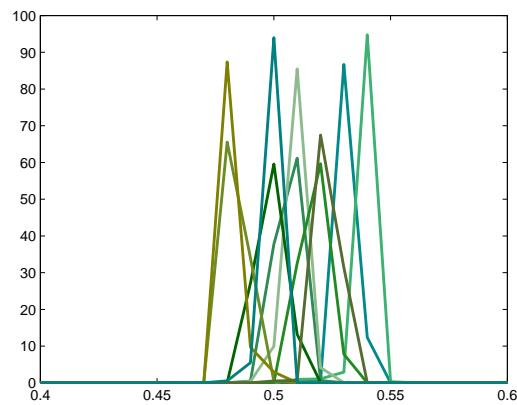


Figure 2: Posterior distribution of the learning parameter β , kernel density estimation, grid exploration with 100 equispaced values in $[.01,1]$, 1,500 observed periods, ten different random seeds. True value is $\beta = 0.55$.

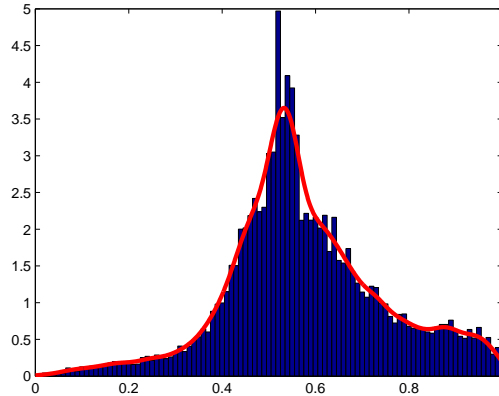


Figure 3: Posterior distribution of the learning parameter β with kernel density estimation of the likelihood and MCMC sampling. 40,000 simulations of 15,000 trading days each are performed. 1,500 pseudo-true trading days are assumed to be observed. The pseudo-true value is $\beta = 0.55$.

As for what concerns performance, it takes 7.2 secs to simulate 1,500 periods on our computer.³⁰ Performing KDE on the 1,500 simulated periods takes an additional 6.2 secs. Simulation time and KDE scale up linearly with the number of periods that are simulated: the time required to simulate the model for 150,000 periods and perform KDE is therefore $100 \times (7.2 + 6.2)$ secs. As grid exploration, for the chosen density of the grid, requires sampling 100 values of the parameter, total computing time is $(100 \times 100) \times (7.2 + 6.2)$ secs = 37.2 hours.

Grid exploration performs very well with a small number of parameters (only one, in our case), as no simulations are ‘wasted’. However, as the number of parameters increases, grid exploration becomes infeasible. We therefore experiment with a random walk Metropolis-Hastings MCMC sampling scheme. We perform 40,000 simulations, each lasting 15,000 periods, using 40 parallel processes each lasting 1,000 simulations, discarding the first 100 simulations as burn-in. For each simulation, KDE is performed. Figure 3 reports the posterior, using the same random seed as in the ‘brute force’ grid exploration.

The posterior is much more dispersed than under ‘brute force’ grid exploration; at the same time, computing time increases by a factor of 40 (from $100 \times 150,000$ periods to $40,000 \times 15,000$ periods), requiring parallelisation. The reason why MCMC sampling is much less efficient when grid exploration is feasible (i.e. with a small number of parameters) is that all simulations are used in grid exploration to produce the density in fig. 1, while most of the simulations end up

tions, is of predictable direction (see Grazzini et al., 2012). Andrieu and Roberts (2009) propose a ‘pseudo-marginal approach’ to stochastic simulation in a MCMC Bayesian framework with an unbiased estimator of the likelihood.

³⁰A Dell Precision R7910 with two 2.5 GHz Intel Xeon CPU E5-2680 v3 processors (each with 12 cores and 24 threads) and 128 GB of RAM, available at the Complexity Lab in Economics at the Department of Economics and Finance at the Catholic University of Milan.

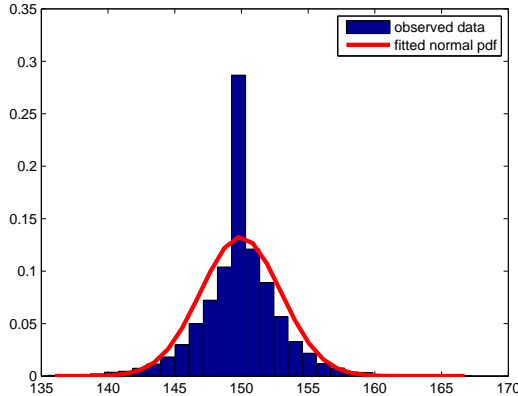


Figure 4: Price distribution in the stationary state, $\beta = 0.55$.

being discarded using MCMC sampling. On the other hand, the number of simulations required to achieve a specific accuracy grows much slower with the number of parameters than in grid exploration.

4.3 Augmentation with external errors

As we have seen, computational costs in Bayesian estimation of AB models come from three sources: (i) the time required for running one simulation, (ii) the time required for obtaining an estimate of the likelihood function, given the simulation outcome, and (iii) the number of times that steps (i) and (ii) have to be repeated. We assume that the coding of the model is already efficient, hence simulation time cannot be reduced.³¹ MCMC techniques, or other efficient sampling methods, are the only feasible options when the number of parameters is not small, but they perform poorly, with respect to ‘brute force’ grid exploration, in low-dimensional spaces: hence, gains have to be searched elsewhere.

We then go parametric and assume Gaussian noise, as described in section 3.2. This assumption finds some support in the simulated data, which however are much more concentrated around the theoretical equilibrium value of $p = 150$ (figure 4).

In spite of this discrepancy, computational experience with the new posterior suggests that it performs very well. In figure 5 we present the new posteriors using grid exploration, which quite understandably are smoother than those obtained under KDE, have only a slightly bigger variance, and are more centered on the true value of the parameter (smaller bias).

³¹An interesting approach, which is related to the normality assumption made in this section but is not further investigated here, is using Gaussian process emulators to approximate the outcomes of the simulation model without actually performing all the simulations (O’Hagan, 2006; Bijak et al., 2013).

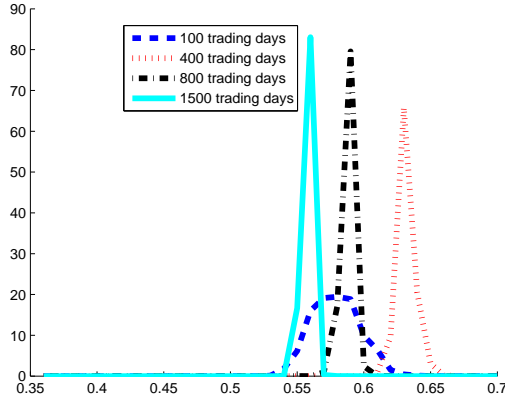


Figure 5: Posterior distribution of the learning parameter β with Gaussian density estimation and grid exploration of the parameter space (100 equispaced values in $[0,1]$). Different numbers of trading days are assumed to be observed. The pseudo-true value is $\beta = 0.55$.

Finally, because the posterior distribution of the learning parameter with MCMC sampling is larger than with grid sampling, the bias in KDE estimation is less noticeable (see figure 3). The posterior obtained assuming a Gaussian density is therefore very similar to the one obtained with KDE, under MCMC sampling, and for the sake of brevity is not reported here.

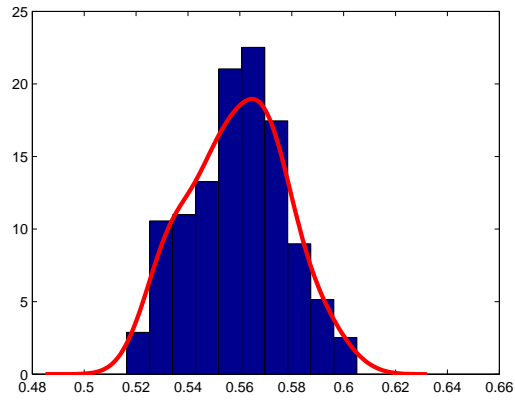
The gains in terms of computing time are, on the other hand, impressive: estimation under the normality assumption is practically instantaneous—it only takes 0.014 secs for simulation (1,500 periods)—compared to 6.2 secs per simulation for KDE. Given that KDE accounted for almost half of total computing times, both with grid exploration and with MCMC sampling, computing costs are therefore also almost halved.

4.4 Approximate Bayesian Computation

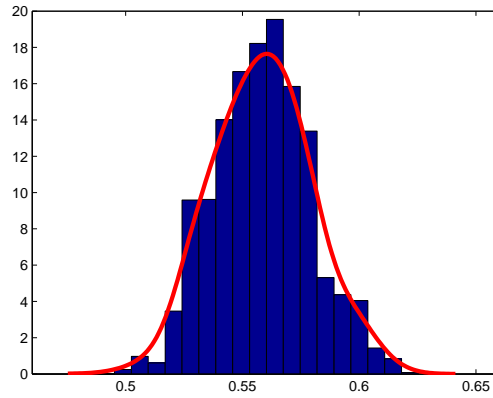
Our last exercise considers ABC techniques. Here, we use as our summary statistics the standard deviation of the price, that we know from Grazzini and Richiardi (2015) discriminates well between different values of the parameter. We apply the most basic rejection sampling ABC algorithm. Figure 6 reports the posterior distributions obtained with 400,000 simulations of 1,500 periods each (parallelised on 40 different cores), for different values of the tolerance threshold h .³²

The number of acceptances, for a fixed computational cost (400,000 runs \times 1,500 periods = 600 millions simulated periods), as theoretically predicted, is inversely proportional to h (the number of dimensions being $q = 1$).

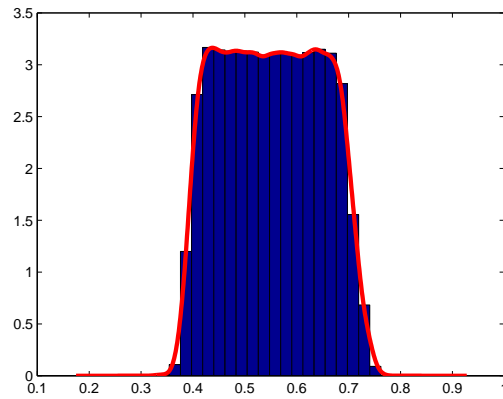
³²This is computationally equivalent to 40,000 simulations of 15,000 periods each, as in our MCMC experiments.



$$h = 0.01, n = 1,331$$



$$h = 0.1, n = 12,988$$



$$h = 1, n = 128,467$$

Figure 6: Posterior distribution of the learning parameter β with likelihood-free ABC estimation and rejection sampling. 400,000 simulations of 1,500 periods each are performed. 1,500 pseudo-true periods are assumed to be observed. Different values of the tolerance threshold are considered. n is the accepted sample size. The standard deviation of the price, used as summary statistics, ranges approximately from 0 to 5, and is equal to 3.0082 in the pseudo-observed series, obtained with $\beta = 0.55$.

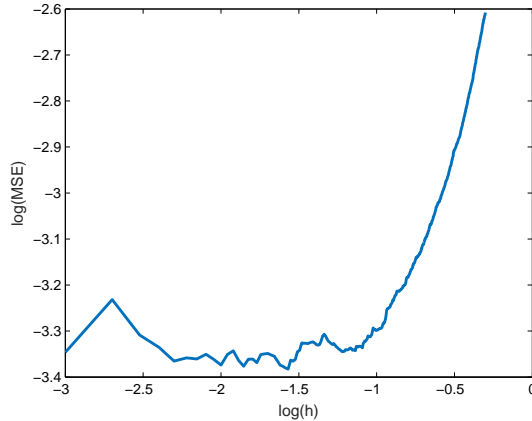


Figure 7: Sensitivity of the mean squared error (MSE) of the estimates with respect to h , log-log scale.

The trade-off between more accuracy but less precision, for small values of h , and less accuracy but more precision, for bigger values of h , is evident. Figure 7 depicts how the empirical mean squared error ($M\hat{S}E$) changes with h , on a log-log scale:

$$M\hat{S}E(\beta(h)) = \frac{\sum_{i=1}^n (\hat{\beta}_i(h) - \beta)^2}{n} = \text{Var}(\hat{\beta}(h)) + \text{bias}(\hat{\beta}(h))^2 \quad (21)$$

The estimated MSE is concave in h .³³ Given that the MSE starts increasing significantly only for values of $h > .1$, a good compromise between the number of ABC samples and the ability to shape a posterior distribution out of the Gaussian prior, in this application, is around this value.

4.5 Performance

Table 1 summarises the performance of the different techniques, in our sample model, together with a qualitative assessment. Grid exploration is much more efficient than other sampling schemes with only one parameter. Parametric approximation of the likelihood function further reduces computing time by almost 50%. ABC achieves good results, but the rejection sampling scheme implemented is highly inefficient, as it implies most runs of the model are discarded. MCMC has the poorest performance, though there are wide margins of improvement of the sampling algorithm.³⁴

³³The precision of the estimate decreases as h gets smaller, due to the reduced number of accepted ABC samples.

³⁴These are explored in the vast literature on MCMC. Given the illustrative purpose of the paper, we adopted a very simple scheme, with little optimisation.

Inferential procedure	Sampling scheme	Simulated periods (millions)	Time per simulated period	Qualitative assessment
Non-parametric KDE	Grid exploration	15	1.86	Very good precision, small bias
Parametric Gaussian	Grid exploration	15	1.00	Very good precision
Non-parametric KDE	MCMC	600	1.86	Poor precision
Parametric Gaussian	MCMC	600	1.00	Poor precision
ABC	Rejection sampling	600	1.00	Good precision

Table 1: Performance of different Bayesian techniques. Simulating one period (trading day) requires 0.0048 secs on our reference machine, normalised to 1 unit of time. Performing KDE requires an additional 0.0041 secs, hence adding 86% of computing time. Gaussian density estimation and ABC require practically no additional costs.

5 Estimation of a macro model with 9 parameters

Of all the techniques analysed in this paper, the non-parametric KDE estimation of the likelihood is the most general, as it does not rely neither on a functional approximation of the likelihood, nor on the choice of specific summary statistics. In this section we apply it to the estimation of the behavioural macroeconomic model of De Grauwe (2012). The model shares with large-scale macro agent-based models the necessity to recur to simulations in order to explore its properties. On the other hand, its simplicity makes it a good choice for an illustrative application. The main difference between small-scale and large-scale agent-based models is the computational time needed for each simulation run. This difference is therefore quantitative, rather than qualitative, and to some extent can be dealt with recourse to computer clusters or cloud computing.³⁵

5.1 The model

This is a standard New Keynesian model, where heterogeneous agents use simple heuristics to forecast the future, while selecting the forecasting rules that have delivered the highest performance in the past. The aggregate behaviour of the model is described by the three standard equations:

$$\text{Aggregate Demand: } y_t = a_1 \tilde{E}_t y_{t+1} + (1 - a_1) y_{t-1} + a_2 (r_t - \tilde{E}_t \pi_{t+1}) + \epsilon_t \quad (22)$$

$$\text{Aggregate Supply: } \pi_t = b_1 \tilde{E}_t \pi_{t+1} + (1 - b_1) \pi_{t-1} + b_2 y_t + \mu_t \quad (23)$$

$$\text{Taylor rule: } r_t = c_1 (\pi_t - \pi^*) + c_2 y_t + c_3 r_{t-1} + u_t \quad (24)$$

³⁵Identification can of course be even more of an issue in complex models spanning many dimensions; however,

where y is the output gap, r is the nominal interest rate, π is the inflation rate, $\pi^* = 0$ is the inflation target, and \tilde{E} is the expectations operator, where the tilde above E refers to expectations that are not formed rationally. Habit formation is usually invoked to justify the presence of lagged output in the aggregate demand equation; a Calvo pricing rule plus some indexation explains the presence of lagged inflation in the New Keynesian Phillips curve. De Grauwe considers two forecasting rules which agents can adopt, a fundamentalist rule, where agents believe that the output gap will revert to its equilibrium value of 0 in the next period, and an adaptive rule, where agents believe that the output gap in the next period will be the same as in the previous period:³⁶

$$\text{Fundamentalist: } \tilde{E}_t^f z_{t+1} = 0 \quad (25)$$

$$\text{Extrapolative: } \tilde{E}_t^e z_{t+1} = z_{t-1}. \quad (26)$$

where $z = \{y, \pi\}$. The model leads to a nice aggregation of expectations, where the market forecast is simply an average of the two forecasts weighted by the fraction of adopters of each rule:

$$\begin{aligned} \tilde{E}_t z_{t+1} &= (1 - p_{z,t}) \tilde{E}_t^f z_{t+1} + p_{z,t} \tilde{E}_t^e z_{t+1} \\ &= p_{z,t} \tilde{E}_t^e z_{t+1} \end{aligned} \quad (27)$$

where p_t is the probability that agents use an extrapolative rule.³⁷

Agents measure the fitness of each rule, $U_{z,t}^f$ and $U_{z,t}^e$, by its mean squared forecasting error (MSFE), with geometrically declining weights for past errors. The probability that an agent will adopt an extrapolative rule is then (Brock and Hommes, 1997)

$$p_{z,t} = \frac{\exp(\gamma U_{z,t}^f)}{\exp(\gamma U_{z,t}^f) + \exp(\gamma U_{z,t}^e)}. \quad (28)$$

The model leads to (i) intermittent effectiveness of inflation targeting, with one regime where the inflation rate remains close to the the inflation target and one regime where it fluctuates

the increased number of parameters is generally associated with a larger set of time series which can be exploited for estimation.

³⁶A similar model is developed in Deák et al. (2015), where individuals and firms use Individual Learning rather than Rational Expectation as a rationality benchmark. In Individual Learning, agents are rational regarding their individual decisions, but they have no macroeconomic model to form expectations of aggregate variables. The model is estimated in a way which is close to the Bayesian approach followed here.

³⁷In the original article the notation is slightly more elaborated.

wildly, and (ii) endogenous business cycles, originated by essentially unpredictable waves of optimism and pessimism. These “animal spirits” are responsible for the the non-normality of the distribution of the output gap (excess kurtosis and fat tails), which occurs even when the model is fed with normally distributed shocks.³⁸

5.2 Estimation

Because the model leads to exact aggregation, it can be simulated using eqs. (22)-(24) and (27)-(28), without having any agent actually doing anything. This permits using the conditional estimation strategy outlined in section 3.3; however, to retain generality we perform an unconditional estimation. The only difference with respect to a model where the aggregate behaviour needs to be computed on the basis of the realised micro-states is therefore, as already noted, a reduced computational time.

The model has 12 parameters (table 2), which in the original paper are calibrated (no details are given on the calibration procedure). For simplicity, we keep the parameters of the Taylor rule (c_1, c_2, c_3) fixed to their calibrated values, the rationale being that they can be separately estimated in the data. We estimate the remaining 9 parameters using KDE (section 3.1), with the MCMC sampling scheme described in section 3.5. Our MCMC implementation involves 6 parallel processes, each consisting of 270,000 draws, where the initial 27,000 draws are discarded (burn-out). For each draw, 5,000 periods are simulated.³⁹

We use data on real GDP and the implicit price deflator of the GDP for the U.S., 1947Q1 to 2016Q3.⁴⁰ The output gap is obtained by applying an HP filter to the log of the GDP series. Whenever possible, the priors are taken from Herbst and Schorfheide (2016). In particular,

³⁸De Grauwe (p. 490) explains the self-fulfilling mechanism which generates the business cycles as follows:

A series of random shocks creates the possibility that one of the two forecasting rules, say the extrapolating one, delivers a higher payoff, i.e. a lower mean squared forecast error (MSFE). This attracts agents that were using the fundamentalist rule. If the successful extrapolation happens to be a positive extrapolation, more agents will start extrapolating the positive output gap. The contagion-effect leads to an increasing use of the optimistic extrapolation of the output-gap, which in turn stimulates aggregate demand. Optimism is therefore self-fulfilling. A boom is created. At some point, negative stochastic shocks and/or the reaction of the central bank through the Taylor rule make a dent in the MSFE of the optimistic forecasts. Fundamentalist forecasts may become attractive again, but it is equally possible that pessimistic extrapolation becomes attractive and therefore fashionable again. The economy turns around.

³⁹Note that systematic grid exploration, which proved very effective in our one-parameter sample application, is infeasible here: sampling all the combinations of 100 different values for each parameter would require 100^9 (1×10^{18}) simulation runs. To give an example, estimating a price learning model with 9 parameters —for instance due to heterogeneity in the learning ability of the agents— would require 3.7×10^{16} hours with systematic grid sampling, under the optimistic assumption that the computational time for running both the model and the Kernel density estimation of the likelihood would remain the same as in the one-parameter case.

⁴⁰The GDPC1 and GDPEF series in the federal Reserve Bank of St. Louis FRED database.

Parameter	Default	Priors	Notes
π^*	0.0	(fixed)	central banks inflation target
a_1	0.5	U(0,3)	coeff. of expected output in AD eq.
$\tau = -1/a_2$	2.0	$\Gamma(2,.5)$	coeff. of real interest rate in AD eq.
b_1	0.5	U(0,1)	coeff. of expected inflation in AS eq.
b_2	0.5	U(0,1)	coeff. of output in AS eq.
c_1	2.0	(fixed)	coeff. of inflation in Taylor eq.
c_2	0.5	(fixed)	coeff. of output in Taylor eq.
c_3	0.5	(fixed)	interest smoothing parameter in Taylor eq.
$\gamma/10$	5.0	B(2,2)	intensity of choice parameter
σ_1	0.1	Inv- $\Gamma(.1,2)$	standard deviation of shocks in output
σ_2	0.1	Inv- $\Gamma(.1,2)$	standard deviation of shocks in inflation
σ_3	0.1	Inv- $\Gamma(.1,2)$	standard deviation of shocks in monetary policy
ρ	0.5	U(0,1)	speed of declining weights in MSFE (memory)

Table 2: Parameters, default values used for the Montecarlo experiment and priors in the behavioural macro model. The parameters of the Γ and Inv-*Gamma* distributions are the mean and standard deviation.

Herbst and Schorfheide use a Gamma distribution with mean 2 and standard deviation 0.5 (shape = 16, scale = .125) for $\tau = -1/a_2$, while they use a Uniform(0,1) for b_2 . Our prior for a_1 , b_1 and ρ are respectively U(0,3), U(0,1) and U(0,1), while for $\gamma/10$ we choose a Beta(2,2) distribution, implying that the support for γ is $[0,10]$, with mean 5 (table 2). For the standard deviation of the shocks, σ_1, σ_2 and σ_3 , we use an Inverse-Gamma with mean .1 and standard deviation 2 (Herbst and Schorfheide (2016) have different parameters of the Inverse-Gamma for the three shocks).

To check the properties of the estimation procedure, we first estimate the model on pseudo-true data, generated by fixing all the parameters to the default values reported in table 2.⁴¹ The length of the pseudo-true data is equal to 270 quarters. The parameters η_j in Eq. (14) are fixed so that the acceptance rate in the estimation on real data is approximately 25%. Figure 8 reports the posterior and the prior distributions for the parameters, for three different seeds for the simulated data.

The parameters a_1, b_2, ρ and γ seem to have a limited effect on the model behaviour: the posterior does not depart significantly from the prior, though for ρ and γ it is at least centred on the true values of the parameters. b_1 —the coefficient of expected inflation in the AS equation— is estimated with error: the posterior is different from the prior, but it is not centred on the true value. On the other hand $a_2(\tau)$ —the coefficient of the real interest rate in the AD equation—

⁴¹The default values differ from the calibrated values in De Grauwe (2012) when the calibrated values are far away from the mean of the priors distributions.

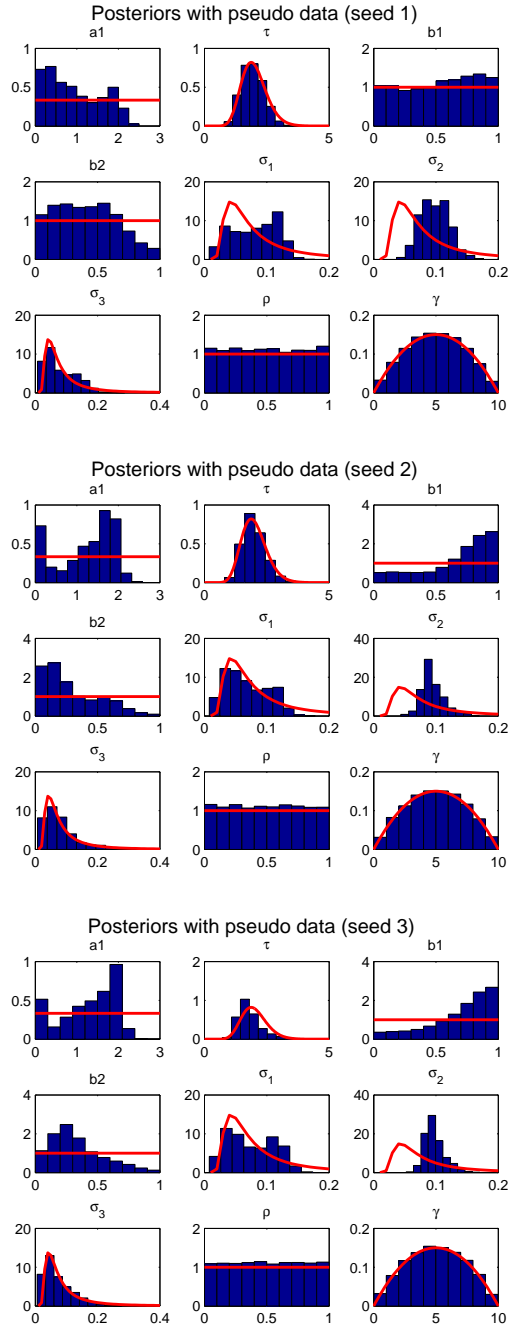


Figure 8: Posterior distributions of the parameters, estimation on pseudo-true data generated with values of the parameters as in table 2, three different seeds. Priors are shown in red.

and the three shock parameters σ_1, σ_2 and σ_3 are more precisely estimated.

We then turn to estimation on real data (figure 9). The posterior distributions of a_1 and b_2 are now significantly different from the priors, and have shrunk considerably. On the other hand, the estimates for b_2 are still extremely volatile. ρ and γ still seem to matter little, while a_2 (τ) and the three shocks parameters are estimated with high precision.

It is also interesting to note that with respect to the priors commonly used in the literature, the estimates point to a much smaller role of the shocks in explaining business cycle dynamics. Their role is picked up by the endogenous swings of optimism and pessimism in expectation formation. To the extent that “shocks are a measure of our ignorance” Abramovitz (1956), the model seems to provide a better explanation of the volatility of the U.S. economy.

6 Conclusions

In this paper we have described how simple Bayesian techniques can be applied to the estimation of AB models, as an alternative to the SMD techniques. Lack of analytical relationships between aggregate variables in AB models implies that the likelihood function has to be estimated by simulation, and cannot be functionally approximated around a steady state. Moreover, because individual behaviour in AB models is not expressed in terms of equilibrium (i.e., rational expectations), the steady state can only be found, generally speaking, by running the model. This implies that the computational burden of these techniques is increased vis-à-vis that of analytical models. This problem is exacerbated by the fact that AB models are generally complicated models with many equations, which run slow.

In our first testbed model with only one parameter, we have compared three approaches to Bayesian estimation, namely non-parametric estimation of the likelihood, parametric (Gaussian) estimation of the likelihood, and likelihood-free Approximate Bayesian Computation. We also experimented with different sampling schemes, from simple ‘brute force’ grid exploration to more sophisticated Markov chain Monte Carlo and rejection sampling algorithms.

Our findings can be summarised as follows. First, simulation time can be a major obstacle to estimating large-scale AB models. Even in our extremely simple model, with one parameter only, simulation time accounts for more than 50% of all estimation time. In this paper we have focused more on the other computational requirements, that is (i) estimating the likelihood, and (ii) sampling the parameters’ space, but simulation time is obviously a constraint. How tight

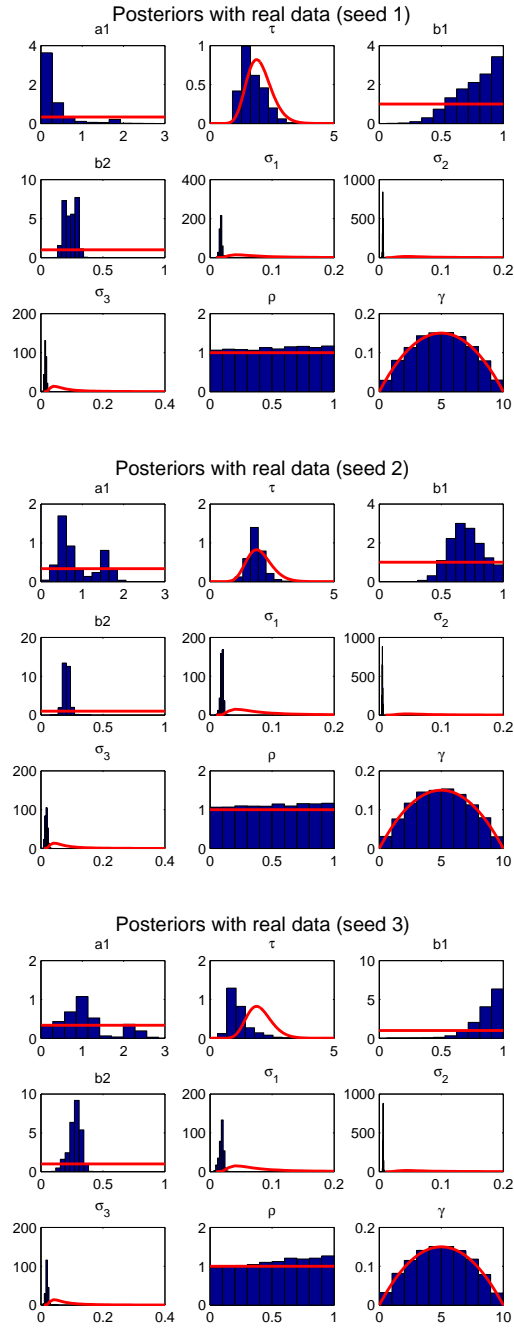


Figure 9: Posterior distributions of the parameters, estimation on observed data, three different seeds. Priors are shown in red.

is this constraint? Very few lessons can be learned, in general: running time depends on the model, on the way the model is coded, on the population size used for running the model and on the computational resources available (e.g. cloud computing), which are rapidly increasing over time.

Second, in simple models grid exploration is by far more efficient than alternative ‘efficient sampling’ schemes that filters proposed values for the parameters to arrive at a sample of values drawn from the posterior distribution, the reason being that grid exploration of the parameters’ space does not ‘throw away’ anything.

Third, both parametric estimation of the likelihood and likelihood-free ABC methods allow for a significant reduction in computing times.

These results can provide insights on common issues in AB modelling, but are obviously specific to the model being tested. In general, systematic grid sampling is not possible, the selection of parametric approximations of the likelihood function requires a lot of data mining and might not yield to satisfactory results, while the choice of summary statistics also requires extensive sensitivity analysis and is, to some extent, arbitrary. Methods that can dispense from making strong assumptions on the behaviour of the model are therefore particularly appealing. In our second application, a macro model with 9 parameters, we have tested one such method, involving a non-parametric estimation of the likelihood via kernel density estimation coupled with a simple MCMC sampling scheme. A full empirical assessment of the approach is hindered by the problem that the data seem not to discriminate well between different values of some parameters, which might be in itself a sign of misspecification, excessive parametrisation, or inherent underidentification (Liu, 1960). These problems are pervasive in complicated, structural non-linear models —including DSGE models (Canova, 2008; Canova and Sala, 2009) and there is unfortunately little that can be done about that, apart from raising awareness and lowering expectations about the estimation results.⁴² However, in our application the method produces precise estimates of most of the parameters. Moreover, the estimated values differ from the priors commonly used for the rational expectations version of the model in ways that can be easily rationalised. Generalisation of these results to more parameters, possibly more complicated models, and extension to other techniques and approaches, is left for future research.

⁴²In the context of DSGE models, Fernández-Villaverde et al. (2016) systematically consider the econometric implications of dealing with misspecified models, together with possible estimation strategies.

Acknowledgements. We thank Christian Robert for introducing us to ABC. We are indebted to seminar participants at the University of Lancaster and to the participants to the workshop “Agent-Based and DSGE Macroeconomic Modelling: Bridging the Gap” at the University of Surrey –to Alessandro Gobbi in particular– for their comments. Jakob Grazzini gratefully acknowledges the support by the European Union, Seventh Framework Programme FP7/2007-2013 Socio-economic Sciences and Humanities under Grant Agreement No. 612796 MACFINROBODS. Matteo Richiardi gratefully acknowledges support by a Marie Curie Intra European Fellowship within the 7th European Community Framework Programme.

References

- Abramovitz, M. (1956). Resource and output trends in the u.s. since 1870. *American Economic Review*, 46(2):5–23.
- Alfarano, S., Lux, F., and Wagner, T. (2007). Empirical validation of stochastic models of interacting agents. *European Physical Journal B - Condensed Matter and Complex Systems*, 55:183–187.
- Alfarano, S., Wagner, T., and Lux, F. (2005). Estimation of agent-based models: The case of an asymmetric herding model. *Computational Economics*, 26:19–49.
- Alfarano, S., Wagner, T., and Lux, F. (2006). Estimation of a simple agent-based model of financial markets: An application to australian stock and foreign exchange data. *Physica A*, 370(1):38–42.
- Andrieu, C., de Freitas, N., Doucet, A., and Jordan, M. (2003). An introduction to mcmc for machine learning. *Machine Learning*, 50(1):5–43.
- Andrieu, C. and Roberts, G. O. (2009). The pseudo-marginal approach for efficient monte carlo computations. *Annals of Statistics*, 37(2):697–725.
- Arulampalam, M., Maskell, S., Gordon, N., and Clapp, T. (2002). A tutorial on particle filters for online nonlinear/non-gaussian bayesian tracking. *IEEE Transactions on Signal Processing*, 50:174–188.
- Banish, S., Lima, R., and Araújo, T. (2012). Aggregation and emergence in agent based models: A markov chain approach. Working Paper WP 25/2012/DE/UECE, School of Economics and Management, Technical University of Lisbon.
- Barber, S., Voss, J., and Webster, M. (2015). The rate of convergence for approximate bayesian computation. *Electronic Journal of Statistics*, 9:80–105.
- Barde, S. (2016). A practical, accurate, information criterion for nth order markov processes. *Computational Economics*, doi:10.1007/s10614-016-9617-9:1–44.
- Beaumont, M. A., Zhang, W., and Balding, D. J. (2002). Approximate bayesian computation in population genetics. *Genetics*, 162(4):2025–2035.

- Bernardo, J. (1997). Noninformative priors do not exist: A discussion. *Journal of Statistical Planning and Inference*, 65(159-189).
- Bijak, J., Hilton, J., Silverman, E., and Cao, V. D. (2013). Reforging the wedding ring: exploring a semi-artificial model of population for the united kingdom with gaussian process emulators. *Demographic Research*, 29(27):729–766.
- Brock, W. A. and Hommes, C. H. (1997). A rational route to randomness. *Econometrica*, 65:1059–1095.
- Brock, W. A. and Hommes, C. H. (1998). Heterogeneous beliefs and routes to chaos in a simple asset pricing model. *Journal of Economic dynamics and Control*, 22(8):1235–1274.
- Canova, F. (2008). How much structure in empirical models? In Mills, T. C. and Patterson, K., editors, *Palgrave Handbook of Econometrics*, volume 2 (Applied Econometrics). Palgrave Macmillan.
- Canova, F. and Sala, L. (2009). Back to square one: Identification issues in dsge models. *Journal of Monetary Economics*, 56:431–449.
- Cliff, D. and Bruten, J. (1997). Minimal-intelligence agents for bargaining behaviors in market based environments. *HP Laboratories Bristol*, (HPL-97-91).
- De Grauwe, P. (2012). Booms and busts in economic activity: A behavioral explanation. *Journal of Economic Behavior & Organization*, 83:484–501.
- Deák, S., Levine, P., and Yang, B. (2015). A new keynesian behavioural model with individual rationality and heterogeneous agents. available at http://www.tinbergen.nl/wp-content/uploads/2015/09/NK_Behavioural_CEF2015_v6.pdf.
- Evensen, G. (2009). *Data Assimilation: The Ensemble Kalman Filter*. Springer, Berlin, 2nd edition.
- Fabretti, A. (2013). On the problem of calibrating an agent based model for financial markets. *Journal of Economic Interaction and Coordination*, 8(2):277–293.
- Fagiolo, G. and Roventini, A. (2012). Macroeconomic policy in dsge and agent-based models. *Revue de l'OFCE*, 124:67–116.

- Fagiolo, G. and Roventini, A. (2016). Macroeconomic policy in dsge and agent-based models redux: New developments and challenges ahead. LEM Working Paper Series 2016/17, Scuola Superiore Sant'Anna.
- Fearnhead, P. and Prangle, D. (2012). Constructing summary statistics for approximate bayesian computation: semi-automatic approximate bayesian computation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 74(3):419–474.
- Fernández-Villaverde, J., Rubio-Ramírez, J., and Schorfheide, F. (2016). Solution and estimation methods for dsge models. In Taylor, J. B. and Uhlig, H., editors, *Handbook of Macroeconomics*, volume 2, chapter 9, pages 527–724. Elsevier.
- Fisher, R. (1922). On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society A*, 222:309–368.
- Fu, Y. and Li, W. (1997). Estimating the age of the common ancestor of a sample of dna sequences. *Molecular biology and evolution*, 14(2):195–199.
- Grazzini, J. (2012). Analysis of the emergent properties: Stationarity and ergodicity. *Journal of Artificial Societies and Social Simulation*, 15(2):7.
- Grazzini, J. and Richiardi, M. (2015). Consistent estimation of agent-based models by simulated minimum distance. *Journal of Economic Dynamics and Control*, 51:148–165.
- Grazzini, J., Richiardi, M., and Sella, L. (2012). Small sample bias in msm estimation of agent-based models. In Andrea Teglio, Simone Alfarano, E. C.-C. M. G.-V., editor, *Managing Market Complexity. The Approach of Artificial Economics.*, Lecture Notes in Economics and Mathematical Systems. Springer.
- Guerini, M. and Moneta, A. (2016). A method for agent-based models validation. Technical report, Scuola Superiore Sant'Anna.
- Hansen, L. P. and Heckman, J. J. (1996). The empirical foundations of calibration. *Journal of Economic Perspectives*, 10(87-104).
- Hartig, F., Calabrese, J. M., Reineking, B., Wiegand, T., and Huth, A. (2011). Statistical inference for stochastic simulation models - theory and application. *Ecology Letters*, 14:816–827.

- Herbst, E. and Schorfheide, F. (2016). *Bayesian Estimation of DSGE Models*. Princeton University Press.
- Holmes, W. R. (2015). A practical guide to the probability density approximation (pda) with improved implementation and error characterization. *Journal of Mathematical Psychology* (forthcoming).
- Kemeny, J. and Snell, J. (1976). *Finite Markov Chains*. Springer, New York.
- Kroese, D. P., Taimre, T., and Botev, Z. I. (2011). *Handbook of Monte Monte Methods*. John Wiley & Sons, Hoboken, NJ.
- Kukacka, J. and Barunik, J. (2016). Estimation of financial agent-based models with simulated maximum likelihood.
- Kydland, F. E. and Prescott, E. C. (1996). The computational experiment: An econometric tool. *Journal of Economic Perspectives*, 10:69–85.
- Lamperti, F. (2015). An information theoretic criterion for empirical validation of time series models. Lem papers series, Sant’Anna School of Advanced Studies, Pisa, Italy.
- Liu, T.-C. (1960). Underidentification, structural estimation, and forecasting. *Econometrica*, 28(4):855–865.
- Macal, C. M. (2016). Everything you need to know about agent-based modelling and simulation. *Journal of Simulation*, 10(2):144–156.
- Marin, J.-M., Pudlo, P., Robert, C. P., and Ryder, R. J. (2011). Approximate bayesian computational methods. *Statistics and Computing*, 22(6):1167–1180.
- Marks, R. E. (2013). Validation and model selection: Three similarity measures compared. *Complexity Economics*, 2:41–61.
- O’Hagan, A. (2006). Bayesian analysis of computer code outputs: A tutorial. *Reliability Engineering & System Safety*, 91(1290-1300).
- Prangle, D. (2014). Lazy abc. *Statistics and Computing*, pages 1–15.
- Recchioni, M. C., Tedeschi, G., and Gallegati, M. (2015). A calibration procedure for analyzing stock price dynamics in an agent-based framework. *Journal of Economic Dynamics & Control*, 60(1):1–25.

- Richiardi, M. (2012). Agent-based computational economics. a short introduction. *The Knowledge Engineering Review*, 27(2):137–149.
- Richiardi, M. G. (2016). The future of agent-based modelling. *Eastern Economic Journal*.
- Rubio-Ramirez, J. F. and Fernández-Villaverde, J. (2005). Estimating dynamic equilibrium economies: linear versus nonlinear likelihood. *Journal of Applied Econometrics*, 20(7):891–910.
- Ruge-Murcia, F. (2007). Methods to estimate dynamic stochastic general equilibrium models. *Journal of Economic Dynamics and Control*, 31(2599-2636).
- Salle, I. and Yildizoglu, M. (2014). Efficient sampling and metamodeling in computational economic models. *Computational Economics*, 44(4):507–553.
- Sani, A., Lamperti, F., Mandel, A., and Roventini, A. (2016). Machine learning surrogates for agent-based models. Centre d’Économie de la Sorbonne, Université Paris 1, Panthéon-Sorbonne, Paris School of Economics.
- Silverman, B. W. (1986). *Density Estimation for Statistics and Data Analysis*. Chapman and Hall.
- Silverman, B. W. (1998). *Density Estimation for Statistics and Data Analysis*. Chapman & Hall/CRC, London.
- Sisson, S. A., Fan, Y., and Beaumont, M., editors (2016). *Handbook of Approximate Bayesian Computation*. Taylor & Francis.
- Sottoriva, A. and Tavaré, S. (2010). Integrating approximate bayesian computation with complex agent-based models for cancer research. In Lechevallier, Y. and Saporta, G., editors, *COMPSTAT 2010. Proceedings in Computational Statistics*. Springer.
- Stordal, A. S., Karlsen, H. A., Nævdal, G., Skaug, H. J., and Vallès, B. (2011). Bridging the ensemble kalman filter and particle filters: the adaptive gaussian mixture filter. *Computational Geosciences*, 15(2):293–305.
- Sun, D. and Berger, J. (2006). Objective bayesian analysis for the multivariate normal model. In *Proc. Valencia / ISBA 8th World Meeting on Bayesian Statistics, Benidorm (Alicante, Spain), June 1st-6th*.

- Tavaré, S., Balding, D. J., Griffiths, R. C., and Donnelly, P. (1997). Inferring coalescent times from dna sequence data. *Genetics*, 145:505–518.
- Turner, B. M. and Zandt, T. V. (2012). A tutorial on approximate bayesian computation. *Journal of Mathematical Psychology*, 56:69–85.
- Ward, J. A., Evans, A. J., and Malleson, N. S. (2016). Dynamic calibration of agent-based models using data assimilation. *Royal Society Open Science*, 3:150703.
- Wood, S. N. (2010). Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466:1102–1107.

A Gaussian density estimation

In this Appendix we derive the expression for the likelihood and the posterior density functions under the assumption of Gaussian errors, either for ϵ (eq. 10) or for \mathbf{u} (eq. 11). Under the Gaussian assumption, disturbances are distributed as $\mathcal{N}_K(\mathbf{0}, \mathbf{\Sigma})$, where K is the dimensionality of \mathbf{y} , the number of aggregate observables that the model is able to reproduce. We can either estimate the free elements of matrix $\mathbf{\Sigma}$ or assume that it is diagonal with diagonal elements σ_{kk} , $k = 1, \dots, K$ which can be thought of as fixed (to some small value) or unknown. In the stationary state, we then have that

$$\mathbf{y}_{t+1} \sim \mathcal{N}_K(\mathbf{g}^*(\boldsymbol{\theta}), \mathbf{\Sigma}(\boldsymbol{\theta})). \quad (29)$$

or, in the case of disturbances added to the data,

$$\mathbf{y}_{t+1}^R \sim \mathcal{N}_K(\mathbf{g}^*(\boldsymbol{\theta}^R), \mathbf{\Sigma}(\boldsymbol{\theta})). \quad (29')$$

The joint distribution of the observables at time t , in turn, is given directly as follows:

$$p(\mathbf{y}_t^R | \boldsymbol{\Xi}_t, \boldsymbol{\theta}, \mathbf{\Sigma}(\boldsymbol{\theta})) \propto |\mathbf{\Sigma}(\boldsymbol{\theta})|^{-1/2} \exp \left\{ -\frac{1}{2} [\mathbf{y}_t^R - \mathbf{g}^*(\boldsymbol{\theta})]' \mathbf{\Sigma}(\boldsymbol{\theta})^{-1} [\mathbf{y}_t^R - \mathbf{g}^*(\boldsymbol{\theta})] \right\}. \quad (30)$$

from which a likelihood function can be easily derived:⁴³

$$\mathcal{L}(\boldsymbol{\theta}, \mathbf{\Sigma}; \mathbf{Y}^R) \propto |\mathbf{\Sigma}|^{-T/2} \exp \left\{ -\frac{1}{2} \text{tr} [\mathbf{A}(\boldsymbol{\theta}) \mathbf{\Sigma}^{-1}] \right\} \quad (31)$$

where $\mathbf{A}(\boldsymbol{\theta}) = \sum_{t=1}^T [\mathbf{y}_t^R - \mathbf{g}^*(\boldsymbol{\theta})] \cdot [\mathbf{y}_t^R - \mathbf{g}^*(\boldsymbol{\theta})]'$.

The likelihood is formed under the assumption that the disturbances are an iid process so that there is, for example, no autocorrelation (the assumption is easy to remove).

Further simplification can be obtained by assuming that the elements of $\mathbf{\Sigma}$ are independent from those of $\boldsymbol{\theta}$:

$$p(\boldsymbol{\theta}, \mathbf{\Sigma}) = p(\boldsymbol{\theta})p(\mathbf{\Sigma}). \quad (32)$$

⁴³We use the ‘trace trick’ for Gaussian likelihoods: $\mathbf{X}'\mathbf{\Sigma}^{-1}\mathbf{X} = \text{tr}(\mathbf{\Sigma}^{-1}\mathbf{X}\mathbf{X}')$. This comes from two properties of the trace: (i) $\text{tr}(AB) = \text{tr}(BA)$ (if all dimensions work out, ie. if AB is a square matrix), (ii) $\text{tr}(c) = c$ if c is a constant. Because $\mathbf{X}'\mathbf{\Sigma}^{-1}\mathbf{X}$ is a scalar, we can then write $\mathbf{X}'\mathbf{\Sigma}^{-1}\mathbf{X} = \text{tr}(\mathbf{X}'\mathbf{\Sigma}^{-1}\mathbf{X}) = \text{tr}(\mathbf{X}\mathbf{X}'\mathbf{\Sigma}^{-1})$.

We then assume diffuse priors:

$$p(\boldsymbol{\theta}) = \text{const.} \quad (33)$$

$$p(\boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{-(K+1)/2} \quad (34)$$

where $p(\boldsymbol{\Sigma})$ is the commonly used independence-Jeffreys prior (Sun and Berger, 2006), that is, invariant to re-parametrisation of $\boldsymbol{\Sigma}$.

This leads to the following posterior:

$$p(\boldsymbol{\theta}, \boldsymbol{\Sigma} | \mathbf{Y}^R) \propto |\boldsymbol{\Sigma}|^{-\frac{T+K+1}{2}} \exp \left\{ -\frac{1}{2} \text{tr} [\mathbf{A}(\boldsymbol{\theta}) \boldsymbol{\Sigma}^{-1}] \right\} \quad (35)$$

Given that $\boldsymbol{\Sigma}$ is generally not an object of interest (though it could be used as measure of how well the model fits the data), we can integrate it out analytically to obtain:

$$p(\boldsymbol{\theta} | \mathbf{Y}^R) \propto |\mathbf{A}(\boldsymbol{\theta})|^{-T/2} \quad (36)$$

In this way, kernel-based estimation is avoided altogether, saving computational time.

Notice that the new posterior does not depend on the different elements of $\boldsymbol{\Sigma}$. The dimensionality of this matrix is the same as the dimensionality of \mathbf{Y}^R which is, in our case, quite low. Generally speaking, it is required that if $\boldsymbol{\Sigma}$ is $K \times K$ we should have $\frac{K(K+1)}{2} \ll T$.