
Inference in differentiable generative models

Matthew Graham¹ Amos Storkey¹

Abstract

Many generative models can be expressed as a differentiable function of random inputs drawn from a known probability distribution. This framework includes both learnt parametric generative models and a large class of procedurally defined simulator models. We present a method for performing efficient Markov chain Monte Carlo (MCMC) inference in such models when conditioning on observations of the model output. For some models this offers an asymptotically exact inference method where Approximate Bayesian Computation might otherwise be employed. We use the intuition that inference corresponds to integrating a density across the manifold corresponding to the set of inputs consistent with the observed outputs. This motivates the use of a constrained variant of Hamiltonian Monte Carlo which leverages the smooth geometry of the manifold to move between inputs exactly consistent with observations.

1. Introduction

There has been a long interest in probabilistic models which are defined *implicitly* (4; 7; 12) - that is where we can generate random values for the latent $\mathbf{z} \in \mathcal{Z}$ and observed $\mathbf{x} \in \mathcal{X}$ variables in the model, but we cannot tractably evaluate the probability distribution $P_{\mathbf{x},\mathbf{z}}$ of those variables and more specifically its density $p_{\mathbf{x},\mathbf{z}}$. Although implicit models are challenging from an inferential perspective, they are ubiquitous in science and engineering in the form of simulator models of physical systems. Typically simulators are specified procedurally in code with any stochasticity introduced by draws from a random number generator. Recently implicit models have also been the subject of much interest in the machine learning community due to the advent of effective adversarial training methods (10) and the gains in modelling flexibility offered by dropping the requirement of being able to compute a density on model outputs (22; 29).

¹University of Edinburgh, Edinburgh, United Kingdom. Correspondence to: Matthew Graham <m.m.graham@ed.ac.uk>.

Any probabilistic model that we can programmatically generate values from can be expressed in the form of a deterministic function which takes as input a (possibly variable length) vector of random inputs sampled from a known distribution. This observation just corresponds to stating that we can track all of the calls to a random number generator in a program, and that given the values sampled from the random number generator all of the operations then performed by the program are deterministic.

To formalise this intuition we introduce some notation. We define an *input space* \mathcal{U} and an associated vector of input random variables $\mathbf{u} \in \mathcal{U}$ which takes on values in this space, with \mathbf{u} having a known distribution $P_{\mathbf{u}} = P$. We then define *generator functions* $\mathbf{g}_{\mathbf{x}} : \mathcal{U} \rightarrow \mathcal{X}$ and $\mathbf{g}_{\mathbf{z}} : \mathcal{U} \rightarrow \mathcal{Z}$ such that we have $\mathbf{x} = \mathbf{g}_{\mathbf{x}}(\mathbf{u})$ and $\mathbf{z} = \mathbf{g}_{\mathbf{z}}(\mathbf{u})$. A common special case is when the input space is partitioned $\mathcal{U} = \mathcal{U}_1 \times \mathcal{U}_2$ and the unobserved variables \mathbf{z} are generated from a subset of the random inputs $\mathbf{u}_1 \in \mathcal{U}_1$ (e.g. corresponding to sampling from a prior distribution over the parameters of a simulator model), with the observed variables \mathbf{x} then generated from a function $\mathbf{g}_{\mathbf{x}|\mathbf{z}} : \mathcal{Z} \times \mathcal{U}_2 \rightarrow \mathcal{X}$, i.e. $\mathbf{x} = \mathbf{g}_{\mathbf{x}|\mathbf{z}}(\mathbf{g}_{\mathbf{z}}(\mathbf{u}_1), \mathbf{u}_2)$. In Appendix A in the supplementary material we show factor graphs visualising these model structures.

We define *differentiable generative models* as the restricted class of models whereby

1. $\mathcal{U} \subseteq \mathbb{R}^M$ and $\mathcal{X} \subseteq \mathbb{R}^{N_{\mathbf{x}}}$
2. P has density ρ with respect to the Lebesgue measure,
3. $\frac{\partial \rho}{\partial \mathbf{u}}$ exists P -almost everywhere,
4. $\frac{\partial \mathbf{g}_{\mathbf{x}}}{\partial \mathbf{u}}$ exists and is full row-rank P -almost everywhere.

In this paper we outline a MCMC method for performing inference in models meeting these conditions. This is an abridged version of a published conference paper (13).

2. Approximate Bayesian Computation

Approximate Bayesian Computation (ABC) (4; 20) is a popular approach for performing approximate inference generative models where we do not have access to an explicit joint density $p_{\mathbf{x},\mathbf{z}}$ but can generate (\mathbf{x}, \mathbf{z}) samples. In ABC the simulated observed outputs \mathbf{x} are decoupled from the observations $\bar{\mathbf{x}}$ by a noise model or *kernel* $p_{\bar{\mathbf{x}}|\mathbf{x}}(\bar{\mathbf{x}}|\mathbf{x}) = k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x})$ with tolerance parameter ϵ , e.g. $k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}) \propto \mathbb{I}[\|\bar{\mathbf{x}} - \mathbf{x}\| < \epsilon]/\epsilon^{N_{\mathbf{x}}}$ (uniform ball kernel) or $k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}) = \mathcal{N}(\bar{\mathbf{x}}|\mathbf{x}, \epsilon^2 \mathbf{I})$ (Gaussian kernel).

The *ABC posterior* is then defined as

$$p_{\mathbf{z}|\bar{\mathbf{x}}}(z|\bar{\mathbf{x}}; \epsilon) = \frac{\int_{\mathcal{X}} k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}) p_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|z) p_{\mathbf{z}}(z) d\mathbf{x}}{\int_{\mathcal{X}} k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}} \quad (1)$$

and we can express expectations of functions of the latent variables \mathbf{z} with respect to this approximate posterior

$$\begin{aligned} \mathbb{E}[f(\mathbf{z})|\bar{\mathbf{x}} = \bar{\mathbf{x}}; \epsilon] &= \int_{\mathcal{Z}} f(\mathbf{z}) p_{\mathbf{z}|\bar{\mathbf{x}}}(z|\bar{\mathbf{x}}; \epsilon) d\mathbf{z} \quad (2) \\ &= \frac{\int_{\mathcal{Z}} \int_{\mathcal{X}} f(\mathbf{z}) k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}) p_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|z) p_{\mathbf{z}}(z) d\mathbf{x} d\mathbf{z}}{\int_{\mathcal{Z}} \int_{\mathcal{X}} k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}) p_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|z) p_{\mathbf{z}}(z) d\mathbf{x} d\mathbf{z}}. \end{aligned}$$

Various Monte Carlo approximate inference schemes can be used to estimate this expectation. The simplest is to generate a set of independent samples $\{\mathbf{x}^{(s)}, \mathbf{z}^{(s)}\}_{s=1}^S$ from $P_{\mathbf{x}, \mathbf{z}}$ ¹ and use the ratio estimator

$$\mathbb{E}[f(\mathbf{z})|\bar{\mathbf{x}} = \bar{\mathbf{x}}; \epsilon] \approx \frac{\sum_{s=1}^S f(\mathbf{z}^{(s)}) k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}^{(s)})}{\sum_{s=1}^S k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}^{(s)})}. \quad (3)$$

In the case of a uniform ball kernel this corresponds to the standard ABC reject algorithm, with expectations being estimated as averages over the latent variable samples for which the corresponding simulated outputs are within a (Euclidean) distance of ϵ from the observations.

Alternatively a MCMC scheme can be used to estimate the ABC posterior expectation in (2) (21). A Markov chain is constructed with a stationary distribution with a density

$$p_{\mathbf{x}, \mathbf{z}|\bar{\mathbf{x}}}(\mathbf{x}, \mathbf{z}|\bar{\mathbf{x}}) \propto k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}) p_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}) p_{\mathbf{z}}(\mathbf{z}) \quad (4)$$

by proposing a new state $(\mathbf{x}', \mathbf{z}')$ given the current state (\mathbf{x}, \mathbf{z}) by sampling $\mathbf{z}' \sim q(\cdot|\mathbf{z})$ from some perturbative proposal distribution q on \mathcal{Z} and then generating a new $\mathbf{x}' \sim p_{\mathbf{x}|\mathbf{z}}(\cdot|\mathbf{z})$. This is then accepted with probability

$$a(\mathbf{x}', \mathbf{z}'|\mathbf{x}, \mathbf{z}) = \min \left[1, \frac{k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}') q(\mathbf{z}|\mathbf{z}') p_{\mathbf{z}}(\mathbf{z}')}{k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{x}) q(\mathbf{z}'|\mathbf{z}) p_{\mathbf{z}}(\mathbf{z})} \right], \quad (5)$$

the overall transition operator leaving (4) stationary. The samples of the chain state can then be used to compute consistent estimators of (2).

For both the ABC reject and ABC MCMC schemes just described, simulated observations \mathbf{x} are independently generated from the conditional $p_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z})$ given the current latent variables \mathbf{z} . The observed values $\bar{\mathbf{x}}$ are a zero-measure set in \mathcal{X} under non-degenerate $p_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z})$ and so as $\epsilon \rightarrow 0$ the probability of accepting a sample / proposed move becomes zero. Applying ABC with a non-zero ϵ therefore can be seen as a practically motivated relaxation of the constraint that true and simulated data exactly match, and hence the ‘approximate’ in *Approximate Bayesian Computation*.

¹In ABC this usually involves generating \mathbf{z} then simulating \mathbf{x} given \mathbf{z} however more generally we can just sample from the joint.

3. Inference in the input space

Using our definitions $\mathbf{x} = \mathbf{g}_{\mathbf{x}}(\mathbf{u})$ and $\mathbf{z} = \mathbf{g}_{\mathbf{z}}(\mathbf{u})$ we can reparametrise (2) as an integral over the input space to the generative model

$$\mathbb{E}[f(\mathbf{z})|\bar{\mathbf{x}} = \bar{\mathbf{x}}; \epsilon] \propto \int_{\mathcal{U}} f \circ \mathbf{g}_{\mathbf{z}}(\mathbf{u}) k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{g}_{\mathbf{x}}(\mathbf{u})) \rho(\mathbf{u}) d\mathbf{u}.$$

This indicates we can estimate ABC expectations by applying standard MCMC methods in the input space to construct a chain with stationary distribution with density

$$\pi_{\epsilon}(\mathbf{u}) \propto k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{g}_{\mathbf{x}}(\mathbf{u})) \rho(\mathbf{u}). \quad (6)$$

We can also however consider the limit of $\epsilon \rightarrow 0$. We have that the kernel term $k_{\epsilon}(\bar{\mathbf{x}}; \mathbf{g}_{\mathbf{x}}(\mathbf{u})) \rightarrow \delta(\bar{\mathbf{x}} - \mathbf{g}_{\mathbf{x}}(\mathbf{u}))$ and so

$$\begin{aligned} \mathbb{E}[f(\mathbf{z})|\mathbf{x} = \bar{\mathbf{x}}] &= \lim_{\epsilon \rightarrow 0} \{\mathbb{E}[f(\mathbf{z})|\bar{\mathbf{x}} = \bar{\mathbf{x}}; \epsilon]\} \quad (7) \\ &\propto \int_{\mathcal{U}} f \circ \mathbf{g}_{\mathbf{z}}(\mathbf{u}) \delta(\bar{\mathbf{x}} - \mathbf{g}_{\mathbf{x}}(\mathbf{u})) \rho(\mathbf{u}) d\mathbf{u}. \end{aligned}$$

The Dirac delta term restricts the integral across the input space \mathcal{U} to an embedded, $M - N_{\mathbf{x}}$ dimensional, implicitly-defined manifold $\mathbf{g}_{\mathbf{x}}^{-1}[\bar{\mathbf{x}}] \equiv \{\mathbf{u} \in \mathcal{U} : \mathbf{g}_{\mathbf{x}}(\mathbf{u}) = \bar{\mathbf{x}}\}$. By applying the *Co-Area Formula* (9, §3.2.12) the integral with respect to the Lebesgue measure across \mathcal{U} in (7) can be rewritten as integral across the embedded manifold $\mathbf{g}_{\mathbf{x}}^{-1}[\bar{\mathbf{x}}]$ with respect the Hausdorff measure for the manifold

$$\mathbb{E}[f(\mathbf{z})|\mathbf{x} = \bar{\mathbf{x}}] \propto \quad (8)$$

$$\int_{\mathbf{g}_{\mathbf{x}}^{-1}[\bar{\mathbf{x}}]} f \circ \mathbf{g}_{\mathbf{z}}(\mathbf{u}) \left| \frac{\partial \mathbf{g}_{\mathbf{x}}}{\partial \mathbf{u}} \frac{\partial \mathbf{g}_{\mathbf{x}}}{\partial \mathbf{u}}^{\top} \right|^{-\frac{1}{2}} \rho(\mathbf{u}) d\mathcal{H}^{M-N_{\mathbf{x}}}(\mathbf{u}).$$

Therefore if we construct a Markov chain $\{\mathbf{u}^{(s)}\}_{s=1}^S$ restricted to $\mathbf{g}_{\mathbf{x}}^{-1}[\bar{\mathbf{x}}]$ and with an invariant distribution with density with respect to the Hausdorff measure of $\mathbf{g}_{\mathbf{x}}^{-1}[\bar{\mathbf{x}}]$

$$\pi(\mathbf{u}) \propto \left| \frac{\partial \mathbf{g}_{\mathbf{x}}}{\partial \mathbf{u}} \frac{\partial \mathbf{g}_{\mathbf{x}}}{\partial \mathbf{u}}^{\top} \right|^{-\frac{1}{2}} \rho(\mathbf{u}) \quad (9)$$

then if the chain is also aperiodic and irreducible we can form MCMC estimators which converge in the limit $S \rightarrow \infty$

$$\mathbb{E} \left[\frac{1}{S} \sum_{s=1}^S (f \circ \mathbf{g}_{\mathbf{z}}(\mathbf{u}^{(s)})) \right] \rightarrow \mathbb{E}[f(\mathbf{z})|\mathbf{x} = \bar{\mathbf{x}}] \quad (10)$$

Intuitively the determinant term in (9) adjusts for the change in the infinitesimal ‘thickness’ (extent in directions orthogonal to the tangent space) of the manifold when mapping through the generator function $\mathbf{g}_{\mathbf{x}}$. The result (8) is given in a slightly different form in (6).

A general framework for performing asymptotically exact inference in differentiable generative models is therefore to define MCMC updates which explore the target density (9) on the pre-image manifold $\mathbf{g}_{\mathbf{x}}^{-1}[\bar{\mathbf{x}}]$. We propose here to use a method which simulates the dynamics of a constrained mechanical system.

4. Constrained Hamiltonian Monte Carlo

In *Hamiltonian Monte Carlo* (HMC) (8; 26) the vector variable of interest \mathbf{u} is augmented with a momentum variable $\mathbf{p} \in \mathbb{R}^M$. The momenta are defined to be independent of \mathbf{u} with a Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})^2$. The joint distribution then has a density proportional to $\exp[-H(\mathbf{u}, \mathbf{p})]$ where $H(\mathbf{u}, \mathbf{p}) = -\log \pi(\mathbf{u}) + \frac{1}{2}\mathbf{p}^\top \mathbf{p}$ is the Hamiltonian. The canonical Hamiltonian dynamic is described by

$$\frac{d\mathbf{u}}{dt} = \frac{\partial H}{\partial \mathbf{p}} = \mathbf{p}, \quad \frac{d\mathbf{p}}{dt} = -\frac{\partial H}{\partial \mathbf{u}} = -\frac{\partial \log \pi}{\partial \mathbf{u}}. \quad (11)$$

This dynamic is time-reversible, volume-preserving and exactly conserves the Hamiltonian. Symplectic integrators allow approximate integration of the Hamiltonian flow while maintaining the time-reversibility and volume-preservation properties (17). These properties make simulated Hamiltonian dynamics an ideal proposal mechanism for a Metropolis MCMC method. The accept probability for a proposal $(\mathbf{u}_p, \mathbf{p}_p)$ generated by simulating the dynamic forward from (\mathbf{u}, \mathbf{p}) is $\exp(H(\mathbf{u}, \mathbf{p}) - H(\mathbf{u}_p, \mathbf{p}_p))$. Typically the change in the Hamiltonian will be small and so the probability of acceptance high.

In our case the system is subject to a constraint of the form $\mathbf{g}_\mathbf{x}(\mathbf{u}) - \bar{\mathbf{x}} = \mathbf{0}$. By introducing Lagrangian multipliers λ_i for each of the constraints, the Hamiltonian for a constrained system can be written as $H(\mathbf{u}, \mathbf{p}) = -\log \pi(\mathbf{u}) + \frac{1}{2}\mathbf{p}^\top \mathbf{p} + \boldsymbol{\lambda}^\top (\mathbf{g}_\mathbf{x}(\mathbf{u}) - \bar{\mathbf{x}})$, with a corresponding constrained Hamiltonian dynamic

$$\frac{d\mathbf{u}}{dt} = \mathbf{p}, \quad \frac{d\mathbf{p}}{dt} = \frac{\partial \log \pi}{\partial \mathbf{u}} - \frac{\partial \mathbf{g}_\mathbf{x}^\top}{\partial \mathbf{u}} \boldsymbol{\lambda}, \quad (12)$$

$$\text{subject to } \mathbf{g}_\mathbf{x}(\mathbf{u}) = \bar{\mathbf{x}}, \quad \frac{\partial \mathbf{g}_\mathbf{x}}{\partial \mathbf{u}} \mathbf{p} = \mathbf{0}. \quad (13)$$

A popular numerical integrator for simulating constrained Hamiltonian dynamics is RATTLE (2). This a generalisation of the Störmer-Verlet (leapfrog) integrator typically used in standard HMC with additional projection steps in which the Lagrange multipliers $\boldsymbol{\lambda}$ are solved for to satisfy the conditions (13). RATTLE is symplectic on the constraint manifold and is time-reversible and volume-preserving (18).

The use of constrained dynamics in HMC has been proposed several times. In the molecular dynamics literature, both (14) and (19) suggest using a simulated constrained dynamic within a HMC framework. Most relevantly here (5) proposes using a constrained HMC variant to perform inference in distributions defined on implicitly defined manifolds. The authors of (5) give sufficient conditions on H and $\mathbf{g}_\mathbf{x}$ for the transition operator to be irreducible and aperiodic: that H is C^2 continuous, and $\mathbf{g}_\mathbf{x}^{-1}[\bar{\mathbf{x}}]$ is a connected smooth and differentiable manifold and $\frac{\partial \mathbf{g}_\mathbf{x}}{\partial \mathbf{u}}$ has full row-rank everywhere.

²For notation simplicity we assume an identity mass matrix.

5. Method

Our constrained HMC implementation is shown in Appendix C in Algorithm 1. We use a generalisation of the RATTLE scheme to simulate the dynamic. The inner updates of the state to solve for the geodesic motion on the constraint manifold are split into multiple smaller steps. This is a special case of the scheme described in (16) and allows more flexibility in choosing an appropriately small step-size to ensure convergence of the iterative solution of the equations projecting on to the constraint manifold while still allowing a more efficient larger step size for updates to the momentum.

Each inner geodesic time-step involves making an unconstrained update $\tilde{\mathbf{u}} \leftarrow \mathbf{u} - \delta t \mathbf{p}$ and then projecting $\tilde{\mathbf{u}}$ back on to $\mathbf{g}_\mathbf{x}^{-1}[\bar{\mathbf{x}}]$ by solving for $\boldsymbol{\lambda}$ which satisfy $\mathbf{g}_\mathbf{x}(\tilde{\mathbf{u}} - \frac{\partial \mathbf{g}_\mathbf{x}^\top}{\partial \mathbf{u}} \boldsymbol{\lambda}) = \bar{\mathbf{x}}$. Here we use a Newton method for solving the system of equations in the projection step. The true Newton update would involve recalculating the Jacobian and solving a dense linear system within the optimisation loop. Instead we use a symmetric quasi-Newton update as proposed in (3), the Jacobian matrix product $\frac{\partial \mathbf{g}_\mathbf{x}}{\partial \mathbf{u}} \frac{\partial \mathbf{g}_\mathbf{x}^\top}{\partial \mathbf{u}}$ evaluated at the previous state used to condition the moves. This matrix is positive-definite and a Cholesky decomposition can be calculated outside the optimisation loop allowing cheaper quadratic cost solves within the loop. For larger systems, the Cholesky decomposition of the Jacobian matrix product will become a dominant cost, generally scaling cubically with $N_\mathbf{x}$. Conditional independency structure of many directed generative models however allows a significantly reduced quadratically scaling computational cost as explained in the supplementary material in Appendix D.

In the Newton iteration projection step, convergence is signalled when the elementwise maximum absolute difference $\|\mathbf{g}_\mathbf{x}(\mathbf{u}) - \bar{\mathbf{x}}\|_\infty$ is below a tolerance ϵ . This acts analogously to the ϵ parameter in ABC methods, however here we typically set this parameter to close to machine precision ($\epsilon = 10^{-8}$ in our experiments) and so the approximation introduced is comparable to that otherwise incurred for using non-exact arithmetic.

6. Example: Lotka–Volterra model

As an example application of our method, we consider inferring the posterior distribution of the parameters of a stochastic predator-prey model given (simulated) observed populations. In Appendices E and F we also show results for 3D human pose and camera parameter inference given 2D joint position information and in-painting of missing regions of digit images using a generative model trained on MNIST. In all experiments Theano (28) was used to specify the generator function and calculate the required derivatives. Python code for the experiments is available

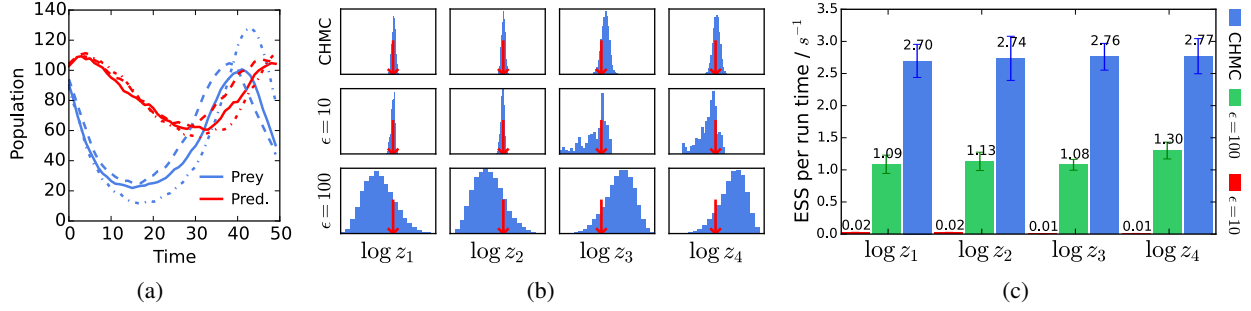


Figure 1. Lotka–Volterra (a) Observed predator-prey populations (solid) and ABC sample trajectories with $\epsilon = 10$ (dashed) and $\epsilon = 100$ (dot-dashed). (b) Marginal empirical histograms for the (logarithm of the) four parameters (columns) from constrained HMC samples (top) and ABC samples with $\epsilon = 10$ (middle) and $\epsilon = 100$ (bottom). Horizontal axes shared across columns. Red arrows indicate true parameter values. (c) Mean ESS normalised by compute time for each of four parameters for ABC with $\epsilon = 10$ (red), $\epsilon = 100$ (green) and our method (blue). Error bars show ± 3 standard errors of mean.

at <https://git.io/dgm>. We considered a stochastic continuous state variant of the Lotka–Volterra model. In particular we consider parameter inference given a simulated solution of the following stochastic differential equations

$$dx_1 = (z_1 x_1 - z_2 x_1 x_2)dt + dn_1, \quad (14)$$

$$dx_2 = (-z_3 x_2 + z_4 x_1 x_2)dt + dn_2, \quad (15)$$

where x_1 represents the prey population, x_2 the predator population, $\{z_i\}_{i=1}^4$ the system parameters and n_1 and n_2 zero-mean, unit variance white noise processes.

The observed data was generated with an Euler-Maruyama discretisation, time-step 1, initial condition $x_1^{(0)} = x_2^{(0)} = 100$ and $z_1 = 0.4$, $z_2 = 0.005$, $z_3 = 0.05$, $z_4 = 0.001$ (chosen to give stable dynamics). The simulation was run for 50 time-steps with the observed outputs defined as the concatenated vector $\mathbf{x} = [x_1^{(1)} x_2^{(1)} \dots x_1^{(50)} x_2^{(50)}]$. Log-normal priors $z_i \sim \log \mathcal{N}(-2, 1)$ were placed on the system parameters.

We compared our method to various ABC approaches (§2) using a uniform ball kernel with radius ϵ . ABC rejection failed catastrophically, with no acceptances in 10^6 samples even with a large $\epsilon = 1000$. ABC MCMC with a Gaussian proposal distribution q also performed very poorly with the dynamic having zero acceptances over multiple runs of 10^5 updates for $\epsilon = 100$ and getting stuck at points in parameter space over many updates for larger $\epsilon = 1000$, even with small proposal steps. Based on a method proposed in the pseudo-marginal literature (25), we tried using alternating elliptical slice sampling (24) updates of the inputs \mathbf{u}_1 used to generate the parameters $\mathbf{z} = \mathbf{g}_z(\mathbf{u}_1)$ and remaining random inputs \mathbf{u}_2 used to generate the observed variables $\mathbf{x} = \mathbf{g}_x(\mathbf{z}, \mathbf{u}_2)$ given \mathbf{z} . The slice sampling updates adapt the size of steps made to ensure a move can always be made. Using this method we were able to obtain reasonable convergence over long runs for both $\epsilon = 100$ and $\epsilon = 10$.

The results are summarised in Figure 1. Figure 1a shows the simulated data used as observations and ABC sample trajec-

tories for $\epsilon = 10$ and $\epsilon = 100$. Though both samples follow the general trends of the observed data there are large discrepancies particularly for $\epsilon = 100$. Our method in contrast samples parameters generating trajectories matching the observations to within $\epsilon = 10^{-8}$ at all points. Figure 1b shows the marginal histograms for the parameters. As would be expected the inferred posterior on the parameters are significantly more tightly distributed about the true values used to generate the observations for our approach and the $\epsilon = 10$ case compared to the results for $\epsilon = 100$. Figure 1c shows the relative sampling efficiency of our approach against the ABC methods, as measured by the *effective sample sizes* (ESS) (computed with R-CODA (27)) normalised by run time averaged across 10 sampling runs for each method. Despite the significantly higher run time per-sample in our method, the reduced autocorrelation due to the much larger moves made by the Hamiltonian dynamic give much higher ESS even over the very approximate $\epsilon = 100$ case.

7. Discussion

We have presented a framework for performing inference in differentiable generative models. Though the constrained HMC updates are computationally costly, the gradient-based exploration of the state space can lead to significantly improved sampling efficiency over simpler methods. Further our approach in some cases allows asymptotically exact inference in differentiable generative models where ABC methods might otherwise be used. In our experiments we were able to condition on high-dimensional observations without the need for dimensionality reduction with summary statistics however if informative summaries are available and correspond to differentiable functions of the observations they can be exploited in our method by adding them to the generator definition. As well as being of practical importance itself, our approach should be useful in providing ‘ground truth’ inferences in test models to assess the affect of the approximations used in ABC methods.

References

- [1] I. Akhter and M. J. Black. Pose-conditioned joint angle limits for 3D human pose reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [2] H. C. Andersen. RATTLE: A velocity version of the SHAKE algorithm for molecular dynamics calculations. *Journal of Computational Physics*, 1983.
- [3] E. Barth, K. Kuczera, B. Leimkuhler, and R. D. Skeel. Algorithms for constrained molecular dynamics. *Journal of computational chemistry*, 1995.
- [4] M. A. Beaumont, W. Zhang, and D. J. Balding. Approximate Bayesian computation in population genetics. *Genetics*, 2002.
- [5] M. A. Brubaker, M. Salzmann, and R. Urtasun. A family of MCMC methods on implicitly defined manifolds. In *International Conference on Artificial Intelligence and Statistics*, 2012.
- [6] P. Diaconis, S. Holmes, and M. Shahshahani. Sampling from a manifold. In *Advances in Modern Statistical Theory and Applications*, pages 102–125. Institute of Mathematical Statistics, 2013.
- [7] P. J. Diggle and R. J. Gratton. Monte Carlo methods of inference for implicit statistical models. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 193–227, 1984.
- [8] S. Duane, A. D. Kennedy, B. J. Pendleton, and D. Roweth. Hybrid Monte Carlo. *Physics Letters B*, 1987.
- [9] H. Federer. *Geometric measure theory*. 1969.
- [10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014.
- [11] C. C. Gordon, T. Churchill, C. E. Clauser, B. Bradtmiller, J. T. McConville, I. Tebbets, and R. A. Walker. Anthropometric survey of US army personell: Final report. Technical report, United States Army, 1988.
- [12] C. Gourieroux, A. Monfort, and E. Renault. Indirect inference. *Journal of applied econometrics*, 8(S1):S85–S118, 1993.
- [13] M. Graham and A. Storkey. Asymptotically exact inference in differentiable generative models. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, pages 499–508, 2017.
- [14] C. Hartmann and C. Schutte. A constrained hybrid Monte-Carlo algorithm and the problem of calculating the free energy in several variables. *ZAMM-Zeitschrift für Angewandte Mathematik und Mechanik*, 2005.
- [15] D. P. Kingma and M. Welling. Auto-encoding variational Bayes. In *Proceedings of the 2nd International Conference on Learning Representations*, 2013.
- [16] B. Leimkuhler and C. Matthews. Efficient molecular dynamics using geodesic integration and solvent-solute splitting. In *Proceedings of the Royal Society A*, 2016.
- [17] B. Leimkuhler and S. Reich. *Simulating Hamiltonian dynamics*. Cambridge University Press, 2004.
- [18] B. J. Leimkuhler and R. D. Skeel. Symplectic numerical integrators in constrained Hamiltonian systems. *Journal of Computational Physics*, 1994.
- [19] T. Lelièvre, M. Rousset, and G. Stoltz. Langevin dynamics with constraints and computation of free energy differences. *Mathematics of computation*, 2012.
- [20] J.-M. Marin, P. Pudlo, C. P. Robert, and R. J. Ryder. Approximate Bayesian computational methods. *Statistics and Computing*, 2012.
- [21] P. Marjoram, J. Molitor, V. Plagnol, and S. Tavaré. Markov chain Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 2003.
- [22] S. Mohamed and B. Lakshminarayanan. Learning in implicit generative models. In *Proceedings of the International Conference on Learning Representations*, 2017.
- [23] I. Murray. Differentiation of the Cholesky decomposition. *arXiv preprint arXiv:1602.07527*, 2016.
- [24] I. Murray, R. P. Adams, and D. J. MacKay. Elliptical slice sampling. In *The Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*, volume 9 of *JMLR: W&CP*, pages 541–548, 2010.
- [25] I. Murray and M. M. Graham. Pseudo-marginal slice sampling. In *International Conference on Artificial Intelligence and Statistics*, 2016.
- [26] R. M. Neal. *MCMC using Hamiltonian dynamics*, pages 113–162. 2011.
- [27] M. Plummer, N. Best, K. Cowles, and K. Vines. CODA: Convergence diagnosis and output analysis for MCMC. *R News*, 6(1):7–11, 2006.
- [28] Theano development team. Theano: A Python framework for fast computation of mathematical expressions. *arXiv e-prints*, abs/1605.02688, 2016.
- [29] D. Tran, R. Ranganath, and D. M. Blei. Deep and hierarchical implicit models. *arXiv preprint arXiv:1702.08896*, 2017.

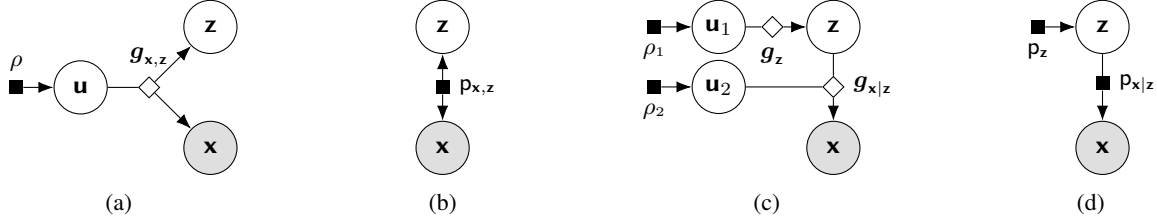


Figure 2. Factor graphs of undirected and directed generative models. Circular nodes represent random variables, filled square nodes probabilistic factors and unfilled diamonds deterministic factors. Shaded circular nodes are observed. Panel (a) shows the more general undirected model case in which observed variables \mathbf{x} and latent variables \mathbf{z} are jointly generated from inputs \mathbf{u} by mapping through a function $g_{\mathbf{x},\mathbf{z}}$, with (b) showing an equivalent factor graph after marginalising out the inputs \mathbf{u} . Panel (c) shows the directed model case in which we first generate the latent variables \mathbf{z} from a subset of the inputs \mathbf{u}_1 then generate the observed variables \mathbf{x} from \mathbf{z} and the remaining inputs \mathbf{u}_2 , with (d) showing a resulting natural directed factorisation of joint distribution when marginalising out \mathbf{u}_1 and \mathbf{u}_2 .

A. Directed and undirected generative models

In Section 1 we described the general case of a generative model in which both the observed and unobserved variables being jointly generated by a function $g_{\mathbf{x},\mathbf{z}}$. This structure is shown as a factor graph in Figure 2a and a corresponding factor graph for just \mathbf{x} and \mathbf{z} with \mathbf{u} marginalised out shown in Figure 2b. As briefly mentioned in Section 1, a common special case is when the input space is partitioned $\mathcal{U} = \mathcal{U}_1 \times \mathcal{U}_2$ and the unobserved variables \mathbf{z} are generated from a subset of the random inputs $\mathbf{u}_1 \in \mathcal{U}_1$, with the observed variables \mathbf{x} then generated from a function $g_{\mathbf{x}|\mathbf{z}}$ which takes as input both the remaining random variables $\mathbf{u}_2 \in \mathcal{U}_2$ and the generated unobserved variables \mathbf{z} , i.e. $\mathbf{x} = g_{\mathbf{x}|\mathbf{z}}(\mathbf{z}, \mathbf{u}_2) = g_{\mathbf{x}|\mathbf{z}}(g_{\mathbf{z}}(\mathbf{u}_1), \mathbf{u}_2)$. This is illustrated as a factor graph in Figure 2c. Again a corresponding factor graph with \mathbf{u} marginalised out is shown in Figure 2d, with in this case the structure of the generator making a *directed* factorisation in terms $p_{\mathbf{z}}$ and $p_{\mathbf{x}|\mathbf{z}}$ natural.

B. Evaluating the target density and its gradient

For the Hamiltonian dynamics we need to be able to evaluate the logarithm of the target density up to an additive constant and its gradient with respect to \mathbf{u} . Dropping the dependence of the Jacobian on \mathbf{u} for brevity we have that

$$\log \pi(\mathbf{u}) = \log \rho(\mathbf{u}) - \frac{1}{2} \log |\mathbf{J}_{g_{\mathbf{x}}} \mathbf{J}_{g_{\mathbf{x}}}^{\top}| - \log Z \quad (16)$$

where Z is the normalising constant for the density which is independent of \mathbf{u} .

In general evaluating the determinant $|\mathbf{J}_{g_{\mathbf{x}}} \mathbf{J}_{g_{\mathbf{x}}}^{\top}|$ has computational cost which scales as $\mathcal{O}(MN_{\mathbf{x}}^2)$. However as part of the constrained dynamics updates the lower-triangular Cholesky decomposition \mathbf{L} of $\mathbf{J}_{g_{\mathbf{x}}} \mathbf{J}_{g_{\mathbf{x}}}^{\top}$ is calculated. Using basic properties of the matrix determinant we have

$$\log \pi(\mathbf{u}) = \log \rho(\mathbf{u}) - \sum_{i=1}^{N_{\mathbf{x}}} \log(L_{ii}) - \log Z. \quad (17)$$

Given the Cholesky factor \mathbf{L} we can therefore evaluate the logarithm of the target density up to an additive constant at a marginal computational cost that scales linearly with $N_{\mathbf{x}}$. For the gradient we can use reverse-mode AD to calculate the gradient of (17) with respect to \mathbf{u} . This requires propagating partial derivatives through the Cholesky decomposition (23); efficient implementations for this are present in many automatic differentiation frameworks.

Alternatively using the standard result for derivative of a log determinant and the invariance of the trace to cyclic permutations we have that the gradient of the log determinant term in (16) can be manipulated in to the form

$$\frac{1}{2} \frac{\partial}{\partial u_i} \log |\mathbf{J}_{g_{\mathbf{x}}} \mathbf{J}_{g_{\mathbf{x}}}^{\top}| = \text{trace} \left((\mathbf{J}_{g_{\mathbf{x}}} \mathbf{J}_{g_{\mathbf{x}}}^{\top})^{-1} \frac{\partial \mathbf{J}_{g_{\mathbf{x}}} \mathbf{J}_{g_{\mathbf{x}}}^{\top}}{\partial u_i} \right) \quad (18)$$

$$= \text{trace} \left(\mathbf{J}_{g_{\mathbf{x}}}^{\top} (\mathbf{J}_{g_{\mathbf{x}}} \mathbf{J}_{g_{\mathbf{x}}}^{\top})^{-1} \frac{\partial \mathbf{J}_{g_{\mathbf{x}}}}{\partial u_i} \right) \quad (19)$$

We denote the matrix vectorisation operator vec such that for a $M \times N$ matrix \mathbf{A} , we have $\text{vec}(\mathbf{A}) = [A_{1,1}, \dots, A_{M,1}, A_{1,2}, \dots, A_{N,M}]^{\top}$. Then as the trace of a matrix product defines an inner product we have that

$\text{trace}(\mathbf{A}\mathbf{B}) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$. We can therefore write the gradient of the log determinant term as

$$\frac{1}{2} \frac{\partial}{\partial \mathbf{u}} \log |\mathbf{J}_{g_x} \mathbf{J}_{g_x}^\top| = \text{vec} \left(\mathbf{J}_{g_x}^\top (\mathbf{J}_{g_x} \mathbf{J}_{g_x}^\top)^{-1} \right)^\top \frac{\partial \text{vec}(\mathbf{J}_{g_x})}{\partial \mathbf{u}} \quad (20)$$

The matrix inside the left vec operator can be computed once by reusing the Cholesky factorisation of $\mathbf{J}_{g_x} \mathbf{J}_{g_x}^\top$ to solving the system of equations by forward and backward substitution. We then have an expression in the form of a vector-Jacobian product which is provided as an efficient primitive in many AD frameworks, e.g. as `Lop` in Theano, and like the gradient (which is actually a special case) can be evaluated at cost which is a constant over head of evaluating the forward function (i.e. the cost of evaluating \mathbf{J}_{g_x} here).

C. Constrained Hamiltonian Monte Carlo implementation

Algorithm 1 Constrained HMC in a differentiable generative model

Require:

g_x : observed variable generator function;
 ϕ : potential energy function $\phi(\mathbf{u}) = -\log \rho(\mathbf{u}) + \frac{1}{2} \log |\mathbf{J}_{g_x}(\mathbf{u}) \mathbf{J}_{g_x}(\mathbf{u})|$;
 \mathbf{x} : observed data values being conditioned on;
 \mathbf{u} : current chain state (model inputs) with $\|g_x(\mathbf{u}) - \mathbf{x}\|_\infty < \epsilon$;
 $(\varphi, \mathbf{J}, \mathbf{L})$: cached values of ϕ , \mathbf{J}_{g_x} and $\text{chol}(\mathbf{J}_{g_x} \mathbf{J}_{g_x}^\top)$ evaluated at \mathbf{u} ;
 ϵ : convergence tolerance for Newton iteration;
 M : number of Newton iterations to try before rejecting for non-convergence;
 δt : integrator time step;
 N_s : number of time steps to simulate;
 N_g : number of geodesic steps per time step.

Ensure:

\mathbf{u}_n : new chain state with $\|g_x(\mathbf{u}_n) - \mathbf{x}\|_\infty < \epsilon$;
 $(\varphi_n, \mathbf{J}_n, \mathbf{L}_n)$: values of ϕ , \mathbf{J}_{g_x} and $\text{chol}(\mathbf{J}_{g_x} \mathbf{J}_{g_x}^\top)$ evaluated at new \mathbf{u}_n .

```

 $n \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
 $\mathbf{p} \leftarrow \text{PROJECTMOM}(n, \mathbf{J}, \mathbf{L})$ 
 $\mathbf{u}_p, \mathbf{p}_p, \mathbf{J}_p, \mathbf{L}_p \leftarrow \text{SIMDYN}(\mathbf{u}, \mathbf{p}, \mathbf{J}, \mathbf{L})$ 
 $\varphi_p \leftarrow \phi(\mathbf{u})$ 
 $r \sim \mathcal{U}(0, 1)$ 
 $p_a \leftarrow \exp(\varphi + \frac{1}{2} \mathbf{p}^\top \mathbf{p} - \varphi_p - \frac{1}{2} \mathbf{p}_p^\top \mathbf{p}_p)$ 
if  $r < p_a$ 
     $\mathbf{u}_n, \varphi_n, \mathbf{J}_n, \mathbf{L}_n \leftarrow \mathbf{u}_p, \varphi_p, \mathbf{J}_p, \mathbf{L}_p$ 
else
     $\mathbf{u}_n, \varphi_n, \mathbf{J}_n, \mathbf{L}_n \leftarrow \mathbf{u}, \varphi, \mathbf{J}, \mathbf{L}$ 
    
```

```

function SIMDYN( $\mathbf{u}, \mathbf{p}, \mathbf{J}, \mathbf{L}$ )
     $\tilde{\mathbf{p}} \leftarrow \mathbf{p} - \frac{\delta t}{2} \nabla \phi(\mathbf{u})$ 
     $\mathbf{p} \leftarrow \text{PROJECTMOM}(\tilde{\mathbf{p}}, \mathbf{J}, \mathbf{L})$ 
     $\mathbf{u}, \mathbf{p}, \mathbf{J}, \mathbf{L} \leftarrow \text{SIMGEO}(\mathbf{u}, \mathbf{p}, \mathbf{J}, \mathbf{L})$ 
    for  $s \in \{2 \dots N_s\}$ 
         $\tilde{\mathbf{p}} \leftarrow \mathbf{p} - \delta t \nabla \phi(\mathbf{u})$ 
         $\mathbf{p} \leftarrow \text{PROJECTMOM}(\tilde{\mathbf{p}}, \mathbf{J}, \mathbf{L})$ 
         $\mathbf{u}, \mathbf{p}, \mathbf{J}, \mathbf{L} \leftarrow \text{SIMGEO}(\mathbf{u}, \mathbf{p}, \mathbf{J}, \mathbf{L})$ 
     $\tilde{\mathbf{p}} \leftarrow \mathbf{p} - \frac{\delta t}{2} \nabla \phi(\mathbf{u})$ 
     $\mathbf{p} \leftarrow \text{PROJECTMOM}(\tilde{\mathbf{p}}, \mathbf{J}, \mathbf{L})$ 
    return  $\mathbf{u}, \mathbf{p}, \mathbf{J}, \mathbf{L}$ 
    
```

```

function PROJECTMOM( $\mathbf{p}, \mathbf{J}, \mathbf{L}$ )
    return  $\mathbf{p} - \mathbf{J}^\top \mathbf{L}^{-\top} \mathbf{L}^{-1} \mathbf{J} \mathbf{p}$ 
    
```

function PROJECTPOS($\mathbf{u}, \mathbf{J}, \mathbf{L}$)

```

     $\delta \leftarrow g_x(\mathbf{u}) - \mathbf{x}$ 
     $i \leftarrow 0$ 
    while  $\|\delta\|_\infty > \epsilon \wedge i < M$ 
         $\mathbf{u} \leftarrow \mathbf{u} - \mathbf{J}^\top \mathbf{L}^{-\top} \mathbf{L}^{-1} \delta$ 
         $\delta \leftarrow g_x(\mathbf{u}) - \mathbf{x}$ 
         $i \leftarrow i + 1$ 
    if  $i = M$ 
        raise REJECTMOVE
    return  $\mathbf{u}$ 
    
```

function SIMGEO($\mathbf{u}, \mathbf{p}, \mathbf{J}, \mathbf{L}$)

```

    for  $i \in \{1 \dots N_g\}$ 
         $\tilde{\mathbf{u}} \leftarrow \mathbf{u} + \frac{\delta t}{N_g} \mathbf{p}$ 
         $\mathbf{u}' \leftarrow \text{PROJECTPOS}(\tilde{\mathbf{u}}, \mathbf{J}, \mathbf{L})$ 
         $\mathbf{J} \leftarrow \mathbf{J}_{g_x}(\mathbf{u}')$ 
         $\mathbf{L} \leftarrow \text{chol}(\mathbf{J} \mathbf{J}^\top)$ 
         $\tilde{\mathbf{p}} \leftarrow \frac{N_g}{\delta t} (\mathbf{u}' - \mathbf{u})$ 
         $\mathbf{p} \leftarrow \text{PROJECTMOM}(\tilde{\mathbf{p}}, \mathbf{J}, \mathbf{L})$ 
         $\mathbf{u}_r \leftarrow \mathbf{u}' - \frac{\delta t}{N_g} \mathbf{p}$ 
         $\mathbf{u}_r \leftarrow \text{PROJECTPOS}(\mathbf{u}_r, \mathbf{J}, \mathbf{L})$ 
        if  $\|\mathbf{u} - \mathbf{u}_r\|_\infty > \sqrt{\epsilon}$ 
            raise REJECTMOVE
         $\mathbf{u} \leftarrow \mathbf{u}'$ 
    return  $\mathbf{u}, \mathbf{p}, \mathbf{J}, \mathbf{L}$ 
    
```

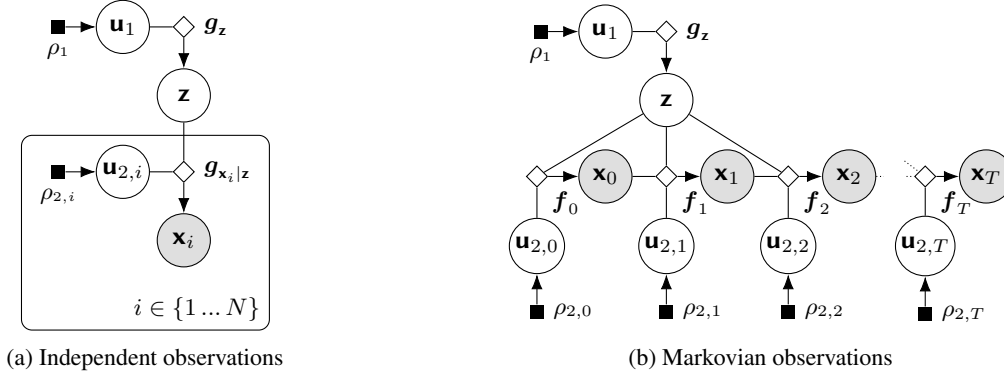


Figure 3. Factor graphs of examples of structured directed generative models.

D. Exploiting model structure

For larger systems, the Cholesky decomposition of the constraint Jacobian matrix product $\mathbf{J}_{g_x} \mathbf{J}_{g_x}^\top$ (line 42) will become a dominant cost, generally scaling cubically with N_x . In many models however conditional independency structure will mean that not all observed variables \mathbf{x} are dependent on all of the input variables \mathbf{u} and so the Jacobian \mathbf{J}_{g_x} has sparsity structure which can be exploited to reduce this worst-case cost.

In particular two common cases are directed generative models in which the observed variables \mathbf{x} can be split into groups $\{\mathbf{x}_g\}_{g=1}^G$ such that all of the \mathbf{x}_i are either conditionally independent given the latent variables $\mathbf{z} = g_z(\mathbf{u}_1)$ (for example independent and identically distributed observations), or each \mathbf{x}_i is conditionally independent of all $\{\mathbf{x}_j\}_{j < i-1}$ given \mathbf{x}_{i-1} and \mathbf{z} (most commonly Markov chains for example from a SDE model though observations with more general tree structured dependencies can also be ordered into this form). Figure 3 shows factor graphs for directed generative models with these two structures, with the conditional independencies corresponding to each \mathbf{x}_i being generated as a function of only a subset $\mathbf{u}_{2,i}$ of the random input variables \mathbf{u}_2 . Equivalently these structures correspond to generator functions g_x which can be expressed in one of the two forms

$$\mathbf{x}_i = g_{\mathbf{x}_i|\mathbf{z}}(\mathbf{z}, \mathbf{u}_{2,i}) \quad (\text{independent}) \quad (21)$$

$$\mathbf{x}_i = f_i(\mathbf{z}, f_{i-1}(\mathbf{z}, f_{i-2}(\dots), \mathbf{u}_{2,i-1}), \mathbf{u}_{2,i}) = g_{\mathbf{x}_i|\mathbf{z}}(\mathbf{z}, \mathbf{u}_{2,\leq i}) \quad (\text{Markov}). \quad (22)$$

For models with these structures the generator Jacobian $\mathbf{J}_{g_x} = [\frac{\partial g_x}{\partial \mathbf{u}_1} \mid \frac{\partial g_x}{\partial \mathbf{u}_2}]$ has a component $\frac{\partial g_x}{\partial \mathbf{u}_2}$ which is either block-diagonal (independent) or block-triangular (Markov). Considering first the simplest case where each $(\mathbf{x}_i, \mathbf{u}_{2,i})$ pair are single dimensional, the Cholesky factor of $\mathbf{J}_{g_x} \mathbf{J}_{g_x}^\top = \frac{\partial g_x}{\partial \mathbf{u}_1} \frac{\partial g_x}{\partial \mathbf{u}_1}^\top + \frac{\partial g_x}{\partial \mathbf{u}_2} \frac{\partial g_x}{\partial \mathbf{u}_2}^\top$ can then be computed by low-rank Cholesky updates of the triangular / diagonal matrix $\frac{\partial g_x}{\partial \mathbf{u}_2}$ with each of the columns of $\frac{\partial g_x}{\partial \mathbf{u}_1}$. As $\dim(\mathbf{u}_1) = L$ is often significantly less than the number of observations being conditioned on N_x , the resulting $\mathcal{O}(LN_x^2)$ cost of the low-rank Cholesky updates is a significant improvement over the original $\mathcal{O}(N_x^3)$. For cases in which each $(\mathbf{x}_i, \mathbf{u}_{2,i})$ pair are both vectors of dimension D (i.e. $N_x = GD$) and so $\frac{\partial g_x}{\partial \mathbf{u}_2}$ is block diagonal / triangular, then the Cholesky factorisation of $\frac{\partial g_x}{\partial \mathbf{u}_2} \frac{\partial g_x}{\partial \mathbf{u}_2}^\top$ can be computed at a cost $\mathcal{O}(GD^3)$ for block diagonal, and $\mathcal{O}(G^2D^3)$ for block triangular $\frac{\partial g_x}{\partial \mathbf{u}_2}$, with then again $\mathcal{O}(LN_x^2)$ cost low-rank updates of this Cholesky factor by the columns of $\frac{\partial g_x}{\partial \mathbf{u}_1}$ performed. When \mathbf{x}_i and $\mathbf{u}_{2,i}$ are vectors of differing dimensions, with generally in this case $\dim(\mathbf{u}_{2,i}) > \dim(\mathbf{x}_i)$ due to the requirement the total number of random inputs M is at least N_x , then though we could choose a subset of each $\mathbf{u}_{2,i}$ of equal dimension to \mathbf{x}_i so as to identify a block-triangular component, generally any gain from exploiting this structure will be minimal and in practice it seems likely to be more efficient to compute the Cholesky of $\mathbf{J}_{g_x} \mathbf{J}_{g_x}^\top$ directly.

E. Human pose and camera model inference

As a second experiment we considered inferring a three-dimensional human pose and camera model from two-dimensional projections of joint positions. We used a 19 joint skeleton model, with a learnt prior distribution over poses parametrised by 47 local joint angles \mathbf{z}_a . The pose model was learnt from the *PosePrior* motion capture data-set (1) with a Gaussian *variational autoencoder* (VAE) (15) trained to match the distribution of the motion capture joint angle data. The circular

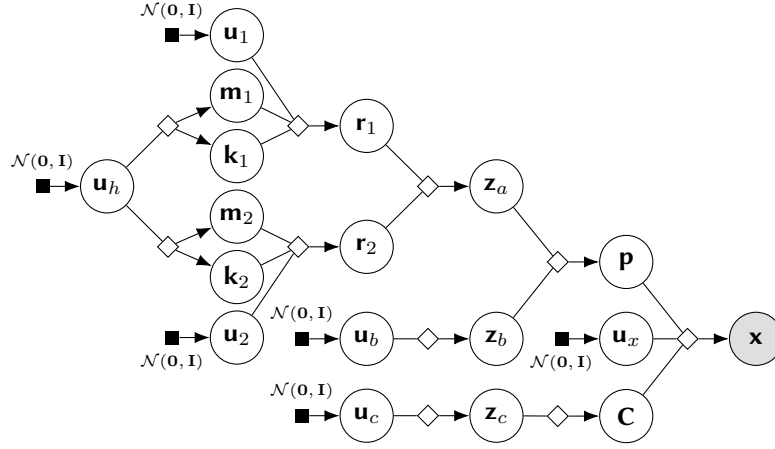


Figure 4. Factor graph of human pose differentiable generative model. The operations corresponding to the deterministic nodes (\diamond) in the graph are described in Algorithm 2.

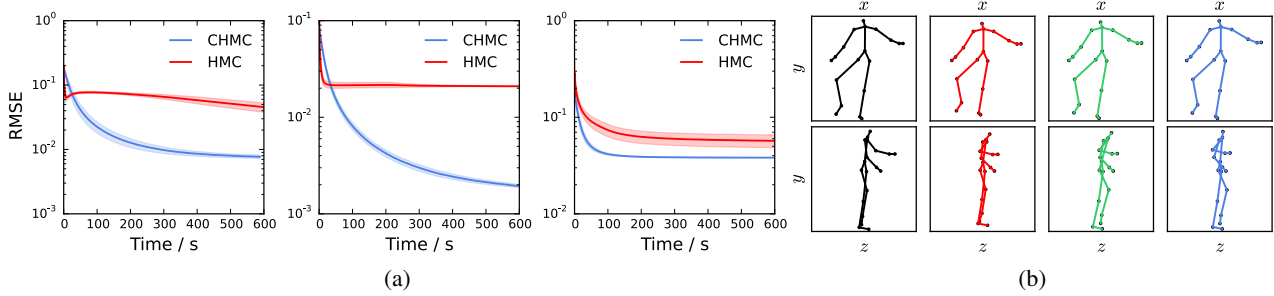


Figure 5. **Human pose** (a) RMSEs of 3D pose posterior mean estimates given binocular projections, using samples from our method (blue) versus running HMC in hierarchical model (red) for three different scenes sampled from the prior. Horizontal axes show computation time to produce number of samples in estimate. Solid curves are average RMSE over 10 runs with different seeds and shaded regions show ± 3 standard errors of mean. (b) Orthographic projections (top: front view, bottom: side view) of 3D poses consistent with monocular projections. Left most pair (black) shows pose used to generate observations, right three show constrained HMC samples.

topology of the angular data is poorly matched by the Euclidean space a Gaussian VAE typically learns a distribution on, and simply ‘unwrapping’ the angles to e.g. $[-\pi, \pi)$ leads to unnatural discontinuities at the $\pm\pi$ cut-point, this both making the initial learning problem challenging (as there is no in-built prior knowledge of continuity across the cut-point) and tending to lead to a learned latent space less amenable to MCMC inference as ‘nearby’ poses with one or more joint angles on opposite sides of the cut-point will likely end up corresponding to points far apart in the latent space.

During training we therefore mapped each vector of 47 joint angles $\mathbf{z}_a^{(i)}$ (corresponding to a single motion capture datapoint) to a pair of 47-dimensional vectors $(\mathbf{r}_1^{(i)}, \mathbf{r}_2^{(i)})$ by sampling a Gaussian random vector $\mathbf{n}^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and then computing $\mathbf{r}_1^{(i)} = \exp \mathbf{n}^{(i)} \odot \cos \mathbf{z}_a^{(i)}$ and $\mathbf{r}_2^{(i)} = \exp \mathbf{n}^{(i)} \odot \sin \mathbf{z}_a^{(i)}$ and training the VAE to maximise (a variational lower bound) on the joint marginal density of the $\{\mathbf{r}_1^{(i)}, \mathbf{r}_2^{(i)}\}_i$ pairs. At the cost of doubling the dimension, this leads to an embedding in a Euclidean space which does not introduce any arbitrary cut-points and empirically seemed to lead to better sample quality from the learned generative model compared to learning the angles directly. Given the trained model we can generate a vector of angles \mathbf{z}_a using the model by sampling a Gaussian code (latent representation) vector \mathbf{u}_h from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ then sampling a pair of 47-dimensional vectors \mathbf{r}_1 and \mathbf{r}_2 from the learnt Gaussian decoder model given \mathbf{u}_h (and further Gaussian random input vectors \mathbf{u}_1 and \mathbf{u}_2), and finally recovering an angle by computing $\mathbf{z}_a = \text{atan2}(\mathbf{r}_2, \mathbf{r}_1)$. The resulting distribution on \mathbf{z}_a is only implicitly defined, but the overall generative model is differentiable with respect to the input vectors $\mathbf{u}_h, \mathbf{u}_1$ and \mathbf{u}_2 .

The *PosePrior* motion capture data includes recordings from only a relatively small number of distinct actors and so limited variation in the ‘bone-lengths’ of the skeleton model. Therefore a separate log-normal model for the bone lengths \mathbf{z}_b was fitted using data from the *ANSUR* anthropometric data-set (11), due to symmetry in the skeleton thirteen independent lengths

Algorithm 2 Human pose model generator functions

Require:

$\{\mathbf{W}_\ell, \mathbf{b}_\ell\}_{\ell=0}^L$: parameters of learnt differentiable network model of pose angles;
 μ_b, Σ : mean and covariance of skeleton bone lengths;
 $\mu_{c,:2}, \sigma_{c,:2}$: means and standard deviations of camera world x, y coordinates;
 $\mu_{c,2}, \sigma_{c,2}$: mean and standard deviation of logarithm of camera world z coordinate;
 ϵ : image joint position observation noise standard deviation;
 JOINTPOSITIONS : maps pose angles and bone lengths to joint positions;
 CAMERAMATRIX : maps camera coordinates to projective camera matrix;
 PROJECT : uses camera matrix to map world coordinates to image coordinates;
 PARTITION : partitions a vector in a specified number of equal length parts;
 FLATTEN : flattens a multidimensional array to a vector.

```

function  $g_z(\mathbf{u}_h; \mathbf{u}_1; \mathbf{u}_2; \mathbf{u}_b; \mathbf{u}_c)$ 
     $\mathbf{h}_L \leftarrow \text{DIFFNET}(\mathbf{u}_h)$ 
     $\mathbf{m}_1, \mathbf{k}_1, \mathbf{m}_2, \mathbf{k}_2 \leftarrow \text{PARTITION}(\mathbf{h}_L, 4)$ 
     $\mathbf{r}_1 \leftarrow \exp(\mathbf{k}_1) \odot \mathbf{u}_1 + \mathbf{m}_1$ 
     $\mathbf{r}_2 \leftarrow \exp(\mathbf{k}_2) \odot \mathbf{u}_2 + \mathbf{m}_2$ 
     $\mathbf{z}_a \leftarrow \text{atan2}(\mathbf{r}_2, \mathbf{r}_1)$ 
     $\mathbf{z}_b \leftarrow \exp(\mu_b + \Sigma_b \mathbf{u}_b)$ 
     $\mathbf{z}_{c,:2} \leftarrow \sigma_{c,:2} \odot \mathbf{u}_{c,:2} + \mu_{c,:2}$ 
     $\mathbf{z}_{c,2} \leftarrow \exp(\sigma_{c,2} \mathbf{u}_{c,2} + \mu_{c,2})$ 
    return  $[\mathbf{z}_a; \mathbf{z}_b; \mathbf{z}_c]$ 
    
```

```

function DIFFNET( $\mathbf{u}_h$ )
     $\mathbf{h}_0 \leftarrow \tanh(\mathbf{W}_0 \mathbf{u}_h + \mathbf{b}_0)$ 
    for  $\ell \in \{1 \dots L-1\}$ 
         $\mathbf{h}_\ell \leftarrow \tanh(\mathbf{W}_\ell \mathbf{h}_{\ell-1} + \mathbf{b}_\ell) + \mathbf{h}_{\ell-1}$ 
    return  $\mathbf{W}_L \mathbf{h}_{L-1} + \mathbf{b}_L$ 

function  $g_{x|z}(\mathbf{z}_a; \mathbf{z}_b; \mathbf{z}_c, \mathbf{u}_x)$ 
     $\mathbf{P} \leftarrow \text{JOINTPOSITIONS}(\mathbf{z}_a, \mathbf{z}_b)$ 
     $\mathbf{C} \leftarrow \text{CAMERAMATRIX}(\mathbf{z}_c)$ 
     $\mathbf{X} \leftarrow \text{PROJECT}(\mathbf{C}, \mathbf{P})$ 
    return FLATTEN( $\mathbf{X}$ ) +  $\epsilon \mathbf{u}_x$ 
    
```

being specified. A simple pin-hole projective camera model with three position parameters \mathbf{z}_c and fixed focal-length was used³. A log-normal prior distribution was placed on the depth co-ordinate $z_{c,2}$ to enforce positivity with normal priors on the other two co-ordinates $z_{c,0}$ and $z_{c,1}$.

Given a generated triplet of joint-angles, bone length and camera parameters $\mathbf{z}_a, \mathbf{z}_b$ and \mathbf{z}_c , a simulated two-dimensional projection of the skeleton \mathbf{x} is produced by first mapping the joint-angles and bone-lengths to a 4×19 matrix of joint positions \mathbf{P} in (homogeneous) world-coordinates by recursing through the skeleton tree. A 3×4 projective camera matrix \mathbf{C} is generated from \mathbf{z}_c and then used to project the world-coordinate joint positions to a 2×19 matrix \mathbf{X} of joint positions in two-dimensional image-coordinates. The projected positions matrix \mathbf{X} is flattened to a vector and a Gaussian vector with standard deviation ϵ added to the projected position vector to give the $19 \times 2 = 38$ dimensional observed vector \mathbf{x} . The noise standard deviation ϵ is chosen so that the noise in the projected joint positions is non-obvious in generated projections. The overall corresponding model generator functions $g_{x|z}$ and g_z are described procedurally in Algorithm ?? and a factor graph summarising the relationships between the variables in the model shown in Figure 4.

Although the Gaussian observed output noise is necessarily not needed to apply our proposed constrained HMC method as the generator without the final additive noise still defines a valid differentiable generative model, using the noisy observation model means that an explicit hierarchical joint density on is defined on $\{\mathbf{x}, \mathbf{u}_h, \mathbf{r}_1, \mathbf{r}_2, \mathbf{z}_b, \mathbf{z}_c\}$ allowing comparison of our constrained HMC method with (non-constrained) HMC as a baseline. Further adding noise to the output ensures the generator Jacobian is full-rank everywhere and also significantly simplifies the process of finding an initial \mathbf{u} such that the generated \mathbf{x} matches observations.

We first considered binocular pose estimation, with the variables defining the three-dimensional scene information $\mathbf{z}_a, \mathbf{z}_b$ and \mathbf{z}_c , inferred given a pair of two-dimensional projections from two simulated cameras with a known offset in their positions (in this case the generator function is adjusted accordingly to output an $19 \times 2 \times 2 = 76$ dimensional observed vector \mathbf{x} corresponding to the concatenation of the flattened projected joint positions from both ‘cameras’). In this binocular case, the disparity in projected joint positions between the two projections gives information about the distances of the correspondings joints from the image plane in the depth direction and so we would expect the posterior distribution on the three-dimensional pose to be tightly distributed around the true values used to generate the observations. We compared our constrained HMC method to running standard HMC on the conditional density of $\{\mathbf{u}_h, \mathbf{r}_1, \mathbf{r}_2, \mathbf{z}_b, \mathbf{z}_c\}$ given \mathbf{x} .

³The camera orientation was assumed fixed to avoid replicating the degrees of freedom specified by the angular orientation of the root joint of the skeleton: only the relative camera–skeleton orientation is important.

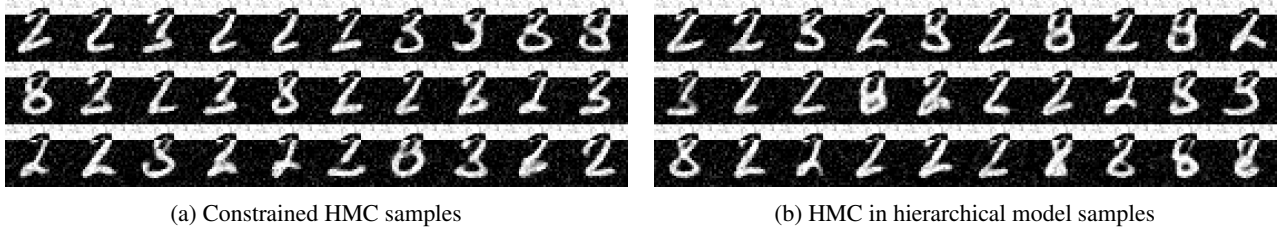


Figure 6. **MNIST** In-painting samples. The top black-on-white quarter of each image is the fixed observed region and the remaining white-on-black region the proposed in-painting. In left-right, top-bottom scan order the images in (a) are consecutive samples from a constrained HMC run; in (b) the images are every 40th sample from a HMC run to account for the $\sim 40\times$ shorter run-time per sample. All images in this run are show in Figure 7.

Figure 5a shows the *root mean squared error* (RMSE) between the posterior mean estimate of the three-dimensional joint positions and the true positions used to generate the observations as the number of samples included in the estimate increases for three test scenes. For both methods the horizontal axis has been scaled by run time. The constrained HMC method (blue curves) tends to give position estimates which converge more quickly to the true position. In this case standard HMC performs relatively poorly despite the significantly cheaper cost of each integrator step compared to the constrained dynamics. This is at least in part due to the small output noise standard deviation ϵ used which requires a small integrator step to be used to maintain reasonable accept rates.

We also considered inferring 3D scene information from a single 2D projection. Monocular projection is inherently information destroying with significant uncertainty to the true pose and camera parameters which generated the observations. Figure 5b shows pairs of orthographic projections of 3D poses: the left most column is the pose used to generate the projection conditioned on and the right three columns are poses sampled using constrained HMC consistent with the observations. The top row shows front x - y views, corresponding to the camera view though with a orthographic rather than perspective projection, the bottom row shows side z - y views with the z axis the depth from the camera. The dynamic is able to move between a range of plausible poses consistent with the observations while reflecting the inherent depth ambiguity from the monocular projection.

F. MNIST in-painting

As a final experiment we considered inferring in-paintings for a missing region of a digit image \mathbf{z} given knowledge of the rest of the pixel values \mathbf{x} . A Gaussian VAE trained on MNIST was used as the generative model, with a 50-dimensional hidden code \mathbf{h} . We compared our method to running HMC in the known conditional distribution on \mathbf{h} given \mathbf{x} (\mathbf{z} can then be directly sampled from its Gaussian conditional distribution given \mathbf{h}).

Example samples are shown in Figure 6. In this case the constrained and standard HMC approaches appear to be performing similarly, with both able to find a range of plausible in-paintings given the observed pixels. Without cost adjustment the standard HMC samples show greater correlation between subsequent updates, however for a fairer comparison the samples shown were thinned to account for the approximately $40\times$ larger run-time per constrained HMC sample. Although the constrained dynamic does not improve efficiency here neither does it seem to hurt it.

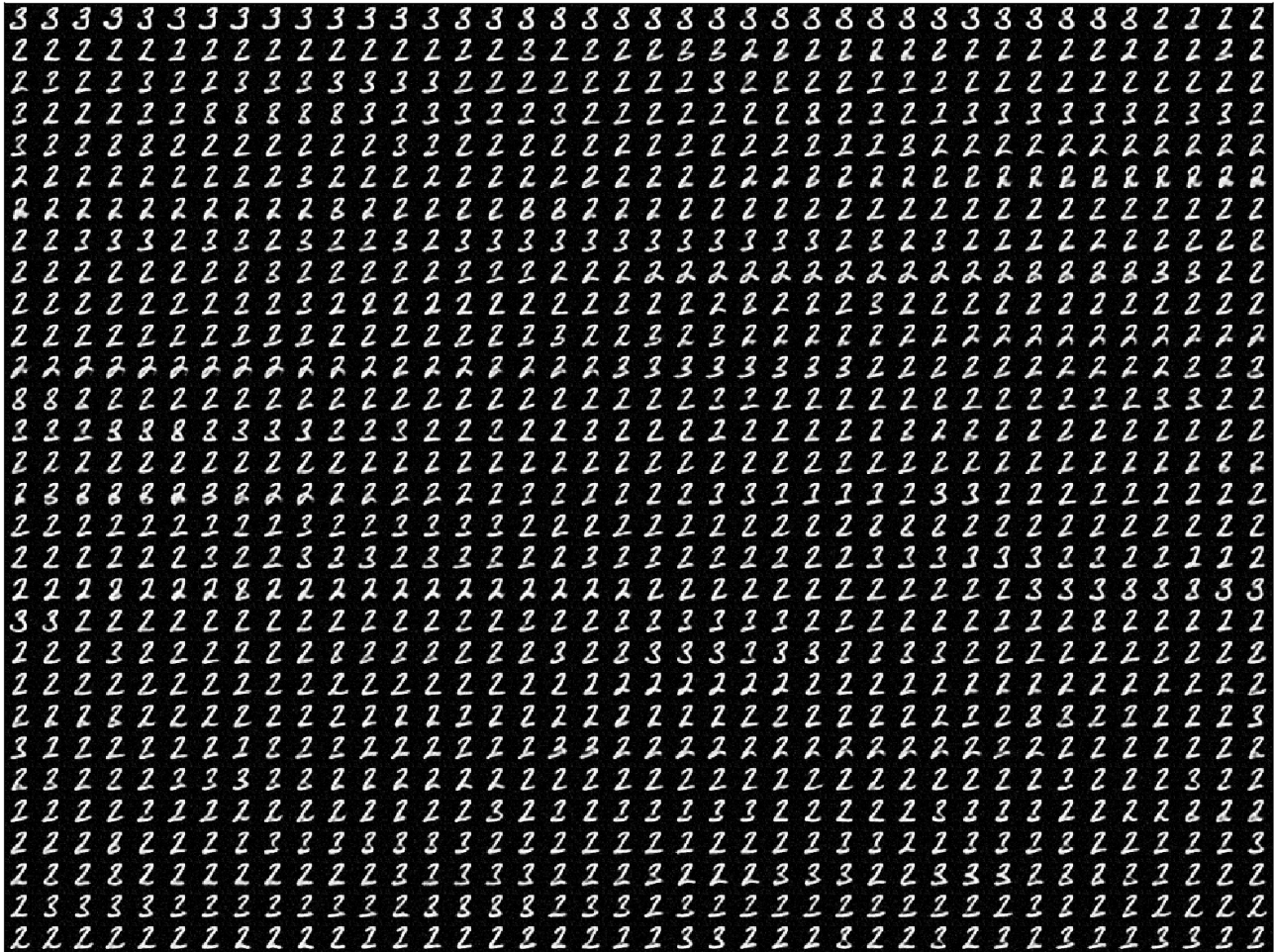


Figure 7. All in-painted MNIST image samples from a single HMC run previously shown thinned by factor 40 in Figure 6b, with samples in chain following left-to-right, top-to-bottom order.