
Hamiltonian ABC

Abstract

Approximate Bayesian computation (ABC) is a powerful and elegant framework for performing inference in simulation-based models. However, due to the difficulty in scaling likelihood estimates, ABC remains useful for relatively low-dimensional problems. [REASONS] We introduce Hamiltonian ABC (HABC), a set of likelihood-free algorithms that apply recent advances in scaling Bayesian learning using stochastic gradients with Hamiltonian dynamics to simulation models. [WHY GOOD: low nbr sims and non-random walk] We find that a small number forward simulations can effectively approximate the true ABC gradient, allowing Hamiltonian dynamics to efficiently traverse parameter spaces.

Models in simulation-based science domains are considered likelihood-free (LF) since they are not defined by a probabilistic model, but by a simulator that maps parameters to pseudo-observations. To perform parameter inference in the models, a powerful set of statistical procedures called Approximate Bayesian Computation (ABC) can be applied. Although there is a wide variety of ABC algorithms, they are dominated by SMC/PMC and MCMC. These have been successful, but their inefficiency of producing samples from the approximate posterior has limited ABC to simulations with relatively few parameters. Hamiltonian Monte Carlo (HMC) provides a framework for extending ABC to high-dimensions, but requires computing the gradient of the log-likelihood. With Hamiltonian ABC, we show how recent advances of stochastic gradient HMC algorithms applied to the big data domain can be also used for efficient sampling for LF models. We find that a small number forward simulations can effectively approximate the true ABC gradient, allowing Hamiltonian dy-

namics to efficiently traverse parameter spaces. We demonstrate HABC on several typical ABC problems, and show that HABC performs comparably to regular Bayesian inference on a high-dimensional problem from machine learning.

1 INTRODUCTION

In simulation-based science, models are defined by simulator and parameters. These are likelihood-free models because, in contrast to probabilistic models, their likelihoods are either intractable or must be approximated by simulations.

Recent advances in Bayesian inference within the *big data* domain have combined the efficiency of stochastic gradient algorithms used for optimization with Hamiltonian dynamics; these include Stochastic Gradient Langevin Dynamics (SGLD) [11], Stochastic Gradient HMC [2], and Stochastic Gradient Thermostats [3]. By using small batches of the full data set these methods are able to quickly compute a stochastic approximation to the true gradient, which with a small enough step-size and with an appropriate amount of injected noise, can produce samples from the true posterior. The Hamiltonian dynamics avoid random walk behavior by providing momentum to the sampler allowing it to reach far away regions of parameter space. Critically for the big data domain, the small step-sizes prevent the trajectories of the sampler from accumulating too much error and therefore they can avoid a Metropolis-Hastings correction step involving the full data set.

Approximate Bayesian Computation (ABC) is a set of inference procedures for so-called *likelihood-free* models [?]. In contrast to typical Bayesian modeling, where the parameters of a probabilistic model of data are the focus, the primary interest in ABC are the parameters of a simulation model. The simulator can generate pseudo or auxiliary data by running the simulator with a set of parameter values, but it cannot compute a likelihood directly. Inference involves running simulations, comparing pseudo data

with observations, and using a distance between the two as an approximate likelihood. Markov chain Monte Carlo (MCMC) or Sequential Monte Carlo (SMC) are the most common inference algorithms. It can be shown that if the statistics used in the likelihood function are sufficient, then these algorithms sample correctly from an approximation to the true posterior.

The approximation to the true posterior can be made arbitrarily small, in principle, by controlling the kernel width parameter ϵ . As $\epsilon \rightarrow 0$, however, simulated data must be closer and closer to the observations (or receive zero likelihood). The acceptance rate of a Markov chain for small ϵ can therefore be low. This inefficiency of ABC sampling is exacerbated by the random walk proposals. As the dimensionality of the parameter space grows, ϵ must grow to prevent low acceptance rates. This has perhaps prevented application of ABC to high-dimensional simulators. Our paper addresses this directly: can we use the sampling efficiency of Hamiltonian dynamics (i.e. high acceptance rates) in high-dimensional simulators?

In this paper we adopt the developments in scaling Bayesian inference to ABC. To do this we use forward simulations to approximate the likelihood-free gradient. The key difference with SG methods mentioned above is that the stochasticity of the gradient does not come from approximating the full data gradient with a batch gradient, but by the stochasticity of the simulator. It is therefore not the expense of the simulator (though this could very well be the case for many interesting simulation-based models, and we address this in Future Work) that requires an approximation to the gradient, but the likelihood-free nature of the problem. For high-dimensional parameter spaces, computing the gradients naively (i.e. by finite differences) can squash the gains brought by the Hamiltonian dynamics; to address this further approximations to the gradient can be used [?].

The approximate gradients used by HABC can be noisy and difficult to manage if care is not taken. To address this we show how a local Gaussian approximation of the simulator output can improve the stability of the gradient estimates over a direct estimate using the ABC kernel function.

A further innovation of this paper is the use of common random numbers (CRN) to improve the efficiency of the Hamiltonian dynamics. By using the same random seed for the finite difference approximation, the first order noise in the gradient is reduced; by keeping the random seeds fixed for several consecutive steps, the second order stochasticity of one gradient estimate to the next is eliminated. We show that this is a valid MCMC procedure.

In Section 3 we review three approaches to stochastic gradient inference using Hamiltonian dynamics: SGLD, SGHMC, and SGNHT. We briefly review ABC-MCMC in Section 2, followed by Hamiltonian ABC in Section 4,

where we will show how to improve the stability of the gradient estimates by using CRNs and local density estimators of the simulator. Extensions to high-dimensional parameter spaces are also discussed. In Section 5 we show how HABC behaves on a simple one-dimensional problem, then in Section 6 we compare HABC with ABC-MCMC for two problems: a low-dimensional model of chaotic population dynamics and a high-dimensional problem.

2 APPROXIMATE BAYESIAN COMPUTATION

Consider the Bayesian inference task of either drawing samples from or learning an approximate model of the following (usually intractable) posterior distribution:

$$\pi(\boldsymbol{\theta}|\mathbf{y}_1, \dots, \mathbf{y}_N) \propto \pi(\boldsymbol{\theta})\pi(\mathbf{y}_1, \dots, \mathbf{y}_N|\boldsymbol{\theta}) \quad (1)$$

where $\pi(\boldsymbol{\theta})$ is a prior distribution over parameters $\boldsymbol{\theta} \in \mathbb{R}^D$ and $\pi(\mathbf{y}_1, \dots, \mathbf{y}_N|\boldsymbol{\theta})$ is the likelihood of N data observations, where $\mathbf{y}_n \in \mathbb{R}^J$. In ABC, the vector of J observations are typically informative statistics of the raw observations. The simulator is treated as generator of random pseudo-observations, i.e. $\mathbf{x} \stackrel{\text{sim}}{\sim} \pi(\mathbf{x}|\boldsymbol{\theta})$ is a draw from the simulator. Discrepancies between the simulator outputs \mathbf{x} and the observations \mathbf{y} are scaled by a closeness parameter ϵ and treated as likelihoods. This is the equivalent to putting an ϵ -kernel around the observations, and using a Monte Carlo estimate of the likelihood using S draws of \mathbf{x} :

$$\pi_\epsilon(\mathbf{y}|\boldsymbol{\theta}) = \int \pi_\epsilon(\mathbf{y}|\mathbf{x})\pi(\mathbf{x}|\boldsymbol{\theta})d\mathbf{x} \approx \frac{1}{S} \sum_{s=1}^S \pi_\epsilon(\mathbf{y}|\mathbf{x}^{(s)}) \quad (2)$$

In ABC Markov chain Monte Carlo (MCMC) [5, 12, 8] the Metropolis-Hastings (MH) proposal distribution is composed of the product of the proposal for the parameters $\boldsymbol{\theta}$ and the proposal for the simulator outputs:

$$q(\boldsymbol{\theta}', \mathbf{x}^{(1)'}, \dots, \mathbf{x}^{(S)'}|\boldsymbol{\theta}) = q(\boldsymbol{\theta}'|\boldsymbol{\theta}) \prod_s \pi(\mathbf{x}^{(s)'}|\boldsymbol{\theta}') \quad (3)$$

Using this form of the proposal distribution, and using the Monte Carlo approximation eq 2, we arrive at the following Metropolis-Hastings accept-reject probability,

$$\alpha = \min \left(1, \frac{\pi(\boldsymbol{\theta}') \sum_{s=1}^S \pi_\epsilon(\mathbf{y}|\mathbf{x}^{(s)'})q(\boldsymbol{\theta}|\boldsymbol{\theta}')}{\pi(\boldsymbol{\theta}) \sum_{s=1}^S \pi_\epsilon(\mathbf{y}|\mathbf{x})q(\boldsymbol{\theta}'|\boldsymbol{\theta})} \right) \quad (4)$$

If the simulations are part of the Markov chain, the algorithm corresponds to the pseudo-marginal (PM) sampler [? 1], otherwise it is a marginal sampler [5, 8]. For this paper we will be interested in the PM sampler because this is equivalent to having the random states that generated the simulation outputs in the state of the Markov chain, which we will use within a valid ABC sampling algorithm in Section 4.

An alternative approach to computing the ABC likelihood is to estimate the parameters of a conditional model $\pi(\mathbf{x}|\boldsymbol{\theta})$, e.g. kernel density estimate [9] or a Gaussian model [14]. While either approach should be adequate and both have their own limits and advantages, for this paper we will use a Gaussian model. In ABC, using a conditional Gaussian model for $\pi(\mathbf{x}|\boldsymbol{\theta})$ is called a *synthetic likelihood* (SL) model [14]. For a SL log-likelihood model, we compute estimators of the first and second moments of $\pi(\mathbf{x}|\boldsymbol{\theta})$. The advantage is that for a Gaussian ϵ -kernel, we can convolve the two densities

$$\pi_{\epsilon}(\mathbf{y}|\boldsymbol{\theta}) = \int \mathcal{N}(\mathbf{y}|\mathbf{x}, \epsilon^2) \mathcal{N}(\mathbf{x}|\mu_{\boldsymbol{\theta}}, \sigma_{\boldsymbol{\theta}}^2) d\mathbf{x} \quad (5)$$

$$= \mathcal{N}(\mathbf{y}|\mu_{\boldsymbol{\theta}}, \sigma_{\boldsymbol{\theta}}^2 + \epsilon^2) \quad (6)$$

Of particular concern to this paper is the behavior of the log-likelihoods for different values of ϵ . In the ϵ -kernel case, the log-likelihood is very sensitive to small values of ϵ :

$$\log \pi_{\epsilon}(\mathbf{y}|\boldsymbol{\theta}) = \log \sum_s \mathcal{N}(\mathbf{y}|\mathbf{x}^{(s)}, \epsilon^2) \quad (7)$$

$$= \log \mathcal{N}(\mathbf{y}|\mathbf{x}^{(s)}, \epsilon^2) + \log(1 + H) \quad (8)$$

$$\approx -\log \epsilon - \frac{1}{2\epsilon^2}(\mathbf{y} - \mathbf{x}^{(m)})^2 \quad (9)$$

where m is the simulation that is closest to \mathbf{y} , H is a sum over terms close to 0. We can see that the log-likelihood can be set arbitrarily small by decreasing ϵ . On the other hand, by using a model of the simulation at $\boldsymbol{\theta}$

$$\log \pi_{\epsilon}(\mathbf{y}|\boldsymbol{\theta}) \approx -\frac{1}{2} \log(\sigma_{\boldsymbol{\theta}}^2 + \epsilon^2) - \frac{(\mathbf{y} - \mu_{\boldsymbol{\theta}})^2}{2(\sigma_{\boldsymbol{\theta}}^2 + \epsilon^2)} \quad (10)$$

For the SL model, ϵ acts as a smoothing term and can be set to small values with little change to the log-likelihood, as long as the SL estimators are fit appropriately. This insensitivity to ϵ will be used in Section 4 for estimating gradients of the ABC likelihood. Before describing HABC in full detail however, we now explain how scaling Hamiltonian dynamics in Bayesian learning can be accomplished using stochastic gradients from batched data.

3 SCALING BAYESIAN INFERENCE USING HAMILTONIAN DYNAMICS

Scaling Bayesian inference algorithms to massive datasets is necessary for their existing relevance in the so-called *big data* era. We now review the role stochastic gradient methods combined with Hamiltonian dynamics have played in recent advances in scaling Bayesian inference. Most importantly, these methods have combined the ability of HMC to explore high-dimensional parameter spaces with the computational efficiency of using stochastic gradients based on small batches of the full dataset. After an

overview of HMC, we will briefly describe stochastic gradient Hamiltonian dynamics (SGHDs), starting with using Langevin dynamics [11], then HMC with friction [2], and finally HMC with thermostats [3]. We will then make the connection between SGHDs and HABC in Section 4.

3.1 Hamiltonian Monte Carlo

Hamiltonian dynamics are often necessary to adequately explore the target distribution of high-dimensional parameter spaces. By proposing parameters that are far from the current location and yet have high acceptance probability, Hamiltonian Monte Carlo [4, 7] can efficiently avoid random walk behavior that can render proposals in high-dimensions inadequate.

HMC simulates the trajectory of a particle along a frictionless surface, using random initial momentum $\boldsymbol{\rho}$ and position $\boldsymbol{\theta}$. The Hamiltonian function computes the energy of the system and the dynamics govern how the momentum and position change over time. The continuous Hamiltonian dynamics can be simulated by discretizing time into small steps η . If η is small, the value of $\boldsymbol{\theta}$ at the end of a simulation can be used as proposals within the Metropolis-Hastings algorithm. Hamiltonian dynamics should propose $\boldsymbol{\theta}$ that are always accepted, but errors due to discretization may require a Metropolis-Hastings correction. It is this correction step that SGHD algorithms want to avoid as it requires computing the log-likelihood over the full data set.

More formally, the Hamiltonian $H(\boldsymbol{\theta}, \boldsymbol{\rho}) = U(\boldsymbol{\theta}) + K(\boldsymbol{\rho})$ is a function of the current potential energy $U(\boldsymbol{\theta})$ and kinetic energy $K(\boldsymbol{\rho}) = \boldsymbol{\rho}^T M^{-1} \boldsymbol{\rho} / 2$ (M is a diagonal matrix of masses which for presentation are set to 1). The potential energy is defined by negative log joint density of the data and prior:

$$U(\boldsymbol{\theta}) = -\log \pi(\boldsymbol{\theta}) - \sum_{i=1}^N \log \pi(\mathbf{y}_i|\boldsymbol{\theta}) \quad (11)$$

The Hamiltonian dynamics follow

$$d\boldsymbol{\theta} = \boldsymbol{\rho} dt \quad d\boldsymbol{\rho} = -\nabla U(\boldsymbol{\theta}) dt \quad (12)$$

in simulation $dt = \eta$.

3.2 Stochastic Gradient Hamiltonian Dynamics

If the log-likelihood over the full data set is replaced with a batch estimate, as is done for the following *stochastic gradient Hamiltonian dynamics* (SGHDs) algorithms, then the error in simulating the Hamiltonian dynamics comes not only from the discretization, but from the variance of the stochastic gradient. As long as this error is controlled, either by using small steps η (SGLD), or adding friction terms A (SGHMC), or using a thermostat ξ (SGNHT), the expensive MH correction step can be avoided and values of

θ from the Hamiltonian dynamics can be used as samples from the posterior.

SGHDs replace the full potential energy and its gradient with a batch approximation:

$$\begin{aligned}\hat{U}(\theta) &= -\log \pi(\theta) - \frac{N}{n} \sum_{i=h_1}^{h_n} \log \pi(\mathbf{y}_i|\theta) \quad (13) \\ \nabla \hat{U}(\theta) &= -\nabla \log \pi(\theta) - \frac{N}{n} \sum_{i=h_1}^{h_n} \nabla \log \pi(\mathbf{y}_i|\theta)\end{aligned} \quad (14)$$

where n is the batch size, and h_i are indices chosen randomly without replacement from $[1, N]$ (i.e. it defined a random batch).

Stochastic gradient Langevin dynamics [11] performs one full leap-frog step of HMC. Starting with a half step for the momentum, the update for θ is

$$\rho_t \sim \mathcal{N}(0, \mathbf{I}_p) \quad (15)$$

$$\rho_{t+\frac{1}{2}} = \rho_t - \eta \nabla \hat{U}(\theta_t)/2 \quad (16)$$

$$\theta_{t+1} = \theta_t + \eta \rho_{t+\frac{1}{2}} \quad (17)$$

It is not necessary to include ρ in the updates since there is only one step:

$$\theta_{t+1} = \theta_t + \eta \mathcal{N}(0, \mathbf{I}_p) - \eta^2 \nabla \hat{U}(\theta_t)/2 \quad (18)$$

One of the potential drawbacks of SGLD is that the momentum term is *refreshed* for every update of the parameters, and since this means the parameter update only uses the current gradient approximation, it limits the benefits of using Hamiltonian dynamics. On the other hand, this also prevents SGLD from accumulating errors in the Hamiltonian dynamics.

Stochastic Gradient HMC (SGHMC) [2] avoids ρ refreshment altogether. By applying HMC directly using the stochastic approximation \hat{U} and $\nabla \hat{U}$, which the authors call *naive SGHMC*, the variance of the gradient will introduce errors that left unaddressed will result in sampling from the incorrect target distribution. Under the assumption that $\nabla \hat{U}(\theta) = \nabla U(\theta) + \mathcal{N}(0, \mathbf{V}_\theta)$, where \mathbf{V}_θ is the covariance of the gradient approximation, and updates $\rho_{t+1} = \rho_t + \Delta \rho_t$ and $\theta_{t+1} = \theta_t + \eta \rho_{t+1}$, the change in momenta $\Delta \rho$ from one full step is

$$-\eta (\nabla U(\theta) + \mathcal{N}(0, \mathbf{V}_\theta)) = -\eta \nabla U(\theta) + \mathcal{N}(0, \eta^2 \mathbf{V}_\theta) \quad (19)$$

By adding a friction term \mathbf{B} to $\Delta \rho$ proportional to \mathbf{V}_θ , the correction step can be avoided

$$\Delta \rho = -\eta \mathbf{B} \rho_t - \eta \nabla U(\theta_t) + \mathcal{N}(0, 2\eta \mathbf{B}) \quad (20)$$

where $\mathbf{B} = \frac{1}{2} \eta \mathbf{V}_\theta$. In practice, since we can only estimate \mathbf{B} by some $\hat{\mathbf{B}}$ and can only compute \hat{U} , a user defined friction term \mathbf{C} is used (with $\mathbf{C} - \hat{\mathbf{B}}$ is semi-positive definite).

Thus the updates used for $\Delta \rho$ for SGHMC:

$$-\eta \mathbf{C} \rho_t - \eta \nabla \hat{U}(\theta_t) + \mathcal{N}(0, 2\eta(\mathbf{C} - \hat{\mathbf{B}})) \quad (21)$$

In our experiments we compute an online estimate $\hat{\mathbf{V}}$ and set $\mathbf{C} = c\mathbf{I}_p + \hat{\mathbf{V}}$.

Stochastic Gradient thermostats (SGT) [3] addresses the difficulty of estimating $\hat{\mathbf{B}}$ by introducing a scalar variable ξ who's addition to the Hamiltonian dynamics maintains the temperature of the system constant, i.e. it acts as a (Nose-Hoover) thermostat [?]. The update equations remain simple: initialize $\xi = \mathbf{C}$ (or c), then for $t = 1 \dots$

$$\rho_{t+1} = \rho_t - \eta \xi_t \rho - \eta \nabla \hat{U}(\theta_t) + \mathcal{N}(0, 2\eta_t \mathbf{C}) \quad (22)$$

$$\theta_{t+1} = \theta_t + \eta \rho_{t+1} \quad (23)$$

$$\xi_{t+1} = \xi_t + \eta (\rho_{t+1}^T \rho_{t+1} / p - 1) \quad (24)$$

In summary, the hyperparameters required for these algorithms are η and \mathbf{C} (for SGHMC and SGNHT only), and in practice, some way of estimating $\hat{\mathbf{V}}$ for SGHMC.

4 HAMILTONIAN ABC

In the application of stochastic gradient Hamiltonian dynamics to ABC, the stochasticity of the gradient due to batch sizes is instead caused by two other factors. The first factor, is due to the finite number of forward simulations that are used to approximate the log-likelihood function. We will explain below how to perform ABC gradient estimates based on finite difference approximations on two log-likelihood formulations: the ϵ -kernel density model and on a Gaussian conditional model, both are based on S random simulations. The second factor is due to the dimension of D . For large D , a finite difference estimate requires $2p$ forward simulations for each random simulation, by using a further stochastic approximation to the gradient, called Simultaneous Perturbation Stochastic Approximation (SPSA) [?], an unbiased estimate of the gradient can be computed with $2q$ forward simulations, where $q \geq 1$ and usually $Q \ll D$.

To begin, we define $\theta_d^+ = \theta + c\Delta_d$ and $\theta_d^- = \theta - c\Delta_d$, where c is a small positive scalar, e.g. $c = 0.001$, and Δ_d is a vector of size D that, for FD estimates is 1 for entry d and 0 otherwise. Our gradient estimate require simulating S times at θ^+ and θ^- , generating sets of pseudo-data $\{\mathbf{x}_{sd}^+\}$ and $\{\mathbf{x}_{sd}^-\}$. These are then plugged into the log-likelihood estimate for \mathbf{y} , which we call \mathcal{L} :

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta_d} \approx \frac{\mathcal{L}(\theta_d^+) - \mathcal{L}(\theta_d^-)}{2c} \quad (25)$$

For the ϵ -kernel density model it is easy to show that the gradient estimate is sensitive to ϵ in this case:

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta_d} \approx \frac{1}{2c\epsilon^2} \left(-\frac{1}{2}(\mathbf{y} - \mathbf{x}_{sd}^+)^2 + \frac{1}{2}(\mathbf{y} - \mathbf{x}_{sd}^-)^2 \right) \quad (26)$$

So by making ϵ arbitrarily small we can make the gradient extremely sensitive to our simulations. Intuitively, we want the gradient to be fairly stable no matter what values of ϵ we use. To do this we will replace log-likelihood with a conditional model. We can see that ϵ now can be set very small as it now acts as a smoother for our estimate, resulting in a robust gradient.

4.1 Gradients in High Dimensions

In the gradient-free setting, Spall [?] provides a stochastic approximate to the true gradient using only 2 forward simulations (function evaluations). This is in contrast to multivariate finite-difference stochastic approximation FDSA [?] requiring $2p$ evaluations.

The gradient estimate is

$$\hat{g}_t(\theta_t) = \begin{bmatrix} \frac{y_t^+ - y_t^-}{2c_{t1}\Delta_{t1}} \\ \vdots \\ \frac{y_t^+ - y_t^-}{2c_{tp}\Delta_{tp}} \end{bmatrix} \quad (27)$$

where c_{tp} is a step-size that is usually constant for all dimensions p , but can be different (as shown in this case); $\Delta_{tp} \in -1, +1$ is a *perturbation mask* (called symmetric Bernouilli variables by Spall), i.e. $\Delta_{tp} \sim 2 * \text{Bernouilli}(0.5) - 1$; and y_t^\pm are function evaluations:

$$y_t^+ = L(\theta + c_t \Delta_t) \quad (28)$$

$$y_t^- = L(\theta - c_t \Delta_t) \quad (29)$$

A way of reducing the noise in the gradient is by averaging over q draws of Δ_t :

$$\hat{g}_t(\theta_t) = \frac{1}{q} \sum_{j=1}^q \hat{g}_t^j(\theta_t) \quad (30)$$

where for each j new perturbation masks are drawn.

For ABC, the gradients we compute are stochastic because we must approximate $\pi(\mathbf{x}|\theta)$ with forward simulations. We can make the approximation arbitrarily good by running more simulations, but this is often unnecessary for sampling, because we can make use of the uncertainty in our current approximation to decide if our sample is correct or not (cite UAI2014?).

- Log likelihood is much noisier in 1, and the 2 there are good reasons why this is a better choice: is no longer depends on *epsilon* (as the true gradient shouldn't), convolving the conditional density is the right thing to do, if possible (though it might not be necessary to perform ABC)
- In both cases, we rewrite the expression conditioning on a set of random seeds ω . Explain the reason

for this. By observing the simulation at different θ but same seeds, we can eliminate much of the noise in the markov chain. By treating ω as part of the state, we can smoothly move along noiseless outputs for several iterations and then step.

4.2 Stochastic ABC Gradients

For low-dimensional problems a finite difference (FD) approximation to ∇U . For each dimension p , forward simulate for each random seed twice at plus c and minus c (only along dimension p): $x_\pm^{(s)} \sim \pi(\mathbf{x}|\theta \pm c\Delta_p, \omega_s)$. These simulations are then plugged into the appropriate estimator (logsumsepx or SL). $\hat{U}_{K\epsilon}$:

$$-\log \pi(\theta \pm c\Delta_p) - \log \sum_s \pi_\epsilon(\mathbf{y}|f(\theta \pm c\Delta_p, \omega_s)) \quad (31)$$

or \hat{U}_{SL}

$$-\log \pi(\theta \pm c\Delta_p) - \log \mathcal{N}(\mathbf{y}|\mu_{\theta \pm c\Delta_p}, \sigma_{\theta \pm c\Delta_p}^2 + \epsilon^2) \quad (32)$$

$$\hat{U}(\theta) = -\log \pi(\theta) - \log \pi_\epsilon(\mathbf{y}|\theta) \quad (33)$$

$$\nabla \hat{U}(\theta) = -\nabla \log \pi(\theta) - \frac{N}{n} \sum_{i=h_1}^{h_n} \nabla \log \pi(\mathbf{y}_i|\theta) \quad (34)$$

4.3 Common Random Numbers

[CITE NEAL How to view an MCMC...] Another variance reduction technique called the method of *common random numbers* (CRN) [?] sets a common seed for the random number generator (RNG) for both calls to the simulator. This technique can remove the effect of the simulator noise in the gradient estimate, leaving on the randomness of the perturbation masks as the source of noise. Consider using SPSA instead of SGD: the CRN technique is equivalent of using the same batch of data vectors to evaluate the log-likelihood which is a sensible approach.

The ABC likelihood is proportional to the convolution between a kernel around the observation statistics \mathbf{y} and the generator of pseudo-statistics \mathbf{x} , i.e. the simulator. Since pseudo-statistics can only be generated, we resort to a Monte Carlo estimate for the ABC likelihood based on S draws from the simulator using the same parameter vector θ but different random seeds:

$$\pi(\mathbf{y}|\theta) = \int \pi_\epsilon(\mathbf{y}|\mathbf{x}) \pi(\mathbf{x}|\theta) d\mathbf{x} \quad (35)$$

$$\approx \frac{1}{S} \sum_{s=1}^S \pi_\epsilon(\mathbf{y}|\mathbf{x}^{(s)}) \quad (36)$$

where $\mathbf{x}^{(s)} = \pi(\mathbf{x}|\theta, \omega_s)$. Note we have passed a unique random seed ω_s into the simulator, making $\mathbf{x}^{(s)}$ a deterministic function of the simulator with parameters $\{\theta, \omega_s\}$. It

will be useful to rewrite this deterministic simulator function as $f: \mathbf{x}^{(s)} = f(\boldsymbol{\theta}, \omega_s)$, allowing us to rewrite the likelihood as

$$\pi(\mathbf{y}|\boldsymbol{\theta}) \approx \frac{1}{S} \sum_{s=1}^S \pi_{\epsilon}(\mathbf{y}|f(\boldsymbol{\theta}, \omega_s)) \quad (37)$$

4.4 Algorithms

5 Demonstration

[FROM UAI] In this illustrative problem we infer the rate of an exponential distribution under a gamma prior with shape α and rate β , having N observations at the true rate θ^* ; this is the *exponential example* in [10]. Let \mathbf{w} be a vector of N draws from an exponential distribution, i.e. $w_n \sim \text{Exp}(\theta)$. The posterior is a gamma distribution with shape $\alpha + N$ and rate $\beta + \sum w_n$. To use this problem with ABC, we use the exponential draws as the simulator and the mean of \mathbf{w} as the statistic y and assume that N is known. The inference problem is therefore to sample from $p(\theta|y, \alpha, \beta, N)$.

5.1 First order CRNs

Using same seed in a gradient estimate at time t (versus not using them). This makes sense if we compare to SG-data where a single batch is used for computing the gradient. If a FD method was used, we would use the same batch.

5.2 Second order CRNs

Using same seed in a gradient estimate at multiple time steps (versus not changing them – the previous example). Analogous to keeping the same batch for consecutive steps in SG-data.

5.3 Comparing Gradient Estimates with SL

- Can estimate ABC posterior by CLT of Gammas (sum of exponential). For large S , the variance of the SL approximates the “true” CLT approximation. Thus we can set ϵ much smaller for the SL case. The kernel log-sumexp suffers when ϵ is small, using only the closest simulation to \mathbf{y} . Instead the SL estimates a density around the simulations and the likelihood is less sensitive to ϵ .
- Better tolerance to low ϵ using SL.

5.4 Using variance of gradients in SGHMC

Compute online \hat{B} then fix.

5.5 Visualizing Noise from Gradient versus Injected Noise

6 Experiments

6.1 Blowfly

[FROM UAI] Adult blowfly populations exhibit dynamic behavior for which several competing population models exist. In this experiment, we use observational data and a simulation model from Wood [14], based on their improvement upon previous population dynamics theory. Population dynamics are modeled using (discretized) differential equations that can produce chaotic behavior for some parameter settings. An example blowfly replicate series is shown in Figure ??, along with times-series generated by a sample from $\pi(\boldsymbol{\theta}|\mathbf{y})$ using GPS-ABC.

In [14] there are several explanations of the population dynamics, corresponding to different simulations and parameters. We concentrate on the equation (1) in section 1.2.3 of the supplementary information, considered “a better alternative model” by the author. The population dynamics equation generates N_1, \dots, N_T using the following update rule:

$$N_{t+1} = P N_{t-\tau} \exp(-N_{t-\tau}/N_0) e_t + N_t \exp(-\delta \epsilon_t)$$

where $e_t \sim \mathcal{G}(1/\sigma_p^2, 1/\sigma_p^2)$ and $\epsilon_t \sim \mathcal{G}(1/\sigma_d^2, 1/\sigma_d^2)$ are sources of noise, and τ is an integer (not to be confused with the τ used as the MH acceptance threshold in our algorithms). In total, there are 6 parameters $\theta = \{\log P, \log \delta, \log N_0, \log \sigma_d, \log \sigma_p, \tau\}$. See [14] for further details about the significance of the parameters. We put Gaussian priors over all the parameters (with Gaussian proposal distributions), except for τ which has a Poisson prior (and a left/right increment proposal). Time-series generated with parameters from this prior distribution produce extremely varied results, some are chaotic, some are degenerate, etc. Modeling this simulator is *very* challenging.

As with any ABC problem the choice of statistics is important as it directly affects the quality of the results. It is also non-trivial and requires careful thought and sometimes trial and error. In total there are 10 statistics: the log of the mean of all 25% quantiles of $N/1000$ (4 statistics), the mean of the 25% quantiles of the first-order differences of $N/1000$ (4 statistics), and the maximal peaks of smoothed N , with 2 different thresholds (2 statistics). With these statistics it is possible to reproduce time-series that appear similar to the observations. Note that these statistics are different from Wood’s, but they capture similar time-series features and are sufficient to produce credible population dynamics. [END FROM UAI]

6.2 Bayesian Autoencoders

We perform Bayesian inference on an autoencoder neural network using MNIST images as observations. We apply two inference algorithms: SGNHT and HABC. For SGNHT, we use the known gradient and for HABC we apply 2SPSA. After 1000 iterations we collect parameter vectors every 100 iterations for a total of 100 sets. (Do we use SGHMC to estimate Bhat in first 1000 iterations?). Use same batches at each iteration for SGNHT and HABC?

results: convergence of training error (batches), convergence in test error using parameters collected (record nbr of forward or backward passes required), show mean and variance of filters (how?), compare with using the last parameter sample only (stochastic optimization).

7 DISCUSSION

8 CONCLUSION

- New set of ABC algorithms for efficient exploration of posterior distribution.
- Key: opens the door to high-dimensional ABC inference.

A further complication, which we do not address on this paper, is the forward simulation might be, and for many interesting problems is, expensive. We would therefore prefer to resort to some sort of surrogate modeling [6, 13?] to avoid simulations whenever possible. We address this scenario in Section ?? . Shown for relatively fast simulators. Extensions to surrogates with GPS (Meeds2014GpsUai, rasmussen:2003, wilkinson:2014) can use derivative of GP for likelihood as GP of gradient. These could be matched by occasional gradient estimates (Osborne).

References

References

- [1] Andrieu, C. and Roberts, G. The pseudo-marginal approach for efficient monte carlo computations. *The Annals of Statistics*, 37(2):697–725, 2009.
- [2] Chen, Tianqi, Fox, Emily B, and Guestrin, Carlos. Stochastic gradient hamiltonian monte carlo. 2014.
- [3] Ding, Nan, Fang, Youhan, Babbush, Ryan, Chen, Changyou, Skeel, Robert D, and Neven, Hartmut. Bayesian sampling using stochastic gradient thermostats. In *Advances in Neural Information Processing Systems*, pp. 3203–3211, 2014.
- [4] Duane, Simon, Kennedy, Anthony D, Pendleton, Brian J, and Roweth, Duncan. Hybrid monte carlo. *Physics letters B*, 195(2):216–222, 1987.
- [5] Marjoram, Paul, Molitor, John, Plagnol, Vincent, and Tavaré, Simon. Markov chain monte carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328, 2003.
- [6] Meeds, Edward and Welling, Max. GPS-ABC: Gaussian process surrogate approximate bayesian computation. *Uncertainty in AI*, 2014.
- [7] Neal, Radford M. Mcmc using hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo*, 2, 2011.
- [8] Sisson, Scott A and Fan, Yanan. Likelihood-free markov chain monte carlo. *Arxiv preprint arXiv:1001.2058*, 2010.
- [9] Turner, Brandon M. and Sederberg, Per B. A generalized, likelihood-free method for posterior estimation. *Psychonomic Bulletin & Review*, 21(2):227–250, 2014.
- [10] Turner, Brandon M and Van Zandt, Trisha. A tutorial on approximate bayesian computation. *Journal of Mathematical Psychology*, 56(2):69–85, 2012.
- [11] Welling, Max and Teh, Yee W. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 681–688, 2011.
- [12] Wilkinson, R. Approximate bayesian computation (ABC) gives exact results under the assumption of model error. *Statistical Applications in Genetics and Molecular Biology*, 12(2):129–142, 2013.
- [13] Wilkinson, R. Accelerating abc methods using gaussian processes. *AISTATS*, 2014.
- [14] Wood, Simon N. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466(7310):1102–1104, 2010.