

Introduction

Objective Designing simple and efficient Bayesian inference algorithms for deep learning models.

Main idea

- Multivariate Stochastic Gradient Nośe-Hoover Thermostat
- (m-SGNHT) with more accurate numerical implementation.
- From Euler integrator to Symmetric Splitting integrator

Illustrations

Double-well Potential Function

- better approximation.
- allows large updates.



Figure: Samples of $\rho(\theta)$ with SSI (1st column) and Euler integrator (2nd column), and the estimated thermostat variable over iterations (3rd column).

Convolutional Neural Networks



Figure: Learning curves of CNN for different step sizes.

High-Order Stochastic Gradient Thermostats for Bayesian Learning of Deep Models

Chunyuan Li, Changyou Chen, Kai Fan and Lawrence Carin Duke University, Durham NC 27708, USA



Algorithms

 θ : model parameter p: momentum ξ : thermostat $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \boldsymbol{p}_t h$ $\{\boldsymbol{p}_{t+1} = \boldsymbol{p}_t - \nabla_{\boldsymbol{\theta}} \tilde{U}_t(\boldsymbol{\theta}_{t+1})h - \mathsf{diag}(\boldsymbol{\xi}_t)\boldsymbol{p}_th + \sqrt{2D}\boldsymbol{\zeta}_{t+1}\}$ $\boldsymbol{\xi}_{t+1} = \boldsymbol{\xi}_t + (\boldsymbol{p}_{t+1} \odot \boldsymbol{p}_{t+1} - 1) h$

Euler Integrator Symmetric Spltting Integrator $A: \theta_{t+1/2} = \theta_t + p_t h/2, \, \xi_{t+1/2} = \xi_t + (p_t \odot p_t - 1) h/2 \to$ $B: \boldsymbol{p}_{t+1/3} = \exp(-\boldsymbol{\xi}_{t+1/2}h/2) \odot \boldsymbol{p}_t \rightarrow$ $O: \boldsymbol{p}_{t+2/3} = \boldsymbol{p}_{t+1/3} - \nabla_{\boldsymbol{\theta}} \tilde{U}_t(\boldsymbol{\theta}_{t+1/2})h + \sqrt{2D}\boldsymbol{\zeta}_{t+1} \rightarrow \boldsymbol{\zeta}_{t+1/2}$ $B: \boldsymbol{p}_{t+1} = \exp(-\boldsymbol{\xi}_{t+1/2}h/2) \odot \boldsymbol{p}_{t+2/3} \rightarrow$ $A: \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_{t+1/2} + \boldsymbol{p}_{t+1}h/2, \, \boldsymbol{\xi}_{t+1} = \boldsymbol{\xi}_{t+1/2} + (\boldsymbol{p}_{t+1} \odot \boldsymbol{p}_{t+1} - 1) h/2$

where h is stepsize, D is diffusion factor, and $\boldsymbol{\zeta}$ is Gaussian noise

Roles of numerical integrators Under certain assumptions, the Bias and MSE of a SG-MCMC algorithm with stepsize h and a Kthorder integrator are:

Bias: $|\mathbb{E}\hat{\phi} - \bar{\phi}| = k$

MSE: $\mathbb{E}\left(\hat{\phi}-\hat{\phi}
ight)^2$ =

where \mathcal{B}_{bias} and \mathcal{B}_{mse} are functions dent of K.

Remark 1 Robustness: that mSGNHT-S is more robust to the stepsizes than mSGNHT-E. $T^{-1/2}$, respectively, indicating mSGNHT-S converges faster. **Remark 3 Measure Accuracy:** In the limit of infinite time $(T \rightarrow$ more accurate than mSGNHT-E.

Advantages for Deep Learning

- It tolerates gradients of various magnitudes, providing a potential solution to mitigate the *vanishing/exploding gradients* problem
- It can estimate model parameters faster and more accurately

Theoretical Justification

$$egin{split} \mathcal{B}_{\mathsf{bias}} + O(h^K) \ &= \mathcal{B}_{\mathsf{mse}} + O(h^{2K}) \ , \ & \mathsf{depending on } (h,T) \ \mathsf{but indepen-substitution} \end{split}$$

The bias and MSE of mSGNHT-S is bounded as: $\mathcal{B}_{bias}+O(h^2)$ and $\mathcal{B}_{mse}+O(h^4)$, compared to $\mathcal{B}_{bias}+O(h)$ and $\mathcal{B}_{mse} + O(h^2)$ for the mSGNHT-E, respectively. This indicates Remark 2 Convergence Rate: The higher order a numerical integrator is, the faster its optimal convergence rate is. Convergence rates in term of *bias* for mSGNHT-S and mSGNHT-E are $T^{-2/3}$ and

 ∞), the terms $\mathcal{B}_{\mathsf{bias}}$ and $\mathcal{B}_{\mathsf{mse}}$ in Lemma 1 vanish, leaving only the $O(h^K)$ terms. This indicates mSGNHT-S is an order of magnitude



A 3-layer Deep Sigmoid Belief Networks (DSBN) is inferred by mSGNHT.

This research was supported by ARO, DARPA, DOE, NGA ONR and NSF.



Acknowledgements