# AMAGOLD:

Amortized Metropolis Adjustment for Efficient Stochastic Gradient MCMC

Ruqi Zhang, **A. Feder Cooper**, Christopher De Sa

June 2020

Cornell University

Markov chain Monte Carlo (MCMC)

Stochastic gradient MCMC (SG-MCMC)

Metropolis-Hastings (M-H) correction

**Markov chain Monte Carlo (MCMC)**

- Are a family of sampling methods popular in Bayesian inference
- Approximate computationally intractable posterior
- Depend on size of inference task's dataset

**Stochastic gradient MCMC (SG-MCMC)**

**Metropolis-Hastings (M-H) correction**

## Markov chain Monte Carlo (MCMC)

- Are a family of sampling methods popular in Bayesian inference
- Approximate computationally intractable posterior
- Depend on size of inference task's dataset

## Stochastic gradient MCMC (SG-MCMC)

- Use subsampling to decouple from dataset size
- Provide speed-ups, but at the cost of introducing bias

## Metropolis-Hastings (M-H) correction

## Markov chain Monte Carlo (MCMC)

- Are a family of sampling methods popular in Bayesian inference
- Approximate computationally intractable posterior
- Depend on size of inference task's dataset

## Stochastic gradient MCMC (SG-MCMC)

- Use subsampling to decouple from dataset size
- Provide speed-ups, but at the cost of introducing bias

## Metropolis-Hastings (M-H) correction

- Removes bias by rejecting fraction of Markov chain's transitions
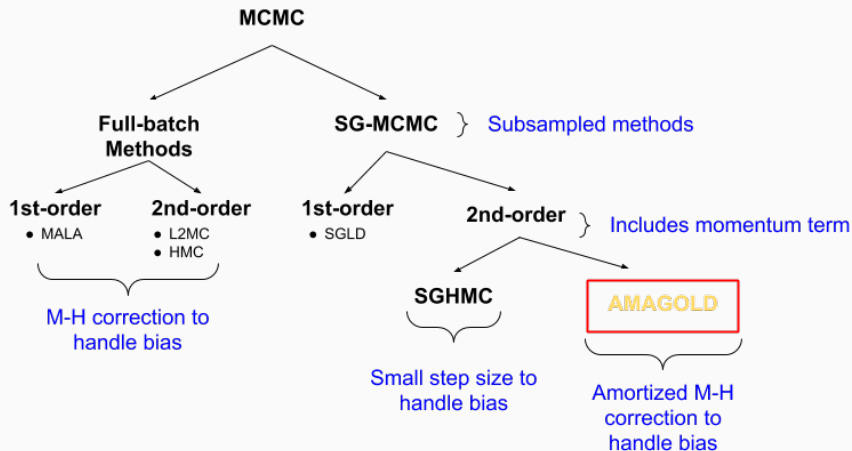- Reintroduces dependency on dataset size

*Question*: Is there a way to construct an unbiased (i.e. exact) SG-MCMC algorithm that retains the efficiency we get from using stochastic gradients?

*Answer*: Yes, which we demonstrate by introducing **AMAGOLD**:

Exact without being prohibitively expensive

Uses M-H correction, amortizing its cost by applying it every *T* algorithm steps

**MCMC**

**Full-batch Methods**

**SG-MCMC** } Subsampled methods

**1st-order**
- MALA

**2nd-order**
- L2MC
- HMC

M-H correction to handle bias

**1st-order**
- SGLD

**2nd-order** } Includes momentum term

**SGHMC**

Small step size to handle bias

AMAGOLD

Amortized M-H correction to handle bias

| Algorithm | Exact? | Stochastic Gradient? |
| --- | --- | --- |
| AMAGOLD | Yes | Yes |
| L2MC | Yes | No |
| HMC | Yes | No |
| SGHMC | No | Yes |

## Bayesian inference

Given: Some dataset $\mathcal{D}$, domain $\Theta$

Sample: From posterior distribution $\pi(\theta) \propto \exp\left(-U(\theta)\right)$ where

$$U(\theta) = -\sum_{x \in \mathcal{D}} \log p(x|\theta) - \log p(\theta).$$

$U(\theta)$: Energy function

$\theta$: Ranges over $\Theta$

$\pi \propto \mu$: $\pi$ is the unique distribution with PDF proportional to $\mu$

## Second-order MCMC

A second-order chain (e.g. HMC, SGHMC, L2MC) *augments* state space with momentum $r$

**Joint distribution**:

$$\pi(\theta, r) \propto \exp(-H(\theta, r)) = \exp\left(-U(\theta) - \frac{1}{2\sigma^2}\|r\|^2\right),$$

$H$ (Hamiltonian): measures total energy of system.

Full-batch energy function (e.g. HMC, L2MC)

$$U(\theta) = -\sum_{x \in \mathcal{D}} \log p(x|\theta) - \log p(\theta).$$

Need M-H correction step to prevent bias due to discretization

Stochastic gradient energy function (e.g. SGHMC, AMAGOLD)

$$\tilde{U}(\theta) \approx -\frac{|\mathcal{D}|}{|\tilde{\mathcal{D}}|} \sum_{x \in \tilde{\mathcal{D}}} \log p(x|\theta) - \log p(\theta).$$

Naively using stochastic gradient estimate can lead to convergence to wrong stationary distribution

**Algorithm 1** SGHMC

1: **given:** Energy $U$, initial state $\theta \in \Theta$, step size $\epsilon$, momentum variance $\sigma^2$, friction $\beta$
2: **loop**
3:     **optionally, resample momentum:**
4:     $r \sim \mathcal{N}(0, \sigma^2)$
5:     **initialize position and momentum:**
6:     $r_{\frac{1}{2}} \leftarrow r, \theta_0 \leftarrow \theta$
7:     **for** $t = 1$ **to** $T$ **do**
8:         **position update:** $\theta_t \leftarrow \theta_{t-1} + \epsilon\sigma^{-2}r_{t-\frac{1}{2}}$
9:         **sample noise** $\eta_t \sim \mathcal{N}(0, 4\epsilon\beta\sigma^2)$
10:        **sample random energy component** $\tilde{U}_t$
11:        **update momentum:**

$$r_{t+\frac{1}{2}} \leftarrow r_{t-\frac{1}{2}} - \epsilon\nabla\tilde{U}_t(\theta_t) - 2\epsilon\beta r_{t-\frac{1}{2}} + \eta_t$$

12:     **end for**
13:     **new values:** $(\theta, r) \leftarrow (\theta_T, r_{T+\frac{1}{2}})$
14:     ▷ no M-H step
15: **end loop**

Second-order MCMC

    Compute posterior distribution using sampling

    Include momentum term

    Full-batch variants (L2MC, HMC) and minibatch

    Minibatch variants (SGHMC, **AMAGOLD**)

## Exact methods

    Guarantee convergence to correct stationary distribution (L2MC, HMC, **AMAGOLD**)

    Use M-H correction to remove bias

## Inexact methods

    Do not have same convergence guarantees (SGHMC)

**Detailed Balance Condition** A Markov chain with transition probability operator *G* is reversible if for any pair of states *x* and *y*

$$\pi(x)G(x,y) = \pi(y)G(y,x).$$

Computing the M-H acceptance probability

$$\tau = \min\left(1, \frac{\pi(y)P(y,x)}{\pi(x)P(x,y)}\right).$$

Given some measure-preserving involution over the state space denoted $x \mapsto x^\perp$, a chain $G$ is *skew-reversible* if $\pi(x) = \pi(x^\perp)$ and

$$\pi(x)G(x, y) = \pi(y^\perp)G(y^\perp, x^\perp).$$

For Hamiltonian dynamics we use the involution that negates the momentum, i.e. $(\theta, r)^\perp = (\theta, -r)$.

AMAGOLD: **A**mortized **M**etropolis-**A**djusted stochastic **G**radient second-**O**rder **L**angevin **D**ynamics

**Convergence rate intuition**:

Essentially equivalent to full-batch L2MC, up to a constant factor

Approaches L2MC's rate as batch size increases or step size decreases

| Algorithm | Exact? | Stochastic Gradient? |
|-----------|--------|----------------------|
| AMAGOLD   | Yes    | Yes                  |
| L2MC      | Yes    | No                   |
| HMC       | Yes    | No                   |
| SGHMC     | No     | Yes                  |

Full-batch
  $\rightarrow$ L2MC with AMA

Full-batch, $\beta = 0$, resample
  $\rightarrow$ HMC

Disable M-H, adjust
hyperparameters
  $\rightarrow$ SGHMC

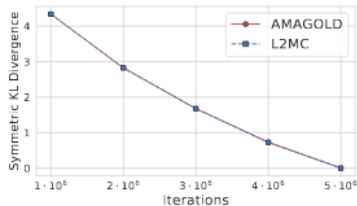(a) SGHMC

(b) AMAGOLD

(c) Tuned AMAGOLD

(d) KL Divergence

(a) Dist1      (b) KL comparison on Dist1