

OPTIMAL TUNING OF SUBSAMPLING HAMILTONIAN MONTE CARLO

KHUE-DUNG DANG, MATIAS QUIROZ

ROBERT KOHN, MINH-NGOC TRAN

MATTIAS VILLANI

DEPARTMENT OF STATISTICS

STOCKHOLM UNIVERSITY

AND

DEPARTMENT OF COMPUTER AND INFORMATION SCIENCE

LINKÖPING UNIVERSITY

- **Hamiltonian Monte Carlo (HMC)**
- **HMC-ECS: HMC with Energy Conserving Subsampling**
- **Block-Poisson estimator of the likelihood**
- **Signed HMC-ECS**
- **Optimal tuning of HMC-ECS**
- **Empirical results**

THE METROPOLIS-HASTINGS (MH) ALGORITHM

■ Bayesian inference

$$\pi(\theta) \propto L(\theta)p(\theta)$$

■ Initialize $\theta^{(0)}$ and iterate for $i = 1, 2, \dots, N$

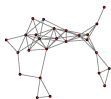
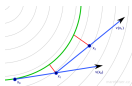
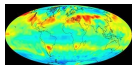
1. Sample $\theta_p \sim q(\cdot | \theta^{(i-1)})$ (the **proposal distribution**)
2. Compute the **acceptance probability**

$$\alpha = \min \left(1, \frac{L(\theta_p)p(\theta_p)}{L(\theta^{(i-1)})p(\theta^{(i-1)})} \frac{q(\theta^{(i-1)} | \theta_p)}{q(\theta_p | \theta^{(i-1)})} \right)$$

3. With probability α set $\theta^{(i)} = \theta_p$ and $\theta^{(i)} = \theta^{(i-1)}$ otherwise.

PROBLEM #1 - LIKELIHOOD EVALUATIONS CAN BE COSTLY.

- **High-dimensional spatio-temporal problems** (GMRFs)
- Models where **numerical methods** are needed for evaluating $p(y_i|\theta)$ (ODEs, optimization, etc)
- **Doubly intractable problems** with costly normalization constants (ERGMs)
- So called **Big data** problems with many observations.



- Hard to make **good proposals** in **high dimensional parameter spaces**.
- **Random walk Metropolis**

$$\theta_p | \theta^{(i-1)} \sim N(\theta^{(i-1)}, c \cdot \Sigma)$$

traverses the parameter space slowly in high dimensions.

- Small c - high acceptance probability, but small steps.
- Large c - large steps, but low acceptance probability.

HAMILTONIAN MONTE CARLO (HMC) - CONTINUOUS TIME

- Augment the posterior with **momentum** variables $\mathbf{r} \in \mathbb{R}^d$.

- **Extended target:**

$$\pi(\theta, \mathbf{r}) \propto \exp(-\mathcal{H}(\theta, \mathbf{r})), \quad \mathcal{H}(\theta, \mathbf{r}) = \mathcal{U}(\theta) + \mathcal{K}(\mathbf{r})$$

$$\mathcal{U}(\theta) = -\log L(\theta) - \log p(\theta) \quad \text{and} \quad \mathcal{K}(\mathbf{r}) = \frac{1}{2} \mathbf{r}^T \mathbf{M}^{-1} \mathbf{r}.$$

- Proposal for (\mathbf{r}, θ) is generated by drawing $\mathbf{r} \sim N(\mathbf{0}, \mathbf{M}^{-1})$ and follow the **Hamiltonian dynamics**:

$$\begin{aligned} \frac{d\theta}{dt} &= \frac{\partial \mathcal{H}(\theta, \mathbf{r})}{\partial \mathbf{r}} = \mathbf{M}^{-1} \mathbf{r} \\ \frac{d\mathbf{r}}{dt} &= -\frac{\partial \mathcal{H}(\theta, \mathbf{r})}{\partial \theta} = -\frac{\partial \mathcal{U}(\theta)}{\partial \theta} \end{aligned}$$

- **Properties** of **Hamiltonian dynamics**:

1. **Reversible** (one-to-one) [leaves target distribution invariant]
2. **Volume preserving** [Jacobian is one]
3. **Energy conserving** [MH acceptance probability is 1]

HAMILTONIAN MONTE CARLO (HMC) - DISCRETE TIME

- **Leap-frog discretization** of the Hamiltonian dynamics.
- Take L **steps** of **size** ϵ .
- MH acceptance probability is no longer 1.
- Problem: each of the L steps require computing posterior gradient $\frac{\partial \log \pi(\theta)}{\partial \theta}$ for each proposal draw. **Costly!**
- Naive solution: simulate the Hamiltonian dynamics on a **subset of the data**.
- Problem: **Energy is not conserved**. MH acceptance probability quickly drops with dimension.
- Our solution: use ideas from MCMC with estimates of the target posterior.

ESTIMATING THE LIKELIHOOD

- **Likelihood estimator**, $\hat{L}_m(\theta, \mathbf{u})$, based on m auxiliary random variables $\mathbf{u} = (u_1, \dots, u_m)$. **Unbiased**:

$$\mathbb{E}_{\mathbf{u}} [\hat{L}_m(\theta, \mathbf{u})] = L(\theta).$$

- **Subsampling**:

- $\mathbf{u} = (u_1, \dots, u_m)$ index sampled observations
- unbiased **log-likelihood estimator** [$\ell(\theta) \equiv \log L(\theta)$]

$$\hat{\ell}(\theta, \mathbf{u}) = \frac{n}{m} \sum_{i \in \mathbf{u}} p(y_i | \theta)$$

- **bias-correction**
- **control variates** to reduce variance
- For (much) more info, see our papers:
 - Speeding Up MCMC by Efficient Data Subsampling, *JASA*. [1]
 - Subsampling MCMC - an Intro for the Survey Statistician, *Sankhya A*.

- Initialize $(\theta^{(0)}, \mathbf{u}^{(0)})$ and iterate for $i = 1, 2, \dots, N$
 1. Sample $\theta_p \sim q(\cdot | \theta^{(i-1)})$ and $\mathbf{u}_p \sim p(\mathbf{u})$ to obtain the **unbiased** estimate $\hat{L}(\theta_p, \mathbf{u}_p)$
 2. Compute the **acceptance probability**

$$\alpha = \min \left(1, \frac{\hat{L}(\theta_p, \mathbf{u}_p) p(\theta_p)}{\hat{L}(\theta^{(i-1)}, \mathbf{u}^{(i-1)}) p(\theta^{(i-1)})} \frac{q(\theta^{(i-1)} | \theta_p)}{q(\theta_p | \theta^{(i-1)})} \right)$$

3. With probability α set $(\theta^{(i)}, \mathbf{u}^{(i)}) = (\theta_p, \mathbf{u}_p)$ and $(\theta^{(i)}, \mathbf{u}^{(i)}) = (\theta^{(i-1)}, \mathbf{u}^{(i-1)})$ otherwise.

- Targets a joint distribution $\tilde{\pi}_m(\theta, \mathbf{u})$ with marginal $\pi(\theta)$ [2]
- .. but **low** $\mathbb{V}(\hat{L}(\theta, \mathbf{u}))$ crucial for **efficient sampling**.

- Pseudo-marginal **extended target**

$$\bar{\pi}_m(\theta, \mathbf{r}, \mathbf{u}) \propto \exp(-\hat{\mathcal{H}}(\theta, \mathbf{r})) p_U(\mathbf{u}), \quad \hat{\mathcal{H}}(\theta, \mathbf{r}) = \hat{\mathcal{U}}(\theta) + \mathcal{K}(\mathbf{r})$$

$$\hat{\mathcal{U}}(\theta) = -\log \hat{L}(\theta) - \log p(\theta) \quad \text{and} \quad \mathcal{K}(\mathbf{r}) = \frac{1}{2} \mathbf{r}^T \mathbf{M}^{-1} \mathbf{r}.$$

- Marginal of θ is still $\pi(\theta)$, if $\hat{L}(\theta)$ is unbiased.

- **HMC-within-Gibbs:**

1. $\theta, \mathbf{r} | \mathbf{u}$ - HMC with energy $\hat{\mathcal{H}}$ given the subsample \mathbf{u}
2. $\mathbf{u} | \theta, \mathbf{r}$ - Metropolis-Hastings update of subsample

- **Crucial:** **the same** $\hat{L}(\theta)$ is used in:

- The Hamiltonian dynamics and
- The MH acceptance probability

- **HMC-ECS.** Subsampling conserves energy.

- **Distant proposals** with **high MH acceptance probability.**

THE BLOCK-POISSON ESTIMATOR

■ The **Block-Poisson estimator** of the likelihood $L(\theta)$:

- Draw $\mathcal{X}_1, \dots, \mathcal{X}_\lambda \stackrel{iid}{\sim} \text{Pois}(1)$.
- For $l = 1, \dots, \lambda$, draw \mathcal{X}_l mini-batches of data of size m .
- Compute unbiased mini-batch estimators

$$\hat{\ell}_m^{(h,l)}, \text{ for } h = 1, \dots, \mathcal{X}_l$$

- Construct likelihood estimate for some constant $a \in \mathbb{R}$

$$\hat{L}_B(\theta) \equiv \prod_{l=1}^{\lambda} \zeta_l \text{ where } \zeta_l \equiv \exp\left(\frac{a + \lambda}{\lambda}\right) \prod_{h=1}^{\mathcal{X}_l} \left(\frac{\hat{\ell}_m^{(h,l)} - a}{\lambda}\right).$$

■ What really matters for MH is the variance of

$$\log \frac{\hat{p}(\mathbf{y}|\theta_p, \mathbf{u}_p)}{\hat{p}(\mathbf{y}|\theta^{(i-1)}, \mathbf{u}^{(i-1)})}$$

■ Product form of $\hat{L}_B(\theta)$: can use **Block Pseudo Marginal (BPM)**.

$$\hat{L}_B(\theta) = \prod_{l=1}^{\lambda} \zeta_l, \text{ where } \zeta_l = \exp\left(\frac{a + \lambda}{\lambda}\right) \prod_{h=1}^{x_l} \left(\frac{\hat{\ell}_m^{(h,l)} - a}{\lambda}\right)$$

- **Unbiased:** $\mathbb{E}(\hat{L}_B(\theta)) = L(\theta)$ for all $\theta \in \Theta$.

Proof parts:

- $\mathbb{E}(\hat{L}_B(\theta)) = \mathbb{E}_{\mathcal{X}_{1:\lambda}} \mathbb{E}_{\mathbf{u}|\mathcal{X}_{1:\lambda}}(\hat{L}_B(\theta))$
- Poisson expectation becomes a power series in ℓ .
- Taylor series expansion of $L(\theta) = \exp(\ell(\theta))$.

- **Positive:** $\hat{L}_B(\theta)$ is almost surely positive only if $\hat{\ell}_m^{(h,l)} \geq a$ almost surely for all h and l .
- For a given λ , $\mathbb{V}(\hat{L}_B(\theta))$ is minimized for $a = \ell - \lambda$.

- Forcing a to be a **lower bound** for all $\hat{\ell}_m^{(h,l)}$ is impractical:
 - Usually need to know ℓ_i for all data points.
 - $a = \ell - \lambda$ implies that λ will be large. Costly!
- **Soft lower bound:** $\Pr(\hat{\ell}_m^{(h,l)} \geq a)$ close to one. More efficient, but $\hat{L}_B(\theta) < 0$ possible.
- **Signed HMC-ECS** [3]
 - **Run HMC-ECS on absolute value** $|\hat{L}_B(\theta)| p(\theta)$
 - **Correct for sign** $s = \text{Sign}(\hat{L}_B(\theta))$ using importance sampling

$$\widehat{\mathbb{E}\psi(\theta)} = \frac{\sum_{i=1}^N \psi(\theta^{(i)}) s^{(i)}}{\sum_{i=1}^N s^{(i)}}.$$

OPTIMAL TUNING OF SIGNED HMC-ECS

- Standard (No U-turn) HMC tuning of:
 - number of **leapfrog** L
 - **step size** ϵ
 - **mass matrix** M .
- **Optimal** λ and m minimizes **Computational Time (CT)**:

$$\text{CT}(\lambda, m) \propto m\lambda \cdot \frac{\text{IF} \left[\sigma^2_{\log|\hat{L}_B|}(\lambda, m) \right]}{(2\tau(\lambda, m) - 1)^2}$$

- Optimal λ and m **balances**
 1. The **cost** of computing \hat{L}_B , which is $m\lambda$ on average
 2. **MH inefficiency**, IF
 3. Probability of a **positive sign** $\tau(\lambda, m) \equiv \Pr(\hat{L}_B \geq 0)$.

OPTIMAL TUNING OF SIGNED HMC-ECS

- To **compute** $\text{CT}(\lambda, m)$, we need expressions for:
 - $\text{IF}(\cdot)$
 - $\sigma^2_{\log|\hat{L}_B|}(\lambda, m)$
 - $\tau(\lambda, m)$
- Need to assume the **distribution of** $\hat{\ell}_m^{(h,l)}$
 - Normal (or CLT)
 - Mixture of Normals (universal approximator)
- The **derivation of IF** is an extension of the theory in [4] to blocked signed PMMH.
- **Guidelines** based on **idealized assumptions**. But accurate in experiments.
- **Conservative guidelines**: $m\lambda$ is not suggested too small.

$$\Pr(\hat{L}_B \geq 0)$$

- Under the minimum variance condition $a = \ell - \lambda$

$$\hat{L}_B(\theta) = \prod_{l=1}^{\lambda} \zeta_l, \text{ where } \zeta_l = \exp\left(\frac{\ell}{\lambda}\right) \prod_{h=1}^{\mathcal{X}_l} \left(\frac{\hat{\ell}_m^{(h,l)} - \ell}{\lambda} + 1\right)$$

- $\hat{L}_B(\theta) > 0$ whenever an even number of ζ_l are negative.
- $\zeta_l > 0$ whenever even number of $\frac{\hat{\ell}_m^{(h,l)} - \ell}{\lambda} + 1$ are negative.
- Applying a result from Feller's first book twice:

$$\Pr(\hat{L}_B \geq 0) = \frac{1}{2} \left[1 + (1 - 2\Psi(m, \lambda))^{\lambda} \right]$$

$$\Psi(m, \lambda) \equiv \Pr(\zeta_l < 0) = \frac{1}{2} \sum_{j=1}^{\infty} \left[1 - (1 - 2\Pr(A_m < 0))^j \right] \Pr(\mathcal{X}_l = j),$$

$$\mathcal{X}_l \stackrel{iid}{\sim} \text{Pois}(1) \text{ and } A_m = \frac{\hat{\ell}_m^{(h,l)} - \ell}{\lambda} + 1.$$

- Under the condition $a = \ell - \lambda$ we have

$$\begin{aligned}\log|\hat{L}| &= \ell + \sum_{l=1}^{\lambda} \sum_{h=1}^{\mathcal{X}_l} \log \left(\left| \frac{\hat{\ell}_m^{(h,l)} - \ell}{\lambda} + 1 \right| \right) \\ &= \ell + \frac{1}{2} \sum_{l=1}^{\lambda} \sum_{h=1}^{\mathcal{X}_l} \log \left(\frac{\hat{\ell}_m^{(h,l)} - \ell}{\lambda} + 1 \right)^2\end{aligned}$$

- $\hat{\ell}_m^{(h,l)} \sim \text{Normal} \Rightarrow \sigma_{\log|\hat{L}_B|}^2(\lambda, m)$ is the variance of a random sum of logs of non-central χ^2 variables.
- Non-central χ^2 is a Poisson mixture of central χ^2 [5]
- Moments of log central χ^2 are known from [6]
- Law of total variance

OPTIMAL TUNING - NORMAL CASE

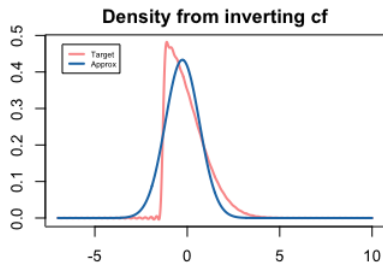
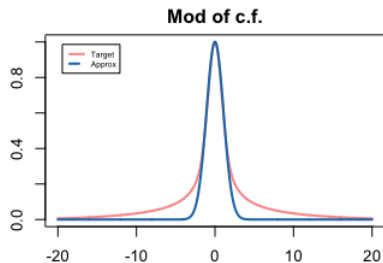
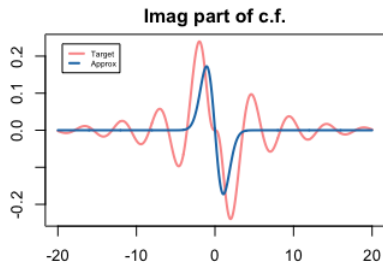
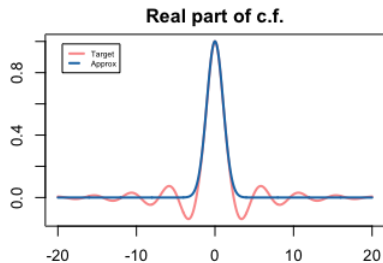
- Assume $\hat{\ell}_m^{(h,l)} \sim \text{Normal}$.
- Both $\Pr(\hat{L}_B \geq 0)$ and $\sigma_{\log|\hat{L}_B|}^2(\lambda, m)$ are functions of the variance of $\hat{\ell}_m^{(h,l)}$

$$\mathbb{V}(\hat{\ell}_m^{(h,l)}(\theta)) = \frac{n^2}{m} \sigma_{\ell_i}^2(\theta)$$

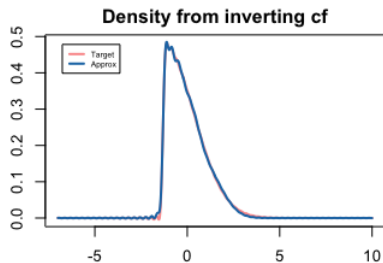
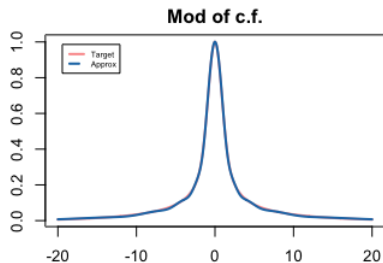
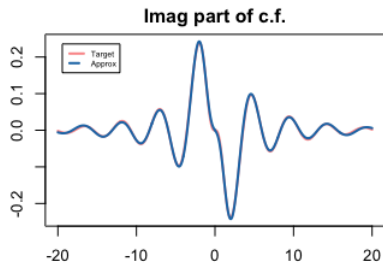
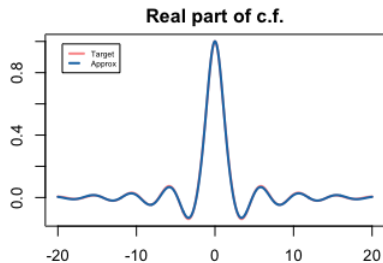
- Optimal tuning therefore depends on $\sigma_{\ell_i}^2(\theta)$.
- Solution: estimate $\sigma_{\ell_i}^2(\theta)$ from a subsample for some selected θ .
- What if $\hat{\ell}_m^{(h,l)}$ are not normal?
- Set $m = 20$ and rely on the CLT. Optimize only λ .
- However, numerical experiments tell us that $m = 1$ is optimal.

- Assume that $\hat{\ell}_m^{(h,l)}$ follows a **mixture of normals**.
- Mixture of normals are **universal approximators**.
- Both $\Pr(\hat{L}_B \geq 0)$ and $\sigma_{\log|\hat{L}_B|}^2$ are still **tractable**.
- ... but estimating $\sigma_{\ell_i}^2(\theta)$ is not enough anymore.
- How to fit a mixture of normals to $\hat{\ell}_m^{(h,l)}$?
- **Matching characteristic functions** (c.f.)
 1. Fit any distribution to a subsample of ℓ_i 's and get the c.f. $\varphi_\ell(t)$.
 2. Compute the c.f. of $\hat{\ell}_m^{(h,l)}$ as $\varphi_{\hat{\ell}_m}(t) = (\varphi_\ell(t/m))^m$.
 3. Approximate the distribution of $\hat{\ell}_m^{(h,l)}$ by a normal mixture by L2-matching of c.f.'s. Plancherel's theorem.

MATCHING A 1-COMPONENT MON TO SKEWED NORMAL



MATCHING A 5-COMPONENT MoN TO SKEWED NORMAL

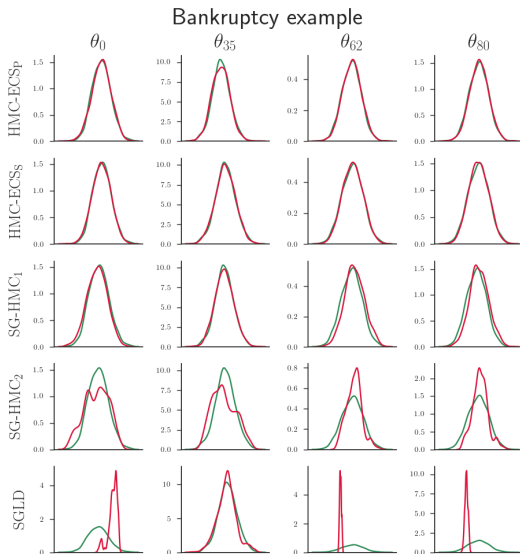


- Performance measure:

$$RCT = \frac{\text{CT of algorithm}}{\text{CT of Perturbed HMC-ECS}}$$

RCT	HMC	HMC-ECS _S	SG-HMC ₁	SG-HMC ₂	SGLD
min	358.5	1	7.0	48.6	53.9
median	466.8	2.2	9.5	100.2	230.1
max	683.6	7.2	538.7	246.4	2784.2






POSTERIOR PERTURBATION (BIAS)



CONCLUSIONS

- **Subsampling** to speed up MCMC and HMC.
- **Block-Poisson** is an **unbiased** and **efficient** estimator of the likelihood.
- **Optimal tuning of Signed HMC-ECS** with Block-Poisson estimator.
- **Very large speed-ups** compared to regular HMC and state-of-the-art subsampling algorithms.
- Can be used to optimally tune Signed HMC in **doubly intractable problems**.

References

-  M. QUIROZ, R. KOHN, M. VILLANI, AND M.-N. TRAN, “SPEEDING UP MCMC BY EFFICIENT DATA SUBSAMPLING,” *JOURNAL OF THE AMERICAN STATISTICAL ASSOCIATION*, NO. FORTHCOMING, PP. 1–35, 2018.
-  C. ANDRIEU AND G. O. ROBERTS, “THE PSEUDO-MARGINAL APPROACH FOR EFFICIENT MONTE CARLO COMPUTATIONS,” *THE ANNALS OF STATISTICS*, PP. 697–725, 2009.
-  A.-M. LYNE, M. GIROLAMI, Y. ATCHADE, H. STRATHMANN, D. SIMPSON, ET AL., “ON RUSSIAN ROULETTE ESTIMATES FOR BAYESIAN INFERENCE WITH DOUBLY-INTRACTABLE LIKELIHOODS,” *STATISTICAL SCIENCE*, VOL. 30, NO. 4, PP. 443–467, 2015.
-  M. K. PITT, R. D. S. SILVA, P. GIORDANI, AND R. KOHN, “ON SOME PROPERTIES OF MARKOV CHAIN MONTE CARLO SIMULATION METHODS BASED ON THE PARTICLE FILTER,” *JOURNAL OF ECONOMETRICS*, VOL. 171, NO. 2, PP. 134–151, 2012.
-  C. WALCK, “HAND-BOOK ON STATISTICAL DISTRIBUTIONS FOR EXPERIMENTALISTS,” TECH. REP., 1996.
[HTTP://INSPIREHEP.NET/RECORD/1389910/FILES/SUF9601.PDF](http://inspirehep.net/record/1389910/files/SUF9601.pdf)



S. E. PAV, “MOMENTS OF THE LOG NON-CENTRAL CHI-SQUARE DISTRIBUTION,” *ARXIV PREPRINT ARXIV:1503.06266*, 2015.