

1

Introduction

Hamiltonian Monte Carlo (HMC): sample from *continuous target distribution* $\pi(q^C) \propto e^{-U(q^C)}$ by introducing auxiliary momentum variables $p^C \in \mathbb{R}^{N_c}$ and simulating *Hamiltonian dynamics*

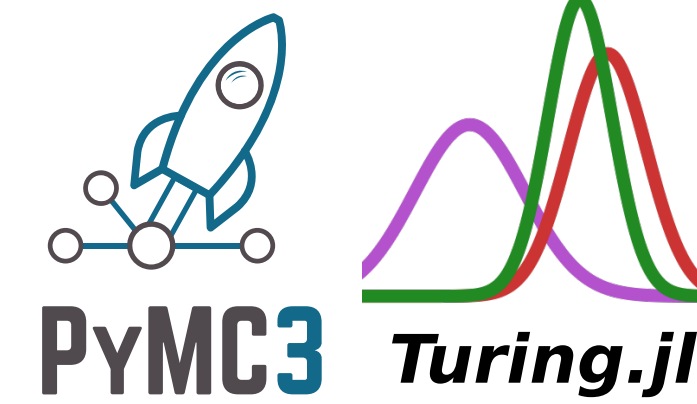
$$\begin{cases} \frac{dq^C(t)}{dt} = \nabla K^C(p^C) \\ \frac{dp^C(t)}{dt} = -\nabla_{q^C} U(q^C) \end{cases}$$

Remarkable empirical success, but can't be applied to distributions with mixed discrete and continuous variables

Our Goal: sample from *target distribution* $\pi(x, q^C) \propto e^{-U(x, q^C)}$ with mixed *discrete* ($x \in \Omega$) and *continuous* ($q^C \in \mathbb{R}^{N_c}$) variables

Existing approaches:

- Integrate out the discrete variables:
 - Only applicable on a small scale
 - Can't always be carried out automatically
- Alternate between continuous HMC and generic discrete updates:
 - Need long HMC trajectory to suppress random walk behavior
 - Discrete updates can only be done infrequently
- Update discrete and continuous variables in tandem:
 - Discontinuous HMC (DHMC)*, *Probabilistic path HMC (PPHMC)*
 - DHMC* is best suited for ordinal paramters, and has inefficient embedding and algorithmic structure
 - PPHMC* only works for phylogenetic trees



2

Mixed Hamiltonian Monte Carlo (M-HMC)

Mixed Hamiltonian Monte Carlo (M-HMC)

- M-HMC also evolves the discrete and continuous variables in tandem
- M-HMC is applicable to any distributions with mixed support
- M-HMC can be efficiently implemented using Laplace momentum

We start with an illustrative example on a 1D Gaussian mixture model (GMM) with 4 mixture components. Use $x \in \{1, 2, 3, 4\}$ to denote the discrete variable, and $q^C \in \mathbb{R}$ to denote the continuous variable. We want to sample from

$$\pi(x, q^C) = \phi_x N(q^C | \mu_x, \Sigma)$$

where

$$\phi_1 = 0.15, \phi_2 = \phi_3 = 0.3, \phi_4 = 0.25, \Sigma = 0.1$$

$$\mu_1 = -2, \mu_2 = 0, \mu_3 = 2, \mu_4 = 4$$

The right panel shows M-HMC in this simple case. The simple change implied by the M-HMC framework leads to a more efficient sampler that is able to correct the bias from naively doing Metropolis-Hastings (MH) updates within HMC.

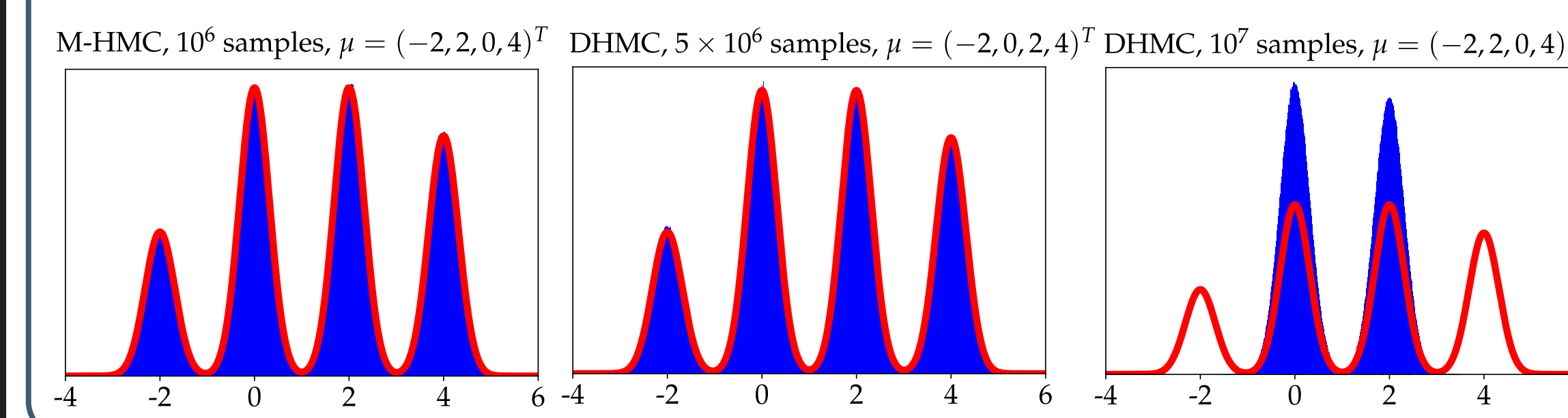
More generally, M-HMC can be easily applied to arbitrary distributions with mixed support, and introduces minimal overhead compared to existing HMC methods. Refer to the paper for more details on the general algorithm.

Comparison with Discontinuous HMC (DHMC):

- DHMC relies on embedding that works best for ordinal parameters

$$x_i = n \iff q_i^D \in (a_n, a_{n+1}], 0 = a_1 \leq a_2 \leq \dots$$

- DHMC needs to update all discrete variables at every step, inefficient



MHwHMC v.s. M-HMC for 1D GMM

Require: U , target potential
 Q , discrete proposal
Input: $x^{(0)}$, current discrete state
 $q^{C(0)}$, current continuous location
 ε , step size; L , # of steps

function M-HMC($x^{(0)}, q^{C(0)}, \varepsilon, L|U, Q$)
 $k^D(0) \sim \text{Exponential}(1), p^C(0) \sim N(0, 1)$
 $x \leftarrow x^{(0)}, k^D \leftarrow k^D(0)$
 $q^C \leftarrow q^{C(0)}, p^C \leftarrow p^C(0)$

for t **from** 1 **to** L **do**
 $q^C, p^C \leftarrow \text{leapfrog}(q^C, p^C, \varepsilon)$

$\tilde{x} \sim Q(\cdot|x), \Delta E \leftarrow \log \frac{e^{-U(x, q^C)} Q(\tilde{x}|x)}{e^{-U(\tilde{x}, q^C)} Q(x|\tilde{x})}$

Naive MH within HMC
if $\text{Exponential}(1) > \Delta E$ **then**
 $x \leftarrow \tilde{x}$
end if

or

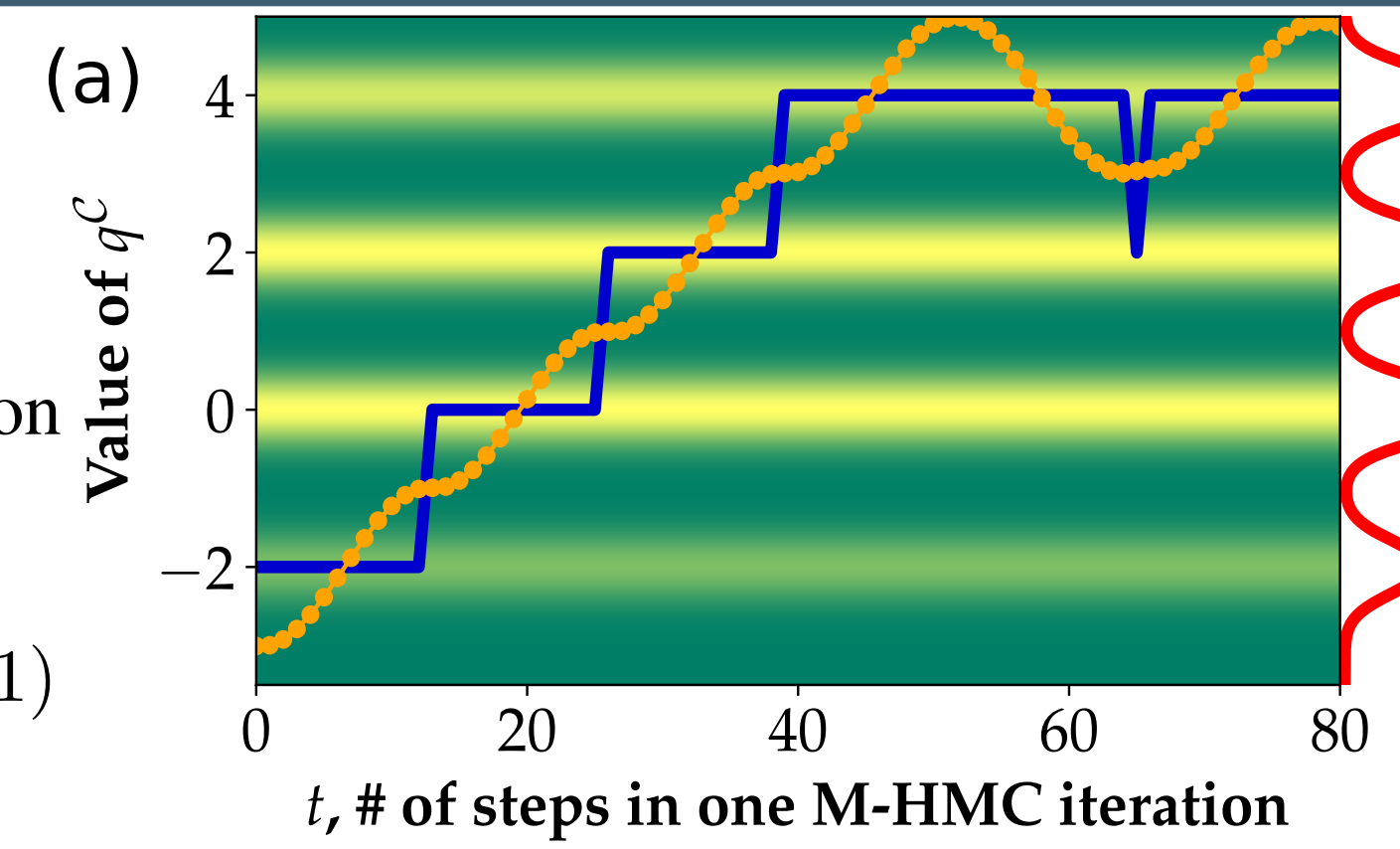
M-HMC
if $k^D > \Delta E$ **then**
 $x \leftarrow \tilde{x}, k^D \leftarrow k^D - \Delta E$
end if

end for

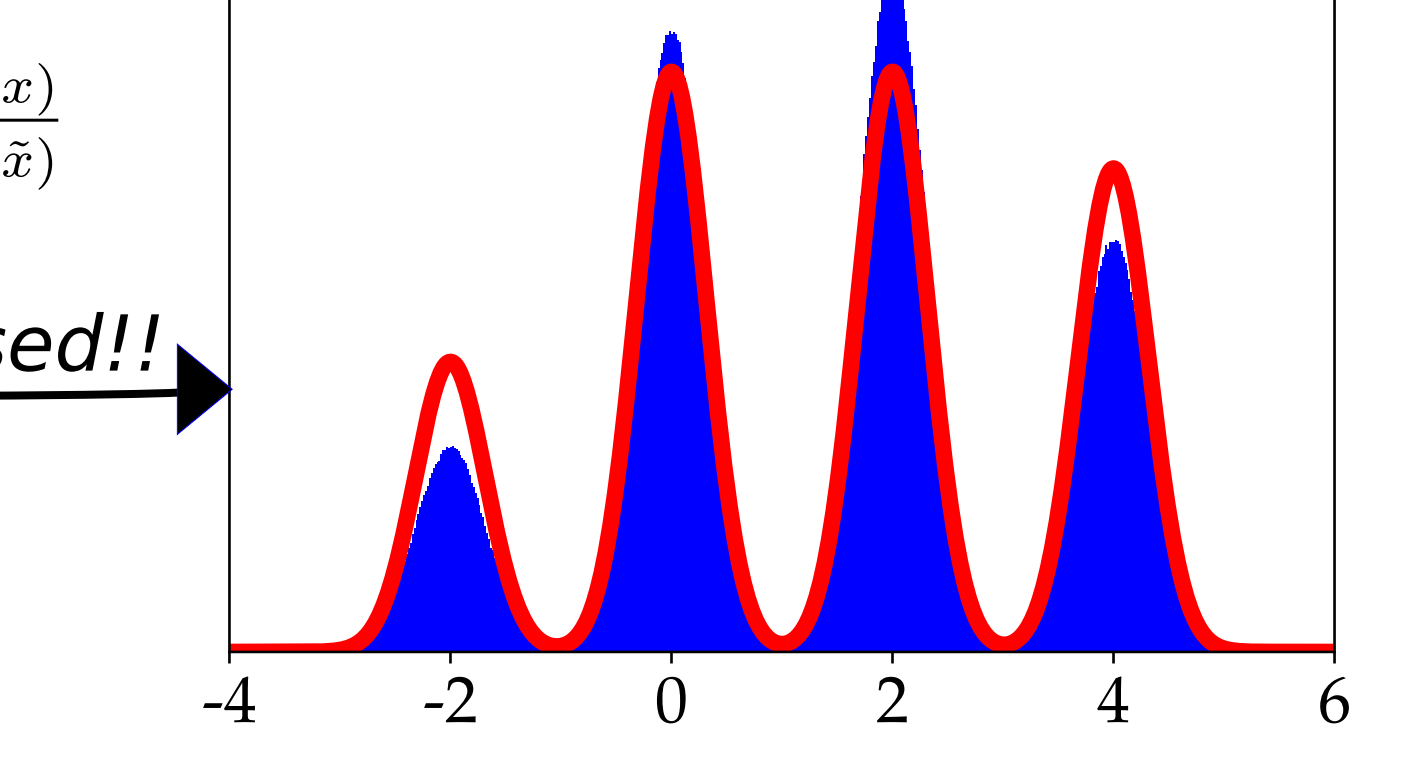
Final MH correction
 $E \leftarrow U(x, q^C) + k^D + K^C(p^C)$
 $E^{(0)} \leftarrow U(x^{(0)}, q^{C(0)}) + k^D(0) + K^C(p^C(0))$
 $x, q^C \leftarrow \text{MHCorrection}(E, E^{(0)})$

return x, q^C

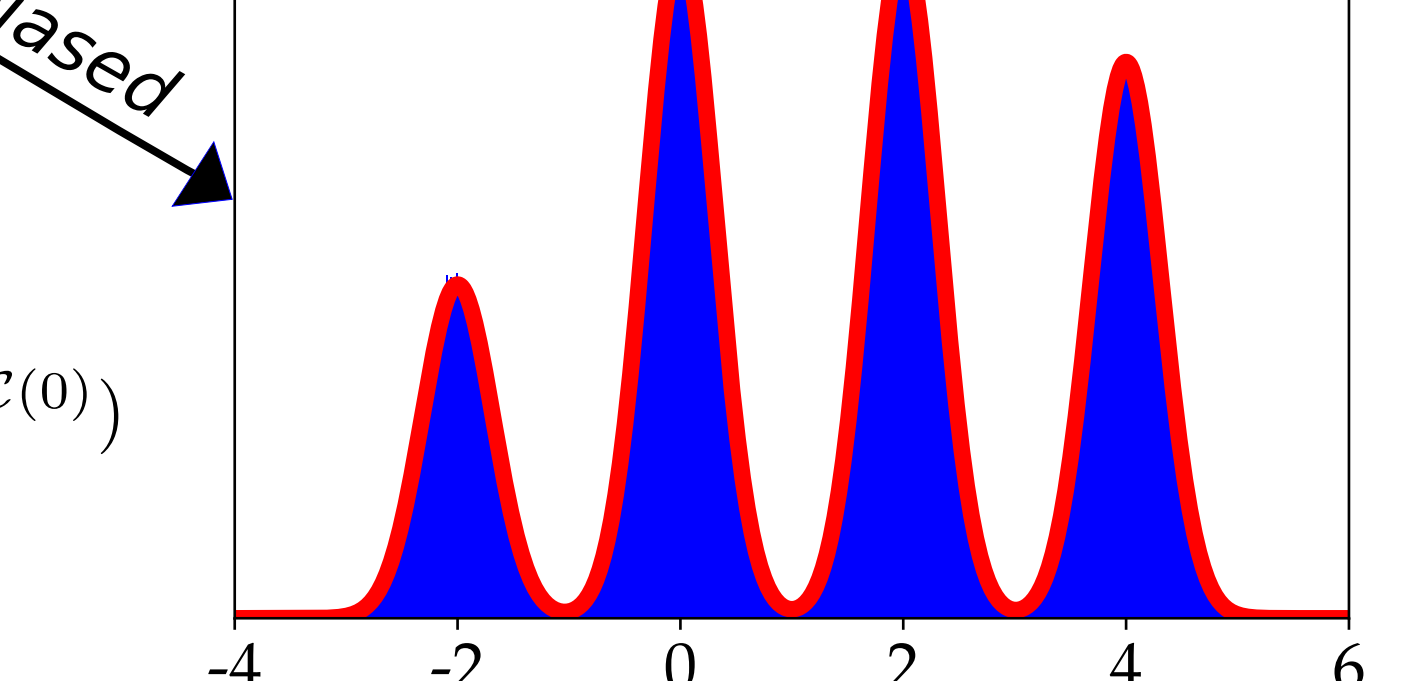
end function



(b) MHwHMC, 10^7 samples, $\mu = (-2, 0, 2, 4)^T$



(c) M-HMC, 10^6 samples, $\mu = (-2, 0, 2, 4)^T$



3

Numerical Experiments

Overview

Methods used:

- Mixed HMC (M-HMC)** : Custom JAX implementation
- Discontinuous HMC (DHMC)** : Custom JAX implementation
- HMC-within-Gibbs (HwG)** : Custom JAX implementation
- No-U-Turn Sampler (NUTS)** : Using NumPyro, for GMMs
- NUTS-within-Gibbs (NwG)** : Using compound step in PyMC3
- Specialized Gibbs samplers** : Custom Numba implementation
 - (Polson et al., 2013) for Bayesian logistic regression
 - (Chen et al., 2013) for Correlated topic models

Performance measure:

- Minimum relative effective sample size (MRESS)
- The minimum ESS over all dimensions
- Normalized by the number of samples
- Estimated using multiple independent chains

Discrete proposals in M-HMC:

$$Q_j(\tilde{x}|x) \propto \pi(\tilde{x}) \rho_j(\tilde{x}|x)$$

$$\text{where } \rho_j(\tilde{x}|x) = \begin{cases} 1 & \text{if } \tilde{x}_j \neq x_j, \tilde{x}_i = x_i, i \neq j \\ 0 & \text{otherwise} \end{cases}$$

24D Gaussian mixture model (GMM)

Basic setup: $\pi(x, q^C) = \phi_x N(q^C | \mu_x, \Sigma)$ where

$$\phi_1 = 0.15, \phi_2 = \phi_3 = 0.3, \phi_4 = 0.25, \Sigma = 3I$$

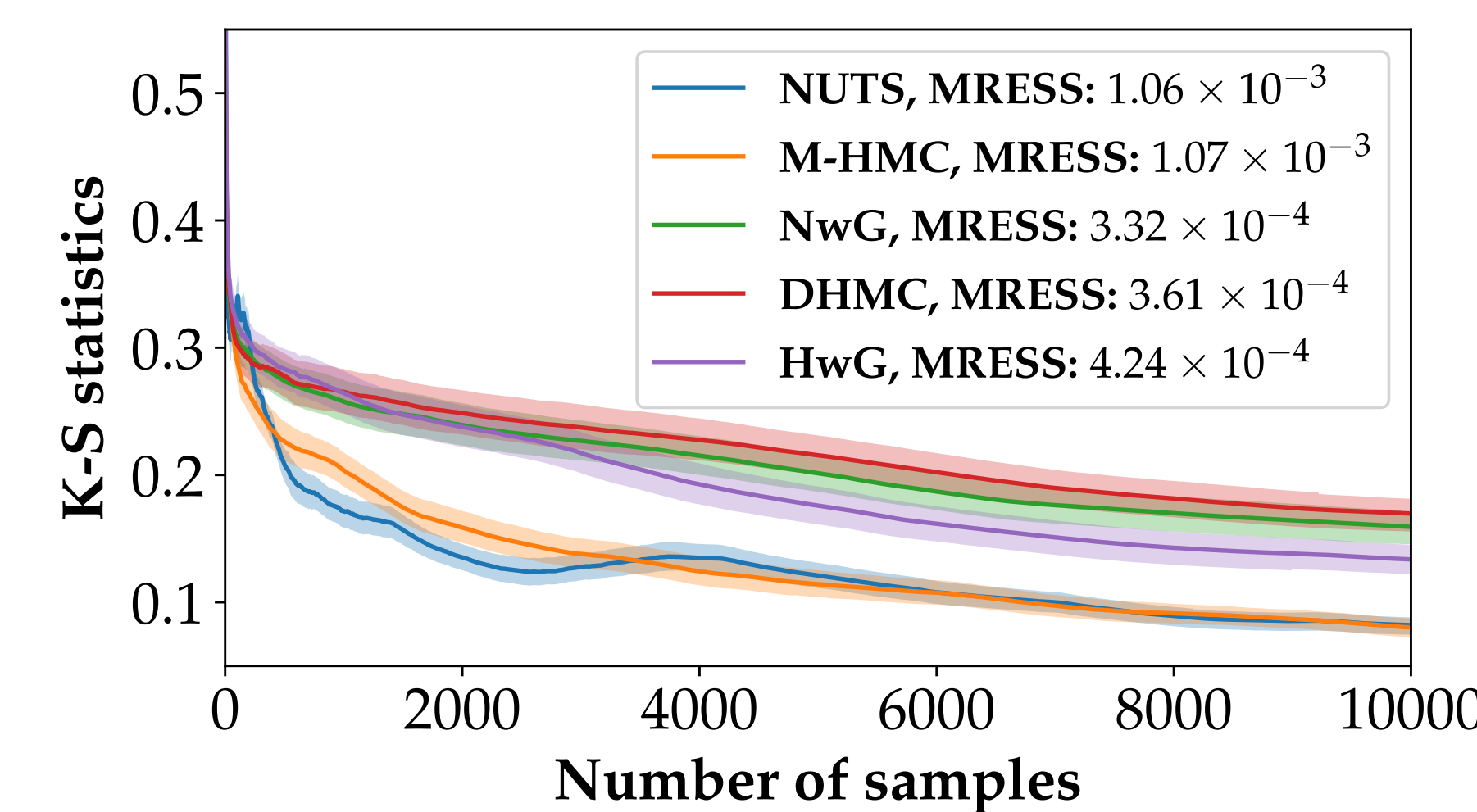
$$\mu_1, \mu_2, \mu_3, \mu_4 \text{ are from 24 permutations of } -2, 0, 2, 4$$

Target distribution: $\pi(x, q^C)$

Methods tested: M-HMC, DHMC, HwG, NwG, NUTS

Results summary:

- All samplers are accurate
- M-HMC is more efficient than DHMC, HwG and NwG
- M-HMC is as efficient as NUTS



Variable selection in Bayesian logistic regression (BLR)

Basic setup: $y_i \sim \text{Bernoulli}(\sigma(X_i^T \beta))$, $i = 1, \dots, 100$
 where $X \in \mathbb{R}^{100 \times 20}$, $\beta \in \mathbb{R}^{20}$, and $\sigma(x) = 1/(1 + e^{-x})$ is the sigmoid function.
 $\gamma_j, j = 1, \dots, 20$ are binary variables indicating the presence of a particular component of β , and $N(0, 25I)$ is an uninformative prior on β .

$$\text{Joint distribution: } p(\beta, \gamma, X, y) = N(\beta|0, 25I) \prod_{i=1}^{100} p_i^{y_i} (1 - p_i)^{1-y_i}$$

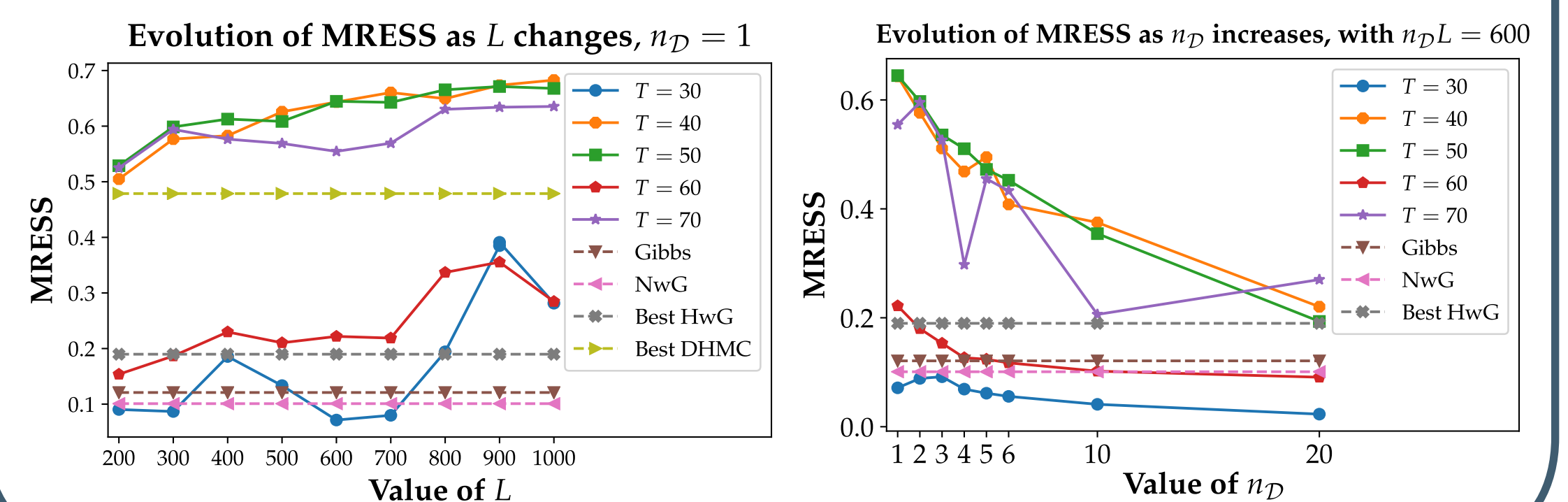
$$\text{where } p_i = \sigma(\sum_{j=1}^{20} X_{ij} \beta_j \gamma_j), i = 1, \dots, 100.$$

Target distribution: $p(\beta, \gamma|X, y)$

Methods tested: M-HMC, DHMC, HwG, NwG, Gibbs (Polson et al., 2013)

Results summary:

- All samplers are accurate
- M-HMC is more efficient than DHMC, HwG, NwG and Gibbs
- M-HMC exhibits U-turn behavior
- M-HMC benefits from distributed/more frequent discrete updates



Correlated topic models (CTMs)

Basic setup:

Given the topics β , a vector $\mu \in \mathbb{R}^K$ and a $K \times K$ covariance matrix Σ , CTM assumes the following generative process for the d th document with N_d words:

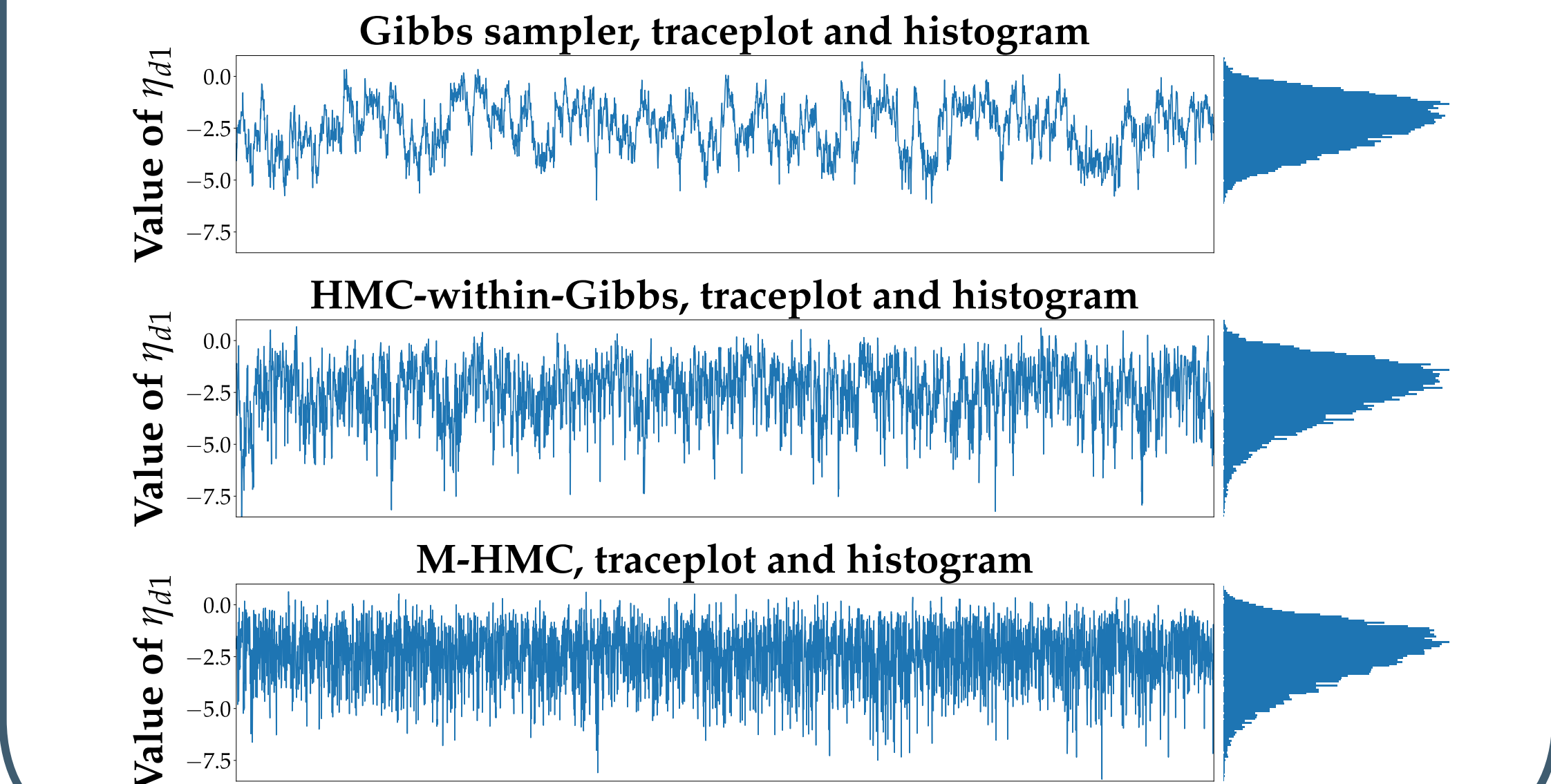
- $\eta_d \sim N(\mu, \Sigma)$
- For each $n \in \{1, \dots, N_d\}$:
 - Draw topic assignment $z_{d,n} | \eta_d \sim \text{Categor}(f(\eta_d))$
 - Draw word $w_{d,n} | z_{d,n}, \beta \sim \text{Categor}(\beta_{z_{d,n}})$

Target distribution: $p(\eta, z|w; \beta, \mu, \Sigma)$

Methods tested: M-HMC, DHMC, HwG, NwG, Gibbs (Chen et al., 2013)

Results summary:

- DHMC fails; Gibbs occasionally fails
- M-HMC is 3x more efficient than HwG/NwG
- M-HMC is 20x more efficient than Gibbs



4

Conclusions

- M-HMC evolves discrete and continuous variables in tandem, and is applicable to any distributions with mixed support.
- M-HMC with Laplace momentum is easy to implement, and introduces minimal overhead when compared with existing HMC methods.
- M-HMC is shown to be more efficient than strong baselines on challenging distributions with mixed support.