

# Sample Adaptive MCMC



Michael H. Zhu (*Stanford University*)

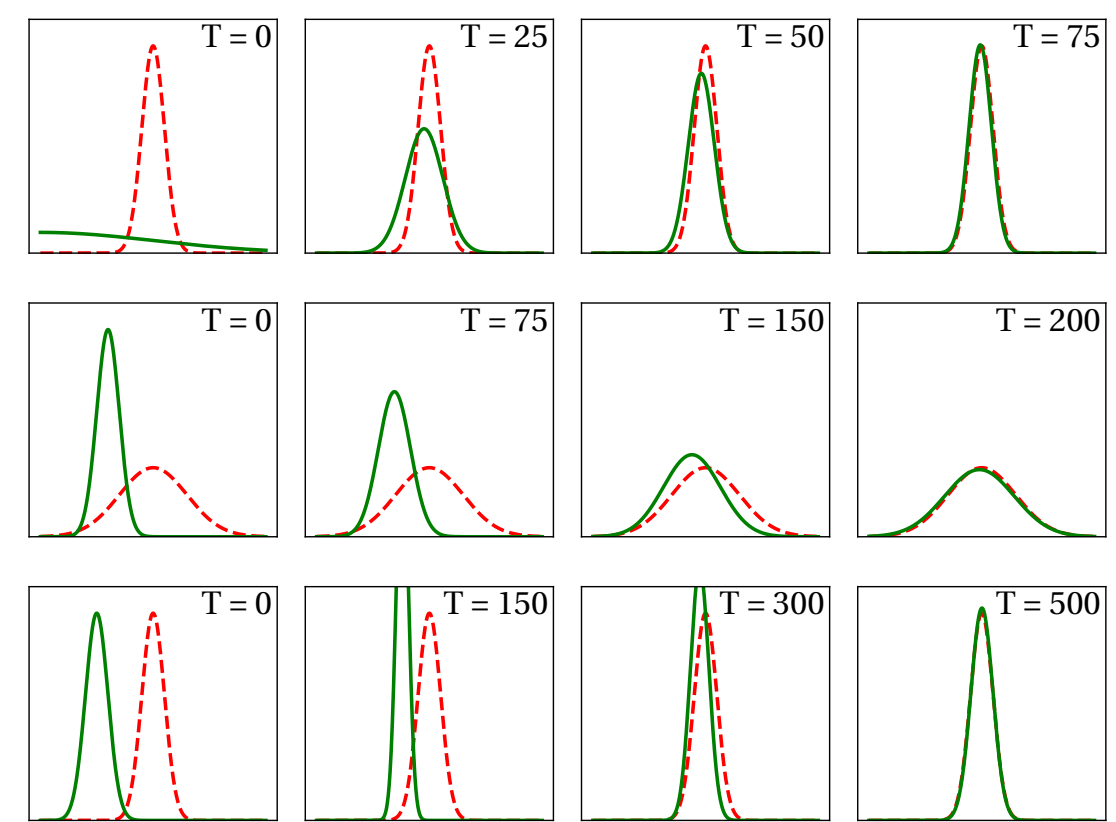
## Abstract

- Sample Adaptive MCMC (SA-MCMC) is a MCMC method based on a reversible Markov chain for  $\pi^{\otimes N}$  which
  - uses an adaptive proposal distribution based on the current state of  $N$  points
  - uses a sequential substitution procedure with one new likelihood evaluation per iteration and at most one updated point each iteration
- The SA-MCMC proposal distribution automatically adapts within its parametric family to best approximate the target distribution.
  - In contrast to many existing MCMC methods, SA-MCMC does not require any tuning of the proposal distribution.
  - SA-MCMC only requires specifying the initial state of  $N$  points, which can often be chosen *a priori*, thereby automating the entire sampling procedure with no tuning required.
- Experimental results demonstrate the fast adaptation and effective sampling of SA-MCMC.

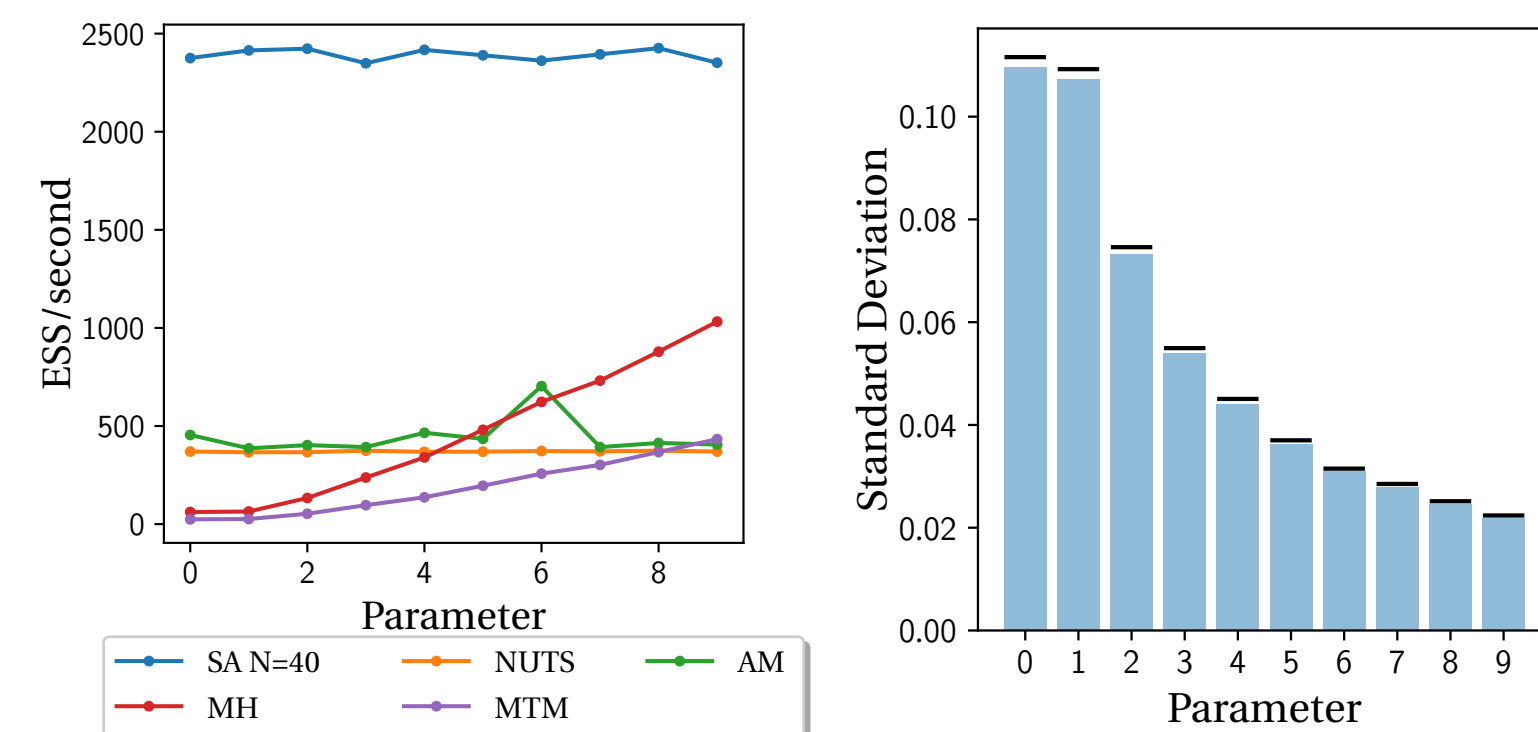
## Background

- For MCMC methods like Metropolis-Hastings, the choice of the proposal distribution  $q(\cdot|\theta^{(k)})$  is important in practice for effective sampling from the target distribution.
- A suboptimal choice for the scale or shape of the proposal can lead to inefficient sampling, yet the design of an optimal proposal distribution is challenging when the properties of the target distribution are unknown, especially in high-dimensional spaces.
- Adaptive MCMC methods such as Adaptive Metropolis continually adapt the proposal distribution based on the entire history of past states.
  - However, the method is no longer based on a valid Markov chain, so the usual MCMC convergence theorems do not apply and the validity of the sampler and an ergodic theorem must be proved for each specific algorithm under certain technical assumptions.
- We use the following notation for Sample Adaptive MCMC:
  - Target distribution  $p(\theta)$
  - Initialization distribution  $q_0(\cdot)$
  - Proposal distribution  $q(\cdot|\mu(S), \Sigma(S))$
  - $N$  points in the state
  - $\kappa$  burn-in iterations,  $K$  estimation iterations

## Simulation examples

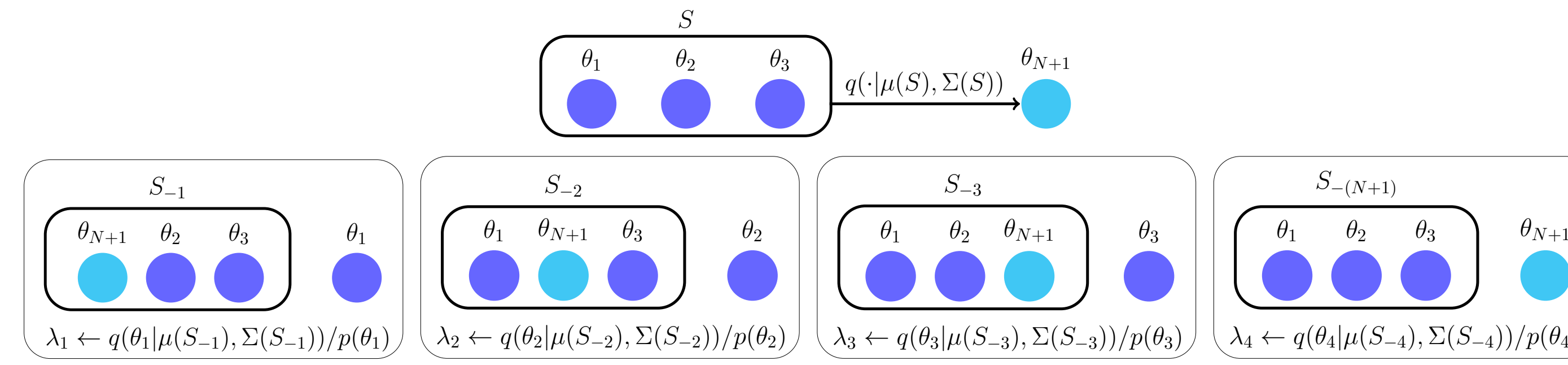


**Figure:** Adaptation of the SA-MCMC proposal distribution (green) to three target distributions (red).



**Figure:** Bayesian linear regression. (*left*) Comparison of ESS/second for each parameter. (*right*) Standard deviation of the SA proposal distribution (blue bar), averaged over iterations, for each parameter compared with the ground truth posterior standard deviation (black line).

## Algorithm



**Figure:** Illustration of one iteration of SA-MCMC for  $N = 3$ . After the proposed point  $\theta_{N+1} \sim q(\cdot|\mu(S), \Sigma(S))$  is sampled, the sets  $S_{-1}, \dots, S_{-(N+1)}$  are used to calculate the substitution probabilities  $\lambda_1, \dots, \lambda_{N+1}$ . One of the sets  $S_{-1}, \dots, S_{-(N+1)}$  is chosen to be the next state with probability proportional to  $\lambda_n$ .

## Algorithm 1 Sample Adaptive MCMC

**Require:**  $p(\theta)$ ,  $q_0(\cdot)$ ,  $q(\cdot|\mu(S), \Sigma(S))$ ,  $N$ ,  $\kappa$ ,  $K$

- 1: Initialize  $S^{(0)} \leftarrow (\theta_1, \dots, \theta_N)$  where  $\theta_n \sim q_0(\cdot)$  for  $n = 1, \dots, N$
- 2: **for**  $k = 1$  **to**  $\kappa + K$  **do**
- 3: Let  $S = (\theta_1, \dots, \theta_N) \leftarrow S^{(k-1)}$
- 4: Sample  $\theta_{N+1} \sim q(\cdot|\mu(S), \Sigma(S))$
- 5: Let  $S_{-n} \leftarrow (S \text{ with } \theta_n \text{ replaced by } \theta_{N+1})$  for  $n = 1, \dots, N$ . Let  $S_{-(N+1)} \leftarrow S$ .
- 6: Let  $\lambda_n \leftarrow q(\theta_n|\mu(S_{-n}), \Sigma(S_{-n}))/p(\theta_n)$  for  $n = 1, \dots, N+1$
- 7: Sample  $j \sim J$  with  $\mathbb{P}[J = n] = \lambda_n / \sum_{i=1}^{N+1} \lambda_i$ ,  $1 \leq n \leq N+1$
- 8: Let  $S^{(k)} \leftarrow S_{-j}$
- 9: **end for**
- 10: Return  $\bigcup_{k=\kappa+1, \dots, \kappa+K} S^{(k)}$

## Theory

- We prove the ergodicity of SA-MCMC under general conditions on the target distribution (the same assumptions as for Metropolis-Hastings) and a family of proposal distributions with diagonal covariance matrices.
- We prove the uniform ergodicity of SA-MCMC under the assumption that  $q(\theta|\gamma)/\pi(\theta)$  is bounded above and below.

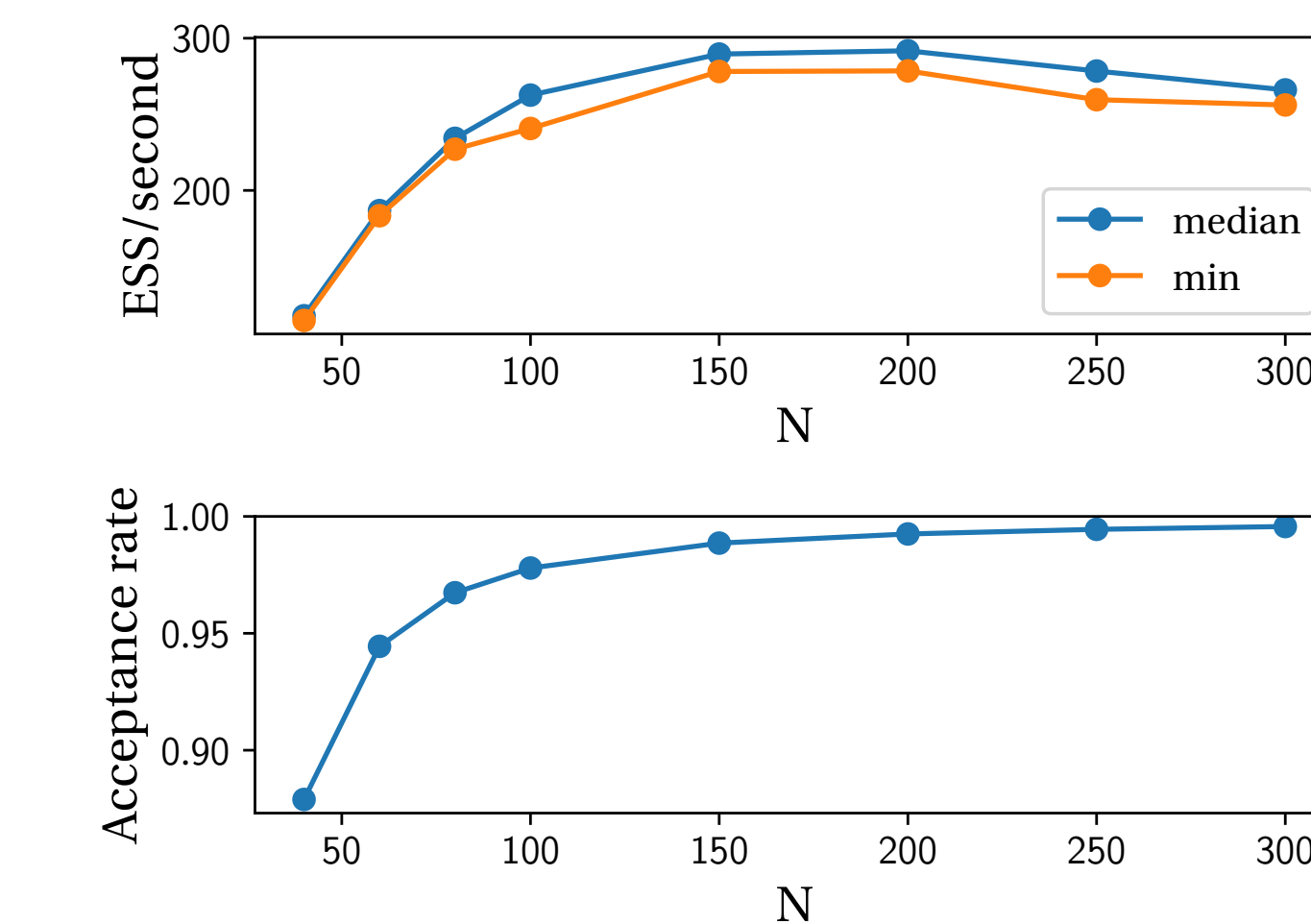
## Experimental setup

- We use the following experimental setup:
  - For Metropolis-Hastings (MH), we use an isotropic normal distribution as the proposal distribution,  $q_{MH}(\cdot|\theta) = \mathcal{N}(\theta, \sigma_{q, MH}^2 \mathbb{I})$ , with scale parameter  $\sigma_{q, MH}$ , and initialize  $\theta^{(0)} \sim q_{0, MH}(\cdot) = \mathcal{N}(0, \sigma_{q, MH}^2 \mathbb{I})$ .
  - For Adaptive Metropolis (AM), we use the optimal MH proposal distribution during the burn-in (non-adaptive) phase and then use the proposal distribution  $q_{AM}(\cdot|\theta^{(1)}, \dots, \theta^{(k-1)}) = \mathcal{N}(\theta^{(k-1)}, s_{AM}^2 \Sigma^{(k-1)})$  at iteration  $k$  with scale parameter  $s_{AM}$  and sample covariance matrix  $\Sigma^{(k-1)}$  of the past samples  $(\theta^{(1)}, \dots, \theta^{(k-1)})$ .
  - For Multiple-Try Metropolis (MTM), we use the optimal MH proposal distribution with 3 tries.
  - For SA-MCMC (SA), we use  $q_{0, SA}(\cdot) = \mathcal{N}(0, \sigma_{q_0, SA}^2 \mathbb{I})$  with scale parameter  $\sigma_{q_0, SA}$  as the distribution for initializing the  $N$  starting points. For the proposal distribution, when using the full covariance matrix, we use the Gaussian family  $q(\cdot|\mu(S), \Sigma(S)) = \mathcal{N}(\cdot|\mu(S), \Sigma(S))$ .

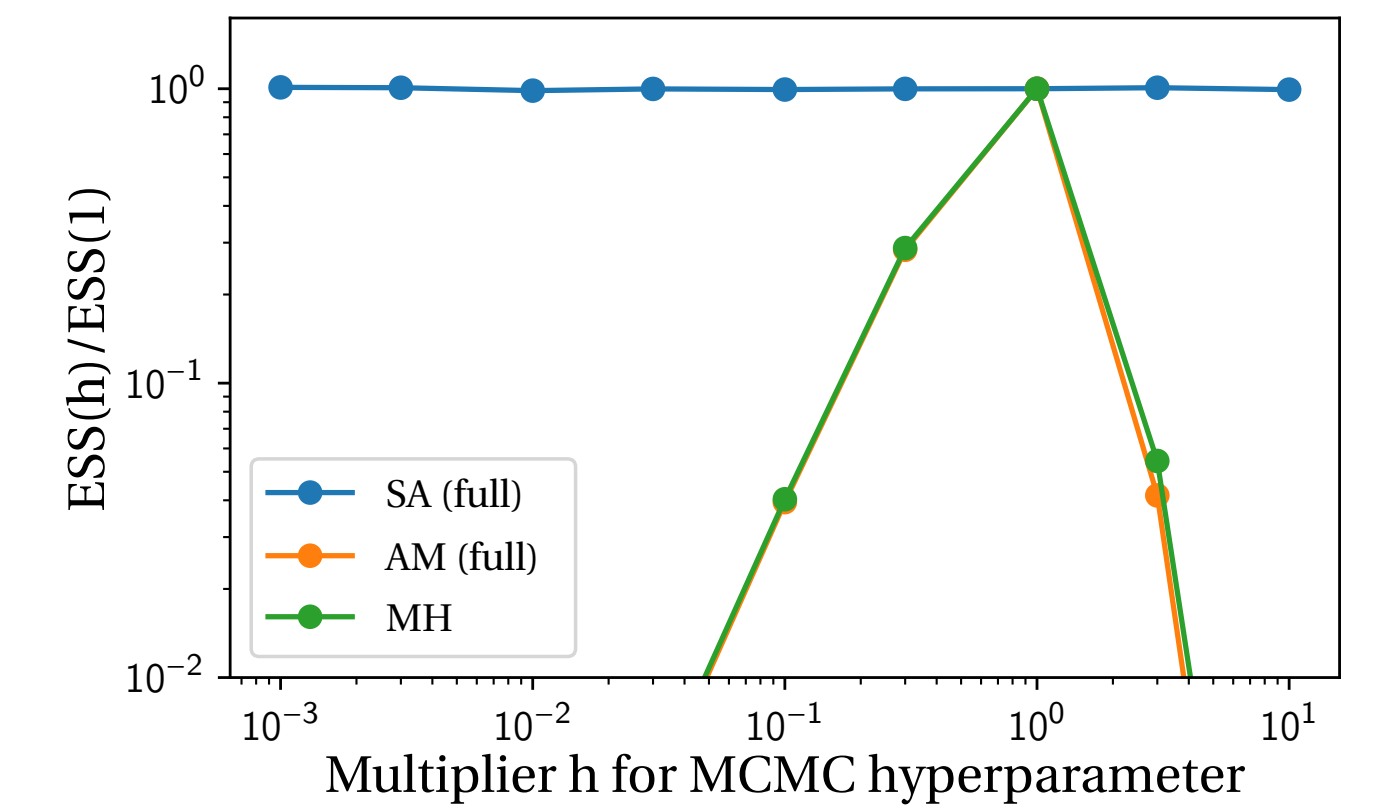
## Experimental results for Bayesian logistic regression

**Table:** Comparison of ESS/second for Bayesian logistic regression on (*top*) 11-dim MNIST 7s vs 9s using 10 features computed with PCA (*bottom*) 7-dim adult census income

	MH	MTM	AM (diag)	AM (full)	SA (diag)	SA (full)	NUTS
min(ESS)/s	13	5	17	37	23	<b>278</b>	54
median(ESS)/s	21	9	23	38	52	<b>290</b>	105
s/chain	733	3651	734	742	782	1112	1160
Hyperparameters	$q=.02$	$q=.02$ $M=3$	$q=.02$ $s=.6$	$q=.02$ $s=.7$	$q_0=1$ $N=40$	$q_0=1$ $N=150$	Stan
Acceptance rate	23%	48%	24%	26%	75%	98.9%	—
min(ESS)/s	1.4	0.6	13	16	67	<b>151</b>	40
median(ESS)/s	17	7	15	17	89	<b>158</b>	49
s/chain	2198	10951	2205	2217	2283	2509	2989
Hyperparameters	$q=.016$	$q=.016$ $M=3$	$q=.016$ $s=.8$	$q=.016$ $s=.85$	$q_0=1$ $N=40$	$q_0=1$ $N=150$	Stan
Acceptance rate	26%	52%	21%	24%	89%	99.2%	—



**Figure:** Plot of ESS/s and acceptance rate for SA-MCMC (full) versus  $N$  on MNIST.



**Figure:** Impact of MCMC hyperparameter on ESS for MNIST. The ratio  $\text{ESS}(h)/\text{ESS}(1)$  measures the drop in ESS using  $0.02h$  for  $q$  in MH,  $0.7h$  for  $s$  in AM, and  $1h$  for  $q_0$  in SA.

**Table:** Comparison of ESS/second for Bayesian logistic regression on (*left*) 55-dim cover type (*right*) 51-dim MiniBooNE between AM (full), SA (full), and NUTS with a dense mass matrix.

	Cover type			MiniBooNE		
	AM	SA	NUTS	AM	SA	NUTS
min(ESS)/s	0.075	<b>2.34</b>	0.099	0.31	<b>3.35</b>	0.023
median(ESS)/s	0.078	<b>2.81</b>	0.114	0.38	<b>6.59</b>	0.039
s/chain	52,469	65,537	25,143	28,178	26,627	33,584
s/chain (burn-in)	4,770	5,958	16,980	1,342	2,421	19,051
s/chain (estimation)	47,699	59,579	8,163	26,836	24,206	14,533
# iter. (burn-in)	100,000	100,000	500	100,000	100,000	500
# iter. (estimation)	1,000,000	1,000,000	2,000	2,000,000	1,000,000	2,000
Hyperparameters	$q=.004$ $s=.32$	$q_0=1$ $N=1,000$	Stan (dense)	$q=.007$ $s=.33$	$q_0=1$ $N=1000$	Stan (dense)
Acceptance rate	25.1%	99.3%	—	25.7%	90.5%	—