

---

# An Efficient Minibatch Acceptance Test for Metropolis-Hastings

---

John Canny, Haoyu Chen, Daniel Seita, Xinlei Pan

University of California, Berkeley

{canny, haoyuchen, seita, xinleipan}@berkeley.edu

## Abstract

Markov chain Monte Carlo (MCMC) methods have many applications in machine learning. We are particularly interested in their application to modeling very large datasets, where it is impractical to perform Metropolis-Hastings tests on the full data. Previous work on reducing the cost of Metropolis-Hastings tests yield variable data consumed per sample, with only constant factor reductions versus using the full dataset for each sample. Here we present a method that can be tuned to provide arbitrarily small batch sizes, by adjusting either proposal step size or temperature. Our approach uses the natural noise present in minibatch likelihood estimates to furnish the randomness in a Metropolis-Hastings test. Our test uses the noise-tolerant Barker acceptance test with a novel additive correction variable. The resulting test can be combined with minibatch proposals to yield updates with the same complexity as a simple SGD update. In this paper we derive the test, analyze its performance, discuss its implementation, and present several experiments.

## 1 Introduction

Markov chain Monte Carlo (MCMC) sampling is a powerful method for computation on intractable distributions. We are interested primarily in large dataset applications, where the goal is to sample a posterior distribution  $p(\theta \mid x_1, \dots, x_N)$  of parameter  $\theta$ , and where the number of data instances  $N$  is large. The Metropolis-Hastings method (M-H) generates sample candidates from a proposal distribution  $q$  which is in general different from the target distribution  $p$ , and decides whether to accept or reject them based on an acceptance test. The acceptance test is usually a Metropolis test [1, 2]. Conventional Metropolis-Hastings requires all  $N$  data instances to generate one posterior sample.

Many state-of-the-art machine learning methods, and deep learning in particular, are based on minibatch updates (such as SGD) to a model. Minibatch updates produce many improvements to the model for each pass over the dataset, and have high sample efficiency. They also map very well onto hardware such as GPUs. In contrast, M-H requires calculations over the full dataset to produce a new sample. Recent results from [3, 4] perform approximate (bounded error) acceptance tests using subsets (minibatches) of the full dataset. The tests depend on minibatch statistics, and on the value of an additional uniform random variable  $u$ . The amount of data consumed for each test varies significantly from one minibatch to the next, and depends on the current sample, the proposed sample, and on the random variable  $u$ . By contrast, [5] performs exact tests but requires a lower bound on parameter distribution across its domain. The amount of data reduction depends on the accuracy of this bound, and such bounds are only available for relatively simple distributions.

Here we derive a new test which incorporates the variability in minibatch statistics as *a natural part of the test*. Because of this, the amount of data required for each test is fixed with high probability and in expected value. We use a Barker test function [6] rather than a Metropolis test, which makes our test naturally error tolerant. The idea of using a noise-tolerant test using Barker’s test function was suggested but not explored empirically in [7] section 6.3. But the asymptotic test statistic CDF and the Barker function are different, which leads to fixed errors for the approach in [7]. Here we

show that the difference between the distributions can be corrected with an additive random variable. This leads to a test which is fast, and whose error can be made arbitrarily small.

Our test is applicable when the variance (over data samples) of the log acceptance probability is small enough (less than 1). It's not clear at first why this quantity should be bounded, but we will show that it is "natural" for well-specified models running Metropolis-Hastings sampling with optimal proposals [8] on a full dataset. When we reduce the amount of data for the test, the variance goes up. We have to reduce variance in one of several ways. Either:

- Increase the temperature of the target distribution. Log likelihoods scale as  $1/T$ , and so the variance of the likelihood ratio will vary as  $1/T^2$ . Our model is no longer well-specified (we are doing inference at a temperature different from that assumed during data generation), but higher temperature can be advantageous for parameter exploration.
- For continuous probability distributions, reduce the proposal step size and variance (for stochastic proposals) compared to an optimal proposal. The variance of the log acceptance probability scales as the square of proposal step size.
- Increase the minibatch size when needed for certain minibatches. Log acceptance variance scales as  $1/k$  vs the minibatch size  $k$ . Our test is adaptive like earlier works, but unlike them, the distribution of minibatch size is Gaussian, not long-tailed. Increased minibatch size also reduces the error rate for the test.

It's worth discussing at this point the typical goals of M-H sampling on very large datasets. By the Bernstein-von Mises Theorem, the posterior distribution of the parameter  $\theta$  for a Bayesian inference task is asymptotically normal, and has variance that scales inversely with the number of data samples  $N$ . This mode is extremely sharp for large datasets, which may contain millions or billions of samples. Simply sampling from this distribution is one application, but an efficient proposal distribution [8] has similar variance to the target distribution and will diffuse to it extremely slowly from an initialization value which is (likely to be) many standard deviations away. If there are any other strong modes, it is very likely for the sampler to find one of them and become trapped in it when run at the normal distribution temperature ( $T = 1$ ). A common solution is to anneal the sampler, running first at high temperature (scaling log likelihoods by  $1/T$ ) which flattens the likelihood landscape. This in turn reduces the variance of the log acceptance probability and allows our acceptance test to be applied.

A second question concerns step size. Once we have fixed temperature, our variance constraint implies that we have to trade-off proposal step size  $s$  and batch size  $b$  ( $b \propto p^2$ ), i.e. we can make many small steps, or one large step, with a given batch of data. One of the primary drivers of this work is our belief in the value of small steps. For applications to neural networks or other models where the posterior is multimodal, posterior inference is arguably a search process. Covering the search space densely with small steps is much more valuable than few sparse steps toward the nearest optimum. In this mode, Metropolis-Hastings can be used in similar fashion to Stochastic Gradient Descent. The goal in SGD is to make gradual progress to a posterior mode with each step, taking small steps so that the cumulative displacement has progressively lower variance. A substantial part of the computational work of MCMC on massive datasets will be similarly in reaching a stationary distribution, which really means finding a deep posterior mode. Taking noisy small steps will nevertheless make steady progress to a posterior mode since their bias is in that direction. We demonstrate this behavior in our logistic regression experiments.

The contributions of this paper are as follows:

- We develop a new, more efficient (in samples per test) minibatch acceptance test with quantifiable error bounds. The test uses a novel additive correction variable to implement a Barker test based on minibatch mean and variance.
- We analyze the test for accuracy and speed.
- We compared performance of our new test and prior approaches on several datasets. We demonstrate the test is several orders of magnitude more efficient than prior work measured as data consumed per test, and that it does not suffer from long-tailed minibatch sizes (up to the dataset size).

## 2 Preliminaries and Related Work

In the Metropolis-Hastings method [9, 10], a difficult-to-compute probability distribution  $p(\theta)$  is sampled using a Markov chain  $\theta_1, \dots, \theta_n$ . The sample  $\theta_{t+1}$  at time  $t + 1$  is generated using a candidate  $\theta'$  from a (simpler) proposal distribution  $q(\theta' | \theta_t)$ , filtered by an acceptance test. The acceptance test is usually a Metropolis test. The Metropolis test has acceptance probability:

$$\alpha(\theta_t, \theta') = \frac{p(\theta')q(\theta_t | \theta')}{p(\theta_t)q(\theta' | \theta_t)} \wedge 1 \quad (1)$$

where  $a \wedge b$  denotes  $\min(a, b)$ . With probability  $\alpha(\theta_t, \theta')$ , we accept  $\theta'$  and set  $\theta_{t+1} = \theta'$ , otherwise set  $\theta_{t+1} = \theta_t$ . The test is often implemented with an auxiliary random variable  $u \sim \mathcal{U}(0, 1)$  with a comparison  $u < \alpha(\theta_t, \theta')$ ; here,  $\mathcal{U}(a, b)$  denotes the uniform distribution on the interval  $[a, b]$ . For notational simplicity, we drop the subscript  $t$  for the current sample  $\theta_t$  and denote it as  $\theta$ .

The acceptance test guarantees detailed balance, which means  $p(\theta)p(\theta' | \theta) = p(\theta')p(\theta | \theta')$ , where  $p(\theta' | \theta)$  is the probability of a transition from state  $\theta$  to  $\theta'$ . Here  $p(\theta' | \theta) = q(\theta' | \theta)\alpha(\theta, \theta')$  so the detailed balance equation becomes

$$p(\theta)q(\theta' | \theta)\alpha(\theta, \theta') = p(\theta')q(\theta | \theta')\alpha(\theta', \theta) \quad (2)$$

This condition, together with ergodicity, guarantees that the Markov chain has a unique stationary distribution  $\pi(\theta) = p(\theta)$ .

For Bayesian inference, we would like to sample from a parameter distribution for  $\theta$  based on some observed data, so the target distribution is  $p(\theta | x_1, \dots, x_N)$ . The acceptance probability is now:

$$\alpha(\theta, \theta') = \frac{p_0(\theta') \prod_{i=1}^N p(x_i | \theta') q(\theta | \theta')}{p_0(\theta) \prod_{i=1}^N p(x_i | \theta) q(\theta' | \theta)} \wedge 1 \quad (3)$$

where  $p_0(\theta)$  is a prior, and  $p(x_i | \theta)$  are the probabilities of the observations. Computing samples this way requires the use of all  $N$  training data points, but this is very expensive for large datasets. To address this challenge, [3, 4] perform approximate Metropolis-Hasting tests using sequential hypothesis testing. During each iteration, they start with a small minibatch of data and test the hypothesis that the sample  $\theta'$  should be accepted based on an approximate version of the test  $u < \alpha(\theta, \theta')$ . If the approximate test cannot make a decision with sufficient confidence, then the minibatch size is increased and the test repeats. This process continues until a decision. The bounds depend on either an asymptotic Central Limit Theorem [3] or a concentration bound [4]. The latter requires direct bounds on the log likelihood ratio, which for general distributions requires knowing  $p(x_i | \theta)$  and  $p(x_i | \theta')$  for all  $N$  samples. In addition, both methods suffer the drawback of resolving small log likelihood ratio differences between the minibatch and full batch versions. In the worst case, all  $N$  data points may be needed, and we soon show that this worst case can occur about  $\Omega(N)$  times during the performance of  $N$  tests.

Following [4], we write the test  $u < \alpha(\theta, \theta')$  in the equivalent form  $\Lambda(\theta, \theta') > \psi(u, \theta, \theta')$ , where<sup>1</sup>

$$\Lambda(\theta, \theta') = \sum_{i=1}^N \log \left( \frac{p(x_i | \theta')}{p(x_i | \theta)} \right) \quad \text{and} \quad \psi(u, \theta, \theta') = \log \left( u \frac{q(\theta' | \theta) p_0(\theta)}{q(\theta | \theta') p_0(\theta')} \right) \quad (4)$$

To reduce computational effort, an unbiased estimate of  $\Lambda(\theta, \theta')$  based on a minibatch can be used:

$$\Lambda^*(\theta, \theta') = \frac{N}{b} \sum_{i=1}^b \log \left( \frac{p(x_i | \theta')}{p(x_i | \theta)} \right) \quad (5)$$

Finally, it will be convenient for our analysis to define the individual components that contribute to the sums above:

$$\Lambda_i(\theta, \theta') = N \log \left( \frac{p(x_i | \theta')}{p(x_i | \theta)} \right) \quad (6)$$

Thus,  $\Lambda(\theta, \theta')$  is the mean of  $\Lambda_i(\theta, \theta')$  over the entire dataset, and  $\Lambda^*(\theta, \theta')$  is the mean of  $\Lambda_i(\theta, \theta')$  over its minibatch.

Since minibatches contains randomly selected samples  $x_i$ , the values  $\Lambda_i$  are independent, identically distributed (iid) random variables<sup>2</sup>. By the Central Limit Theorem, we expect  $\Lambda^*(\theta, \theta')$  to be

<sup>1</sup>Our definitions differ from those in [4] by a factor of  $N$  to simplify our analysis later.

<sup>2</sup>The analysis assumes sampling with replacement although implementations on typical large datasets will approximate this by sampling without replacement.

approximately Gaussian. The acceptance test then becomes a statistical test of the hypothesis that  $\Lambda(\theta, \theta') > \psi(u, \theta, \theta')$  by establishing that  $\Lambda^*(\theta, \theta')$  is substantially larger than  $\psi(u, \theta, \theta')$ . In [3] an asymptotic central limit argument was used to derive this gap, while in [4] a concentration bound was used. In both cases, the resulting tests were shown to give useful reductions in number of samples required over using the full dataset, but there were no worst-case bounds given on the average number of samples needed per iteration.

We next show that for a simple distribution, the lower bound of average number of data instances consumed for one iteration of [3] and [4] is  $\Omega(N)$ , where  $N$  is the number of data points.

## 2.1 A Worst-Case Gaussian Example

Consider a Gaussian model where  $x_1, \dots, x_N \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta, 1)$  with a known variance  $\sigma^2 = 1$ , true mean  $\theta = 0.5$ , and a uniform prior on  $\theta$ . The log likelihood ratio is

$$\Lambda^*(\theta, \theta') = \frac{N}{b} \sum_{i=1}^b \log \frac{p(x_i | \theta')}{p(x_i | \theta)} = N(\theta' - \theta) \left( \frac{1}{b} \sum_{i=1}^b x_i - \theta - \frac{\theta' - \theta}{2} \right) \quad (7)$$

which is normally distributed over selection of the Normal samples  $x_i$ . Since the  $x_i$  have unit variance, their mean has variance  $1/b$ , and the variance of  $\Lambda^*(\theta, \theta')$  is  $\sigma^2(\Lambda^*) = (\theta' - \theta)^2 N^2 / b$ . In order to pass a hypothesis test that  $\Lambda > \psi$ , there needs to be a large enough gap (several  $\sigma(\Lambda^*)$ ) between  $\Lambda^*(\theta, \theta')$  and  $\psi(u, \theta, \theta')$ .

Our goal is to sample from the posterior  $p(\theta | x_1, \dots, x_N)$ , which is a normal distribution  $\mathcal{N}(\mu, 1/N)$  centered on the sample mean  $\mu$ , and with variance  $1/N$ . In one dimension, an efficient proposal distribution has the same variance as the target distribution [8], so we use the proposal  $\mathcal{N}(\theta, 1/N)$ , which is implemented as  $q(\theta' | \theta) = \phi((\theta' - \theta)N)$ , where  $\phi(x)$  is the Normal density function with zero mean and unit variance. This proposal is symmetric  $q(\theta' | \theta) = q(\theta | \theta')$ , and since we assumed a uniform prior on  $\theta$ , we see that  $\psi(u, \theta, \theta')$  reduces to  $\log u$ . Our worst-case scenario is specified in Lemma 1.

**Lemma 1.** *For the model in Section 2.1, there exists a fixed (independent of  $N$ ) constant  $c$  such that with probability  $\geq c$  over the joint distribution of  $(\theta, \theta', u)$ , the tests from [3, 4] consume all  $N$  samples.*

*Proof.* See Appendix A. □

Similar results can be shown for other distributions and proposals by identifying regions in product space  $(\theta, \theta', u)$  such that the hypothesis test needs to separate nearly-equal values. It follows that the accelerated M-H tests in [3, 4] require at least a constant fraction  $\geq c$  in the amount of data consumed per test compared to full-dataset tests, so their speed-up is at most  $1/c$ .

The problem is that these methods use tail bounds to separate  $\Delta$  away from zero, but for certain input/random  $u$  combinations,  $\Delta$  can be arbitrarily close to zero. We will avoid this by using the *approximately normal* variation in  $\Delta^*$  to *replace* the variation due to  $u$ .

## 2.2 MCMC Posterior Inference

There is a separate line of MCMC work drawing principles from statistical physics. By viewing random variables as particles in a system, one can apply Hamiltonian Monte Carlo (HMC) [11] methods which generate high acceptance *and* distant proposals when run on full batches of data. Recently Langevin Dynamics [12, 13] has been applied to Bayesian estimation on minibatches of data. This simplified dynamics uses local proposals and avoids MH tests by using small proposal steps whose acceptance approaches 1 in the limit. However, the constraint on proposal step size is severe, and the state space exploration reduces to a random walk. Full minibatch HMC for minibatches was recently described in [14] which allows momentum-augmented proposals with larger step sizes. However, step sizes are still limited by the need to run accurately without MH tests. By providing an M-H test with similar cost to standard gradient steps, our work opens the door to applying those methods with much more aggressive step sizes without loss of accuracy.

### 3 A New Metropolis-Hastings Acceptance Test

#### 3.1 Log-Likelihood Ratios

For our new M-H test, we denote the exact and approximate log likelihood ratios as  $\Delta$  and  $\Delta^*$ . First,  $\Delta$  is defined as

$$\Delta(\theta, \theta') = \log \frac{p_0(\theta') \prod_{i=1}^N p(x_i | \theta') q(\theta | \theta')}{p_0(\theta) \prod_{i=1}^N p(x_i | \theta) q(\theta' | \theta)}, \quad (8)$$

where  $p_0, p$ , and  $q$  match the corresponding functions within Equation 3. We separate out terms dependent and independent of the data  $x$  as:

$$\Delta(\theta, \theta') = \sum_{i=1}^N \log \frac{p(x_i | \theta')}{p(x_i | \theta)} - \psi(1, \theta, \theta') = \Lambda(\theta, \theta') - \psi(1, \theta, \theta'). \quad (9)$$

A minibatch estimator of  $\Delta$ , denoted as  $\Delta^*$ , is

$$\Delta^*(\theta, \theta') = \frac{N}{b} \sum_{i=1}^b \log \frac{p(x_i | \theta')}{p(x_i | \theta)} - \psi(1, \theta, \theta') = \Lambda^*(\theta, \theta') - \psi(1, \theta, \theta') \quad (10)$$

Note that  $\Delta$  and  $\Delta^*$  are evaluated on the full dataset and a minibatch of size  $b$  respectively. The scaling term  $N/b$  ensures that  $\Delta^*(\theta, \theta')$  is an unbiased estimator of  $\Delta(\theta, \theta')$ .

The key to our test is a smooth acceptance function. We consider functions other than the classical Metropolis test that satisfy the detailed balance (see Equation 2) condition needed for accurate posterior estimation. A class of functions leading to detailed balance is specified as follows:

**Lemma 2.** *If  $g(s)$  is any function such that  $g(s) = \exp(s)g(-s)$ , then the acceptance function  $\alpha(\theta, \theta') \triangleq g(\Delta(\theta, \theta'))$  satisfies detailed balance.*

This result is used in [6] to define the Barker acceptance test. As a sanity check, choosing  $g(s) = \exp(s) \wedge 1$  — a function satisfying the requirement of Lemma 2 — produces the classical Metropolis acceptance test  $\alpha(\theta, \theta') = g(\Delta(\theta, \theta')) = \frac{p(\theta')q(\theta|\theta')}{p(\theta)q(\theta'|\theta)} \wedge 1$ . In fact,  $g(s) = \exp(s) \wedge 1$  is the optimal acceptance function in terms of acceptance rate, since it accepts with probability 1 for  $\Delta > 0$ .

#### 3.2 Barker (Logistic) Acceptance Function

For our new MH test we use the Barker logistic [6] function:  $g(s) = (1 + \exp(-s))^{-1}$ . Straightforward arithmetic shows that it satisfies the condition in Lemma 2. While it is slightly less efficient than the Metropolis test when used on the full dataset, we will see that its smoothness allows it to naturally tolerate substantial variance in its input argument. This in turn will lead to a much more efficient test on subsets of data.

Assume we begin with the current sample  $\theta$  and a candidate sample  $\theta'$ , and that  $V \sim \mathcal{U}(0, 1)$  is a uniform random variable. We accept  $\theta'$  if  $g(\Delta(\theta, \theta')) > V$ , and reject otherwise. Since  $g(s)$  is monotonically increasing, its inverse  $g^{-1}(s)$  is well-defined and unique. So an equivalent test is to accept  $\theta'$  iff

$$\Delta(\theta, \theta') > X = g^{-1}(V) \quad (11)$$

where  $X$  is a random variable with the logistic distribution (its CDF is the logistic function). To see this notice that  $\frac{dV}{dX} = g'$ , that  $g'$  is the density corresponding to a logistic CDF, and finally that  $\frac{dV}{dX}$  is the density of  $X$ . The density of  $X$  is symmetric, so we can equivalently test whether

$$\Delta(\theta, \theta') + X > 0 \quad (12)$$

for a logistic random variable  $X$ .

#### 3.3 A Minibatch Acceptance Test

We now describe acceptance testing using the minibatch estimator  $\Delta^*(\theta, \theta')$ . From Equation 10,  $\Delta^*(\theta, \theta')$  can be represented as a constant term plus the mean of  $b$  IID terms  $\Lambda_i(\theta, \theta')$  of the form  $N \log \frac{p(x_i | \theta')}{p(x_i | \theta)}$ . As  $b$  increases,  $\Delta^*(\theta, \theta')$  therefore has a distribution which approaches a normal

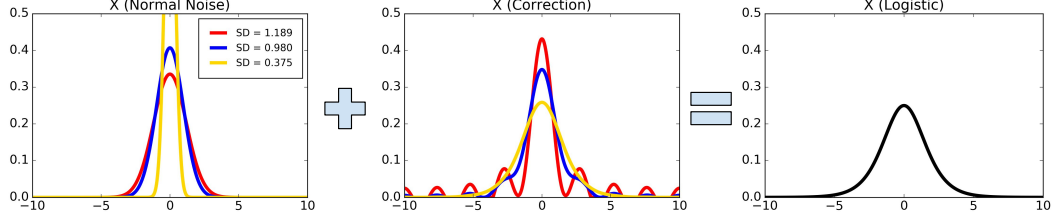


Figure 1: Three examples of  $X_{\text{norm}}$  and  $X_{\text{corr}}$  distributions that convolve to form the standard logistic distribution. These were obtained based on our original linear programming formulation, as discussed in Section 4. We use three standard deviation values of  $X_{\text{norm}}$ . The two red curves convolve to form the logistic, etc. The  $y$ -axis is capped at 0.5 for readability. This figure must be viewed in color.

distribution by the Central Limit Theorem. We now describe this using an asymptotic argument and defer specific bounds between the CDFs of  $\Delta^*(\theta, \theta')$  and a Gaussian to Section 5.

In the limit, since  $\Delta^*$  is normally distributed about its mean  $\Delta$ , we can write

$$\Delta^* = \Delta + X_{\text{norm}}, \quad X_{\text{norm}} \sim \tilde{\mathcal{N}}(0, \sigma^2(\Delta^*)), \quad (13)$$

where  $\tilde{\mathcal{N}}(0, \sigma^2(\Delta^*))$  denotes a distribution which is approximately normal with variance  $\sigma^2(\Delta^*)$ . But to perform the test in Equation 12 we want  $\Delta + X$  for a logistic random variable  $X$  (call it  $X_{\text{log}}$  from now on). In [7] it was proposed to use  $\Delta^*$  in a Barker test anyway and tolerate the fixed error caused by this approximation.

Our approach is to instead decompose  $X_{\text{log}}$  as

$$X_{\text{log}} = X_{\text{norm}} + X_{\text{corr}}, \quad (14)$$

where we assume  $X_{\text{norm}} \sim \mathcal{N}(0, \sigma^2)$  and that  $X_{\text{corr}}$  is a zero-mean “correction” variable with density  $C_\sigma(X)$ . The two variables are added (i.e., their distributions convolve) to form  $X_{\text{log}}$ . This decomposition requires an appropriate  $X_{\text{corr}}$  distribution. We show in Section 4 that it is possible to use deconvolution to derive a highly accurate representation of  $C_\sigma(X)$ .

Using  $X_{\text{corr}}$  samples from  $C_\sigma(X)$ , the acceptance test is now

$$\Delta + X_{\text{log}} = (\Delta + X_{\text{norm}}) + X_{\text{corr}} = \Delta^* + X_{\text{corr}} > 0. \quad (15)$$

Therefore, assuming the variance of  $\Delta^*$  is small enough, if we have an estimate of  $\Delta^*$  from the current data minibatch, we test acceptance by adding a random variable  $X_{\text{corr}}$  and then accept  $\theta'$  if the result is positive (and reject otherwise).

If  $\tilde{\mathcal{N}}(0, \sigma^2(\Delta^*))$  is exactly  $\mathcal{N}(0, \sigma^2(\Delta^*))$ , the above test is exact as well, and as we show in Section 5, if there is a maximum error  $\epsilon$  between the CDF of  $\tilde{\mathcal{N}}(0, \sigma^2(\Delta^*))$  and the CDF of  $\mathcal{N}(0, \sigma^2(\Delta^*))$ , then the acceptance test has an error of at most  $\epsilon$  relative to the full batch version.

## 4 Computing the Correction Distribution

Our proposed test in Equation 15 requires knowing the distribution of the correction variable  $X_{\text{corr}}$  such that  $X_{\text{norm}} + X_{\text{corr}} = X_{\text{log}}$ , where  $X_{\text{norm}} \sim \mathcal{N}(0, \sigma^2)$  and  $X_{\text{log}}$  has a standard logistic CDF,  $(1 + \exp(-X))^{-1}$ . In Section 5, we show that the accuracy of the test depends on the absolute error between the CDFs of  $X_{\text{norm}} + X_{\text{corr}}$  and  $X_{\text{log}}$ . Consequently, we need to minimize this in our construction of  $X_{\text{corr}}$ . More formally, let  $\Phi_{s_X} = \Phi(X/s_X)$  where  $\Phi$  is the standard normal CDF<sup>3</sup>,  $S(X)$  be the logistic function, and  $C_\sigma(X)$  be the *density* of the correction  $X_{\text{corr}}$  distribution. Our goal, based on Equation 14, is to solve the following optimization problem:

$$C_\sigma^* = \arg \min_{C_\sigma} |\Phi_{s_X} * C_\sigma - S| \quad (16)$$

where  $*$  denotes convolution.

<sup>3</sup>Hence,  $\Phi_{s_X}$  is the CDF of a zero-mean Gaussian with standard deviation  $s_X$ .

For computation of  $C_\sigma$ , we assume that its input  $Y$  and another variable  $X$  lie in the intervals  $[-V, V]$  and  $[-2V, 2V]$ , respectively. In continuous form, the optimization problem is therefore

$$C_\sigma^* = \arg \min_{C_\sigma} \sup_{x \in [-2V, 2V]} \left| \int_{-V}^V \Phi_\sigma(x - y) C_\sigma(y) dy - S(x) \right| \quad (17)$$

To get  $C_\sigma^*$  in a computable form, we discretize  $X$  and  $Y$  into  $4N + 1$  and  $2N + 1$  values respectively. If  $i \in \{-2N, \dots, 2N\}$  and  $j \in \{-N, \dots, N\}$ , then we can write  $X_i = i(V/N)$  and  $Y_j = j(V/N)$ , and the condition in Equation 17 can be re-expressed as:

$$C_\sigma^* = \arg \min_{C_\sigma} \max_{i \in \{-2N, \dots, 2N\}} \left| \sum_{j=-N}^N \Phi_\sigma(X_i - Y_j) C_\sigma(Y_j) - S(X_i) \right|. \quad (18)$$

To make the problem easier to understand, we can define a matrix  $M$  and vectors  $u$  and  $v$  such that  $M_{ij} = \Phi_\sigma(X_i - Y_j)$ ,  $u_j = C_\sigma(Y_j)$ , and  $v_i = S(X_i)$ , where the indices  $i$  and  $j$  are appropriately translated to be non-negative indices for  $M$ ,  $u$ , and  $v$ . Thus, Equation 18 is equivalent to minimizing the  $L_\infty$  norm:

$$u^* = \arg \min_u \|Mu - v\|_\infty. \quad (19)$$

We have the additional constraint that  $u_j > 0$  for all  $j$ , since  $u$  represents a density. We first explored optimizing this problem with linear programming to find a  $u$  minimizing  $\|Mu - v\|_\infty$  subject to  $u \geq 0$ . Figure 1 shows three examples of  $X_{\text{norm}}$  and  $X_{\text{corr}}$  densities that sum to  $X_{\text{log}}$  from this process.

This was tractable up to a few hundred dimensions for  $u$ . However, the discretization error is bounded by the curvature of the underlying functions times the squared quantization step. The curvatures are here slightly less than one, i.e. the errors are of the order of  $(V/N)^2$ . Here we chose  $V = 20$  to provide adequate containment of the distributions (the CDFs are extremely close to either zero or one outside this range). So with 200 points, we have a discretization error of approximately 0.01, which is relatively large. To yield higher resolution and lower error, we switched to a least squares solution:

$$u^* = \arg \min_u \|Mu - v\|_2^2 + \lambda \|u\|_2^2, \quad (20)$$

where we add the standard regularization factor  $\lambda$  for stability in high dimensions. The solution is well-known from the normal equations:  $u^* = (M^T M + \lambda I)^{-1} M^T v$ , and in practice, it yields an acceptable  $L_\infty$  norm for Equation 19.

With this approach, there is no guarantee that  $u^* \geq 0$ . However, we have some flexibility in the choice of  $\sigma$  in Equation 17. As we decrease the variance of  $X_{\text{norm}}$ , the variance of  $X_{\text{corr}}$  grows by the same amount and is in fact the result of convolution with a Gaussian whose variance is the difference. Thus as  $\sigma$  decreases,  $C_\sigma(X)$  grows and approaches the derivative of a logistic function at  $\sigma = 0$ . It retains some very weak negative values for  $\sigma > 0$  but removal of those values leads to very small error.

As can be seen from Table 1, the errors between  $X_{\text{norm}} + X_{\text{corr}}$  and  $X_{\text{log}}$  can be made very small, approaching single floating precision (about  $10^{-7}$ ).

The complete procedure is given in Algorithm 1. A few points:

- It uses an adaptive step size so as to use the smallest possible average minibatch size. Unlike previous work however (and as we show in Section 6) the size distribution is short-tailed.
- An additional normal variable  $X_{\text{nc}}$  is added to  $\Delta^*$  to produce a variable with unit variance. This is not mathematically necessary, but allows us to use a single correction distribution  $C_1$  with  $\sigma = 1$  for  $X_{\text{corr}}$ , saving on memory footprint.
- The sample variance  $s_{\Delta^*}^2$  is proportional to  $\|\theta' - \theta\|_2^2$  whose distribution for Normal proposals is the square of a normal variable.

## 5 Analysis

We now derive error bounds for our M-H test, and for the approximate target distribution that it generates. From Table 1, we know that it is possible to generate the correction samples  $X_{\text{corr}}$  with a

**Input** : Number of samples  $T$ , minibatch size  $m$ , error bound  $\delta$ , pre-computed correction  $C_1(X)$  distribution, initial sample  $\theta_1$ .

**Output** : A chain of  $T$  samples  $\{\theta_1, \dots, \theta_T\}$  from  $p(\theta)$ ;

**for**  $t = \{1, \dots, T\}$  **do**

Propose a candidate  $\theta'$  from proposal distribution  $q(\theta' | \theta_t)$ ;

Draw a minibatch of  $m$  points  $x_i$ , compute  $\Delta^*(\theta_t, \theta')$  and sample variance  $s_{\Delta^*}^2$ ;

Estimate moments  $E|\Lambda_i - \Lambda|$  and  $E|\Lambda_i - \Lambda|^3$  from the sample, and error  $\epsilon$  from Equation 27;

**while**  $s_{\Delta^*}^2 \geq 1$  **or**  $\epsilon > \delta$  **do**

Draw  $m$  more samples to augment the minibatch, update  $\Delta^*$ ,  $s_{\Delta^*}^2$  and  $\epsilon$  estimates;

**end**

Draw  $X_{\text{nc}} \sim \mathcal{N}(0, 1 - s_{\Delta^*}^2)$  and  $X_{\text{corr}}$  from the correction distribution  $C_1(X)$ ;

**if**  $\Delta^* + X_{\text{nc}} + X_{\text{corr}} > 0$  **then**

Accept the candidate,  $\theta_{t+1} = \theta'$ ;

**else**

Reject and re-use the old sample,  $\theta_{t+1} = \theta_t$ ;

**end**

**end**

**Algorithm 1:** A description of M-H sampling with our acceptance test.

Table 1: Error in  $X_{\text{norm}} + X_{\text{corr}}$  versus  $X_{\text{log}}$

$N$	$\sigma$	$\lambda$	$L_\infty$ error
4000	0.9	1	1.0e-4
4000	0.8	0.03	5.0e-6

CDF error approaching single-precision floating point error. We therefore treat the  $X_{\text{corr}}$  variable as a sample from the exact correction distribution and we will not analyze its errors.

In the most similar prior works, [3] uses asymptotic arguments based on the Central Limit Theorem to argue that its approximate acceptance test error tends to zero as batch size increases. But no quantitative bounds are given. In [4], explicit bounds are given, but they depend on bounding

$$C_{\theta, \theta'} = \max_{1 \leq i \leq N} |\log p(x_i | \theta') - \log p(x_i | \theta)|. \quad (21)$$

Such bounds can be derived efficiently for models such as logistic regression, but it is unclear how to derive them for a complex model such as a neural network. In general, since a new  $\theta'$  value is obtained each iteration, one would need to iterate through all the  $p(x_i | \theta')$  terms, defeating the purpose of minibatch MCMC.

In this paper, we use quantitative forms of the Central Limit Theorem which rely on measurable statistics from a single minibatch. Thus a sampler using our approach does not need to “see” data beyond the current minibatch. This supports use of the sampler on very large datasets, and in online mode where the dataset is presented as a stream.

For the analysis, in Section 5.1, we first present bounds on the absolute and relative error (in terms of the CDFs) of the distribution of  $\Delta^*$  vs a Gaussian. We then show in Section 5.2 that these error bounds are preserved after the addition of other random variables, in particular  $X_{\text{nc}}$  and  $X_{\text{corr}}$ . From this it follows that the acceptance test has the same error bound.

### 5.1 Bounding the Distribution of $\Delta^*$ Versus a Gaussian

To start, we use the following quantitative central-limit result:

**Lemma 3.** *Let  $X_1, \dots, X_n$  be a set of zero-mean, independent, identically-distributed random variables with sample mean  $\bar{X}$  and sample variance  $s_X^2$  where:*

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{and} \quad s_X = \frac{1}{n} \left( \sum_{i=1}^n (X_i - \bar{X})^2 \right)^{\frac{1}{2}}. \quad (22)$$



This implies  $t = \bar{X}/s_X$  has an approximate Student's distribution which approaches a normal distribution in the limit. Then

$$\sup_x |\Pr(t < x) - \Phi(x)| \leq \frac{6.4E|X|^3 + 2E|X|}{\sqrt{n}}. \quad (23)$$

*Proof.* The following bound is given immediately after Corollary 2 in from [15]:

$$-6.4E|X|^3 - 2E|X| \leq \sup_x |\Pr(t < x) - \Phi(x)|\sqrt{n} \leq 1.36E|X|^3 \quad (24)$$

This bound applies to  $x \geq 0$ . Applying the bound to  $-x$  when  $x < 0$  and combining with  $x > 0$ , we obtain the weaker but unqualified bound in Equation 23.  $\square$

Lemma 3 demonstrates that as long as we know the first and third absolute moments  $E|X|$  and  $E|X|^3$ , we can bound the error of the normal approximation, which decays as  $O(n^{-\frac{1}{2}})$ . Making the change of variables  $y = xs_X$ , Equation 23 becomes

$$\sup_y \left| \Pr(\bar{X} < y) - \Phi\left(\frac{y}{s_X}\right) \right| \leq \frac{6.4E|X|^3 + 2E|X|}{\sqrt{n}}, \quad (25)$$

showing that the distribution of  $\bar{X}$  approaches the normal distribution  $\mathcal{N}(0, s_X)$  whose variance is  $s_X$ , measured from the sample.

To apply this to our test, let  $X_i = \Lambda_i(\theta, \theta') - \Lambda(\theta, \theta')$ , so that the  $X_i$  are zero-mean, i.i.d. variables. If instead of all  $n$  samples, we only extract a subset of  $b$  samples corresponding to our minibatch, we can connect  $\bar{X}$  with our  $\Delta^*$  term:

$$\bar{X} = \Delta^*(\theta, \theta') - \Delta(\theta, \theta'), \quad (26)$$

so that  $s_X = s_{\Delta^*}$ . This results in the following Corollary:

**Corollary 1.** *We can now substitute into Equation 25 and displace by the mean, giving:*

$$\sup_y \left| \Pr(\Delta^* < y) - \Phi\left(\frac{y - \Delta}{s_{\Delta^*}}\right) \right| \leq \frac{6.4E|X|^3 + 2E|X|}{\sqrt{b}} = \epsilon(\theta, \theta', b) \quad (27)$$

Corollary 1 shows that the distribution of  $\Delta^*$  approximates a Normal distribution with mean  $\Delta$  and variance  $s_{\Delta^*}^2$ . Furthermore, it bounds the error with *estimable quantities*: both  $E|X|$  and  $E|X|^3$  can be estimated as means of  $|\Lambda_i - \Lambda|$  and  $|\Lambda_i - \Lambda|^3$ , respectively, on each minibatch. We expect this will often be accurate enough on minibatches with hundreds or thousands of points, but otherwise bootstrap CIs can be computed from those sequences. Since the bounds are monotone in  $E|X|$  and  $E|X|^3$ , using upper bootstrap CI limits will provide high-confidence error bounds.

## 5.2 Bounds are Preserved After Adding Random Variables

We next relate the CDFs of distributions and show that bounds are preserved after adding random variables.

**Lemma 4.** *Let  $P(x)$  and  $Q(x)$  be two cumulative distributions satisfying  $\sup_x |P(x) - Q(x)| \leq \epsilon$  with  $x$  in some real range. Let  $R(y)$  be the density of another random variable  $y$ . Let  $P'$  be the convolution  $P * R$  and  $Q'$  be the convolution  $Q * R$ . Then  $P'(z)$  (resp.  $Q'(z)$ ) is the CDF of sum  $z = x + y$  of independent random variables  $x$  with CDF  $P(x)$  (resp.  $Q(x)$ ) and  $y$  with density  $R(y)$ . Then*

$$\sup_x |P'(x) - Q'(x)| \leq \epsilon \quad (28)$$

*Proof.* We first observe that

$$P'(z) - Q'(z) = \int_{-\infty}^{+\infty} (P(z - x) - Q(z - x))R(x)dx, \quad (29)$$

and since  $\sup_x |P(x) - Q(x)| \leq \epsilon$  it follows that for all  $z$ :

$$-\epsilon = \int_{-\infty}^{+\infty} -\epsilon R(x)dx \leq \int_{-\infty}^{+\infty} (P(z - x) - Q(z - x))R(x)dx \leq \int_{-\infty}^{+\infty} \epsilon R(x)dx = \epsilon \quad (30)$$

$\square$

Relating Lemma 4 to our acceptance test, we observe that:

**Corollary 2.** *If  $\sup_y |\Pr(\Delta^* < y) - \Phi(\frac{y-\Delta}{s_{\Delta^*}})| \leq \epsilon(\theta, \theta', b)$ , then*

$$\sup_y |\Pr(\Delta^* + X_{\text{nc}} + X_{\text{corr}} < y) - S(y - \Delta)| \leq \epsilon(\theta, \theta', b) \quad (31)$$

where  $S(x)$  is the standard logistic function, and  $X_{\text{nc}}$  and  $X_{\text{corr}}$  are generated as per Algorithm 1.

*Proof.* We apply Lemma 4 twice. First take:

$$P(y) = \Pr(\Delta^* < y) \quad \text{and} \quad Q(y) = \Phi\left(\frac{y - \Delta}{s_{\Delta^*}}\right) \quad (32)$$

and convolve with the distribution of  $X_n$  which has density  $\phi(X/\sigma_n)$  where  $\sigma_n^2 = 1 - s_{\Delta^*}^2$ . This yields the next iteration of  $P$  and  $Q$ :

$$P'(y) = \Pr(\Delta^* + X_{\text{nc}} < y) \quad \text{and} \quad Q'(y) = \Phi(y - \Delta) \quad (33)$$

Now we convolve with the distribution of  $X_{\text{corr}}$  which gives:

$$P''(y) = \Pr(\Delta^* + X_{\text{nc}} + X_{\text{corr}} < y) \quad \text{and} \quad Q''(y) = S(y - \Delta) \quad (34)$$

Both steps preserve the error bound  $\epsilon(\theta, \theta', b)$ . Finally  $S(y - \Delta)$  is a logistic CDF centered at  $\Delta$ , and so  $S(y - \Delta) = \Pr(\Delta + X_{\text{log}} < y)$  for a logistic random  $X_{\text{log}}$ . We conclude that the probability of acceptance for the actual test  $\Pr(\Delta^* + X_{\text{nc}} + X_{\text{corr}} > 0)$  differs from the exact test  $\Pr(\Delta + X_{\text{log}} > 0)$  by at most  $\epsilon$ .  $\square$

As our discussion in Section 3 showed, as the distribution of  $\Delta^*$  approaches a Gaussian, our new MH test becomes more accurate. Corollary 2 shows that the bounds from Section 5.1 are preserved after the addition of the random variables we use, showing that our test should remain accurate.

In fact we can do better than this approximation (showing the error decreases as  $O(n^{-1})$ ) by using a more precise limit distribution under an additional assumption. Let  $\mu_i$  denote the  $i^{\text{th}}$  moment, and  $b_i$  denote the  $i^{\text{th}}$  absolute moment of  $X$ . If Cramer's condition holds:

$$\lim_{t \rightarrow \infty} \sup |E(\exp(itX))| < 1, \quad (35)$$

then Equation 2.2 in Bentkus et al.'s work on Edgeworth expansions [16] provides the following:

**Lemma 5.** *Let  $X_1, \dots, X_n$  be a set of zero-mean, independent, identically-distributed random variables with sample mean  $\bar{X}$  and with  $t$  defined as in Lemma 3. If  $X$  satisfies Cramer's condition, then*

$$\sup_x \left| \Pr(t < x) - G\left(x, \frac{\mu_3}{b_2^{3/2}}\right) \right| \leq \frac{c(\epsilon, b_2, b_3, b_4, b_{4+\epsilon})}{n} \quad (36)$$

where

$$G_n(x, y) = \Phi(x) + \frac{y(2x^2 + 1)}{6\sqrt{n}} \Phi'(x). \quad (37)$$

Lemma 5 shows that the average of the  $X_i$  has a more precise, skewed CDF limit  $G_n(x, y)$  where the skew term has weight proportional to a certain measure of skew derived from the moments:  $\frac{\mu_3}{b_2^{3/2}}$ .

Note that if the  $X_i$  are symmetric, the weight of the correction term is zero, and the CDF of the average of the  $X_i$  converges to  $\Phi(x)$  at a rate of  $O(n^{-1})$ .

Here the limit  $G_n(x, y)$  is a normal CDF plus a correction term that decays as  $n^{1/2}$ . Importantly, since  $\phi''(x) = x^2\phi(x) - \phi(x)$  where  $\phi(x) = \Phi'(x)$ , the correction term can be rewritten giving:

$$G_n(x, y) = \Phi(x) + \frac{y}{6\sqrt{n}} (2\phi''(x) + 3\phi(x)) \quad (38)$$

From which we see that  $G_n(x, y)$  is a linear combination of  $\Phi(x)$ ,  $\phi(x)$  and  $\phi''(x)$ . This is quite fortuitous. In Algorithm 1, we correct for the difference in  $\sigma$  between  $\Delta^*$  and the variance needed by  $X_{\text{corr}}$  using  $X_{\text{nc}}$ . This same method works when we wish to estimate the error in  $\Delta^*$  vs  $G_n(x, y)$ . Since all of the component functions of  $G_n(x, y)$  are derivatives of a (unit variance)  $\Phi(x)$ , adding a

normal variable with variance  $\sigma'$  increases the variance of all three functions to  $1 + \sigma'$ . Thus we add  $X_{nc}$  as per Algorithm 1 preserving the limit in Equation 38.

The deconvolution approach can be used to construct a correction variable  $X_{corr}$  between  $G_n(x, y)$  and  $S(x)$  the standard logistic function. An additional complexity is that  $G_n(x, y)$  has additional parameters  $y$  and  $n$ . Since these act as a single multiplier  $\frac{y}{6\sqrt{n}}$  in Equation 38, its enough to consider a function  $g(x, y')$  parametrized by  $y' = \frac{y}{6\sqrt{n}}$ . This function can be computed and saved offline. As we have shown above, errors in the “limit” function propagate directly through as errors in the acceptance test. To achieve a test error of  $10^{-6}$  (close to single floating point precision), we need a  $y'$  spacing of  $10^{-6}$ . It should not be necessary to tabulate values all the way to  $y' = 1$ , since  $y'$  is scaled inversely by the square root of minibatch size. Assuming a max  $y'$  of 0.1 requires us to tabulate about 100,000. Since our  $x$  resolution is 10,000, this leads to a table with about 1 billion values, which can comfortably be stored in memory. However, if  $g(x, y)$  is moderately smooth in  $y$ , it should be possible to achieve similar accuracy with a much smaller table. We leave further analysis and experiments with  $g(x, y)$  as future work.

## 6 Experiments

We conduct two sets of experiments to explore the benefits of our minibatch MH test and to benchmark it with previous work. In Section 6.1, we show that our test enables samples to converge to the posterior distribution of a heated Gaussian mixture model (additional details are in Appendix B). In Section 6.2, we analyze its efficiency on logistic regression.

### 6.1 Mixture of Gaussians

We start with a simple Gaussian mixture model, borrowing an experiment from [12]. The parameter is 2-D,  $\theta = \langle \theta_1, \theta_2 \rangle$ , and the parameter/data generation process is

$$\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \right); \quad x_i \sim 0.5 \cdot \mathcal{N}(\theta_1, \sigma_x^2) + 0.5 \cdot \mathcal{N}(\theta_1 + \theta_2, \sigma_x^2). \quad (39)$$

We set  $\sigma_1^2 = 10$ ,  $\sigma_2^2 = 1$  and  $\sigma_x^2 = 2$ . We fix  $\theta = \langle 0, 1 \rangle$ . The original paper sampled 100 data points from this distribution, and estimated the posterior in the parameters. We are interested in performance on larger problems and so sampled 1,000,000 points to form the posterior of  $p(\theta) \prod_{i=1}^{1,000,000} p(x_i | \theta)$  with the same prior from Equation 39. The larger number of samples produces a much sharper posterior with two very narrow peaks. Our goal is to reproduce the original posterior, so we adjust the temperature to  $T = 10,000$ . Taking logs, we get the target as shown in the far left of Figure 2.

We run MCMC with our MH test using minibatch size  $m = 100$ . We also run this using the methods from [3] (“Korattikara 2014”) and [4] (“Bardenet 2014”), also with  $m = 100$ . For [3], we increment minibatches by 100 as recommended, and for [4], we increase the minibatch size geometrically with a ratio of  $\gamma = 1.5$  as recommended. These tests use different minibatch size updates because of their different approaches to repeated testing. [4] uses straightforward discounting, allocating a fraction of the test error  $\delta$  to each minibatch test (call these subtests). Its important to bound the number of subtests therefore, so that the tolerable error per subtest is not too small. [3] uses a more aggressive repeated testing approach, explicitly modeling correlation between tests and allowing many more tests for a given  $\delta$ .

The tolerance for making a decision in [3] is  $\epsilon = 0.005$ , and the error bound control parameters in [4] are  $p = 2$  and  $\delta = 0.01$ . All methods (including ours) use the same random walk proposer with covariance  $\Sigma = \text{diag}(0.3, 0.3)$ . Since the proposer generates a random walk, all shaping of distribution of the samples is due to the MH test. All methods are run 5000 times to collect 5000 samples.

Figure 2 shows scatter plots of the resulting  $\theta$  samples for the three methods, with darker regions indicating a greater density of points. There are no obvious differences, so to better understand any differences, we measure the similarity between each set of samples and the actual posterior.

To effect the comparison, we first discretize the posterior coordinates into bins (non-overlapping ranges) with respect to the two components of the parameter  $\theta$ . The probability of a sample falling into one of these bins is the integral of the true posterior probability over the area of that bin. Let

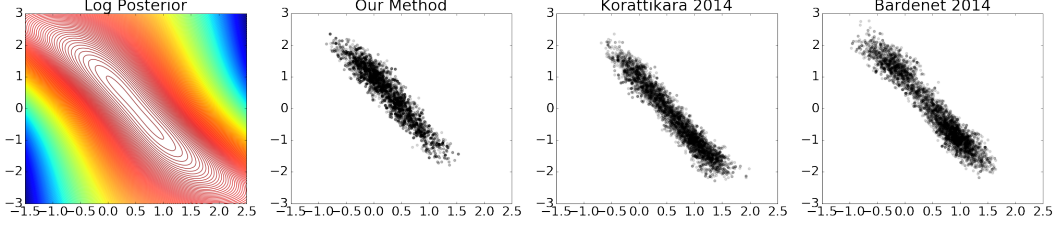


Figure 2: The log posterior contours and scatter plots of sampled  $\theta$  values using different methods.

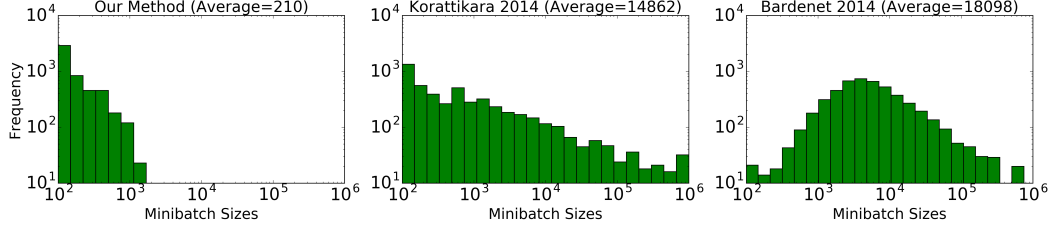


Figure 3: Counts of minibatch sizes that the three methods used in the experiment of Section 6.1. To make comparisons easier, the axes have the same range (using a log-log scale).

this probability be denoted  $P_i$  where  $i$  is the bin number. A single sample from any of the MH methods should therefore be multinomial with distribution  $P$ , and a set of  $n$  (ideally independent) samples should follow a  $\text{Multinomial}(P, n)$  distribution. Since the ideal distribution is simple, we can use it to measure the quality of the sample distributions rather than use general purpose tests like KL-divergence or likelihood-ratio. The latter in any case are problematic with zero counts in some bins as we have here.

In practice, it is difficult to compute multinomial probabilities for large  $n$ , in that case the per-bin distributions are well-approximated by Poisson distributions with parameter  $\lambda_i = P_i n$ . Given a set of samples  $\Theta = \{\theta_1, \dots, \theta_T\}$ , let  $c_j$  denote the number of individual samples  $\theta_i$  that fall in bin  $j$ . Then the log probability for the set is approximated as

$$\log p(\Theta) \approx \frac{1}{N_{\text{bins}}} \sum_{j=1}^{N_{\text{bins}}} c_j \log(nP_j) - nP_j - \log(\Gamma(c_j + 1)). \quad (40)$$

The results for each model are shown in Table 2 ... [Daniel: IN PROGRESS](#)

While the Poisson results provide general guidance on the relative accuracy of the models it is difficult interpret the quantitative scores. We therefore perform significance tests to test the null hypothesis that the ideal distribution really generated the samples. We perform two tests:

1. An exact test using simulation to generate an approximate distribution of scores given the null hypothesis.
2. An approximate test using the chi-squared distribution as the test statistic.

The results were... [Daniel: IN PROGRESS](#)

Figure 3 shows a histogram of the final minibatch sizes used by the three methods in each iteration. Our method consumes significantly less data; most minibatch sizes are smaller than 1000, and the average size is 210. The other methods occasionally need to consume a large proportion of the entire data set of one million elements; the average minibatch sizes are 14862 and 18098 for Korattikara 2014 and Bardenet 2014, respectively. The average minibatch sizes accurately predict the running times of these methods since all methods have a running time proportional to the total data consumed, with one caveat: the method of [4] requires a global bound on the term  $C_{\theta, \theta'}$  (see Equation 21) over the entire the data set at every iteration in order to calculate the error bound. This bound must be either computed analytically using properties of the specific dataset, or it must be checked at each sample if a black-box test is required. We found that the sample code provided by the authors of [4]

Table 2: [Daniel: TODO Poisson Probabilities](#)

Model	Ours	Korattikara 2014	Bardenet 2014
Poisson Value (Eq. 40)	TODO	TODO	TODO

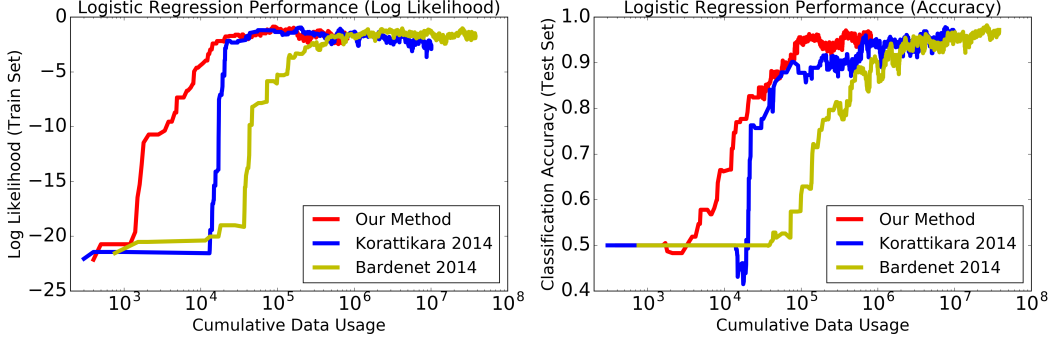


Figure 4: Logistic regression performance (accuracy/log likelihood) based on cumulative data usage.

in fact computed this bound explicitly for each sample generated, i.e. the implementation traversed the entire dataset for each sample, providing no performance advantage over the complete test. For this reason, we omitted the time to compute  $C_{\theta, \theta'}$  from our measurements of the running time of the sample implementation of [4], and consider only the time to process the minibatches.

## 6.2 Logistic Regression

We next use logistic regression for the binary classification of 1s versus 7s in the MNIST dataset [17]. The data has 12007 training and 1000 testing elements (we used a random subset of the test data). The pixel values are scaled to be in  $(0, 1)$ . We impose a flat prior on  $\theta$ . For the proposer, we again use a random walk, this time with covariance matrix  $0.05I$  for the  $784 \times 784$  identity matrix  $I$ . The posterior temperature is set at  $T = 1000$ . We run our MH test and compare with Korattikara 2014 and Bardenet 2014; minibatch sizes and other parameters are set to be the same value as with the Gaussian Mixtures experiment from Section 6.1.

Figure 4 shows the training log likelihood and the testing prediction accuracy, each as a function of the cumulative training data points processed.<sup>4</sup> To generate the curves, for each of the 5000 sampled vectors  $\theta_t$ ,  $t \in \{1, \dots, 5000\}$ , we use  $\theta_t$  as the parameter for logistic regression. We see that our minibatch MH test is more efficient compared to the other two methods, because it achieves similar or superior performance while consuming fewer data points.

Figure 5, in a similar manner as Figure 3, shows the histogram of minibatch sizes for all three methods on a log-log scale. With an initial minibatch size of 100, Korattikara 2014 and Bardenet 2014 achieve average minibatch size of 2026 and 7612 respectively over the MNIST classification task, while our method achieves that of 163, showing significantly better performance than the benchmark test methods.

## 7 Conclusions

In this paper, we have derived a new MH test for minibatch MCMC methods. We demonstrated how a simple deconvolution process allows us to use a minibatch approximation to the full data tests. We experimentally show the benefits of our test on Gaussian mixtures and logistic regression. Straightforward directions for future work include running more experiments with a particular focus on investigation of the variance precondition. More elaborate extensions include testing on deep neural networks and combining our results with Hamiltonian Monte Carlo methods, providing

<sup>4</sup>Note that the curves do not span the same length over the x-axis, because our test consumes fewer samples throughout the MCMC procedure, so the corresponding curve will “end” before the other two.

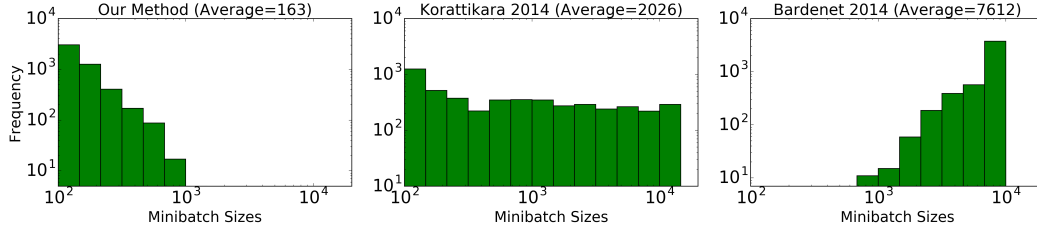


Figure 5: Counts of minibatch sizes that the three methods used in the MNIST logistic regression experiment from Section 6.2. It is analogous to Figure 3.

a recipe for how to use our algorithm (following the framework of [18]), or integrating parallel MCMC [19, 20] concepts.

## References

- [1] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, “Equation of state calculations by fast computing machines,” *The Journal of Chemical Physics*, vol. 21, 1953.
- [2] W. K. Hastings, “Monte carlo sampling methods using markov chains and their applications,” *Biometrika*, vol. 57, pp. 97–109, 1970.
- [3] A. Korattikara, Y. Chen, and M. Welling, “Austerity in MCMC land: Cutting the metropolis-hastings budget,” in *Proceedings of the 31st International Conference on Machine Learning (ICML)*, 2014.
- [4] R. Bardenet, A. Doucet, and C. Holmes, “Towards scaling up markov chain monte carlo: an adaptive subsampling approach,” in *Proceedings of the 31st International Conference on Machine Learning (ICML)*, 2014.
- [5] D. Maclaurin and R. P. Adams, “Firefly monte carlo: Exact MCMC with subsets of data,” in *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence (UAI)*, 2014.
- [6] A. A. Barker, “Monte-carlo calculations of the radial distribution functions for a proton-electron plasma,” *Australian Journal of Physics*, vol. 18, pp. 119–133, 1965.
- [7] R. Bardenet, A. Doucet, and C. Holmes, “On markov chain monte carlo methods for tall data,” *arXiv preprint arXiv:1505.02827v1*, 2015.
- [8] G. O. Roberts and J. S. Rosenthal, “Optimal scaling for various metropolis–hastings algorithms,” *Statistical Science*, vol. 16, no. 4, p. 351–367, 2001.
- [9] W. Gilks and D. Spiegelhalter, *Markov chain Monte Carlo in practice*. Chapman & Hall/CRC, 1996.
- [10] S. Brooks, A. Gelman, G. Jones, and X.-L. Meng, *Handbook of Markov Chain Monte Carlo*. CRC press, 2011.
- [11] R. M. Neal, “MCMC using Hamiltonian dynamics,” *Handbook of Markov Chain Monte Carlo*, vol. 54, pp. 113–162, 2010.
- [12] M. Welling and Y. W. Teh, “Bayesian learning via stochastic gradient langevin dynamics,” in *Proceedings of the 28th International Conference on Machine Learning (ICML)*, 2011.
- [13] S. Ahn, A. K. Balan, and M. Welling, “Bayesian posterior sampling via stochastic gradient fisher scoring,” in *Proceedings of the 29th International Conference on Machine Learning (ICML)*, 2012.
- [14] T. Chen, E. Fox, and C. Guestrin, “Stochastic gradient Hamiltonian Monte Carlo,” in *Proceedings of the 31st International Conference on Machine Learning (ICML)*, 2014.
- [15] Y. Novak, “On self-normalized sums and student’s statistic,” *Theory of Probability and its Applications*, vol. 49, no. 2, pp. 336–344, 2005.
- [16] V. Bentkus, F. Gotze, and W.R.vanZwet, “An edgeworth expansion for symmetric statistics,” *Annals of Statistics*, vol. 25, no. 2, 1997.
- [17] Y. LeCun and C. Cortes, “MNIST handwritten digit database,”

- [18] Y. Ma, T. Chen, and E. Fox, “A complete recipe for stochastic gradient mcmc,” in *Advances in Neural Information Processing Systems* 28, 2015.
- [19] E. Angelino, E. Kohler, A. Waterland, M. Seltzer, and R. P. Adams, “Accelerating MCMC via parallel predictive prefetching,” in *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence, (UAI)*, 2014.
- [20] S. Ahn, B. Shahbaba, and M. Welling, “Distributed stochastic gradient MCMC,” in *Proceedings of the 31st International Conference on Machine Learning, (ICML)*, 2014.

## Outline of Appendix

In this appendix, we describe the following topics:

- The proof of Lemma 1.
- More information on Section 6.1.

### A Proof of Lemma 1

Daniel: changes I made:

1. (Somewhat important) I added slightly more detail on why if we assume  $Y$  has variance  $1/N$ , then  $Y \in \frac{1}{\sqrt{N}}[0.5, 1]$  has probability  $\Phi(1) - \Phi(0.5)$ ; note that I'm assuming the positive interval only, not the negative one. However, you also made  $p_0$  double the value of  $p_1$ . I struggled to understand why, until I realized that you used the phrase “with matching sign.” So I *think* you meant that  $\theta' - \theta$  could be either positive or negative, but that  $\theta - 0.5$  was restricted in its sign so could be only half its value. Is my thinking right?
2. Changed all occurrences of  $\Lambda^*(\theta', \theta)$  to  $\Lambda^*(\theta, \theta')$ , i.e., switched the input ordering.
3. Originally,  $\log u_0$  was defined as  $N(\Lambda^*(\theta, \theta') - \psi(1, \theta, \theta'))$  but I changed the starting definition so it does *not* have a factor of  $N$  in its definition (because we later get  $N$  to appear in a subsequent step).
4. In addition,  $\psi(1, \theta, \theta') = 0$  due to uniform priors and symmetric proposals and there didn't seem to be a good enough reason to include it in the definition so I deleted it (sometimes it's insightful to have a nice form of zero inside, but I didn't see it here).
5. I added in an “ $E$ ” because it seemed like you wanted the “ $u$ ” that would make  $\psi(u, \theta, \theta') = \log u = E[\Lambda^*(\theta, \theta')]$ .
6. Added more detail for the  $\Lambda^*(\theta, \theta') \approx \psi(u, \theta, \theta')$  part.
7. (TODO later) After Equation 45, we say that “the resulting test with our chosen  $u_0$  will never correctly succeed [...]”, but I think we should say it will not succeed with high probability, because we could select a minibatch of samples such that the mean *actually is* different from 0.5 (it depends on the distribution of the variable that corresponds to  $Y = \frac{1}{b} \sum_{i=1}^b x_i$ ).
8. (TODO later) Unfortunately, I am unable to understand why “the size of the range in  $u$  is at least  $\exp([-2, -1.125])$ ”, unless the real range is  $\exp([-2, 0])$  (as I think it should be) and you were simply cropping off the extra range.
9. (TODO later) My main confusion now is that I am unable to derive Equation 46, and we may want to clarify the notation used when we put two intervals side by side.
10. (TODO later) I do not understand why after Equation 48, we say that the standard deviation of the LHS “given the interval constraints is at least  $0.5/\sqrt{b}$ .” We already know that the  $X_i$  are iid  $\mathcal{N}(\theta = 0.5, 1)$  variables with a fixed  $\theta = 0.5$  value, so we already know the standard deviation is supposed to be  $\sqrt{(\text{Var}(x_1) + \dots + \text{Var}(x_b))/b^2} = 1/\sqrt{b}$ .
11. I added Equation 49, which I think is how you determined the range. Note that it doesn't seem to matter whether it is  $[-1, 0.75]$  or  $[-0.75, 1]$  since the main thing that matters is that we have maximum absolute value of one.
12. At the end, it said the joint distribution of  $(\theta, \theta' - \theta, u)$ , but I changed it to  $(\theta, \theta', u)$ , based on your change to the text of Lemma 1. But here's a question: how is the sample of  $\theta'$  (or even  $\theta' - \theta$ , as it was written earlier) independent of  $\theta$ ?

Assume  $(\theta' - \theta) \in \pm \frac{1}{\sqrt{N}}[0.5, 1]$  and  $\theta - 0.5 \in \pm \frac{1}{\sqrt{N}}[0.5, 1]$  with matching sign (i.e., first establish  $(\theta' - \theta)$ , then “force”  $\theta - 0.5$  to be the same sign). As discussed in Lemma 1, both  $q(\theta' | \theta)$  and  $p(\theta | x_1, \dots, x_N)$  have variance  $1/N$ . If we denote  $\Phi$  as the CDF of the standard normal distribution, then the former event occurs with probability  $p_0 = 2(\Phi(\sqrt{N} \frac{1}{\sqrt{N}}) - \Phi(\sqrt{N} \frac{0.5}{\sqrt{N}})) = 2(\Phi(1) - \Phi(0.5)) \approx 0.2997$ . The latter event, because we restrict signs, occurs with probability  $p_1 = \Phi(\sqrt{N} \frac{1}{\sqrt{N}}) - \Phi(\sqrt{N} \frac{0.5}{\sqrt{N}}) = \Phi(1) - \Phi(0.5) \approx 0.14988$ .



These events together guarantee that  $\Lambda^*(\theta, \theta')$  is negative, intuitively because the next choice of  $\theta'$  must be “further away” from 0.5 than  $\theta$ . This fact means that we can find a  $u \in (0, 1)$  so that  $\psi(u, \theta, \theta') = \log u < 0$  equals  $E[\Lambda^*(\theta, \theta')]$ . Specifically, choose  $u_0$  to satisfy  $\log u_0 = E[\Lambda^*(\theta, \theta')]$ . Using  $E[x_i] = 0.5$  and Equation 7, we see that

$$\log u_0 = N(\theta' - \theta) \frac{1}{b} \cdot E \left[ \sum_{i=1}^b x_i - \theta - \frac{\theta' - \theta}{2} \right] \quad (41)$$

$$= -N(\theta' - \theta) \left( \theta - 0.5 + \frac{\theta' - \theta}{2} \right). \quad (42)$$

Next, consider the minibatch acceptance test  $\Lambda^*(\theta, \theta') \not\approx \psi(u, \theta, \theta')$  required by these methods, where  $\not\approx$  informally means “significantly different from” under the distribution over samples of  $x_i$ . This turns out to be

$$\Lambda^*(\theta, \theta') \not\approx \psi(u_0, \theta, \theta') \iff N(\theta' - \theta) \cdot \frac{1}{b} \sum_{i=1}^b x_i - \theta - \frac{\theta' - \theta}{2} \not\approx \log u_0 \quad (43)$$

$$\iff \frac{1}{b} \sum_{i=1}^b x_i - \left( \theta + \frac{\theta' - \theta}{2} + \frac{\log u_0}{N(\theta' - \theta)} \right) \not\approx 0 \quad (44)$$

$$\iff \frac{1}{b} \sum_{i=1}^b x_i - 0.5 \not\approx 0. \quad (45)$$

Since the  $x_i$  have mean 0.5, the resulting test with our chosen  $u_0$  will never correctly succeed and must use all  $N$  data points. Furthermore, if we sample values of  $u$  near enough to  $u_0$ , the terms in parenthesis will not be sufficiently different from 0.5 to allow the test to succeed.

The choices above for  $\theta$  and  $\theta'$  guarantee that

$$\log u_0 \in -[0.5, 1][0.75, 1.5] = [-1.5, -0.375]. \quad (46)$$

Next, consider the range of  $u$  values near  $u_0$ :

$$\log u \in \log u_0 + [-0.5, 0.375]. \quad (47)$$

The size of the range in  $u$  is at least  $\exp([-2, -1.125]) \approx [0.13534, 0.32465]$  and occurs with probability at least  $p_2 = 0.18932$ . With  $u$  in this range, we rewrite the test as:

$$\frac{1}{b} \sum_{i=1}^b x_i - 0.5 \not\approx \frac{\log u / u_0}{N(\theta' - \theta)} \quad (48)$$

so that, as in Equation 45, the LHS has expected value zero. Given our choice of intervals for the variables, we can compute the range for the right hand side (RHS) assuming<sup>5</sup> that  $\theta' - \theta > 0$ :

$$\min\{\text{RHS}\} = \frac{-0.5}{\sqrt{N} \cdot 0.5} = -\frac{1}{\sqrt{N}} \quad \text{and} \quad \max\{\text{RHS}\} = \frac{0.375}{\sqrt{N} \cdot 0.5} = \frac{0.75}{\sqrt{N}} \quad (49)$$

Thus, the RHS is in  $\frac{1}{\sqrt{N}}[-1, 0.75]$ . The standard deviation of the LHS given the interval constraints is at least  $0.5/\sqrt{b}$ . Consequently, the gap between the LHS and RHS in Equation 48 is at most  $2\sqrt{b/N}$  standard deviations, limiting the range in which the test will be able to “succeed” without requiring more samples.

The samples  $\theta$ ,  $\theta'$  and  $u$  are drawn independently and so the probability of the conjunction of these events is  $c = p_0 p_1 p_2 = 0.0085$ .

## B Gaussian Mixture Experiment Details

In this section, we provide more details on Section 6.1. Our parameter is a 2-D vector  $\theta = \langle \theta_1, \theta_2 \rangle$ , where

$$\theta_1 \sim \mathcal{N}(0, \sigma_1^2) \quad \text{and} \quad \theta_2 \sim \mathcal{N}(0, \sigma_2^2) \quad (50)$$

<sup>5</sup>If  $\theta' - \theta < 0$ , then the range would be  $\frac{1}{\sqrt{N}}[-0.75, 1]$  but this does not matter for the purposes of our analysis.

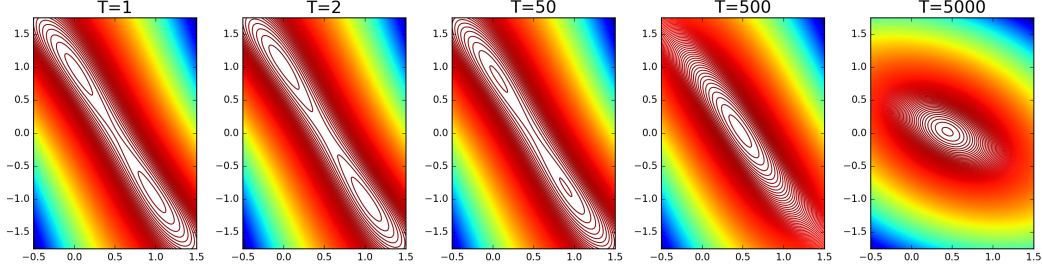


Figure 6: The posterior distribution, with  $N = 10000$  and temperature  $T$  varying.

and where  $\mathcal{N}$  indicates the normal distribution (more generally, the multivariate normal). We consider the above as our prior. Following [12], we set  $\sigma_1^2 = 10$  and  $\sigma_2^2 = 1$ , so the covariance matrix of  $\theta$  is  $\Sigma = \text{diag}(10, 1)$ . Therefore, the log prior probability we endow on  $\theta$  is

$$\log p(\theta) = \log \left( \frac{1}{2\pi\sqrt{10}} \right) - \frac{1}{2}\theta^T \Sigma^{-1} \theta. \quad (51)$$

To generate the data, we use the following Gaussian mixture with tied means:

$$x_i \sim \frac{1}{2}\mathcal{N}(\theta_1, \sigma_x^2) + \frac{1}{2}\mathcal{N}(\theta_1 + \theta_2, \sigma_x^2) \quad (52)$$

where, again following [12], we set  $\sigma_x^2 = 2$ . This means the log likelihood of a single data instance is

$$\log p(x_i | \theta) = \log \frac{1}{4\sqrt{\pi}} + \log \left( \exp \left( -\frac{1}{4}(x_i - \theta_1)^2 \right) + \exp \left( -\frac{1}{4}(x_i - (\theta_1 + \theta_2))^2 \right) \right) \quad (53)$$

The official problem statement is as follows: given some number of conditionally independent data points  $x_1, x_2, \dots, x_N$  generated according to (52), determine the posterior distribution of  $\theta$ :

$$\log p(\theta | x_1, \dots, x_N) = \log p(\theta) + \sum_{i=1}^N \log p(x_i | \theta). \quad (54)$$

Alternatively, depending on the application, we may opt to instead pick a point estimate<sup>6</sup> of  $\theta$ , such as the MAP estimate. (If  $N$  is extremely large, it will cause the posterior to peak sharply at its modes, reducing distribution estimates to point estimates.)

One technique we use to make our test more applicable is to add a temperature to the posterior. In general, we will add  $T > 1$  so that the posterior is  $p(\theta)(\prod_{i=1}^N p(x_i | \theta))^{1/T}$ , resulting in the log posterior of

$$\log p(\theta | x_1, \dots, x_N) \approx \log p(\theta) + \frac{1}{T} \sum_{i=1}^N \log p(x_i | \theta). \quad (55)$$

which has the extra  $1/T$  to decrease the scale factor, so the variance scales as  $1/T^2$ .

To gain intuition on the shape of the posterior, Figure 6 shows simulated contour plots of the posterior based on  $N = 10000$  data points, and temperature varying. The (non-tempered) posterior is a multimodal distribution with modes at  $(0, 1)$  and  $(1, -1)$ . A larger  $T$  implies a flatter posterior, which can be useful for parameter exploration so that the samples do not get trapped in a mode.

<sup>6</sup>We do this for the logistic regression experiment, described in Section 6.2.