
A Simple Minibatch Acceptance Test for MCMC

Anonymous Author(s)

Affiliation
Address
email

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. To do this, we describe a novel acceptance test which is applied to minibatches of data. Previous work on minibatch Metropolis-Hastings use adaptive batch sizes, which may use far more than one minibatch of data and therefore yield few samples. Here we describe a simpler test on fixed-sized minibatches. Instead of adaptive batch sizes, our test requires a precondition on the variance of the log proposal probability ratio. Given a measurement of this variance, a minibatch size can be chosen to support the test in one step, which yields a number of samples proportional to the batch/minibatch ratio. Alternatively, the variance condition can be used to adjust the temperature of the sampling distribution to a “natural” value to provide regular samples at a given minibatch size. The resulting proposal/test is only slightly more complex than a simple SGD update. In this paper we derive the test, discuss its implementation, and present several experiments.

1 Introduction

Markov chain Monte Carlo (MCMC) sampling is a powerful method for computation on intractable distributions. MCMC algorithms propose samples from a proposal distribution q which is in general different from the target distribution p , and decide whether to accept or reject them based on an acceptance test. The acceptance test is usually a Metropolis-Hastings test [1, 2]. For example, for Bayesian posterior inference, the posterior distribution $p(\theta | x_1, \dots, x_N) \propto p(\theta) \prod_{i=1}^N p(x_i | \theta)$ of parameter θ based on conditionally independent data $\{x_i\}_{i=1}^N$ is the target distribution which one aims to approximate by sampling θ values. When N is large, it is very expensive to compute the exact posterior because the likelihood must be evaluated on every x_i .

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 therefore have high sample efficiency. They also map very well onto hardware such as GPUs. In contrast, MCMC methods typically require calculations over the full dataset to produce a new sample. A recent result from [3] attempted to bridge the gap between minibatch optimization and MCMC by proposing an *adaptive minibatch* method. Their method uses a subset of the N data points for the MH test each iteration. However the number of points is not known a priori, and is adjusted dynamically during every test. This undermines many of the practical advantages of minibatch calculations: GPU kernels work best for certain minibatch sizes, memory blocks can be reused when minibatch size is fixed, etc.

In this paper, we develop an alternative minibatch acceptance test for MCMC methods which uses only one fixed-size minibatch per test. Instead of dynamically adjusting the batch size, we add a prerequisite for the test based on the variance of the log acceptance probability. If this variance is

38 small enough, the test requires only a single step. This variance changes slowly on typical problems,
39 and we can estimate it before performing the test. The minibatch variance decreases as the minibatch
40 size grows, so the variance measurement can be used to set the minimum minibatch size.

41 There is another advantage when using our test for Bayesian posterior inference. As the dataset size
42 grows, the posterior distribution sharpens and often concentrates to a few sharp peaks. Samples
43 from this distribution will come from those peaks, and will easily become stuck in one of them.
44 The samples are not exploring the posterior parameter space, and are unlikely to concentrate at a
45 strong local optimum. They therefore fail to accurately represent a sample of the posterior because of
46 inadequate search (i.e., mixing). We argue that it is much more natural to sample the distribution at a
47 higher temperature determined by the minibatch size rather than the dataset size. Starting at a high
48 temperature (small minibatch size) which mixes easily, one can reduce the temperature (increasing
49 minibatch size) until one reaches the desired target distribution temperature. Although this approach
50 also involves changing minibatch size, the changes increase minibatch size monotonically, and only a
51 few changes are performed in the course of a sampling session. Annealing can also be used for MAP
52 and ML estimation because changing temperature preserves posterior modes.

53 To be precise, the contributions of this paper are as follows:

- 54 1. We develop a new minibatch acceptance test which satisfies detailed balance.
55 2. We compare performance of our new test and the adaptive sampling method.
56 3. We experiment using the approach for posterior estimation.

57 2 Preliminaries and Related Work

58 In standard MCMC methods [4, 5] for Bayesian inference with parameter θ and conditionally
59 independent data, the goal is to compute the distribution $p(\theta | x_1, \dots, x_N)$. To do so, one generates
60 a chain of (correlated) samples $\theta_1, \dots, \theta_T$ for large T and approximates p using the sample counts.
61 Each iteration t has a current θ_t , and a *proposal distribution* $q(\theta' | \theta_t)$ determines a new candidate θ' .
62 With probability P_a , the sample is accepted, so $\theta_{t+1} = \theta'$. Otherwise, $\theta_{t+1} = \theta_t$. This is done by
63 drawing a uniform random variable $u \sim \text{Unif}[0, 1]$ and accepting if $u < P_a$. Traditionally, P_a is

$$P_a = \min \left\{ 1, \frac{f(\theta')q(\theta_t | \theta')}{f(\theta_t)q(\theta' | \theta_t)} \right\} = \min \left\{ 1, \frac{p(\theta') \prod_{i=1}^N p(x_i | \theta') q(\theta_t | \theta')}{p(\theta_t) \prod_{i=1}^N p(x_i | \theta_t) q(\theta' | \theta_t)} \right\}, \quad (1)$$

64 where $f(\theta_t) = p(\theta_t | x_1, \dots, x_N)$. This P_a satisfies detailed balance, so if one samples long enough,
65 one will arrive at a stationary distribution matching the posterior, though a burn-in period and/or only
66 using samples at regular intervals is often done in practice.

67 Unfortunately, computing f requires the use of all N training data points. Moreover, it is difficult to
68 design tests using substantially fewer than N points that also satisfy detailed balance. To reconcile
69 these competing objectives, [3] proposes an adaptive minibatch MCMC algorithm which uses a
70 sequential hypothesis test. During each iteration, the algorithm starts with a small minibatch of data
71 and tests the hypothesis that the sample θ' should be accepted or rejected based on a probability ratio.
72 If the test cannot make a decision over a certain confidence threshold, then the minibatch size is
73 increased and the test repeats. This process continues until a decision. The downside of this algorithm
74 is the need to keep incrementing the minibatch size; in the worst case, all N data points may be
75 needed to decide on θ' . A similar approach from [6] has slightly more robust theoretical properties
76 based on concentration inequalities, but comes at the cost of more computational time.

77 The work of [7] presents another minibatch MCMC method which gets the *exact* posterior, so long
78 as there exists a (cheap) lower bound on the likelihood for each data point. Their approach is only
79 a starting point, however, and it can be difficult to derive these lower bounds in general cases. The
80 extra auxiliary variables they introduce also hurt mixing.

81 There is a separate line of MCMC work on simulating the physics of a probability distribution.
82 By viewing random variables as particles in a system, one can apply Hamiltonian Monte Carlo
83 (HMC) [8] methods which generate the best case scenario for proposals: high quality *and* distant.
84 HMC methods require a full gradient computation of the posterior, which means they also face
85 the problem of computing the likelihood at each point. Stochastic versions of HMC [9, 10] use
86 a minibatch of data each iteration but require controlling the minibatch noise to avoid divergence.

87 Langevin Dynamics [11, 12] is a similar approach which does not use the momentum terms in
 88 HMC. Minibatch SGHMC and Langevin Dynamics are orthogonal to our work and are other ways of
 89 attempting to solve the same problem of MCMC methods on big data. We can combine our methods
 90 with these results, as we demonstrate in Section 5.3.

91 3 A New Metropolis-Hastings Test

92 For our new MH test, we use two key values, Δ and Δ' :

$$\Delta = \log \left(\frac{p(\theta') \prod_{i=1}^N p(x_i | \theta') q(\theta_t | \theta')}{p(\theta_t) \prod_{i=1}^N p(x_i | \theta_t) q(\theta' | \theta_t)} \right); \quad \Delta' = \log \left(\frac{p(\theta') (\prod_{i=1}^n p(x_i | \theta'))^{\frac{N}{n}} q(\theta_t | \theta')}{p(\theta_t) (\prod_{i=1}^n p(x_i | \theta_t))^{\frac{N}{n}} q(\theta' | \theta_t)} \right) \quad (2)$$

93 where p and q have similar definitions as in Equation 1. Thus, Δ and Δ' are probability ratios in
 94 log space. Note that the only difference between Δ and Δ' is that they use N and n data terms x_i ,
 95 respectively, and that the latter uses an N/n scaling term. We assume $n \ll N$ (so Δ' is substantially
 96 faster to compute) and that the n data points x_i represent a random minibatch without replacement.

97 3.1 The Full Data Test

98 To motivate the use of Δ , we turn to Lemma 1.

99 **Lemma 1.** *An acceptance function g such that $g(\Delta) = \exp(\Delta)g(-\Delta)$ satisfies detailed balance.*
 100 *That is, $f(\theta_t)p(\theta' | \theta_t) = f(\theta')p(\theta_t | \theta')$, where $p(\theta_y | \theta_x)$ is the probability of jumping from θ_x to*
 101 *θ_y in our chain.*

102 *Proof.* We begin by deriving $p(\theta' | \theta_t)$. This is equivalent to the probability of proposing θ' and then
 103 accepting it, so $p(\theta' | \theta_t) = q(\theta' | \theta_t)g(\Delta)$. Similarly, $p(\theta_t | \theta') = q(\theta_t | \theta')g(-\Delta)$. Notice that the
 104 probability of accepting a transition from θ' to θ_t is $g(-\Delta)$ because this inverts the fraction inside the
 105 logarithm term of Δ . By assumption, we can expand $g(\Delta) = \exp(\Delta)g(-\Delta)$ in $p(\theta' | \theta_t)$. Doing
 106 this, and combining the result with the definition of $p(\theta_t | \theta')$, we get

$$g(-\Delta) = \frac{p(\theta' | \theta_t)}{q(\theta' | \theta_t) \exp(\Delta)} = \frac{p(\theta_t | \theta')}{q(\theta_t | \theta')} \quad (3)$$

107 Rearranging terms and expanding $\exp(\Delta)$, we have

$$\frac{p(\theta' | \theta_t)f(\theta_t)q(\theta' | \theta_t)}{q(\theta' | \theta_t)f(\theta')q(\theta_t | \theta')} = \frac{p(\theta_t | \theta')}{q(\theta_t | \theta')} \quad (4)$$

108 Cancellations result in $f(\theta')p(\theta_t | \theta') = f(\theta_t)p(\theta' | \theta_t)$. Thus, detailed balance is satisfied. \square

109 As a sanity check, the standard Metropolis-Hastings acceptance function $g(\Delta) = \min\{1, e^\Delta\} =$
 110 $\min\left\{1, \frac{f(\theta')q(\theta_t | \theta')}{f(\theta_t)q(\theta' | \theta_t)}\right\}$ satisfies the condition $g(\Delta) = \exp(\Delta)g(-\Delta)$.

111 For our MH test, the key is that we use a different g , the logistic function: $g(\Delta) = (1 + \exp(-\Delta))^{-1}$.
 112 Straightforward arithmetic shows that it satisfies the condition in Lemma 1. The logistic function
 113 is nice because we can easily sample from it using the following procedure. At any iteration with
 114 a current parameter θ_t and a candidate sample θ' , we can compute Δ . Let u be a uniform random
 115 variable $u \sim \text{Unif}[0, 1]$. We accept θ' if $g(\Delta) > u$, and reject otherwise. This process is equivalent to
 116 sampling a random variable X with cumulative distribution function $F_X(x) = g(x)$, and accepting if
 117 $\Delta > X$ and rejecting otherwise. We can define $X = g^{-1}(u)$ so that its CDF is the logistic function,
 118 which follows because for arbitrary $X = x$, we have

$$F_X(x) = \Pr(X \leq x) = \Pr(g^{-1}(u) \leq x) = \Pr(u \leq g(x)) = \int_0^{g(x)} 1 dx = g(x),$$

119 as the density of u is one. Thus, the criteria to accept the candidate θ' is equivalent to whether $\Delta > X$.
 120 Moreover, the density of X (a standard logistic random variable) is symmetric about zero, so the
 121 acceptance criteria can also be expressed as $\Delta + X > 0$.

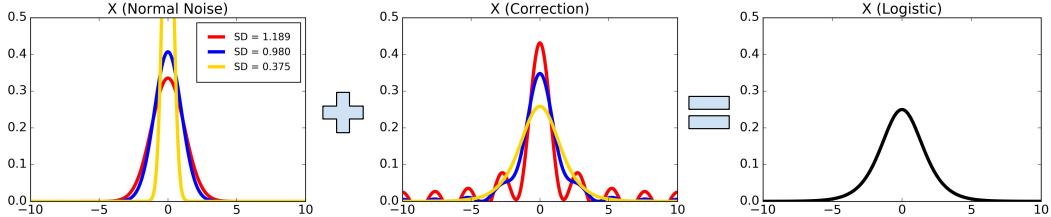


Figure 1: Three examples of X_{norm} and X_{corr} distributions that convolve to form the standard logistic distribution. We use three standard deviation values of X_{norm} . The two red (respectively, blue, green) curves form the logistic. The y -axis is capped at 0.5. This figure must be viewed in color.

122 3.2 The Minibatch Version

123 With large datasets, Δ is intractable to compute. An intuitive fix is to replace Δ with the minibatch
124 version, Δ' , for our MH test. The following lemma characterizes the distribution of Δ' .

125 **Lemma 2.** *If the minibatch data is chosen randomly without replacement and has sufficiently many
126 elements, then the distribution of Δ' is approximately Gaussian.*

127 *Proof.* By expanding the definition of Δ' from Equation 2, we get

$$\Delta' = \log p(\theta') - \log p(\theta_t) + \log q(\theta_t \mid \theta') - \log q(\theta' \mid \theta_t) + \frac{N}{n} \sum_{i=1}^n (\log p(x_i \mid \theta') - \log p(x_i \mid \theta_t)).$$

128 During a given iteration, θ_t and θ' are fixed. Therefore, the randomness in Δ' comes only from
129 the minibatch. Since the minibatch is chosen randomly, the distribution of $\sum_{i=1}^n (\log p(x_i \mid \theta') -$
130 $\log p(x_i \mid \theta_t))$ converges to a Gaussian for sufficiently large N by the Central Limit Theorem. \square

131 From Lemma 2, since Δ' is a noisy approximation of Δ , the relationship is expressed as

$$\Delta' = \Delta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2), \quad (5)$$

132 for a Gaussian noise term ε . The noise means we cannot just replace $\Delta + X$ with $\Delta' + X$ for our
133 test, so we need to slightly change our acceptance criteria. Our insight is to decompose X as

$$X = X_{\text{norm}} + X_{\text{corr}}, \quad (6)$$

134 where we assume X_{norm} is a zero-mean Gaussian and X_{corr} is a zero-mean ‘‘correction’’ term. These
135 two add (i.e., their distributions convolve) to form X . The criteria to accept is now

$$\Delta + X = \Delta + X_{\text{norm}} + X_{\text{corr}} \approx \Delta' + X_{\text{corr}} > 0. \quad (7)$$

136 The approximation exists because X_{norm} is an estimate of ε . If X_{norm} is exact, then we get the same
137 MH test with Δ' as we do with the full data version Δ .

138 It is important to understand why we use Equations 6 and 7. The deconvolution process allows us
139 to use $\Delta' + X_{\text{corr}}$ as a proxy for $\Delta + X$. It works because if we know the distribution¹ of ε (from
140 Equation 5), then we can pretend that our Δ' is in fact Δ . To ‘‘insert’’ ε into our test, we deconvolve
141 X so that one of its components, X_{norm} , plays the part of ε . In light of our revised MH test, the
142 previous discussion raises the question: in an iteration of MCMC with current θ_t and proposed θ' ,
143 how do we know the distribution of X_{corr} ? We need the correct distribution to sample the *values* of
144 X_{corr} appropriately. This procedure consists of two major steps.

- 145 1. We first need to estimate the standard deviation of X_{norm} . We generate K values of Δ' ,
146 each using a different random minibatch of n data points, but each using the same θ_t and θ' .
147 In our experiments (see Section 5) we make K small, around 5-10, for speed purposes, and
148 n is usually 50-200. Using these K values, we compute the empirical standard deviation
149 $\text{std}(\Delta')$.

¹Computing the distributions of ε and X_{norm} reduces to finding their standard deviations.

Input : Number of samples T , number of Δ' estimates K , minibatch size m , pre-computed X_{corr} distributions, $\{x_i\}_{i=1}^N$ values from a target distribution, and initial sample θ_0 .

Output : A chain of T samples $\{\theta_1, \dots, \theta_T\}$ for either posterior inference or MAP/ML estimation.

```

for  $t = \{1, \dots, T\}$  do
    Propose a candidate  $\theta'$  and create an empty list for  $\Delta'$  estimates,  $d = []$ ;
    for  $k = \{1, \dots, K\}$  do
        | Sample a random minibatch of size  $m$  from the complete data;
        | Compute  $\Delta'$  from this minibatch (see Equation 2), and add  $\Delta'$  to  $d$ ;
    end
    Compute  $\text{std}(d)$ , the standard deviation estimate of  $\Delta'$  from list  $d$ ;
    if  $\text{std}(d) \geq 1.2$  then
        | Skip loop (set  $\theta_{t+1} = \theta_t$ ) and apply variance preconditioning fix (if desired);
    else
        | Choose the closest  $X_{\text{corr}}$  distribution, and sample a value  $X_c$  from it;
        | Using another random minibatch of  $m$  data points, compute a new  $\Delta'$  (call it  $\Delta'_{\text{real}}$ );
        | if  $X_c + \Delta'_{\text{real}} > 0$  then
            | | Accept the candidate,  $\theta_{t+1} = \theta'$ ;
        | else
            | | Reject and re-use the old one,  $\theta_{t+1} = \theta_t$ ;
        | end
    end
end
```

Algorithm 1: A description of our MH test within the MCMC algorithm.

150 2. Once we have $\text{std}(\Delta')$, we can determine the distribution of X_{corr} from standard deconvolution
151 techniques. These rely on the well-known fact that the Fourier transform of a
152 convolution is the product of the individual Fourier transforms. In practice, we pre-compute
153 distributions of X_{corr} for different $\text{std}(\Delta')$ values to facilitate the sampling process.

154 Algorithm 1 describes our MH test within the MCMC algorithm.

155 3.3 Discussions and Practical Considerations of the New MH Test

156 **The Convolution.** Figure 1 demonstrates Equation 6 and provides three (color-coded) examples of
157 densities for X_{norm} and X_{corr} . The density of X , a logistic random variable with mean $\mu = 0$ and
158 scale $s = 1$, is known, as are the X_{norm} densities for different standard deviations. The deconvolution
159 provides us with the X_{corr} distributions. Note, however, that as the standard deviation of X_{norm}
160 increases, the X_{corr} distribution becomes increasingly unstable and “bumpier”, because the logistic
161 curve has fatter tails than Gaussians. In addition, we cap the standard deviation possibilities of X_{norm}
162 at ≈ 1.2 in our experiments. In our experiments, we discretize the possible standard deviations of
163 X_{norm} , and we also limit the densities to consider the range $[-10, 10]$ as shown in Figure 1.

164 **Variance Preconditioning.** Returning to the discussion from Section 1, our test requires a variance
165 check. Equation 6 means that the variance of X , which is $\pi^2/3 \approx 3.29$ is an upper bound on the
166 variances of X_{norm} . Therefore, the test cannot work if the variance of X_{norm} is too large. The
167 standard deviation of X is about $\sqrt{3.29} \approx 1.81$, so this bounds $\text{std}(X_{\text{norm}})$ (as well as $\text{std}(\Delta')$).
168 From the discussion above, we bound the standard deviation at 1.2 so that X_{corr} can handle the
169 remaining noise. Recall from Section 3.2 that we estimate the $\text{std}(X_{\text{norm}})$ each iteration. If the
170 variance/standard deviation test fails, then we can (a) skip the iteration, (b) increase the minibatch
171 size, (c) change the proposal to decrease step sizes, or (d) increase temperature, as we now discuss.

172 **Temperature.** One way to satisfy the variance precondition is to increase the temperature of the
173 target distribution. For a temperature $T > 1$, the augmented target distribution would become

$$\log p(\theta | x_1, \dots, x_N) \approx \log p(\theta) + \frac{1}{T} \frac{N}{n} \sum_{i=1}^n \log p(x_i | \theta). \quad (8)$$

174 with an extra $1/T$ term augmented. As the temperature increases, the values of Δ' get smaller in
175 absolute value, thus reducing variance. The flatter posterior results in easier mixing.

176 **4 Theoretical Results**

177 In this section, we explore the convergence properties of our MH test. Due to space constraints, all
 178 our proofs are in Appendix A.

179 We first introduce some notation. Let our true and approximated acceptance probabilities be $P_a =$
 180 $g(\Delta) = \Pr(\Delta + X > 0)$ and $P'_a = \Pr(\Delta' + X_{\text{corr}} > 0)$. Define $X_\xi = \varepsilon - X_{\text{norm}}$ which means,
 181 invoking Equations 5 and 6, that X_ξ is the random variable characterizing the accuracy of our X_{norm} .
 182 Thus, $\Delta' = \Delta + X_{\text{norm}} + X_\xi$. We assume the standard deviation of X_{norm} is less than $\pi/\sqrt{3}$, so
 183 that we can use Equation 6. We denote $F_X(x)$ as the CDF of X .

184 We use some terms from standard Markov chain theory [13]. Let the *transition kernels* of MCMC
 185 on complete and minibatch data at step i be P_i and P'_i , respectively. When these kernels (e.g., P_i)
 186 are applied on a probability distribution D_1 , they generate a new distribution D_2 , which we write as
 187 $P_i \circ D_1 = D_2$. Define π to be the stationary distribution obtained by P_i , and π'_i to be the stationary
 188 distribution obtained by P'_i (π' needs a subscript in this case). Finally, we indicate the *total variation
 189 distance* between two kernels using $\|\cdot\|$.

190 Our first result shows that the distance $\|P_i \circ D - P'_i \circ D\| \leq 4\zeta\ell$ for all distributions D .

191 **Lemma 3.** *If $|X_\xi| < \zeta$ and $|\nabla F_X(x)| < \ell$ for all x , then $\|P_i \circ D - P'_i \circ D\| \leq 4\zeta\ell \quad \forall D$.*

192 Thus, the accuracy of our transition kernel P'_i is related to the quality of our estimated X_{norm} . We
 193 next show that $\text{Var}(\Delta')$ is roughly proportional to the step size of our proposer.

194 **Lemma 4.** *For one step sampling from θ_t to θ' , if the jumping step $\|\theta_t - \theta'\|_2 < \epsilon$, and the gradient
 195 of the log likelihood is bounded by a constant factor, i.e. $\|\nabla(\log \Pr(x_i | \theta))\|_2 < k$, the variance of
 196 Δ' is bounded by $4\epsilon^2 k^2(m - \frac{m(m-1)}{N-1})$, where m is the minibatch size, and N is the total data size.*

197 Since $\text{Var}(\Delta') < 4\epsilon^2 k^2(m - \frac{m(m-1)}{N-1})$, we can assume the maximum value of X_ξ is proportional to
 198 the standard deviation of Δ' , i.e. ϵ . Thus, we can assume $\zeta = \epsilon C$, where C is a constant factor. We
 199 finally introduce Theorem 1 to characterize the convergence of our samples.

200 **Theorem 1.** *Assume π satisfies $\|P_i \circ D_0 - \pi\| \leq \eta \|D_0 - \pi\|$, where $\eta \in [0, 1]$ for all distributions
 201 D_0 . Also assume there exists $0 < \rho_t < 1$ such that $\|P'_i \circ D - \pi'_i\| \leq 2\rho_i \quad \forall D$. We can get
 202 $\|P'_t \circ P'_{t-1} \circ \dots \circ P'_1 \circ D_0 - \pi\| \leq \sum_{s=1}^t \{\prod_{u=s+1}^t \rho_u(1 - \alpha_u)\} \rho_s \alpha_s + \alpha_t$, where $\alpha_i = \frac{\epsilon_i C}{1-\eta}$ and
 203 ϵ_i is the step size at iteration i .*

204 Theorem 1 indicates that the difference between the target and sampled distributions consists of two
 205 components: ρ_i , which is determined by the efficiency of P'_i , and the α_i s, which are determined by
 206 the error of our transition kernels. These involve competing objectives: the larger the step size, the
 207 more efficient the P'_i but its error increases. The reverse happens with a step size too small.

208 **5 Experiments**

209 We conduct three sets of experiments to explore the benefits of our minibatch MH test and to
 210 benchmark it with previous work. In Section 5.1, we show that our test enables samples to converge
 211 to the posterior distribution of a heated Gaussian mixture model. In Section 5.2, we analyze its
 212 efficiency on logistic regression. In Section 5.3, we apply the test for deep learning. Appendices B, C,
 213 and D contain more detailed information on these respective experiments.

214 **5.1 Mixture of Gaussians**

215 We start with a simple Gaussian mixture model, borrowing an experiment from [11]. The parameter
 216 is 2-D, $\theta = (\theta_1, \theta_2)$, and the parameter/data generation process is

$$(\theta_1, \theta_2) \sim \mathcal{N}((0, 0), \text{diag}(\sigma_1^2, \sigma_2^2)); \quad x_i \sim 0.5 \cdot \mathcal{N}(\theta_1, \sigma_x^2) + 0.5 \cdot \mathcal{N}(\theta_1 + \theta_2, \sigma_x^2). \quad (9)$$

217 We set $\sigma_1^2 = 10$, $\sigma_2^2 = 1$ and $\sigma_x^2 = 2$. Fixing $\theta = (0, 1)$, we draw 10000 data points so that the target
 218 distribution is $p(\theta) \prod_{i=1}^{10000} p(x_i | \theta)$, with the prior based on the θ generation process in Equation 9.
 219 This results in rather sharp posterior modes and high $\text{std}(\Delta')$ estimates, so we apply a temperature
 220 $T = 110$ to reduce $\text{std}(\Delta')$. Taking logs, we get the target as shown in the far left of Figure 2.

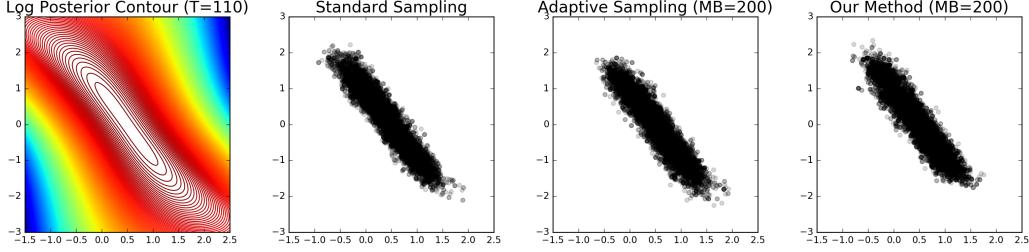


Figure 2: The log posterior contours (temperature 110) and three scatter plots of sampled θ values.

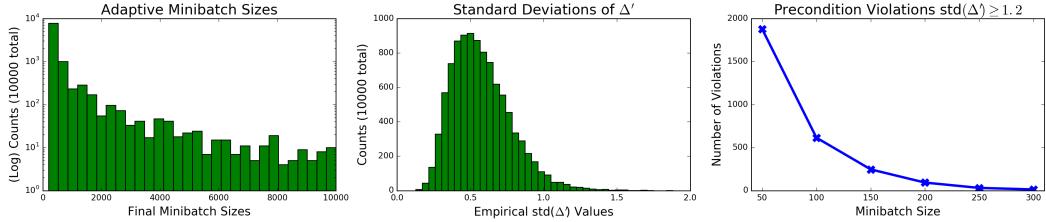


Figure 3: Histograms of minibatch sizes and $\text{std}(\Delta')$ values, along with “ $\text{std}(\Delta')$ violations” (right).

221 We run MCMC with our MH test using minibatch size $m = 200$. We also run this using standard full-
 222 batch MCMC (“standard sampling”) with the standard test from Equation 1, and the method from [3]
 223 (“adaptive sampling”). For the latter, we also use $m = 200$ and increment minibatches by that amount
 224 within an iteration if necessary. The tolerance for making a decision is $\epsilon = 0.1$. To make comparisons
 225 easier, all three use the same random walk proposer with covariance $\Sigma = \text{diag}(0.03, 0.03)$. This is a
 226 poor proposer, but it is sufficient for our purposes as the quality of the samples will be primarily due
 227 to the MH test. All methods are run 10000 times to collect 10000 samples.

228 Figure 2 shows scatter plots of the resulting θ samples for the three methods, with darker regions
 229 indicating a greater density of points. All three methods obtain the same rough form of the posterior,
 230 so our MH test can indeed result in the same posterior as the other two methods.

231 Figure 3 (to the left) shows a histogram of the final minibatch sizes used by the adaptive subsampling
 232 method each iteration (on a log scale for readability). Usually, the method can make a decision
 233 with 200 samples. Occasionally, however, it must use all 10000 points, resulting in lots of wasted
 234 computation. For this particular experiment, that happened eight times. In contrast, our method keeps
 235 the minibatch size fixed at 200, but requires that $\text{std}(\Delta') < 1.2$. Figure 3 (in the middle) shows a
 236 histogram of the estimated $\text{std}(\Delta')$ values each iteration. Only 75 iterations (out of 10000) resulted in
 237 $\text{std}(\Delta') \geq 1.2$ and for these, we skipped the rest of the iteration (setting $\theta_{t+1} = \theta_t$). To smooth the
 238 $\text{std}(\Delta')$ values, we use moving average updates. The third plot in Figure 3 investigates the number
 239 of times $\text{std}(\Delta') \geq 1.2$ based on minibatch sizes. For six sizes (50, 100, 150, 200, 250, and 300),
 240 we ran MCMC with our MH test five times and averaged the number of “violations.” As the size
 241 increases, the number of violations decreases, as expected.

242 5.2 Logistic Regression

243 We next use logistic regression for the binary classification of 1s versus 7s in the MNIST dataset [14].
 244 The data has 12007 and 1000 training and testing points, respectively (we used a random subset of
 245 the test data). The proposer is again a random walk with covariance matrix $0.1I$ for the 784×784
 246 identity matrix I . We set the posterior temperature at $T = 3000$. We set the minibatch size $m = 50$
 247 and compare with adaptive sampling with tolerance 0.02.

248 Figure 4 shows the prediction accuracy and log likelihood on the test set as a function of the
 249 cumulative training data points processed. Our test increments the cumulative data by a fixed amount
 250 per iteration, but the adaptive method may require more data per iteration. We see that our minibatch
 251 MH test is more efficient; it has similar or better performance while consuming fewer data points.

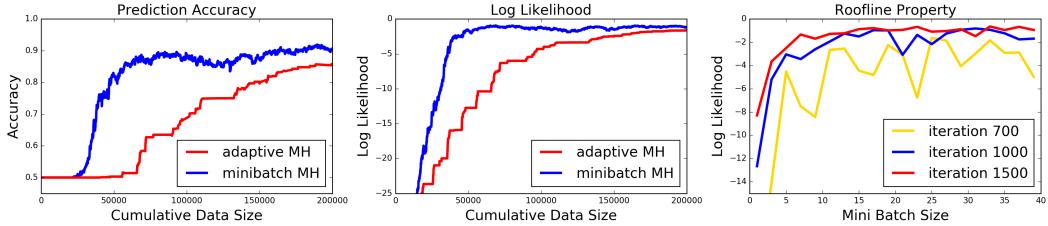


Figure 4: Logistic regression performance (accuracy/log likelihood) and minibatch size analysis.

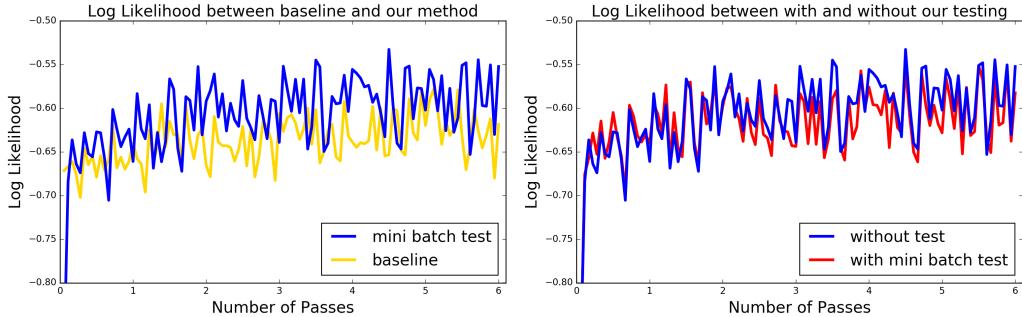


Figure 5: Log likelihood versus number of passes for optimization of neural network.

252 The third plot of Figure 4 shows performance based on minibatch size. The parameter is evaluated on
 253 the test set after training iterations 700, 1000, and 1500. For those three cases, when the minibatch
 254 size increases beyond about 10, the performance of the sampled data will not improve with a large
 255 minibatch size. This makes sense because the error of our MH test is proportional to our estimation
 256 error of $\text{std}(\Delta')$, and as the minibatch increases, $\text{std}(\Delta')$ decreases. Beyond a certain standard
 257 deviation (or minibatch) threshold, the error in our minibatch MH test is negligible.

258 5.3 Neural Network Optimization

259 In this experiment, we apply our minibatch MH test with a Stochastic Gradient Hamiltonian Monte
 260 Carlo (SGHMC) [9] proposer to sample the complete parameters for a fully connected neural network.
 261 We use the Higgs data set from the UCI dataset [15], which is a binary classification task with one
 262 million instances, each of which are 28-dimensional. The network we use has four layers: the first is
 263 the input, the last is the softmax, and the intermediate ones apply the sigmoid activation unit. We
 264 also employ dropout and batch normalization [16]. In order to control $\text{std}(\Delta')$, we initialize the
 265 temperature at 1000, and adjust it at iteration i according to $T_i = \max\{1, 1000/(i + 1)^{0.5}\}$.

266 For our comparison baseline, we train the same network architecture using the Adaptive Gradient
 267 Descent method [17]. Since the frequencies of the features are quite different, we use the same
 268 strategy as in Adaptive Gradient to rescale the gradient before applying SGHMC. We implement our
 269 MH test and SGHMC in the open-source BIDMach project [18].

270 Figure 5 illustrates our experiment results and reveals that SGHMC with our minibatch MH test
 271 achieves higher log likelihood than the baseline. By comparing SGHMC with and without our
 272 minibatch MH test, we see that our MH test has limited benefit. Because we decrease the step size
 273 quickly, the acceptance rate of SGHMC is already high, so there is little need to test.

274 6 Conclusions

275 In this paper, we have derived a new MH test for minibatch MCMC methods. We demonstrated
 276 how a simple deconvolution process allows us to use a minibatch approximation to the full data
 277 tests. We experimentally show the benefits of our test on Gaussian mixtures, logistic regression, and
 278 deep learning. Straightforward directions for future work include running more experiments with
 279 a particular focus on investigation of the variance precondition. More elaborate extensions include
 280 combining our results with Hamiltonian Monte Carlo methods, providing a recipe for how to use our
 281 algorithm (following the framework of [19]), or integrating parallel MCMC [20, 21] concepts.

282 **References**

- 283 [1] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, “Equation of state
284 calculations by fast computing machines,” *The Journal of Chemical Physics*, vol. 21, 1953.
- 285 [2] W. K. Hastings, “Monte carlo sampling methods using markov chains and their applications,” *Biometrika*,
286 vol. 57, pp. 97–109, 1970.
- 287 [3] A. K. Balan, Y. Chen, and M. Welling, “Austerity in MCMC land: Cutting the metropolis-hastings budget,”
288 in *Proceedings of the 31st International Conference on Machine Learning (ICML)*, 2014.
- 289 [4] W. Gilks and D. Spiegelhalter, *Markov chain Monte Carlo in practice*. Chapman & Hall/CRC, 1996.
- 290 [5] S. Brooks, A. Gelman, G. Jones, and X.-L. Meng, *Handbook of Markov Chain Monte Carlo*. CRC press,
291 2011.
- 292 [6] R. Bardenet, A. Doucet, and C. Holmes, “Towards scaling up markov chain monte carlo: an adaptive
293 subsampling approach,” in *Proceedings of the 31st International Conference on Machine Learning (ICML)*,
294 2014.
- 295 [7] D. Maclaurin and R. P. Adams, “Firefly monte carlo: Exact MCMC with subsets of data,” in *Proceedings
296 of the Thirtieth Conference on Uncertainty in Artificial Intelligence, (UAI)*, 2014.
- 297 [8] R. M. Neal, “MCMC using Hamiltonian dynamics,” *Handbook of Markov Chain Monte Carlo*, vol. 54,
298 pp. 113–162, 2010.
- 299 [9] T. Chen, E. Fox, and C. Guestrin, “Stochastic gradient Hamiltonian Monte Carlo,” in *Proceedings of the
300 31st International Conference on Machine Learning (ICML)*, 2014.
- 301 [10] N. Ding, Y. Fang, R. Babbush, C. Chen, R. D. Skeel, and H. Neven, “Bayesian sampling using stochastic
302 gradient thermostats,” in *Advances in Neural Information Processing Systems 27*, 2014.
- 303 [11] M. Welling and Y. W. Teh, “Bayesian learning via stochastic gradient langevin dynamics,” in *Proceedings
304 of the 28th International Conference on Machine Learning (ICML)*, 2011.
- 305 [12] S. Ahn, A. K. Balan, and M. Welling, “Bayesian posterior sampling via stochastic gradient fisher scoring.,”
306 in *Proceedings of the 29th International Conference on Machine Learning (ICML)*, 2012.
- 307 [13] S. Meyn and R. L. Tweedie, *Markov Chains and Stochastic Stability*. Cambridge University Press, 2009.
- 308 [14] Y. LeCun and C. Cortes, “MNIST handwritten digit database,”
- 309 [15] M. Lichman, “UCI machine learning repository,” 2013.
- 310 [16] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal
311 covariate shift,” in *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, 2015.
- 312 [17] P. Bartleet, E. Hazan, and A. Rakhlin, “Adaptive online gradient descent,” in *Advances in Neural Informa-
313 tion Processing Systems 20*, 2007.
- 314 [18] J. Canny and H. Zhao, “Bidmach: Large-scale learning with zero memory allocation,” *BIGLearn NIPS
315 Workshop*, 2013.
- 316 [19] Y. Ma, T. Chen, and E. Fox, “A complete recipe for stochastic gradient mcmc,” in *Advances in Neural
317 Information Processing Systems 28*, 2015.
- 318 [20] E. Angelino, E. Kohler, A. Waterland, M. Seltzer, and R. P. Adams, “Accelerating MCMC via parallel
319 predictive prefetching,” in *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence,
320 (UAI)*, 2014.
- 321 [21] S. Ahn, B. Shahbaba, and M. Welling, “Distributed stochastic gradient MCMC,” in *Proceedings of the 31st
322 International Conference on Machine Learning, (ICML)*, 2014.
- 323 [22] Y. Yang and D. B. Dunson, “Sequential markov chain monte carlo,” *arXiv preprint arXiv:1308.3861*, 2013.

Outline of Appendix

325 In this appendix, we describe the following topics, with an emphasis on clarity and understanding:

- 326 • Proofs for material in Section 4.
- 327 • More information on Section 5.1.
- 328 • More information on Section 5.2.
- 329 • More information on Section 5.3.

330 **A Proofs**

331 **A.1 Proof of Lemma 3**

332 We first prove that the difference in our acceptance rates is bounded: $|P_a - P'_a| \leq 2\zeta\ell$ for some small
 333 constant $\zeta > 0$. Given $X_\xi = \xi$ the distance $|P_a - P'_a|$ is:

$$\begin{aligned} |P_a - P'_a| &= |\Pr(\Delta + X > 0) - \Pr(\Delta' + X_{\text{corr}} > 0)| \\ &= |\Pr(\Delta + X > 0) - \Pr(\Delta + X_{\text{norm}} + \xi + X_{\text{corr}} > 0)| \\ &= |\Pr(\Delta + X > 0) - \Pr(\Delta + X + \xi > 0)| \\ &= |F_X(-\Delta) - F_X(-\Delta - \xi)| \\ &= |\nabla F(-\Delta)\xi + o(\xi^2)| \\ &\leq |\nabla F(-\Delta)\xi| + |o(\xi^2)| \leq 2|\nabla F(-\Delta)||\xi| \leq 2\ell\xi, \end{aligned}$$

334 where we use Taylor's Theorem ($o(\xi^2)$ represents higher-order terms that have smaller absolute value)
 335 and the Triangle Inequality. Since $|X_\xi| < \zeta$, $|P_a - P'_a| < 2\zeta\ell$.

336 Next, we show the distance of new distributions obtained by applying P_i and P'_i on an arbitrary
 337 distribution D once is bounded if we use the same proposer, i.e. $\|P_i \circ D - P'_i \circ D\| \leq 4\zeta\ell \forall D$. We
 338 write $P_i(\theta_t \mid \theta') = P_a(\theta_t, \theta')q(\theta' \mid \theta_t) + (1 - P_a(\theta_t, \theta'))\delta_D(\theta' - \theta_t)$, where δ_D is the Dirac delta
 339 function. The total variation distance between two generated distributions is

$$\begin{aligned} \|P_i \circ D - P'_i \circ D\| &= \int_{\theta'} \left| \int_{\theta_t} ((P_a - P'_a)q(\theta' \mid \theta_t) + (1 - P_a - 1 + P'_a)\delta_D(\theta_t - \theta'))dD(\theta_t) \right| d\theta' \\ &= \int_{\theta'} \left| \int_{\theta_t} (P_a - P'_a)(q(\theta' \mid \theta_t) - \delta_D(\theta_t - \theta'))dD(\theta_t) \right| d\theta' \\ &\leq 2\zeta\ell \int_{\theta'} \left| \int_{\theta_t} (q(\theta' \mid \theta_t) + \delta_D(\theta_t - \theta'))dD(\theta_t) \right| d\theta' \\ &\leq 2\zeta\ell \int_{\theta'} (D'(\theta') + D(\theta'))d\theta' \\ &= 4\zeta\ell, \end{aligned}$$

340 where D' is the distribution generated if we always accept the proposed point. (In the last line, we
 341 integrate two probability density functions, resulting in a value of two.) This matches our desired
 342 bound.

343 **A.2 Proof of Lemma 4**

344 From Lemma 2, only the $\sum_{i=1}^n (\log p(x_i \mid \theta') - \log p(x_i \mid \theta_t))$ term brings randomness into Δ' . To
 345 simplify the subsequent notation, we define $g_i(\theta) = \log p(x_i \mid \theta)$, and $q_i = g_i(\theta') - g_i(\theta_t)$. Then by
 346 Taylor's Theorem:

$$|q| = |g(\theta') - g(\theta_t)| = |\nabla g(\theta_t)^T(\theta' - \theta_t) + o(\theta' - \theta_t)^2| \leq 2\epsilon k.$$

³⁴⁷ Thus, given θ and θ' , the variance of Δ' can be written as:

$$\begin{aligned}\text{Var}(\Delta') &= \text{Var} \left[\sum_{i=1}^m g_i(\theta') - g_i(\theta_t) \right] = \text{Var} \left(\sum_{i=1}^m q_i \right) \\ &= m\text{Var}(q_i) + m(m-1)\text{Cov}_{i \neq j}(q_i, q_j) \\ &= m\text{Var}(q_i) - \frac{m(m-1)}{N-1}\text{Var}(q_i) \\ &\leq 4 \left(m - \frac{m(m-1)}{N-1} \right) k^2 \epsilon^2,\end{aligned}$$

³⁴⁸ as desired.

³⁴⁹ A.3 Proof of Theorem 1

³⁵⁰ First, by Theorem 1 in [3], $\|P_i \circ D_0 - \pi\| \leq \eta \|D_0 - \pi\|$ and $\|P_i \circ D - P'_i \circ D\| \leq \epsilon_i C$ for all D .
³⁵¹ The norm between the stationary distributions is $\|\pi - \pi'_i\| \leq \frac{\epsilon_i C}{1-\eta}$.

³⁵² Second, by Theorem 3.6 in [22], we have

$$\|P'_t \circ \dots \circ P'_1 \circ D_0 - \pi_t\| \leq \sum_{s=1}^t \left\{ \prod_{u=s+1}^t \rho_u(1-\alpha_u) \right\} \rho_s \alpha_s,$$

³⁵³ where $\alpha = \frac{\epsilon_i C}{1-\eta}$. Thus, we get:

$$\begin{aligned}\|P'_t \circ \dots \circ P'_1 \circ D_0 - \pi_0\| &\leq \|P'_t \circ \dots \circ P'_1 \circ D_0 - \pi_t\| + \|\pi_t - \pi_0\| \\ &\leq \sum_{s=1}^t \left\{ \prod_{u=s+1}^t \rho_u(1-\alpha_u) \right\} \rho_s \alpha_s + \alpha_t,\end{aligned}$$

³⁵⁴ as desired.

³⁵⁵ B Gaussian Mixture Experiment Details

³⁵⁶ In this section, we go over the math details on the Gaussian mixture model problem borrowed
³⁵⁷ from [11]. Our parameter is a 2-D vector $\theta = (\theta_1, \theta_2)$, where

$$\theta_1 \sim \mathcal{N}(0, \sigma_1^2) \quad \text{and} \quad \theta_2 \sim \mathcal{N}(0, \sigma_2^2) \quad (10)$$

³⁵⁸ where \mathcal{N} indicates the normal distribution (more generally, the multivariate normal). We consider
³⁵⁹ the above as our prior. Following [11], we set $\sigma_1^2 = 10$ and $\sigma_2^2 = 1$, so the covariance matrix of θ is
³⁶⁰ $\Sigma = \text{diag}(10, 1)$. Therefore, the log prior probability we endow on θ is

$$\log p(\theta) = \log \left(\frac{1}{2\pi\sqrt{10}} \right) - \frac{1}{2} \theta^T \Sigma^{-1} \theta. \quad (11)$$

³⁶¹ To generate the data, we use the following Gaussian mixture with tied means:

$$x_i \sim \frac{1}{2} \mathcal{N}(\theta_1, \sigma_x^2) + \frac{1}{2} \mathcal{N}(\theta_1 + \theta_2, \sigma_x^2) \quad (12)$$

³⁶² where, again following [11], we set $\sigma_x^2 = 2$. This means the log likelihood of a single data instance is

$$\log p(x_i | \theta) = \log \left(\frac{1}{4\sqrt{\pi}} \right) + \log \left(\exp \left(-\frac{1}{4}(x_i - \theta_1)^2 \right) + \exp \left(-\frac{1}{4}(x_i - (\theta_1 + \theta_2))^2 \right) \right) \quad (13)$$

³⁶³ Here is the problem statement: given some number of conditionally independent data points
³⁶⁴ x_1, x_2, \dots, x_N generated according to (12), determine the posterior distribution of θ :

$$\log p(\theta | x_1, \dots, x_N) = \log p(\theta) + \sum_{i=1}^N \log p(x_i | \theta). \quad (14)$$

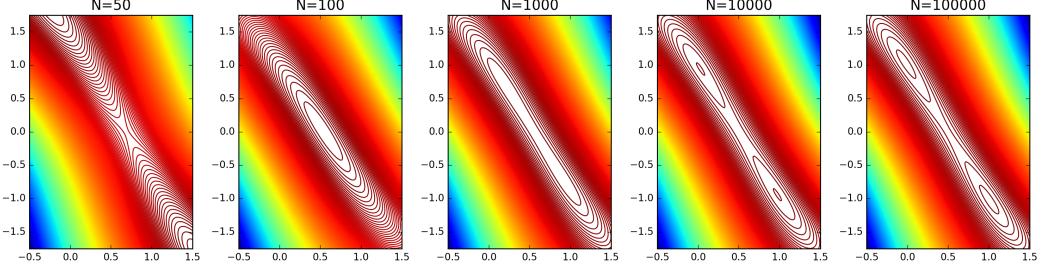


Figure 6: The posterior distribution, from 50 to 100k samples, with temperature set at 1.

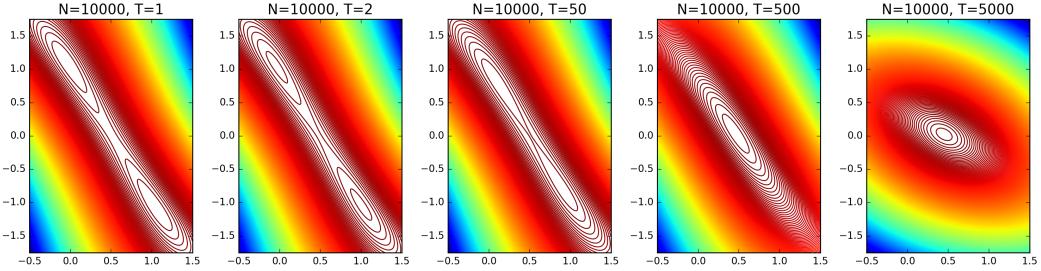


Figure 7: The posterior distribution, with $N = 10000$ but with temperature T varying.

365 Alternatively, if there are too many data points, we may opt to instead pick a point estimate of θ ,
 366 generally the MAP estimate. (If N is extremely large, it will cause the posterior to peak sharply at its
 367 modes, reducing distribution estimates to point estimates.) Note that in many cases, we will need to
 368 take a *minibatch estimate* of (14). In that case, the literature generally uses

$$\log p(\theta | x_1, \dots, x_N) \approx \log p(\theta) + \frac{N}{n} \sum_{i=1}^n \log p(x_i | \theta). \quad (15)$$

369 where we only use $n \ll N$ samples, but we must scale up the likelihood contribution by N/n . If we
 370 didn't add this scaling factor, then the contribution of the likelihood terms would be weaker.

371 One technique we use is adding a *temperature* to our distribution. In general, we will want to add
 372 $T > 1$ so that our posterior is $p(\theta)((\prod_{i=1}^n p(x_i | \theta))^{N/n})^{1/T}$, resulting in the log posterior of

$$\log p(\theta | x_1, \dots, x_N) \approx \log p(\theta) + \frac{1}{T} \frac{N}{n} \sum_{i=1}^n \log p(x_i | \theta). \quad (16)$$

373 which has the extra $1/T$ to decrease the scale factor. Equation (16) is what we use for our experiments,
 374 because warmer distributions help us satisfy our $\text{std}(\Delta') < 1.2$ requirement.

375 To gain some intuition on what the posterior looks like, Figure 6 shows simulated contour plots of
 376 the posterior based on varying numbers of data points N , with the temperature set at $T = 1$. Note
 377 that because we are using all N points here, the scale factor $N/n = 1$. As N increases, the posterior
 378 converges to a multimodal distribution with modes at $(0, 1)$ and $(1, -1)$. Figure 7 is similar, except
 379 this time we fix the number of samples at $N = 10000$, but show how changing the temperature T
 380 affects the distribution. A larger T implies a flatter posterior, one that (weakly) peaks in between the
 381 two true modes.

382 C Logistic Regression Experiment Details

383 In this section, we go over some details of our logistic regression experiment. The feature vector for
 384 an image consists of its pixel values, normalized between 0 and 1. For simplicity, we only consider
 385 the binary classification case, so we only use digits one (denoted as output $y = 1$) and seven (denoted

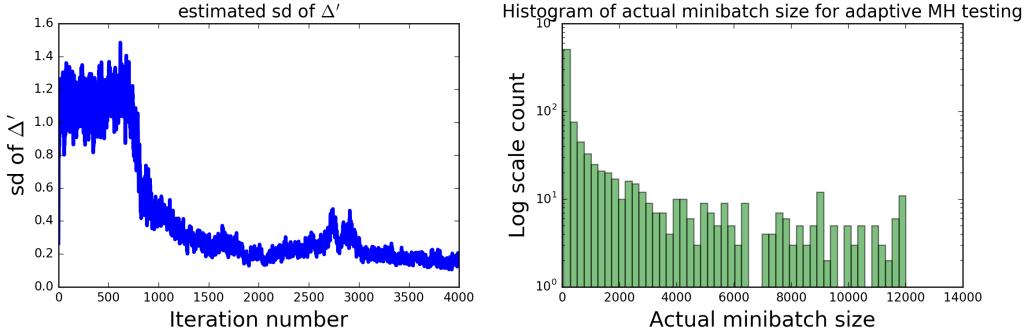


Figure 8: Additional results for the logistic regression experiment.

386 as $y = -1$). The probability of the i^{th} output $y_i \in \{-1, +1\}$ with the input vector x_i is

$$p(y_i | x_i) = \frac{1}{1 + \exp(-y_i \theta^T x_i)}, \quad (17)$$

387 where θ is the 784-length parameter vector.

388 For our experiment, we impose a uniform prior to represent our lack of knowledge about θ . We use
389 a random walk proposer, which can be modeled as $\theta' = \theta_i + \mathcal{N}(0, \sigma^2 I)$, where θ_i is the current
390 sample, θ' is the proposed sample, and we choose the variance to be a constant $\sigma^2 = 0.01$ for all
391 components. We initialize θ_0 to be a vector of all ones, and set our minibatch size as $m = 50$.

392 For our minibatch MH test, in order to enforce the $\text{std}(\Delta') < 1.2$ condition, we use a constant
393 temperature $T = 3000$. If our estimated $\text{std}(\Delta') \geq 1.2$, we ignore the current iteration. Figure 8
394 plots our estimated $\text{std}(\Delta')$ values versus iteration count.

395 For adaptive MH testing, our experimental settings are the same as with our MH test, except we do
396 not impose a temperature. The minibatch size of adaptive MH testing is also initialized as 50, but it
397 may increase by that amount each iteration. Figure 8 shows the histogram of the actual minibatch
398 size at the end of each iteration in adaptive MH testing.

399 D Neural Network Experiment Details

400 For our neural network experiment, we use the architecture discussed in Section 5.3. We use the
401 SGHMC as the proposer for our minibatch MH testing. For the baseline, we use a tuned adaptive
402 gradient descent optimizer, whose step size changes by $0.01/(i+1)^{0.4}$. Both methods have minibatch
403 sizes set at $m = 200$. There are one million total data instances x_i .

404 For SGHMC, we use the simplified update equations [9]:

$$\Delta\theta = v, \quad \Delta v = -\eta \nabla U(\theta) - \alpha v + \mathcal{N}(0, 2(\alpha - \hat{\beta})\eta), \quad (18)$$

405 where v represents auxiliary momentum variables, and $\nabla U(\theta)$ is the gradient of the system. We
406 set the hyperparameters to be $\eta_i = 0.01/(i+1)^{0.4}$, and $\alpha = 0.1$. We use the empirical Fisher [12]
407 information $V(\theta)$ to estimate the value of $\hat{\beta}$, so that $\hat{\beta} = \frac{1}{2}\eta V(\theta)$. In order to control $\text{std}(\Delta')$, we
408 initialize the temperature at 1000, and adjust it at iteration i according to $T_i = \max\{1, 1000/(i+1)^{0.5}\}$.
409