

# Measuring Sample Quality with Stein’s Method

**Authored by:**

Jackson Gorham  
Lester Mackey

## **Abstract**

To improve the efficiency of Monte Carlo estimation, practitioners are turning to biased Markov chain Monte Carlo procedures that trade off asymptotic exactness for computational speed. The reasoning is sound: a reduction in variance due to more rapid sampling can outweigh the bias introduced. However, the inexactness creates new challenges for sampler and parameter selection, since standard measures of sample quality like effective sample size do not account for asymptotic bias. To address these challenges, we introduce a new computable quality measure based on Stein’s method that bounds the discrepancy between sample and target expectations over a large class of test functions. We use our tool to compare exact, biased, and deterministic sample sequences and illustrate applications to hyperparameter selection, convergence rate assessment, and quantifying bias-variance tradeoffs in posterior inference.

## **1 Paper Body**

When faced with a complex target distribution, one often turns to RMarkov chain Monte Carlo (MCMC) [1] to approximate intractable expectations  $EP[h(Z)] = \int h(x)p(x)dx$  with asymptotically exact sample estimates  $EQ[h(X)] = \frac{1}{n} \sum_{i=1}^n h(x_i)$ . These complex targets commonly arise as posterior distributions in Bayesian inference and as candidate distributions in maximum likelihood estimation [2]. In recent years, researchers [e.g., 3, 4, 5] have introduced asymptotic bias into MCMC procedures to trade off asymptotic correctness for improved sampling speed. The rationale is that more rapid sampling can reduce the variance of a Monte Carlo estimate and hence outweigh the bias introduced. However, the added flexibility introduces new challenges for sampler and parameter selection, since standard sample quality measures, like effective sample size, asymptotic variance, trace and mean plots, and pooled and within-chain variance diagnostics, presume eventual convergence to the target [1] and hence do not account for asymptotic bias. To address this shortcoming, we develop a new measure of sample quality suitable for comparing asymptotically exact, asymptotically biased, and even deterministic sample sequences. The quality

measure is based on Stein's method and is attainable by solving a linear program. After outlining our design criteria in Section 2, we relate the convergence of the quality measure to that of standard probability metrics in Section 3, develop a streamlined implementation based on geometric spanners in Section 4, and illustrate applications to hyperparameter selection, convergence rate assessment, and the quantification of bias-variance tradeoffs in posterior inference in Section 5. We discuss related work in Section 6 and defer all proofs to the appendix. Notation We denote the  $\ell_2$ ,  $\ell_1$ , and  $\ell_\infty$  norms on  $\mathbb{R}^d$  by  $\|\cdot\|_2$ ,  $\|\cdot\|_1$ , and  $\|\cdot\|_\infty$  respectively. We will often refer to a generic norm  $\|\cdot\|$  on  $\mathbb{R}^d$  with associated dual norms  $\|\cdot\|_*$ ,  $\sup_{\|w\|=1} \langle v, w \rangle$  for vectors  $w \in \mathbb{R}^d$ ,  $\|M\|_*$ ,  $\sup_{\|v\|=1} \langle Mv, v \rangle$  for matrices  $M \in \mathbb{R}^{d \times d}$ , and  $\|T\|_*$ ,  $\sup_{\|v\|=1} \langle T[v], v \rangle$  for tensors  $T \in \mathbb{R}^{d \times d \times d}$ . We denote the  $j$ -th standard basis vector by  $e_j$ , the  $\partial$  partial derivative  $\partial x_k$  by  $r_k$ , and the gradient of any  $\mathbb{R}^d$ -valued function  $g$  by  $\nabla g$  with components  $(\nabla g(x))_j = r_j g(x)$ .

## 2

### Quality Measures for Samples

Consider a target distribution  $P$  with open convex support  $X \subseteq \mathbb{R}^d$  and continuously differentiable density  $p$ . We assume that  $p$  is known up to its normalizing constant and that exact integration under  $P$  is intractable for most functions of interest. We will approximate expectations under  $P$  with the aid of a weighted sample, a collection of distinct sample points  $x_1, \dots, x_n \in X$  with weights  $q(x_i)$  encoded in a probability mass function  $q$ . The probability mass function  $q$  induces a discrete distribution  $Q$  and an approximation  $\mathbb{E}_Q[h(X)] = \sum_{i=1}^n q(x_i) h(x_i)$  for any target expectation  $\mathbb{E}_P[h(Z)]$ . We make no assumption about the provenance of the sample points; they may arise as random draws from a Markov chain or even be deterministically selected. Our goal is to compare the fidelity of different samples approximating a common target distribution. That is, we seek to quantify the discrepancy between  $\mathbb{E}_Q$  and  $\mathbb{E}_P$  in a manner that (i) detects when a sequence of samples is converging to the target, (ii) detects when a sequence of samples is not converging to the target, and (iii) is computationally feasible. A natural starting point is to consider the maximum deviation between sample and target expectations over a class of real-valued test functions  $H$ ,  $d_H(Q, P) = \sup_{h \in H} |\mathbb{E}_Q[h(X)] - \mathbb{E}_P[h(Z)]|$ .

$$(1) \quad \mathbb{E}_P[h(Z)] - \mathbb{E}_Q[h(X)]$$

When the class of test functions is sufficiently large, the convergence of  $d_H(Q_m, P)$  to zero implies that the sequence of sample measures  $(Q_m)_{m=1}^\infty$  converges weakly to  $P$ . In this case, the expression (1) is termed an integral probability metric (IPM) [6]. By varying the class of test functions  $H$ , we can recover many well-known probability metrics as IPMs, including the total variation distance, generated by  $H = \{h : X \rightarrow \mathbb{R} \mid \sup_{x \in X} |h(x)| \leq 1\}$ , and the Wasserstein distance (also known as the Kantorovich-Rubinstein or earth mover's distance),  $d_{Wk}^k$ , generated by  $H = \{h : X \rightarrow \mathbb{R} \mid \sup_{x \in X} |h(x) - h(y)| \leq k \|x - y\|\}$ .

The primary impediment to adopting an IPM as a sample quality measure is

that exact computation is typically infeasible when generic integration under  $P$  is intractable. However, we could skirt this intractability by focusing on classes of test functions with known expectation under  $P$ . For example, if we consider only test functions  $h$  for which  $EP[h(Z)] = 0$ , then the IPM value  $dH(Q, P)$  is the solution of an optimization problem depending on  $Q$  alone. This, at a high level, is our strategy, but many questions remain. How do we select the class of test functions  $h$ ? How do we know that the resulting IPM will track convergence and non-convergence of a sample sequence (Desiderata (i) and (ii))? How do we solve the resulting optimization problem in practice (Desideratum (iii))? To address the first two of these questions, we draw upon tools from Charles Stein's method of characterizing distributional convergence. We return to the third question in Section 4.

3

Stein's Method

Stein's method [7] for characterizing convergence in distribution classically proceeds in three steps: 1. Identify a real-valued operator  $T$  acting on a set  $G$  of  $\mathbb{R}^d$ -valued functions of  $X$  for which  $EP[(Tg)(Z)] = 0$  for all  $g \in G$ .

(2)

Together,  $T$  and  $G$  define the Stein discrepancy,

$$S(Q, T, G) = \sup_{g \in G} |EQ[(Tg)(X)]| = \sup_{g \in G} |EQ[(Tg)(X)] - EP[(Tg)(Z)]|$$

$$EP[(Tg)(Z)] = dT G(Q, P),$$

an IPM-type quality measure with no explicit integration under  $P$ . 2. Lower bound the Stein discrepancy by a familiar convergence-determining IPM  $dH$ . This step can be performed once, in advance, for large classes of target distributions and ensures that, for any sequence of probability measures  $(\mu_m)_{m=1}^\infty$ ,  $S(\mu_m, T, G)$  converges to zero only if  $dH(\mu_m, P)$  does (Desideratum (ii)). 3. One commonly considers real-valued functions  $g$  when applying Stein's method, but we will find it more convenient to work with vector-valued  $g$ .

2

3. Upper bound the Stein discrepancy by any means necessary to demonstrate convergence to zero under suitable conditions (Desideratum (i)). In our case, the universal bound established in Section 3.3 will suffice. While Stein's method is typically employed as an analytical tool, we view the Stein discrepancy as a promising candidate for a practical sample quality measure. Indeed, in Section 4, we will adopt an optimization perspective and develop efficient procedures to compute the Stein discrepancy for any sample measure  $Q$  and appropriate choices of  $T$  and  $G$ . First, we assess the convergence properties of an equivalent Stein discrepancy in the subsections to follow. 3.1

Identifying a Stein Operator

The generator method of Barbour [8] provides a convenient and general means of constructing operators  $T$  which produce mean-zero functions under  $P$  (2). Let  $(Z_t)_{t \geq 0}$  represent a Markov process with unique stationary distribution  $P$ . Then the infinitesimal generator  $A$  of  $(Z_t)_{t \geq 0}$ , defined by  $(Au)(x) = \lim_{t \downarrow 0} (E[u(Z_t) - u(Z_0) | Z_0 = x]) / t$

for

$$u : \mathbb{R}^d \rightarrow \mathbb{R},$$

satisfies  $\mathbb{E}P[(Au)(Z)] = 0$  under mild conditions on  $A$  and  $u$ . Hence, a candidate operator  $T$  can be constructed from any infinitesimal generator. For example, the overdamped Langevin diffusion, defined by the stochastic differential equation  $dZ_t = \frac{1}{2} \nabla \log p(Z_t) dt + dW_t$  for  $(W_t)_{t \geq 0}$  a Wiener process, gives rise to the generator  $\frac{1}{2} \Delta (APu)(x) = \frac{1}{2} \nabla u(x) \cdot \nabla \log p(x) + \frac{1}{2} \Delta u(x)$ . (3) After substituting  $g$  for  $\frac{1}{2} \nabla u$ , we obtain the associated Stein operator

$(TPg)(x) = \frac{1}{2} \nabla g(x) \cdot \nabla \log p(x) + \frac{1}{2} \Delta g(x)$ . (4) The Stein operator  $TP$  is particularly well-suited to our setting as it depends on  $P$  only through the derivative of its log density and hence is computable even when the normalizing constant of  $p$  is not. If we let  $\partial X$  denote the boundary of  $X$  (an empty set when  $X = \mathbb{R}^d$ ) and  $n(x)$  represent the outward unit normal vector to the boundary at  $x$ , then we may define the classical Stein set  $\{x \in \mathbb{R}^d : \nabla g(x) \cdot n(x) \leq k, g : X \rightarrow \mathbb{R} \text{ s.t. } \sup_{x \in \partial X} \max\{k(x), |\nabla g(x) \cdot n(x)|\} \leq 1 \text{ and } \int_X \nabla g(x) \cdot n(x) dx = 0\}$  with  $n(x)$  defined of sufficiently smooth functions satisfying a Neumann-type boundary condition. The following proposition is a consequence of integration by parts and shows that  $G_k$  is a suitable domain for  $TP$ . Proposition 1. If  $\mathbb{E}P[\nabla \log p(Z) \cdot k] \leq 1$ , then  $\mathbb{E}P[(TPg)(Z)] = 0$  for all  $g \in G_k$ . Together,  $TP$  and  $G_k$  form the classical Stein discrepancy  $S(Q, TP, G_k)$ , our chief object of study. 3.2

### Lower Bounding the Classical Stein Discrepancy

In the univariate setting ( $d = 1$ ), it is known for a wide variety of targets  $P$  that the classical Stein discrepancy  $S(\mu_m, TP, G_k)$  converges to zero only if the Wasserstein distance  $d_{Wk}(\mu_m, P)$  does [9, 10]. In the multivariate setting, analogous statements are available for multivariate Gaussian targets [11, 12, 13], but few other target distributions have been analyzed. To extend the reach of the multivariate literature, we show in Theorem 2 that the classical Stein discrepancy also determines Wasserstein convergence for a large class of strongly log-concave densities, including the Bayesian logistic regression posterior under Gaussian priors. Theorem 2 (Stein Discrepancy Lower Bound for Strongly Log-concave Densities). If  $X = \mathbb{R}^d$ , and  $\log p$  is strongly concave with third and fourth derivatives bounded and continuous, then, for any probability measures  $(\mu_m)_{m \geq 1}$ ,  $S(\mu_m, TP, G_k) \rightarrow 0$  only if  $d_{Wk}(\mu_m, P) \rightarrow 0$ . We emphasize that the sufficient conditions in Theorem 2 are certainly not necessary for lower bounding the classical Stein discrepancy. We hope that the theorem and its proof will provide a template for lower bounding  $S(Q, TP, G_k)$  for other large classes of multivariate target distributions. 3

### 3.3

### Upper Bounding the Classical Stein Discrepancy

We next establish sufficient conditions for the convergence of the classical Stein discrepancy to zero. Proposition 3 (Stein Discrepancy Upper Bound). If  $X \subseteq \mathbb{R}^d$  and  $Z \sim P$  with  $\nabla \log p(Z)$  integrable,  $S(Q, TP, G_k) \leq \mathbb{E}[kX \cdot Z] + \mathbb{E}[k \nabla \log p(X) \cdot \nabla \log p(Z)] + \mathbb{E}[\nabla \log p(Z) \cdot (X - Z)] + \mathbb{E}[kX \cdot Z] + \mathbb{E}[k \nabla \log p(X) \cdot \nabla \log p(Z)] + \mathbb{E}[k \nabla \log p(Z) \cdot (X - Z)]$ . One implication of Proposition 3 is that  $S(\mu_m, TP, G_k)$  converges to zero whenever  $\mu_m \rightarrow P$ .



$\ell_1$ ,

the family of functions which satisfy the classical constraints and certain implied Taylor compatibility constraints at pairs of points in  $E$ . Remarkably, if the graph  $G_1$  consists of edges between all distinct sample points  $x_i$ , then the associated complete graph Stein discrepancy  $S(Q, \mathcal{T}_P, \|\cdot\|_k, Q, G_1)$  is equivalent to the classical Stein discrepancy in the following strong sense. 4

**Proposition 5 (Equivalence of Classical and Complete Graph Stein Discrepancies).** If  $X = \mathbb{R}^d$ , and  $G_1 = (\text{supp}(Q), E_1)$  with  $E_1 = \{(x_i, x_l) \in \text{supp}(Q)^2 : x_i \neq x_l\}$ , then  $S(Q, \mathcal{T}_P, \|\cdot\|_k) \leq S(Q, \mathcal{T}_P, \|\cdot\|_k, Q, G_1) \leq c_d S(Q, \mathcal{T}_P, \|\cdot\|_k)$ , where  $c_d$  is a constant, independent of  $(Q, P)$ , depending only on the dimension  $d$  and norm  $\|\cdot\|_k$ .

Proposition 5 follows from the Whitney-Glaeser extension theorem for smooth functions [14, 15] and implies that the complete graph Stein discrepancy inherits all of the desirable convergence properties of the classical discrepancy. However, the complete graph also introduces order  $n^2$  constraints, rendering computation infeasible for large samples. To achieve the same form of equivalence while enforcing only  $O(n)$  constraints, we will make use of sparse geometric spanner subgraphs. 4.2

#### Geometric Spanners

For a given dilation factor  $t \geq 1$ , a  $t$ -spanner [16, 17] is a graph  $G = (V, E)$  with weight  $w_{xy}$  on each edge  $(x, y) \in E$  and a path between each pair  $x_0 \neq y_0 \in V$  with total weight no larger than  $t w_{x_0 y_0}$ . The next proposition shows that spanner Stein discrepancies enjoy the same convergence properties as the complete graph Stein discrepancy. **Proposition 6 (Equivalence of Spanner and Complete Graph Stein Discrepancies).** If  $X = \mathbb{R}^d$ ,  $G_t = (\text{supp}(Q), E)$  is a  $t$ -spanner, and  $G_1 = (\text{supp}(Q), \{(x_i, x_l) \in \text{supp}(Q)^2 : x_i \neq x_l\})$ , then  $S(Q, \mathcal{T}_P, \|\cdot\|_k, Q, G_1) \leq S(Q, \mathcal{T}_P, \|\cdot\|_k, Q, G_t) \leq 2t S(Q, \mathcal{T}_P, \|\cdot\|_k, Q, G_1)$ .

Moreover, for any  $p$  norm, a 2-spanner with  $O(dn)$  edges can be computed in  $O(dn \log(n))$  expected time for  $d$  a constant depending only on  $d$  and  $\|\cdot\|_k$  [18]. As a result, we will adopt a 2-spanner Stein discrepancy,  $S(Q, \mathcal{T}_P, \|\cdot\|_k, Q, G_2)$ , as our standard quality measure. 4.3

#### Decoupled Linear Programs

The final unspecified component of our Stein discrepancy is the choice of norm  $\|\cdot\|_k$ . We recommend the  $\ell_1$  norm, as the resulting optimization problem decouples into  $d$  independent finite-dimensional linear programs (LPs) that can be solved in parallel. More precisely,  $S(Q, \mathcal{T}_P, \|\cdot\|_1, Q, (V, E))$  equals  $P_d$

$$\begin{aligned} & \min_{\mathbf{V}} \sum_{j=1}^d \sum_{i=1}^n q(v_i) \left( \sum_{j=1}^d \log p(v_i) + \sum_{j=1}^d \sum_{l=1}^n \mathbb{1}_{(v_i, v_l) \in E} \left( \sum_{j=1}^d (e_{ij} - e_{lj}) \right) \right) \\ & \text{s.t. } \sum_{j=1}^d v_j = 1, \sum_{j=1}^d v_j = 1, \text{ and } \sum_{i=1}^n (v_i, v_l) \in E, \sum_{j=1}^d (e_{ij} - e_{lj}) \leq 1 \\ & \quad \sum_{j=1}^d v_j = 1, \sum_{j=1}^d v_j = 1, \text{ and } \sum_{i=1}^n (v_i, v_l) \in E, \sum_{j=1}^d (e_{ij} - e_{lj}) \leq 1 \end{aligned}$$

l 1

We have arbitrarily numbered the elements  $v_i$  of the vertex set  $V$  so that value  $g_j(v_i)$ , and  $jki$  represents the gradient value  $rk\ g_j(v_i)$ . 4.4

$j_l\ h\ j\ e_l, v_i\ 2\ 1\ 2\ k v_i\ v_l\ k_l$

$j_i$

$v_l\ i$ —

?

? 1.

represents the function

Constrained Domains

A small modification to the unconstrained formulation (7) extends our tractable Stein discrepancy computation to any domain defined by coordinate boundary constraints, that is, to  $X = (\tau_1, 1) \dots (\tau_d, d)$  with  $1 \leq \tau_j \leq 1$  for all  $j$ . Specifically, for each dimension  $j$ , we augment the  $j$ -th coordinate linear program of (7) with the boundary compatibility constraints  $\dots$   $(v_b) = \max_{i,j} -v_{ij} j_i b_j, -v_{ijk} i b_j, j_i (v_{jji} b_j)^2 j \leq 1$ , for each  $i, b_j \in \{j, j\} \subseteq R$ , and  $k \in \{j\}$ . (8) 2

$i_j$

$j$

These additional constraints ensure that our candidate function and gradient values can be extended to a smooth function satisfying the boundary conditions  $hg(z), n(z)i = 0$  on  $\partial X$ . Proposition 15 in the appendix shows that the spanner Stein discrepancy so computed is strongly equivalent to the classical Stein discrepancy on  $X$ . Algorithm 1 summarizes the complete solution for computing our recommended, parameter-free spanner Stein discrepancy in the multivariate setting. Notably, the spanner step is unnecessary in the univariate setting, as the complete graph Stein discrepancy  $S(Q, TP, G_k, Q, G_1)$  can be computed directly by sorting the sample and boundary points and only enforcing constraints between consecutive points in this ordering. Thus, the complete graph Stein discrepancy is our recommended quality measure when  $d = 1$ , and a recipe for its computation is given in Algorithm 2. 5

Algorithm 1 Multivariate Spanner Stein Discrepancy input:  $Q$ , coordinate bounds  $(\tau_1, 1), \dots, (\tau_d, d)$  with  $1 \leq \tau_j \leq 1$  for all  $j$  G2 Compute sparse 2-spanner of  $\text{supp}(Q)$  for  $j = 1$  to  $d$  do (in parallel)  $r_j$  Solve  $j$ -th coordinate linear program (7) with graph  $G_2$  and boundary constraints (8) Pd return  $j=1$   $r_j$

Algorithm 2 Univariate Complete Graph Stein Discrepancy

input:  $Q$ , bounds  $(\tau, 1)$  with  $1 \leq \tau \leq 1$   $(x(1), \dots, x(n_0))$   $SORT(\{x_1, \dots, x_n, \tau, \tau\} \subseteq R)$   $P \subseteq [n_0, d]$  return  $\sup_{i=1}^{2Rn_0} \frac{1}{2Rn_0} \sum_{i=1}^{2Rn_0} q(x(i)) \left( \int dx \log p(x(i)) + i \right) \leq \tau$  s.t.  $k \leq 1, \delta_i \leq n_0, -i \leq \tau \leq 1 \leq x(i) \leq 1$ , and,  $\delta_i \leq n_0, \tau \leq i \leq i+1 \leq x(i) \leq x(i+1) \leq -i \leq i+1 \leq i+1 \leq \max_{x \in (i+1), 1 \leq x(i), x(i+1)} x(i), (x \leq x)^2 (x \leq i+1)$

2

5

(i)

2







```

h = TP g
?
n = 3000
Stein discrepancy
?
Scaled Student's t
Gaussian
? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?
? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?
? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?
n = 300
?
1.0 0.5 0.0 ?0.5 ?1.0 1.0 0.5 0.0 ?0.5 ?1.0 1.0 0.5 0.0 ?0.5 ?1.0
n = 300
?
0.10
Scaled Student's t
Gaussian
?
6 ?6 ?3 0
3 6
x
Figure 1: Left: Complete graph Stein discrepancy for a N (0, 1) target.
Middle / right: Optimal Stein functions g and discriminating test functions h
= TP g recovered by the Stein program. 6
seed = 8
seed = 9 ?
? ?
?
? ?
? ?
?? ? ?
?
?
? ?
? ?
?
?
? ?
0.03 0.01
?
? ? ?
? ? ??? ?
? ???
?
?

```

? ? ?  
 ?? ? ?  
 ?  
 ? ? ?  
 ? ???  
 ?  
 ?  
 ? ?  
 ???  
 ? ?  
 ?  
 Gaussian  
 0.030 0.010  
 ? ?  
 ? ?  
 ?  
 ?  
 ?  
 ??  
 ?  
 ? ?  
 ? ? ? ??  
 0.003 0.001  
 ?  
 ? ? ? ?  
 ?  
 ? ???  
 ?  
 ?  
 ?  
 ?  
 ? ? ? ??? ?  
 ?  
 ?  
 ?  
 Discrepancy ?  
 Classical Stein Wasserstein  
 ?  
 ?  
 ? ?  
 ?  
 ?  
 ?  
 ? ?  
 ? ??? ?  
 Uniform

Discrepancy value  
 seed = 7 0.30 0.10  
 ??  
 ???  
 Complete graph Stein  
 ?  
 ?  
 ?  
 100  
 1000  
 10000  
 100  
 1000  
 10000  
 100  
 1000  
 10000  
 Number of sample points, n

Figure 2: Comparison of discrepancy measures for sample sequences drawn i.i.d. from their targets. which is computable for simple univariate target distributions [22] and provably lower bounds the non-uniform Stein discrepancies (5) with  $c1:3 = (0.5, 0.5, 1)$  for  $P = \text{Unif}(0, 1)$  and  $c1:3 = (1, 4, 2)$  for  $P = N(0, 1)$  [9, 23]. For  $N(0, 1)$  and  $\text{Unif}(0, 1)$  targets and several random number generator seeds, we generate a sequence of sample points i.i.d. from the target distribution and plot the nonuniform classical and complete graph Stein discrepancies and the Wasserstein distance as functions of the first n sample points in Figure 2. Two apparent trends are that the graph Stein discrepancy very closely approximates the classical and that both Stein discrepancies track the fluctuations in Wasserstein distance even when a magnitude separation exists. In the  $\text{Unif}(0, 1)$  case, the Wasserstein distance in fact equals the classical Stein discrepancy because  $TP g = g 0$  is a Lipschitz function. 5.3

#### Selecting Sampler Hyperparameters

Stochastic Gradient Langevin Dynamics (SGLD) [3] with constant step size  $\gamma$  is a biased MCMC procedure designed for scalable inference. It approximates the overdamped Langevin diffusion, but, because no Metropolis-Hastings (MH) correction is used, the stationary distribution of SGLD deviates increasingly from its target as  $\gamma$  grows. If  $\gamma$  is too small, however, SGLD explores the sample space too slowly. Hence, an appropriate choice of  $\gamma$  is critical for accurate posterior inference. To illustrate the value of the Stein diagnostic for this task, we adopt the bimodal Gaussian mixture model (GMM) posterior of [3] as our target. For a range of step sizes  $\gamma$ , we use SGLD with minibatch size 5 to draw 50 independent sequences of length  $n = 1000$ , and we select the value of  $\gamma$  with the highest median quality  $\hat{Q}$  either the maximum effective sample size (ESS, a standard diagnostic based on autocorrelation [1]) or the minimum spanner Stein discrepancy  $\hat{S}$  across these sequences. The average discrepancy computation consumes 0.4s for spanner construction and 1.4s per coordinate

linear program. As seen in Figure 3a, ESS, which does not detect distributional bias, selects the largest step size presented to it, while the Stein discrepancy prefers an intermediate value. The rightmost plot of Figure 3b shows that a representative SGLD sample of size  $n$  using the  $\gamma$  selected by ESS is greatly overdispersed; the leftmost is greatly underdispersed due to slow mixing. The middle sample, with  $\gamma$  selected by the Stein diagnostic, most closely resembles the true posterior.

#### 5.4 Quantifying a Bias-Variance Trade-off

The approximate random walk MH (ARWMH) sampler [5] is a second biased MCMC procedure designed for scalable posterior inference. Its tolerance parameter  $\gamma$  controls the number of datapoint likelihood evaluations used to approximate the standard MH correction step. Qualitatively, a larger  $\gamma$  implies fewer likelihood computations, more rapid sampling, and a more rapid reduction of variance. A smaller  $\gamma$  yields a closer approximation to the MH correction and less bias in the sampler stationary distribution. We will use the Stein discrepancy to explicitly quantify this bias-variance trade-off. We analyze a dataset of 53 prostate cancer patients with six binary predictors and a binary outcome indicating whether cancer has spread to surrounding lymph nodes [24]. Our target is the Bayesian logistic regression posterior [1] under a  $N(0, I)$  prior on the parameters. We run RWMH ( $\gamma = 0$ ) and ARWMH ( $\gamma = 0.1$  and batch size = 2) for 105 likelihood evaluations, discard the points from the first 103 evaluations, and thin the remaining points to sequences of length 1000. The discrepancy computation time for 1000 points averages 1.3s for the spanner and 12s for a coordinate LP. Figure 4 displays the spanner Stein discrepancy applied to the first  $n$  points in each sequence as a function of the likelihood evaluation count. We see that the approximate sample is of higher Stein quality for smaller computational budgets but is eventually overtaken by the asymptotically exact sequence.

```

diagnostic = ESS
Step size,  $\gamma = 5e+05$ 
Step size,  $\gamma = 5e+03$ 
Step size,  $\gamma = 5e+02$ 
? ? ?
2.0
?
4
?
3
1.0
2 1
?
?
?
diagnostic = Spanner Stein
? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?
? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?

```

?  
 ?2 ?3 ?4 1e?04  
 1e?03  
 1e?02  
 ?2 ?1  
 Step size, ?  
 (a) Step size selection criteria  
 0  
 Spanner Stein discrepancy  
 0.3  
 ?  
 0.2 ?  
 16  
 ? ?  
 ??  
 ?? ?? ?? ??  
 0.1  
 3e+03 1e+04 3e+04 1e+05  
 3  
 ?2 ?1  
 0  
 x1  
 1  
 2  
 3  
 ?2 ?1  
 0  
 1  
 2  
 ; Stein discrepancy minimized at ? = 5 ? 10  
 ?  
 3  
 1.0  
 ?  
 2.0 ?  
 1.5  
 ? ? ?  
 ??? ? ? ?? ? ?? ?? ?? ??? ????? ? ?? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?  
 ? ??? ?????  
 0.5  
 3e+03 1e+04 3e+04 1e+05  
 ??  
 ? ?? ??? ? ?? ? ? ? ????? ?? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?  
 ?? ????? ???  
 .  
 ? ?

(b) 1000 SGLD sample points with equidensity contours of  $p$  overlaid

Figure 3: (a) ESS maximized at  $\gamma = 5$   $\gamma = 10$  Discrepancy

[illegible]

?? ? ??? ? ? ?? ? ? ? ??? ? ?? ? ? ? ?? ? ?? ? ??? ?? ? ? ? ? ? ?? ?? ??  
? ? ? ??? ????? ?? ????? ? ?? ? ? ? ? ?? ??? ?? ????? ?? ?? ? ? ? ? ??? ?  
? ?? ? ? ? ? ? ? ? ??? ??  
? ? ??? ? ? ??  
?? ?? ? ? ? ? ? ??  
?  
???? ?? ??  
?  
? ?? ??? ?????  
? ??  
????  
?  
? ? ? ? ? ?

?  
1.5  
x2  
Log median diagnostic  
2.5  
1.0 0.5  
3e+03 1e+04 3e+04 1e+05  
Hyperparameter  
? ? ?  
??  
? ? ??  
?? ????? ?

?  
?=0 ? = 0.1  
3e+03 1e+04 3e+04 1e+05  
Number of likelihood evaluations  
Figure 4: Bias-variance trade-off curves for Bayesian logistic regression with approximate RWMH. To corroborate our result, we use a Metropolis-adjusted Langevin chain [25] of length 107 as a surrogate  $Q$  for the target and compute several error measures for each sample  $Q$ : normalized probability  $\max -E[X Z]$  error  $\max_l -E[(hX, w_l i) (hZ, w_l i)] - /k w_l k_1$ , mean error  $\max_{jj} -EQ?_j [Z_{jj}] -$ , and second moment  $\max -E[X X Z Z] -$   
 $j k j k j, k$  error  $\max$  for  $X ? Q, Z ? Q?$ ,  $(t)$ ,  $1+e1 t$ , and  $w_l$  the  $l$ -th datapoint covariate  $j, k -EQ? [Z_j Z_k] -$  vector. The measures, also found in Figure 4, accord with the Stein discrepancy quantification.

5.5  
Assessing Convergence Rates  
The Stein discrepancy can also be used to assess the quality of deterministic sample sequences. In Figure 5 in the appendix, for  $P = \text{Unif}(0, 1)$ , we plot the complete graph Stein discrepancies of the first  $n$  points of an i.i.d.  $\text{Unif}(0, 1)$  sample, a deterministic Sobol sequence [26], and a deterministic R1 kernel



herding sequence [27] defined by the norm  $\|h\|_H = \int_0^1 (h(x))^2 dx$ . We use the median value over 50 sequences in the i.i.d. case and estimate the convergence rate for each sampler using the slope of the best least squares affine fit to each log-log plot. The discrepancy computation time averages 0.08s for  $n = 200$  points, and the recovered rates of  $n^{-0.49}$  and  $n^{-1}$  for the i.i.d. and Sobol sequences accord with expected  $O(1/n)$  and  $O(\log(n)/n)$  bounds from the literature [28, 26]. As 0.96 witnessed also in other outpaces its best known bound of  $p$  metrics [29], the herding rate of  $n^{-dH(Q_n, P)} = O(1/n)$ , suggesting an opportunity for sharper analysis.

6

#### Discussion of Related Work

We have developed a quality measure suitable for comparing biased, exact, and deterministic sample sequences by exploiting an infinite class of known target functionals. The diagnostics of [30, 31] also account for asymptotic bias but lose discriminating power by considering only a finite collection of functionals. For example, for a  $N(0, 1)$  target, the score statistic of [31] cannot distinguish two samples with equal first and second moments. Maximum mean discrepancy (MMD) on a characteristic Hilbert space [32] takes full distributional bias into account but is only viable when the expected kernel evaluations are easily computed under the target. One can approximate MMD, but this requires access to a separate trustworthy ground-truth sample from the target. Acknowledgments The authors thank Madeleine Udell, Andreas Eberle, and Jessica Hwang for their pointers and feedback and Quirijn Bouts, Kevin Buchin, and Francis Bach for sharing their code and counsel. 8

## 2 References

- [1] S. Brooks, A. Gelman, G. Jones, and X.-L. Meng. Handbook of Markov chain monte carlo. CRC press, 2011.
- [2] C. J. Geyer. Markov chain monte carlo maximum likelihood. Computer Science and Statistics: Proc. 23rd Symp. Interface, pages 156–163, 1991.
- [3] M. Welling and Y.-W. Teh. Bayesian learning via stochastic gradient Langevin dynamics. In Proceedings of the 28th International Conference on Machine Learning, pages 681–688, 2011.
- [4] S. Ahn, A. Korattikara, and M. Welling. Bayesian posterior sampling via stochastic gradient fisher scoring. In Proceeding of 29th International Conference on Machine Learning (ICML’12), 2012.
- [5] A. Korattikara, Y. Chen, and M. Welling. Austerity in MCMC land: Cutting the Metropolis-Hastings budget. In Proceeding of 31th International Conference on Machine Learning (ICML’14), 2014.
- [6] A. Müller. Integral probability metrics and their generating classes of functions. Advances in Applied Probability, 29(2):pp. 429–443, 1997.
- [7] C. Stein. A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 2: Probability Theory, pages 583–602, Berkeley, CA, 1972. University of California Press.
- [8] A. D. Barbour. Stein’s method and Poisson process convergence. J. Appl. Probab.,

(Special Vol. 25A): 175–184, 1988. A celebration of applied probability. [9] L. HY. Chen, L. Goldstein, and Q.-M. Shao. Normal approximation by Steins method. Springer Science & Business Media, 2010. [10] S. Chatterjee and Q.-M. Shao. Nonnormal approximation by Steins method of exchangeable pairs with application to the Curie-Weiss model. *Annals of Applied Probability*, 21(2):464–483, 2011. [11] G. Reinert and A. Röllin. Multivariate normal approximation with Steins method of exchangeable pairs under a general linearity condition. *Annals of Probability*, 37(6):2150–2173, 2009. [12] S. Chatterjee and E. Meckes. Multivariate normal approximation using exchange-able pairs. *Alea*, 4: 257–283, 2008. [13] E. Meckes. On Steins method for multivariate normal approximation. In *High dimensional probability V: The Luminy volume*, pages 153–178. Institute of Mathematical Statistics, 2009. [14] G. Glaeser. Etude de quelques algèbres tayloriennes. *J. Analyse Math.*, 6:1–124; erratum, insert to 6 (1958), no. 2, 1958. [15] P. Shvartsman. The Whitney extension problem and Lipschitz selections of set-valued mappings in jetspaces. *Transactions of the American Mathematical Society*, 360(10):5529–5550, 2008. [16] P. Chew. There is a planar graph almost as good as the complete graph. In *Proceedings of the Second Annual Symposium on Computational Geometry, SCG ’86*, pages 169–177, New York, NY, 1986. ACM. [17] D. Peleg and A. A. Schaffer. Graph spanners. *Journal of Graph Theory*, 13(1):99–116, 1989. [18] S. Har-Peled and M. Mendel. Fast construction of nets in low-dimensional metrics and their applications. *SIAM Journal on Computing*, 35(5):1148–1184, 2006. [19] Q. W. Bouts, A. P. ten Brink, and K. Buchin. A framework for computing the greedy spanner. In *Proceedings of the Thirtieth Annual Symposium on Computational Geometry, SOCG’14*, pages 11:11–11:19, New York, NY, 2014. ACM. [20] M. Lubin and I. Dunning. Computing in operations research using Julia. *INFORMS Journal on Computing*, 27(2):238–248, 2015. [21] Gurobi Optimization. Gurobi optimizer reference manual, 2015. URL <http://www.gurobi.com>. [22] S. S. Vallender. Calculation of the Wasserstein distance between probability distributions on the line. *Theory of Probability & Its Applications*, 18(4):784–786, 1974. [23] C. Döbler. Stein’s method of exchangeable pairs for the Beta distribution and generalizations. arXiv:1411.4477, 2014. [24] A. Canty and B. D. Ripley. *boot: Bootstrap R (S-Plus) Functions*, 2015. R package version 1.3-15. [25] G. O. Roberts and R. L. Tweedie. Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli*, pages 341–363, 1996. [26] R. E. Caflisch. Monte carlo and quasi-monte carlo methods. *Acta numerica*, 7:1–49, 1998. [27] Y. Chen, M. Welling, and A. Smola. Super-samples from kernel herding. In *Proceeding of 26th Uncertainty in Artificial Intelligence (UAI’10)*, 2010. [28] E. del Barrio, E. Gin, and C. Matrn. Central limit theorems for the Wasserstein distance between the empirical and the true distributions. *Ann. Probab.*, 27(2):1009–1071, 04 1999. [29] F. Bach, S. Lacoste-Julien, and G. Obozinski. On the equivalence between herding and conditional gradient algorithms. In *Proceeding of 29th International Conference on Machine Learning (ICML’12)*, 2012. [30] A. Zellner and C.-K. Min. Gibbs sampler convergence criteria. *Journal of the American Statistical Association*, 90(431):921–927, 1995. [31] Y. Fan, S. P. Brooks, and A. Gelman.

Output assessment for monte carlo simulations via the score statistic. *Journal of Computational and Graphical Statistics*, 15(1), 2006. [32] A. Gretton, K. M. Borgwardt, M. Rasch, B. Schölkopf, and A. J. Smola. A kernel method for the twosample-problem. In *Advances in Neural Information Processing Systems*, pages 513–520, 2006.