# Supplement

This supplement contains supporting material for the paper *Stein $\Pi$-Importance Sampling*. The mathematical background on Stein kernels is contained in Appendix A. The proof of Theorem 1 is contained in Appendix B. For implementation of Stein $\Pi$-Importance Sampling without the aid of automatic differentiation, various explicit derivatives are required; the relevant calculations can be found in Appendix C. The empirical protocols and additional empirical results are presented in Appendix D.

## A  Mathematical Background

This appendix contains mathematical background on reproducing kernels and Stein kernels, as used in the main text. Appendix A.1 introduces matrix-valued reproducing kernels, while Appendix A.2 specialises to Stein kernels by application of a Stein operator to a matrix-valued kernel. A selection of useful Stein kernels are presented in Appendix A.3.

### A.1  Matrix-Valued Reproducing Kernels

A *matrix-valued kernel* is a function $K : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}^{d \times d}$, that is both

1. symmetric; $K(x, y) = K(y, x)$ for all $x, y \in \mathbb{R}^d$, and

2. positive semi-definite; $\sum_{i=1}^{n} \sum_{j=1}^{n} \langle c_i, K(x_i, x_j) c_j \rangle \geq 0$ for all $x_1, \dots, x_n \in \mathbb{R}^d$ and $c_1, \dots, c_n \in \mathbb{R}^d$.

Let $K_x = K(\cdot, x)$. For vector-valued functions $g, g' : \mathbb{R}^d \to \mathbb{R}^d$, defined by $g = \sum_{i=1}^{n} K_{x_i} c_i$ and $g' = \sum_{j=1}^{m} K_{x'_j} c'_i$, define an inner product

$$\langle g, g' \rangle_{\mathcal{H}(K)} = \sum_{i=1}^{n} \sum_{j=1}^{m} \langle c_i, K(x_i, x'_j) c'_j \rangle. \tag{9}$$

There is a unique Hilbert space of such vector-valued functions associated to $K$, denoted $\mathcal{H}(K)$; see Proposition 2.1 of Carmeli et al. (2006). This space is characterised as

$$\mathcal{H}(K) = \overline{\text{span}}\{K_x c : x, c \in \mathbb{R}^d\}$$

where here the closure is taken with respect to the inner product in (9). It can be shown that $\mathcal{H}(K)$ is in fact a reproducing kernel Hilbert space (RKHS) which satisfies the *reproducing property*

$$\langle g, K_x c \rangle_{\mathcal{H}(K)} = \langle g(x), c \rangle$$

for all $g \in \mathcal{H}(K)$ and $x, c \in \mathbb{R}^d$. Matrix-valued kernels are the natural starting point for construction of KSDs, as described next.

### A.2  Stein Kernels

A general construction for Stein kernels is to first identify a matrix-valued RKHS $\mathcal{H}(K)$ and an operator $S_P : \mathcal{H}(K) \to L^1(P)$ for which $\int S_p h \, dP = 0$ for all $h \in \mathcal{H}(K)$. Such an operator will be called a *Stein operator*. The collection $\{S_p h : h \in \mathcal{H}(K)\}$ inherits the structure of an RKHS, whose reproducing kernel

$$k_P(x, y) = \langle S_P K_x, S_P K_y \rangle_{\mathcal{H}(K)} \tag{10}$$

is a Stein kernel, meaning that $\mu_P = 0$ where $\mu_P$ is the kernel mean embedding from (1); see Barp et al. (2022b). Explicit calculations for the Stein kernels considered in this work can be found in Appendix C.

For univariate distributions, Barbour (1988) proposed to obtain Stein operators from infinitessimal generators of $P$-invariant continuous-time Markov processes; see also Barbour (1990); Gotze (1991).

13

The approach was extended to multivariate distributions in Gorham and Mackey (2015). The starting point is the $P$-invariant Itô diffusion

$$\mathrm{d}X_t = \frac{1}{2}\frac{1}{p(X_t)}\nabla \cdot [p(X_t)M(X_t)]\mathrm{d}t + M(X_t)^{1/2}\mathrm{d}W_t, \tag{11}$$

where $p$ is the density of $P$, assumed to be positive, $M : \mathbb{R}^d \to \mathbb{R}^{d\times d}$ is a symmetric matrix called the *diffusion matrix*, and $W_t$ is a standard Wiener process (Kent, 1978; Roberts and Stramer, 2002). Here the notation $[\nabla \cdot A]_i = \nabla \cdot (A_{i,\cdot}^\top)$ indicates the divergence operator applied to each row of the matrix $A(x) \in \mathbb{R}^{d\times d}$. The infinitessimal generator is

$$(A_P u)(x) = \frac{1}{2}\frac{1}{p(x)}\nabla \cdot [p(x)M(x)\nabla u(x)].$$

Substituting $h(x)$ for $\frac{1}{2}\nabla u(x)$, we obtain a Stein operator

$$(S_P h)(x) = \frac{1}{p(x)}\nabla \cdot [p(x)M(x)h(x)] \tag{12}$$

called the *diffusion Stein operator* (Gorham et al., 2019). This is indeed a Stein operator, since under mild integrability conditions on $K$, the divergence theorem gives that $\int S_p h \, \mathrm{d}P = 0$ for all $h \in \mathcal{H}(K)$; for full details and a proof see Barp et al. (2022b).

## A.3 Selecting a Stein Kernel

There are several choices for a Stein kernel, and which we should use depends on what form of convergence we hope to control (Gorham and Mackey, 2017; Gorham et al., 2019; Hodgkinson et al., 2020; Barp et al., 2022b; Kanagawa et al., 2022). Appendix A.3.1 describes the Langevin–Stein kernel for weak convergence control, Appendix A.3.2 describes the KGM–Stein kernels for additional control over moments, and Appendix A.3.3 presents the Riemann–Stein kernel, whose convergence properties have to-date been less well-studied.

All of the kernels that we consider have length scale parameters that need to be specified, and some also have location parameters to be specified. As a reasonably automatic default we define

$$x_\star \in \arg\max p(x), \qquad \Sigma^{-1} = -\nabla^2 \log p(x_\star)$$

as a location and a matrix of characteristic length scales for $P$ that will be used throughout. These values can typically be obtained using gradient-based optimisation, which is usually cheaper to perform compared to full approximation of $P$. It is assumed that $\nabla^2 \log p(x_\star)$ is positive definite in the sequel.

### A.3.1 Weak Convergence Control with Langevin–Stein Kernels

The first kernel we consider, which we called the Langevin–Stein kernel in the main text, was introduced by Gorham and Mackey (2017). This Stein kernel was developed for the purpose of controlling the weak convergence of a sequence $(Q_n)_{n\in\mathbb{N}} \subset \mathcal{P}(\mathbb{R}^d)$ to $P$. Recall that a sequence $(Q_n)_{n\in\mathbb{N}}$ is said to *converge weakly* (or *in distribution*) to $P$ if $\int f \mathrm{d}Q_n \to \int f \mathrm{d}P$ for all continuous bounded functions $f : \mathbb{R}^d \to \mathbb{R}$. This convergence is denoted $Q_n \overset{\mathrm{d}}{\to} P$ in shorthand.

The problem considered in Gorham and Mackey (2017) was how to select a combination of matrix-valued kernel $K$ (and, implicitly, a diffusion matrix $M$) such that the Stein kernel $k_P$ in (10) generates a KSD $D_P(Q)$ in (4) for which $D_P(Q_n) \to 0$ implies $Q_n \overset{\mathrm{d}}{\to} P$. Their solution was to combine the inverse multi-quadric kernel with an identity diffusion matrix;

$$K(x,y) = (1 + \|x - y\|_\Sigma^2)^{-\beta} I, \qquad M(x) = I$$

for $\beta \in (0,1)$. Provided that $P$ has a density $p$ for which $\nabla \log p(x)$ is Lipschitz, and that $P$ is *distantly dissipative* (see Definition 4 of Gorham and Mackey, 2017), the associated KSD enjoys weak convergence control. Technically, the results in Gorham and Mackey (2017) apply only when $\Sigma = I$, but Theorem 4 in Chen et al. (2019) demonstrated that they hold also for any positive definite $\Sigma$. Following the recommendation of several previous authors, including Chen et al. (2018, 2019); Riabiz et al. (2022), we take $\beta = \frac{1}{2}$ throughout.

### A.3.2 Moment Convergence Control with KGM–Stein Kernels

Despite its many elegant properties, weak convergence can be insufficient for applications where we are interested in integrals $\int f \, dP$ for which the integrand $f : \mathbb{R}^d \to \mathbb{R}$ is unbounded. In particular, this is the case for moments of the form $f(x) = x_1^{\alpha_1} \dots x_d^{\alpha_d}$, $0 \neq \alpha \in \mathbb{N}_0^d$. In such situations, we may seek also the stronger property of *moment convergence control*. The development of KSDs for moment convergence control was recently considered by Kanagawa et al. (2022), and we refered to their construction as the KGM–Stein kernels in the main text. (For convenience, we have adopted the initials of the authors in naming the KGM–Stein kernel.)

A sequence $(Q_n)_{n \in \mathbb{N}} \subset \mathcal{P}(\mathbb{R}^d)$ is said to converge to $P$ *in the sth order moment* if $\int \|x\|^s dQ_n(x) \to \int \|x\|^s dP(x)$. To establish convergence of moments, we need an additional condition on top of weak convergence control: uniform integrability control. A sequence of measures $(Q_n)_{n \in \mathbb{N}}$ is said to have *uniformly integrable sth moments* if for any $\varepsilon > 0$, we can take $r > 0$ such that

$$\sup_{n \in \mathbb{N}} \int_{\|x\| > r} \|x\|^s \, dQ_n(x) < \varepsilon.$$

This condition essentially states that the tail decay of the measures is well-controlled (so that it has a convergent moment). The KSD convergence $D_P(Q_n) \to 0$ implies uniform integrability if for any $\varepsilon > 0$, we can take $r_\varepsilon > 0$ and $f_\varepsilon \in \mathcal{H}(K)$ such that

$$S_P f_\varepsilon(x) \geq \|x\|^s 1\{\|x\| > r_\varepsilon\} - \varepsilon, \tag{13}$$

i.e., the Stein-modified RKHS can approximate the (norm-weighted) indicator function arbitrarily well. Such a function $f_\varepsilon$ can be explicitly constructed (while not guaranteed to be a member of the RKHS). Specifically, the choice $f_\varepsilon = (1 - \iota_\varepsilon) g$ satisfies (13) under an appropriate dissipativity condition, where $\iota_\varepsilon$ is a differentiable indicator function vanishing outside a ball, and $g(x) = -x/\sqrt{1 + \|x\|^2}$. This motivated Kanagawa et al. (2022) to introduce the *sth order KGM–Stein kernel*, which is based on the matrix-valued kernel and diffusion matrix

$$K(x,y) = [\phi(\|x - y\|_\Sigma) + \kappa_{\lin}(x,y)] I, \qquad M(x) = (1 + \|x - x_\star\|_\Sigma^2)^{\frac{s-1}{2}} I,$$

where $(x,y) \mapsto \phi(\|x - y\|_\Sigma)$ is a $C_0^1$ universal kernel (see Barp et al., 2022b, Theorem 4.8). For comparability of our results, we take $\phi$ to be the inverse multi-quadric $\phi(r) = (1 + r^2)^{-1/2}$, and

$$\kappa_{\lin}(x,y) = \frac{1 + (x - x_\star)^\top \Sigma^{-1}(y - x_\star)}{\sqrt{1 + \|x - x_\star\|_\Sigma^2}\sqrt{1 + \|y - x_\star\|_\Sigma^2}}.$$

Here the normalised linear kernel $\kappa_{\lin}$ ensures $g \in \mathcal{H}(K)$, while the $C_0^1$ universal kernel $\phi$ allows approximation of $S_P \iota_\varepsilon g$; see Kanagawa et al. (2022).

### A.3.3 Exploiting Geometry with Riemann–Langevin–Stein Kernels

For academic interest only, here we describe the *Riemann–Stein kernel* that featured in Figure 2 of the main text. This Stein kernel is motivated by the analysis of Gorham et al. (2019), who argued that the use of rapidly mixing Itô diffusions in Stein operators can lead to sharper convergence control. The Riemann–Stein kernel is based on the class of so-called *Riemannian* diffusions considered in Girolami and Calderhead (2011), who proposed to take the diffusion matrix $M$ in (11) to be $M = (\mathcal{I}_{\text{prior}} + \mathcal{I}_{\text{Fisher}})^{-1}$, the inverse of the Fisher information matrix, $\mathcal{I}_{\text{Fisher}}$, regularised using the Hessian of the negative log-prior, $\mathcal{I}_{\text{prior}}$. For the two-dimensional illustration in Section 3.2, this leads to the diffusion matrix

$$M(x) = \left( I + \sum_{i=1}^n [\nabla f_i(x)][\nabla f_i(x)]^\top \right)^{-1},$$

where we recall that $y_i = f_i(x) + \epsilon_i$, where the $\epsilon_i$ are independent with $\epsilon_i \sim \mathcal{N}(0,1)$, and the prior is $x \sim \mathcal{N}(0,1)$. For the presented experiment we paired the above diffusion matrix with the inverse multi-quadric kernel $K(x,y) = (1 + \|x - y\|_\Sigma^2)^{-\beta}$ for $\beta = \frac{1}{2}$. The Riemann–Stein kernel extends naturally to distributions $P$ defined on Riemannian manifolds $\mathcal{X}$; see Barp et al. (2022a) and Example 1 of Hodgkinson et al. (2020).

Unfortunately, the Riemann–Stein kernel is prohibitively expensive in most real applications, since each evaluation of $M$ requires a full scan through the size-$n$ dataset. The computational complexity

of Stein Π-Thinning with the Riemann–Stein kernel is therefore $O(m^2 n^2)$, which is unfavourable compared to the $O(m^2 n)$ complexity in the case where the Stein kernel is not data-dependent. Furthermore, the convergence control properties of the Riemann–Stein kernel have yet to be established. For these reasons we included the Riemann–Stein kernel for illustration only; further groundwork will be required before the Riemann-Stein kernel can be practically used.

# B  Proof of Theorem 1

This appendix is devoted to the proof of Theorem 1. The proof is based on the recent work of Durmus and Moulines (2022), on the geometric convergence of MALA, and on the analysis of sparse (greedy) approximation of kernel discrepancies performed in Riabiz et al. (2022); these existing results are recalled in Appendix B.1. An additional technical result on preconditioned MALA is contained in Appendix B.2. The proof of Theorem 1 itself is contained in Appendix B.3.

## B.1  Auxiliary Results

To precisely describe the results on which our analysis is based, we first need to introduce some notation and terminology. Let $V : \mathcal{X} \to [1, \infty)$ and, for a function $f : \mathcal{X} \to \mathbb{R}$ and a measure $\mu$ on $\mathcal{X}$, let

$$\|f\|_V := \sup_{x \in \mathcal{X}} \frac{|f(x)|}{V(x)}, \qquad \|\mu\|_V := \sup_{\|f\|_V \leq 1} \left| \int_{\mathcal{X}} f \, \mathrm{d}\mu \right|.$$

Recall that a $Q$-invariant Markov chain $(x_i)_{i \in \mathbb{N}} \subset \mathcal{X}$ with $n^{\text{th}}$ step transition kernel $Q^n$ is *V-uniformly ergodic* (see Theorem 16.0.1 of Meyn and Tweedie, 2012) if and only if $\exists R \in [0, \infty), \rho \in (0, 1)$ such that

$$\|Q^n(x, \cdot) - Q\|_V \leq R \rho^n V(x) \tag{14}$$

for all initial states $x \in \mathcal{X}$ and all $n \in \mathbb{N}$.

Although MALA (Algorithm 1) is classical (Roberts and Stramer, 2002), until recently explicit sufficient conditions for ergodicity of MALA had not been obtained. The first result we will need is due Durmus and Moulines (2022), who presented the first explicit conditions for $V$-uniform convergence of MALA. It applies only to *standard* MALA, meaning that the preconditioning matrix $M$ appearing in Algorithm 1 is the identity matrix. The extension of this result to preconditioned MALA will be handled in Appendix B.2.

**Theorem 2.** *Let $Q \in \mathcal{P}(\mathbb{R}^d)$ admit a density, $q$, such that*

(DM1) *there exists $x_0$ with $\nabla \log q(x_0) = 0$*

(DM2) *$q$ is twice continuously differentiable with $\sup_{x \in \mathbb{R}^d} \|\nabla^2 \log q(x - x_0)\| < \infty$*

(DM3) *there exists $b > 0$ and $B \geq 0$ such that $-\nabla^2 \log q(x - x_0) \succeq bI$ for all $\|x - x_0\| \geq B$.*

*Then there exists $\epsilon_0 > 0$ such that for all step sizes $\epsilon \in (0, \epsilon_0)$, standard $Q$-invariant MALA (i.e. with $M = I$) is $V$-uniformly ergodic for $V(x) = \exp\left(\frac{b}{16} \|x - x_0\|^2\right)$.*

*Proof.* This is Theorem 1 of Durmus and Moulines (2022). $\qquad \square$

The next result that we will need establishes consistency of the greedy algorithm applied to samples from a Markov chain that is $Q$-invariant.

**Theorem 3.** *Let $P, Q \in \mathcal{P}(\mathcal{X})$ with $P \ll Q$. Let $k_P : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ be a Stein kernel and let $D_P : \mathcal{X} \times \mathcal{X} \to [0, \infty]$ denote the associated KSD. Consider a $Q$-invariant, time-homogeneous Markov chain $(x_i)_{i \in \mathbb{N}} \subset \mathcal{X}$ such that*

($R^+1$) *$(x_i)_{i \in \mathbb{N}}$ is $V$-uniformly ergodic, such that $V(x) \geq \frac{\mathrm{d}P}{\mathrm{d}Q}(x) \sqrt{k_P(x)}$*

($R^+2$) *$\sup_{i \in \mathbb{N}} \mathbb{E}\left[ \frac{\mathrm{d}P}{\mathrm{d}Q}(x_i) \sqrt{k_P(x_i)} V(x_i) \right] < \infty$*

16

534    (R$^+$3) *there exists $\gamma > 0$ such that $\sup_{i \in \mathbb{N}} \mathbb{E}\left[\exp\left\{\gamma \max\left(1, \frac{\mathrm{d}P}{\mathrm{d}Q}(x_i)^2\right) k_P(x_i)\right\}\right] < \infty.$*

535    *Let $P_{n,m}$ be the result of running the greedy algorithm in (5). If $m \leq n$ and $\log(n) = O(m^{\gamma/2})$ for*
536    *some $\gamma < 1$, then $D_P(P_{n,m}) \to 0$ almost surely as $m, n \to \infty$.*

537    *Proof.* This is Theorem 3 of Riabiz et al. (2022).    $\square$

## B.2   Preconditioned MALA

539    In addition to the auxiliary results in Appendix B.1, which concern standard MALA (i.e. with
540    $M = I$), we require an elementary fact about MALA, namely that preconditioned MALA is
541    equivalent to standard MALA under a linear transformation of the state variable. Recall that the
542    $M$-preconditioned MALA algorithm is a Metropolis–Hastings algorithm whose proposal is the
543    Euler–Maruyama discretisation of the Itô diffusion (11).

544    **Proposition 1.** *Let $M(x) \equiv M$ for a symmetric positive definite and position-independent matrix*
545    *$M \in \mathbb{R}^{d \times d}$. Let $Q \in \mathcal{P}(\mathbb{R}^d)$ admit a probability density function (PDF) $q$ for which the $Q$-invariant*
546    *diffusion $(X_t)_{t \geq 0}$, given by setting $p = q$ in (11), is well-defined. Then under the change of variables*
547    *$Y_t := M^{1/2} X_t$,*

$$\mathrm{d}Y_t = \frac{1}{2}(\nabla \log \tilde{q})(Y_t)\mathrm{d}t + \mathrm{d}W_t, \tag{15}$$

548    *where $\tilde{q}(x) \propto q(M^{-1/2}x)$ for all $x \in \mathbb{R}^d$.*

549    *Proof.* From the chain rule,

$$(\nabla \log \tilde{q})(y) = \nabla_y \log q(M^{-1/2}y) = M^{-1/2}(\nabla \log q)(M^{-1/2}y),$$

550    and thus, substituting $Y_t = M^{1/2} X_t$, (15) is equal to

$$\mathrm{d}X_t = M^{-1/2}\left[\frac{1}{2}M^{-1/2}(\nabla \log q)(M^{-1/2}M^{1/2}X_t) + \mathrm{d}W_t\right]$$
$$= \frac{1}{2}M^{-1}(\nabla \log q)(X_t) + M^{-1/2}\mathrm{d}W_t,$$

551    which is identical to (11) in the case where $M(x) = M$ is constant.    $\square$

552    Let $Q$ and $\tilde{Q}$ be the distributions referred to in Proposition 1, whose PDFs are respectively $q(x)$
553    and $\tilde{q}(x) \propto q(M^{-1/2}x)$. Proposition 1 then implies that the $M$-preconditioned MALA algorithm
554    applied to $Q$ (i.e. Algorithm 1 for $\Pi = Q$) is equivalent to the standard MALA algorithm (i.e.
555    $M = I$) applied to $\tilde{Q}$. This fact allows us to generalise the result of Theorem 2 as follows:

556    **Corollary 1.** *Consider a symmetric positive definite matrix $M \in \mathbb{R}^{d \times d}$. Assume that conditions*
557    *(DM1-3) in Theorem 2 are satisfied. Then there exists $\epsilon_0' > 0$ and $b' > 0$ such that for all step*
558    *sizes $\epsilon \in (0, \epsilon_0')$, the $M$-preconditioned $Q$-invariant MALA is $V$-uniformly ergodic for $V(x) =$*
559    *$\exp\left(\frac{b'}{16}\|x - x_0\|^2\right)$.*

560    *Proof.* From Theorem 2 and Proposition 1, the result follows if we can establish (DM1-3) for $\tilde{Q}$, since
561    $M$-preconditioned MALA is equivalent to standard MALA applied to $\tilde{Q}$. For a matrix $A \in \mathbb{R}^{d \times d}$,
562    let $\lambda_{\min}(A)$ and $\lambda_{\max}(A)$ respectively denote the minimum and maximum eigenvalues of $A$. For
563    (DM1) we set $y_0 = M^{1/2}x_0$ and observe that

$$(\nabla \log \tilde{q})(y_0) = M^{-1/2}(\nabla \log q)(x_0) = 0.$$

564    For (DM2) we have that

$$\sup_{y \in \mathbb{R}^d} \|\nabla^2(\log \tilde{q})(y - y_0)\| = \sup_{y \in \mathbb{R}^d} \|M^{-1/2}(\nabla^2 \log q)(M^{-1/2}(y - y_0))M^{-1/2}\|$$
$$\leq \lambda_{\min}(M)^{-1} \sup_{x \in \mathbb{R}^d} \|(\nabla^2 \log q)(x - x_0)\| < \infty.$$

For (DM3) we have that

$$-(\nabla^2 \log \tilde{q})(y - y_0) = -M^{-1/2}(\nabla^2 \log q)(M^{-1/2}(y - y_0))M^{-1/2}$$
$$= -M^{-1/2}(\nabla^2 \log q)(x - x_0)M^{-1/2} \succeq M^{1/2}(bI)M^{1/2} = bM^{-1} \succeq b'I$$

where $b' = b\lambda_{\max}(M)^{-1}$, which holds for all $\|x - x_0\| \geq B$, and in particular for all $\|y - y_0\| \geq B'$ where $B' = B\lambda_{\max}(M)^{1/2}$. Thus (DM1-3) are established for $\tilde{Q}$. $\qquad\square$

**Remark 1.** *The choice $M = \Sigma^{-1}$, which sets the preconditioner matrix $M$ equal to the inverse of the length scale matrix $\Sigma$ used in the specification of the kernel $K$ (c.f. Appendix A.3), leads to the elegant interpretation that Stein $\Pi$-Importance Sampling applied to $M$-preconditioned MALA is equivalent to the Stein $\Pi$-Importance Sampling applied to standard MALA (i.e. with $M = I$) for the whitened target $\tilde{P}$ with PDF $\tilde{p}(x) \propto p(M^{-1/2}x)$. For our experiments, however, the preconditioner matrix $M$ was learned during a warm-up phase of MALA, since in general the curvature of $P$ (captured by $\Sigma$) and the curvature of $\Pi$ (captured by $M^{-1}$) may be different.*

### B.3 Proof of Theorem 1

The route to establishing Theorem 1 has three parts. First, we establish (DM1-3) of Theorem 2 with $Q = \Pi$, to deduce from Corollary 1 that $\Pi$-invariant $M$-preconditioned MALA is $V$-uniformly ergodic. This in turn enables us to establish conditions (R$^+$1-3) of Theorem 3, again for $Q = \Pi$, from which the strong consistency $D_P(P_{n,m}) \xrightarrow{\text{a.s.}} 0$ of S$\Pi$T-MALA is established. Finally, we note that $0 \leq D_P(P_n^\star) \leq D_P(P_{n,m})$, since the support of $P_{n,m}$ is contained in the support of $P_n^\star$, and the latter is optimally weighted, whence also the strong consistency of S$\Pi$IS-MALA.

**Establish (DM1-3)** First we establish (DM1-3) for $Q = \Pi$. Fix $x_0 \in \mathbb{R}^d$. For (DM2), first recall that the range of $k_P$ is $[C_1^2, \infty)$ where $C_1 > 0$, from Assumption 1. Since $\log(\cdot)$ has bounded second derivatives on $[C_1^2, \infty)$, there is a constant $C > 0$ such that

$$\forall x \in \mathbb{R}^d, \qquad \|\nabla^2 \log k_P(x)\| \leq C\|\nabla^2 k_P(x)\|.$$

Thus, using compactness of the set $\{x : \|x - x_0\| \leq B_2\}$,

$$\sup_{x \in \mathbb{R}^d} \|\nabla^2 \log k_P(x)\| \leq C \max \left( \underbrace{\sup_{\|x - x_0\| \leq B_2} \|\nabla^2 k_P(x)\|}_{<\infty \text{ by (A3)}}, \underbrace{\sup_{\|x - x_0\| \geq B_2} \|\nabla^2 k_P(x)\|}_{<b_2\|I\| \text{ by (A4)}} \right) < \infty. \quad (16)$$

Now, $\pi$ is twice differentiable as it is the product of twice differentiable functions $p$ and $k_P^{1/2}$ from (A1) and (A3), and moreover

$$\sup_{x \in \mathbb{R}^d} \|\nabla^2 \log \pi(x - x_0)\| \leq \underbrace{\sup_{x \in \mathbb{R}^d} \|\nabla^2 \log p(x)\|}_{<\infty \text{ by (A1)}} + \frac{1}{2}\underbrace{\sup_{x \in \mathbb{R}^d} \|\nabla^2 \log k_P(x)\|}_{<\infty \text{ by (16)}} < \infty,$$

so (DM2) is satisfied. For (DM3), first note from the chain and product rules that for all $\|x\| \geq B_2$

$$\nabla^2 \log k_P(x - x_0) = \underbrace{\frac{\nabla^2 k_P(x - x_0)}{k_P(x - x_0)}}_{\preceq(b_2/C_1^2)I \text{ by (A4)}} - \underbrace{\frac{[\nabla k_P(x - x_0)][\nabla k_P(x - x_0)]^\top}{k_P(x - x_0)^2}}_{\succeq 0} \preceq \frac{b_2}{C_1^2}I. \quad (17)$$

Thus, for all $\|x - x_0\| \geq B := \|x_0\| + \max(B_1, B_2)$,

$$-\nabla^2 \log \pi(x - x_0) = \underbrace{-\nabla^2 \log p(x - x_0)}_{\succeq b_1 I \text{ by (A2)}} - \frac{1}{2}\underbrace{\nabla^2 \log k_P(x - x_0)}_{\preceq(b_2/C_1^2)I \text{ by (17)}} \succeq \underbrace{\left(b_1 - \frac{b_2}{2C_1^2}\right)I}_{=:b>0} \quad (18)$$

as required. The same argument establishes (DM1); from (18) we have $\lim_{\|x\|\to\infty} \pi(x) = 0$, and since $\pi$ is a continuously differentiable density there must exist an $x_0$ at which $\pi$ is locally minimised. Thus we have established (DM1-3) for $Q = \Pi$ and we may conclude from Corollary 1 that there is an $\epsilon_0' > 0$ and $b' > 0$ such that, for all $\epsilon \in (0, \epsilon_0')$, the $\Pi$-invariant $M$-preconditioned MALA chain $(x_i)_{i\in\mathbb{N}}$ is $V$-uniformly ergodic for $V(x) = C_2 \exp\left(\frac{b'}{16}\|x - x_0\|^2\right)$ (since if a Markov chain is $V$-uniformly ergodic, then it is also $CV$-uniformly ergodic).

18

**Establish ($R^+$1-3)** The aim is now to establish conditions ($R^+$1-3) of Theorem 3 for $Q = \Pi$. By construction $dP/d\Pi = C_2/\sqrt{k_P(x)} < C_2/C_1 < \infty$, where $C_1$ and $C_2$ were defined in Assumption 1, so that $P \ll \Pi$. It has already been established that $(x_i)_{i \in \mathbb{N}}$ is $V$-uniformly ergodic, and further

$$V(x) = C_2 \exp\left(\frac{b'}{16}\|x - x_0\|^2\right) \geq C_2 = \frac{dP}{d\Pi}(x)\sqrt{k_P(x)}$$

for all $x$, which establishes ($R^+$1). Let $R$ and $\rho$ denote constants for which the $V$-uniform ergodicity property (14) is satisfied. From $V$-uniform ergodicity, the integral $\int V \, d\Pi$ exists and

$$\left|\mathbb{E}\left[\frac{dP}{d\Pi}(x_i)\sqrt{k_P(x_i)}V(x_i)\right] - C_2 \int V \, d\Pi\right| = C_2 \left|\mathbb{E}[V(x_i)] - \int V \, d\Pi\right|$$
$$\leq C_2 R\rho^n V(x_0) \to 0$$

which establishes ($R^+$2). Fix $\gamma > 0$. By construction $dP/d\Pi \leq C_2/C_1$, and thus

$$\exp\left\{\gamma \max\left(1, \frac{dP}{d\Pi}(x)^2\right)k_P(x)\right\} < \exp\{\tilde{\gamma}k_P(x)\}$$

where $\tilde{\gamma} = \max(1, C_2/C_1)\gamma$. Since we have assumed that $k_P$ is continuous with, from (A4),

$$b_3 := \limsup_{\|x\| \to \infty} \frac{k_P(x)}{\|x\|^2} < \infty,$$

we may take $\gamma$ such that $\tilde{\gamma}b_3 < b'/16$, so that $\|x \mapsto \exp\{\tilde{\gamma}k_P(x)\}\|_V < \infty$ and in particular

$$\left|\mathbb{E}\left[\exp\{\tilde{\gamma}k_P(x_i)\}\right] - \int \exp\{\tilde{\gamma}k_P(x)\} \, d\Pi(x)\right| \leq \|x \mapsto \exp\{\tilde{\gamma}k_P(x)\}\|_V \times R\rho^n V(x_0) \to 0$$

which establishes ($R^+$3). Thus we have established ($R^+$1-3) for $Q = \Pi$, so from Theorem 3 we have strong consistency of SΠT-MALA (i.e. $D_P(P_{n,m}) \overset{\text{a.s.}}{\to} 0$) provided that $m \leq n$ with $\log(n) = O(m^{\gamma/2})$ for some $\gamma < 1$. The latter condition is equivalent to $m = \Omega((\log n)^\delta)$ for some $\delta > 2$, which we used for the statement. Since $0 \leq D_P(P_n^\star) \leq D_P(P_{n,m})$, the strong consistency of SΠIS-MALA is also established.

# C   Explicit Calculation of Stein Kernels

This appendix contains explicit calculations for the Langevin–Stein and KGM–Stein kernels $k_P$, which are sufficient to implement Stein $\Pi$-Importance Sampling and Stein $\Pi$-Thinning. These calculations can also be performed using automatic differentiation, but comparison to the analytic expressions is an important step in validation of computer code.

To proceed, we observe that the diffusion Stein operator $S_P$ in (12) applied to a matrix-valued kernel $K$ is equivalent to the Langevin–Stein operator applied to the kernel $C(x, y) = M(x)K(x, y)M(y)^\top$. In the case of the Langevin–Stein and KGM–Stein kernels we have $K(x, y) = \kappa(x, y)I$ for some $\kappa(x, y)$ and $M(x) = (1 + \|x - x_\star\|_\Sigma^2)^{(s-1)/2}I$ for some $s \in \{0, 1, 2, \dots\}$. Thus $C(x, y) = c(x, y)I$ where

$$c(x, y) := (1 + \|x - x_\star\|_\Sigma^2)^{(s-1)/2}(1 + \|y - x_\star\|_\Sigma^2)^{(s-1)/2}\kappa(x, y)$$

and

$$k_P(x, y) = \nabla_x \cdot \nabla_y c(x, y) + [\nabla_x c(x, y)] \cdot [\nabla_y \log p(y)] + [\nabla_y c(x, y)] \cdot [\nabla_x \log p(x)]$$
$$+ c(x, y)[\nabla_x \log p(x)] \cdot [\nabla_y \log p(y)],$$

following the calculations in Oates et al. (2017). To evaluate the terms in this formula we start by differentiating $c(x, y)$, to obtain

$$\nabla_x c(x, y) = (1 + \|x - x_\star\|_\Sigma^2)^{(s-1)/2}(1 + \|y - x_\star\|_\Sigma^2)^{(s-1)/2}$$
$$\times \left[ \frac{(s-1)\kappa(x,y)\Sigma^{-1}(x - x_\star)}{1 + \|x - x_\star\|_\Sigma^2} + \nabla_x \kappa(x, y) \right]$$

$$\nabla_y c(x, y) = (1 + \|x - x_\star\|_\Sigma^2)^{(s-1)/2}(1 + \|y - x_\star\|_\Sigma^2)^{(s-1)/2}$$
$$\times \left[ \frac{(s-1)\kappa(x,y)\Sigma^{-1}(y - x_\star)}{1 + \|y - x_\star\|_\Sigma^2} + \nabla_y \kappa(x, y) \right]$$

$$\nabla_x \cdot \nabla_y c(x, y) = (1 + \|x - x_\star\|_\Sigma^2)^{(s-1)/2}(1 + \|y - x_\star\|_\Sigma^2)^{(s-1)/2}$$
$$\times \left[ \frac{(s-1)^2\kappa(x,y)(x - x_\star)^\top \Sigma^{-2}(y - x_\star)}{(1 + \|x - x_\star\|_\Sigma^2)(1 + \|y - x_\star\|_\Sigma^2)} + \frac{(s-1)(y - x_\star)^\top \Sigma^{-1}\nabla_x \kappa(x,y)}{(1 + \|y - x_\star\|_\Sigma^2)} \right.$$
$$\left. + \frac{(s-1)(x - x_\star)^\top \Sigma^{-1}\nabla_y \kappa(x,y)}{(1 + \|x - x_\star\|_\Sigma^2)} + \nabla_x \cdot \nabla_y \kappa(x, y) \right].$$

These expressions involve gradients of $\kappa(x, y)$, and explicit formulae for these are presented for the choice of $\kappa(x, y)$ corresponding to the Langevin–Stein kernel in Appendix C.1, and to the KGM–Stein kernel in Appendix C.2.

To implement Stein $\Pi$-Thinning we require access to both $k_P(x)$ and $\nabla k_P(x)$, the latter for use in the proposal distribution and acceptance probability in MALA. These quantities will now be calculated. In what follows we assume that $\kappa(x, y)$ is continuously differentiable, so that partial derivatives with respect to $x$ and $y$ can be interchanged. Then

$$c_0(x) := c(x, x)$$
$$= (1 + \|x - x_\star\|_\Sigma^2)^{s-1}\kappa(x, x)$$
$$c_1(x) := \nabla_x c(x, y)|_{y \to x}$$
$$= (1 + \|x - x_\star\|_\Sigma^2)^{s-1}\left[ \frac{(s-1)\kappa(x,x)\Sigma^{-1}(x - x_\star)}{(1 + \|x - x_\star\|_\Sigma^2)} + \nabla_x \kappa(x, y)|_{y \to x} \right]$$
$$c_2(x) := \nabla_x \cdot \nabla_y c(x, y)|_{y \to x}$$
$$= (1 + \|x - x_\star\|_\Sigma^2)^{s-1}\left[ \frac{(s-1)^2\kappa(x,x)(x - x_\star)^\top \Sigma^{-2}(x - x_\star)}{(1 + \|x - x_\star\|_\Sigma^2)^2} \right.$$
$$\left. + \frac{2(s-1)(x - x_\star)^\top \Sigma^{-1}\nabla_x \kappa(x,y)|_{y \to x}}{(1 + \|x - x_\star\|_\Sigma^2)} + \nabla_x \cdot \nabla_y \kappa(x, y)|_{y \to x} \right]$$

so that

$$k_P(x) := k_P(x, x) = c_2(x) + 2c_1(x) \cdot \nabla_x \log p(x) + c_0(x)\|\nabla_x \log p(x)\|^2. \tag{19}$$

Let $[\nabla_x c_1(x)]_{i,j} = \partial_{x_i}[c_1(x)]_j$ and $[\nabla_x^2 \log p(x)]_{i,j} = \partial_{x_i}\partial_{x_j} \log p(x)$. Now we can differentiate (19) to get

$$\nabla_x k_P(x) = \nabla_x c_2(x) + 2[\nabla_x c_1(x)][\nabla_x \log p(x)] + 2[\nabla_x^2 \log p(x)]c_1(x)$$
$$+ [\nabla_x c_0(x)]\|\nabla_x \log p(x)\|^2 + 2c_0(x)[\nabla_x^2 \log p(x)][\nabla_x \log p(x)]. \tag{20}$$

In what follows we also derive explicit formulae for $c_0(x)$, $c_1(x)$ and $c_2(x)$, and hence for $\nabla_x c_0(x)$, $\nabla_x c_1(x)$ and $\nabla_x c_2(x)$, for the case of the Langevin–Stein kernel in Appendix C.1, and the KGM–Stein kernel in Appendix C.2.

20

### C.1 Explicit Formulae for the Langevin–Stein Kernel

The Langevin–Stein kernel from Appendix A.3.1 corresponds to the choice $s = 1$ and $\kappa(x, y)$ the inverse multi-quadric kernel, so that

$$\kappa(x, y) = (1 + \|x - y\|_\Sigma^2)^{-\beta}$$
$$\nabla_x \kappa(x, y) = -2\beta(1 + \|x - y\|_\Sigma^2)^{-\beta-1}\Sigma^{-1}(x - y)$$
$$\nabla_y \kappa(x, y) = 2\beta(1 + \|x - y\|_\Sigma^2)^{-\beta-1}\Sigma^{-1}(x - y)$$
$$\nabla_x \cdot \nabla_y \kappa(x, y) = -4\beta(\beta + 1)(1 + \|x - y\|_\Sigma^2)^{-\beta-2}(x - y)^\top \Sigma^{-2}(x - y)$$
$$+ 2\beta\text{tr}(\Sigma^{-1})(1 + \|x - y\|_\Sigma^2)^{-\beta-1}.$$

Evaluating on the diagonal:

$$\kappa(x, x) = 1$$
$$\nabla_x \kappa(x, y)|_{y \to x} = \nabla_y \kappa(x, y)|_{y \to x} = 0$$
$$\nabla_x \cdot \nabla_y \kappa(x, y)|_{y \to x} = 2\beta\text{tr}(\Sigma^{-1}),$$

so that $c_0(x) = 1$, $c_1(x) = 0$, $c_2(x) = 2\beta\text{tr}(\Sigma^{-1})$. Differentiating these formulae, $\nabla_x c_0(x) = 0$, $\nabla_x c_1(x) = 0$, $\nabla_x c_2(x) = 0$.

### C.2 Explicit Formulae for the KGM–Stein Kernel

The KGM kernel of order $s$ from Appendix A.3.2 corresponds to the choice

$$\kappa(x, y) = (1 + \|x - y\|_\Sigma^2)^{-\beta} + \frac{1 + (x - x_\star)^\top \Sigma^{-1}(y - x_\star)}{(1 + \|x - x_\star\|_\Sigma^2)^{s/2}(1 + \|y - x_\star\|_\Sigma^2)^{s/2}},$$

for which we have

$$\nabla_x \kappa(x, y) = -2\beta(1 + \|x - y\|_\Sigma^2)^{-\beta-1}\Sigma^{-1}(x - y)$$
$$+ \frac{\Sigma^{-1}(y - x_\star) - s[1 + (x - x_\star)^\top \Sigma^{-1}(y - x_\star)]\Sigma^{-1}(x - x_\star)(1 + \|x - x_\star\|_\Sigma^2)^{-1}}{(1 + \|x - x_\star\|_\Sigma^2)^{s/2}(1 + \|y - x_\star\|_\Sigma^2)^{s/2}}$$
$$\nabla_y \kappa(x, y) = 2\beta(1 + \|x - y\|_\Sigma^2)^{-\beta-1}\Sigma^{-1}(x - y)$$
$$+ \frac{\Sigma^{-1}(x - x_\star) - s[1 + (x - x_\star)^\top \Sigma^{-1}(y - x_\star)]\Sigma^{-1}(y - x_\star)(1 + \|y - x_\star\|_\Sigma^2)^{-1}}{(1 + \|x - x_\star\|_\Sigma^2)^{s/2}(1 + \|y - x_\star\|_\Sigma^2)^{s/2}}$$
$$\nabla_x \cdot \nabla_y \kappa(x, y) = -4\beta(\beta + 1)(1 + \|x - y\|_\Sigma^2)^{-\beta-2}(x - y)^\top \Sigma^{-2}(x - y) + 2\beta\text{tr}(\Sigma^{-1})(1 + \|x - y\|_\Sigma^2)^{-\beta-1}$$
$$+ \frac{\begin{bmatrix} \text{tr}(\Sigma^{-1}) - s(1 + \|x - x_\star\|_\Sigma^2)^{-1}(x - x_\star)^\top \Sigma^{-2}(x - x_\star) \\ -s(1 + \|y - x_\star\|_\Sigma^2)^{-1}(y - x_\star)^\top \Sigma^{-2}(y - x_\star) \\ +s^2[1 + (x - x_\star)^\top \Sigma^{-1}(y - x_\star)](1 + \|x - x_\star\|_\Sigma^2)^{-1}(1 + \|y - x_\star\|_\Sigma^2)^{-1} \\ \times (x - x_\star)^\top \Sigma^{-2}(y - x_\star) \end{bmatrix}}{(1 + \|x - x_\star\|_\Sigma^2)^{s/2}(1 + \|y - x_\star\|_\Sigma^2)^{s/2}}.$$

Evaluating on the diagonal:

$$\kappa(x, x) = 1 + (1 + \|x - x_\star\|_\Sigma^2)^{-s+1}$$
$$\nabla_x \kappa(x, y)|_{y \to x} = \nabla_y \kappa(x, y)|_{y \to x} = -(s - 1)\Sigma^{-1}(x - x_\star)(1 + \|x - x_\star\|_\Sigma^2)^{-s}$$
$$\nabla_x \cdot \nabla_y \kappa(x, y)|_{y \to x} = 2\beta\text{tr}(\Sigma^{-1}) + \text{tr}(\Sigma^{-1})(1 + \|x - x_\star\|_\Sigma^2)^{-s}$$
$$+ s(s - 2)(1 + \|x - x_\star\|_\Sigma^2)^{-s-1}(x - x_\star)^\top \Sigma^{-2}(x - x_\star)$$

so that

$$c_0(x) = 1 + (1 + \|x - x_\star\|_\Sigma^2)^{s-1}$$
$$c_1(x) = (s - 1)(1 + \|x - x_\star\|_\Sigma^2)^{s-2}\Sigma^{-1}(x - x_\star)$$
$$c_2(x) = \frac{[(s - 1)^2(1 + \|x - x_\star\|_\Sigma^2)^{s-1} - 1](x - x_\star)^\top \Sigma^{-2}(x - x_\star)}{(1 + \|x - x_\star\|_\Sigma^2)^2} + \frac{\text{tr}(\Sigma^{-1})[1 + 2\beta(1 + \|x - x_\star\|_\Sigma^2)^s]}{(1 + \|x - x_\star\|_\Sigma^2)}.$$

645   Differentiating these formulae:

$$\nabla_x c_0(x) = 2(s-1)(1+\|x-x_\star\|_\Sigma^2)^{s-2}\Sigma^{-1}(x-x_\star)$$

$$\nabla_x c_1(x) = 2(s-1)(s-2)(1+\|x-x_\star\|_\Sigma^2)^{s-3}[\Sigma^{-1}(x-x_\star)][\Sigma^{-1}(x-x_\star)]^\top$$
$$+ (s-1)(1+\|x-x_\star\|_\Sigma^2)^{s-2}\Sigma^{-1}$$

$$\nabla_x c_2(x) = 2(s-1)^2(s-3)(1+\|x-x_\star\|_\Sigma^2)^{s-4}[(x-x_\star)^\top\Sigma^{-2}(x-x_\star)]\Sigma^{-1}(x-x_\star)$$
$$+ 2(s-1)^2(1+\|x-x_\star\|_\Sigma^2)^{s-3}\Sigma^{-2}(x-x_\star)$$
$$+ 4\beta\mathrm{tr}(\Sigma^{-1})(s-1)(1+\|x-x_\star\|_\Sigma^2)^{s-2}\Sigma^{-1}(x-x_\star)$$
$$- 2(1+\|x-x_\star\|_\Sigma^2)^{-2}[\Sigma^{-2}(x-x_\star)+\mathrm{tr}(\Sigma^{-1})\Sigma^{-1}(x-x_\star)]$$
$$+ 4(1+\|x-x_\star\|_\Sigma^2)^{-3}[(x-x_\star)^\top\Sigma^{-2}(x-x_\star)]\Sigma^{-1}(x-x_\star).$$

646   These complete the analytic calculations necessary to compute the Stein kernel $k_P$ and its gradient.

647 ## D   Empirical Assessment

648   This appendix contains full details of the empirical protocols that were employed and the additional
649   empirical results described in the main text. Appendix D.1 discusses the effect of dimension on
650   our proposed $\Pi$. Additional illuatrative results from Section 3.2 are contained in Appendix D.2.
651   The full details for how MALA was implemented are contained in Appendix D.3. An additional
652   illustration using a generalised auto-regressive moving average (GARCH) model is presented in
653   Appendix D.4. The full results for S$\Pi$IS-MALA are contained in Appendix D.5, and in Appendix D.6
654   the convergence of the sparse approximation provided by S$\Pi$T-MALA to the optimal weighted
655   approximation is investigated. Finally, the performance of KSDs is quantified using the 1-Wasserstein
656   divergence in Appendix D.7.

657 ### D.1   The Effect of Dimension on $\Pi$

658   The improvement of Stein $\Pi$-Importance Sampling over the default Stein importance sampling
659   algorithm (i.e. $\Pi = P$) can be expected to reduce as the dimension $d$ of the target $P$ is increased. To
660   see this, consider the Langevin–Stein kernel

$$k_P(x) = c_1 + c_2\|\nabla\log p(x)\|_\Sigma^2 \tag{21}$$

661   for some $c_1, c_2 > 0$; see Appendix C. Taking $P = \mathcal{N}(0, I_{d\times d})$, for which the length scale matrix $\Sigma$
662   appearing in Appendix A.3 is $\Sigma = I_{d\times d}$, we obtain

$$k_P(x) = c_1 + c_2\|x\|^2.$$

663   However, the sampling distribution $\Pi$ defined in (8) depends on $k_P$ only up to an unspecified
664   normalisation constant; we may therefore equally consider the asymptotic behaviour of $\tilde{k}_P(x) :=$
665   $k_P(x)/d$. Let $X \sim P$. Then $\mathbb{E}[\tilde{k}_P(X)] = c_2$ is a $d$-independent constant, and

$$\left\|\tilde{k}_P - \mathbb{E}[\tilde{k}_P(X)]\right\|_{L^2(P)}^2 = \int\left[\frac{k_P(x)-(c_1+c_2 d)}{d}\right]^2 \mathrm{d}P(x) = \frac{2c_2^2}{d} \to 0$$

666   as $d \to \infty$. This shows that $\tilde{k}_P$ converges to a constant function in $L^2(P)$, and thus for "typical"
667   values of $x$ in the effective support of $P$,

$$\pi(x) \propto p(x)\sqrt{\tilde{k}_P(x)} \overset{\approx}{\propto} p(x),$$

668   so that $\Pi \approx P$ in the $d \to \infty$ limit. This intuition is borne out in simulations involving both the
669   Langevin–Stein kernel (as just discussed) and also the KGM3–Stein kernel. Indeed, Figure S1 shows
670   that as the dimension $d$ is increased, the marginal distributions of $\Pi$ become increasingly similar to
671   those of $P$.

Figure S1: The effect of dimension on $\Pi$: Here $P$ was taken to be the standard Gaussian distribution $\mathcal{N}(0, I_{d \times d})$ in $\mathbb{R}^d$ and the proposed distribution $\Pi$ was computed. The marginal distribution of the first component of $\Pi$ is plotted for $d \in \{1, 2, 10\}$, for both (a) the Langevin–Stein kernel and (b) the KGM3–Stein kernel.



Figure S2: Assessing the performance of the sampling distributions $\Pi$ shown in Figure 2. The mean kernel Stein discrepancy (KSD) for computation performed using the Langevin–Stein kernel (purple), the KGM3–Stein kernel (blue), and the Riemann–Stein kernel (red); in each case, KSD was computed using the same Stein kernel used to construct $\Pi$. Solid lines indicate the baseline case of sampling from $P$, while dashed lines indicate the proposed approach of sampling from $\Pi$. (The experiment was repeated 10 times and standard error bars are plotted.)

## D.2    2D Illustration from the Main Text

Section 3.2 of the main text contained a 2-dimensional illustration of Stein $\Pi$-Importance Sampling and presented the distributions $\Pi$ corresponding to different choices of Stein kernel. Here, in Figure S2, we present the mean KSDs for Stein $\Pi$-Importance Sampling performed using the Langevin–Stein kernel (purple), the KGM3–Stein kernel (blue), and the Riemann–Stein kernel (red), corresponding to the sampling distributions $\Pi$ displayed in Figure 2 of the main text.

For this experiment, exact sampling from both $P$ and $\Pi$ was performed using a fine grid on which all probabilities were calculated and appropriately normalised. Results are in broad agreement with the 1-dimensional illustration contained in the main text, in the sense that in all cases Stein $\Pi$-Importance Sampling provides a significant improvement over the default Stein importance sampling method with $\Pi$ equal to $P$.

23

**Algorithm 4** Adaptive MALA

---

**Require:** $x_{0,0}$ (initial state), $\epsilon_0$ (initial step size), $M_0$ (initial preconditioner matrix), $\{n_i\}_{i=0}^{h-1}$ (epoch lengths),
$\quad$ $\{\alpha_i\}_{i=1}^{h-1}$ (learning schedule), $h$ (number of epochs), $k_P$ (Stein kernel)
1: $\{x_{0,1} \ldots, x_{0,n_0}\} \leftarrow \texttt{MALA}(x_{0,0}, \epsilon_0, M_0, n_0, k_P)$
2: **for** $i = 1, \ldots, h-1$ **do**
3: $\quad$ $x_{i,0} \leftarrow x_{i-1,n_{i-1}}$
4: $\quad$ $\rho_{i-1} \leftarrow \frac{1}{n_{i-1}} \sum_{j=1}^{n_{i-1}} 1_{x_{i-1,j} \neq x_{i-1,j-1}}$ $\qquad\qquad$ ▷ Average acceptance rate for chain $i$
5: $\quad$ $\epsilon_i \leftarrow \epsilon_{i-1} \exp(\rho_{i-1} - 0.57)$ $\qquad\qquad\qquad\qquad$ ▷ Update step size
6: $\quad$ $M_i \leftarrow \alpha_i M_i + (1 - \alpha_i)\mathrm{cov}(\{x_{i-1,1} \ldots, x_{i-1,n_{i-1}}\})$ $\qquad$ ▷ Update preconditioner matrix
7: $\quad$ $\{x_{i,1} \ldots, x_{i,n_i}\} \leftarrow \texttt{MALA}(x_{i,0}, \epsilon_i, M_i, n_i, k_P)$
8: **end for**

---

## D.3 Implementation of MALA

For implementation of MALA in Algorithm 4 we are required to specify a step size $\epsilon$ and a preconditioner matrix $M$. In general, suitable values for both of these parameters will be problem-dependent. Standard practice is to perform some form of manual or automated tuning to arrive at parameter values for which the average acceptance rate is close to 0.57, motivated by the asymptotic analysis of Roberts and Rosenthal (1998). Adaptive MCMC algorithms, which seek to optimise the parameters of MCMC algorithms such as MALA during the warm-up period, provide an appealing solution, and was the approach taken in this work.

The adaptive MALA algorithm which we used is contained in Algorithm 4, where we have let $\texttt{MALA}(x, \epsilon, M, n, k_P)$ denote the output from the preconditioned MALA with initial state $x$, step size $\epsilon$, preconditioner matrix $M$, and chain length $n$, described in Algorithm 1. In Algorithm 4, we use $\mathrm{cov}(\cdot)$ to denote the sample covariance matrix. The algorithm monitors the average acceptance rate and increases or decreases it according to whether it is below or above, respectively, the 0.57 target. For the preconditioner matrix, the sample covariance matrix of samples obtained from the penultimate tuning run of MALA is used. For all experiments that we report using MALA, we set $\epsilon_0 = 1$, $M_0 = I_d$, $h = 10$, and $\alpha_1 = \cdots = \alpha_9 = 0.3$. The warm-up epoch lengths were $n_0 = \cdots = n_8 = 1,000$ and the final epoch length was $n_9 = 10^5$. The samples $\{x_{h-1,1}, \ldots, x_{h-1,n_{i-1}}\}$ from the final epoch are returned, and constituted output from MALA for our experimental assessment.

To sample from $P$ instead of $\Pi$, we used Algorithm 4 we formally set $k_P(x) = 1$ for all $x \in \mathbb{R}^d$, which recovers $\Pi = P$ as the target.

## D.4 Illustration on a GARCH Model

This appendix contains an additional illustrative experiment, concerning a GARCH model that is a particular instance of a model from the $\texttt{PosteriorDB}$ database discussed in Section 4. The purpose of this illustration is to facilitate an empirical investigation in a slightly higher dimension ($d = 4$) and to explore the effect of changing the order $s$ of the KGM–Stein kernel defined in Appendix A.3.2.

First we describe the GARCH model that was used. These models are widely-used in econometrics to describe time series data $\{y_t\}_{t=1}^n$ in settings where the volatility process is assumed to be time-varying (but stationary). In particular, we consider the GARCH(1,1) model

$$y_t = \phi_1 + a_t,$$
$$a_t = \sigma_t \epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, 1),$$
$$\sigma_t^2 = \phi_2 + \phi_3 a_{t-1}^2 + \phi_4 \sigma_{t-1}^2,$$

where $\phi_2 > 0$, $\phi_3 > 0$, $\phi_4 > 0$, and $\phi_3 + \phi_4 < 1$ are the model parameters, constrained to a subset of $\mathbb{R}^4$. For ease of sampling, a change of variables $\tau : (\phi_1, \phi_2, \phi_3, \phi_4) \mapsto \theta$ is performed in such a way that the parameter $\theta \in \mathbb{R}^4$ is unconstrained. Assuming an improper flat prior on $\theta$, the log-posterior density for $\theta$ is given up to an additive constant by

$$\log p(\theta \mid y_1, \ldots, y_n) \stackrel{+C}{=} \sum_{t=1}^n \left[ -\frac{1}{2} \log(\sigma_t^2) - \frac{y_t^2}{2\sigma_t^2} \right] + \log|J_{\tau^{-1}}(\theta)|,$$

where $|J_{\tau^{-1}}(\theta)|$ is the Jacobian determinant of $\tau^{-1}$.

Figure S3: Illustrating the shape of $\Pi$ based on the KGM$s$–Stein kernel for a GARCH(1,1) model, controlling convergence of moments up to order $s \in \{2, 3, 4\}$. The marginal density functions of each distribution were approximated using one-million samples obtained using MCMC.

For this illustration, real data were provided within the model description of `PosteriorDB`, for which the estimated *maximum a posteriori* parameter is $\hat{\phi} = (5.04, 1.36, 0.53, 0.31)$. The marginal distributions of $\Pi$ corresponding to the KGM–Stein kernels of orders $s \in \{2, 3, 4\}$ are compared to the marginals of $P$ in Figure S3. It can be seen that higher orders $s$ correspond to greater over-dispersion of $\Pi$; this makes intuitive sense since larger $s$ corresponds to a more stringent KSD (controlling the convergence of moments up to order $s$) which places greater emphasis on how the tails of $P$ are approximated. Further, for the final skewed marginal of $P$, we note that the distribution $\Pi$ exaggerates the skew, placing more of its mass in the tail of the direction which is positively skewed. Further discussion of skewed targets is contained in Appendix D.8.

### D.5  Stein $\Pi$-Importance Sampling for `PosteriorDB`

To introduce objectivity into our assessment, we exploited the `PosteriorDB` benchmark (Magnusson et al., 2022). This ongoing project is an attempt toward standardised benchmarking, consisting of a collection of posteriors to be numerically approximated. The test problems in `PosteriorDB` are defined in the `Stan` probabilistic programming language, and so `BridgeStan` (Roualdes et al., 2023) was used to directly access posterior densities and their gradients as required. The ambition of `PosteriorDB` is to provide an extensive set of benchmark tasks; at the time we conducted our research, `PosteriorDB` was at Version 0.4.0 and contained 149 models, of which 47 came equipped with a gold-standard sample of size $n = 10^3$, generated from a long run of Hamiltonian Monte Carlo (the No-U-Turn sampler in `Stan`). Of these 47 models, a subset of 40 were found to be compatible with `BridgeStan`, which was at Version 1.0.2 at the time this research was performed. The version of `Stan` that we used was `Stanc3` Version 2.31.0 (Unix). Thus we used a total of 40 test problems for our empirical assessment.

For each test problem, a total of 10 replicate experiments were performed and standard errors were computed. A sampling method was defined as being *significantly better* for approximation of a given target, compared to all other methods considered, if had lower mean KSD *and* the standard error bar did not overlap with the standard error bar of any other method. Table 1 in the main text summarises the performance of S$\Pi$IS-MALA, fixing the number of samples to be $n = 3 \times 10^3$. In this appendix, full empirical results are provided.

25

For sampling from MALA, we used the adaptive algorithm described in Appendix D.3 with a final epoch of length $n_{\max} = 10^5$. Then, whenever a set of $n \ll n_{\max}$ consecutive samples from MALA are required for our experimental assessment, these were obtained by selecting at random a consecutive sequence of length $n$ from the total chain of length $10^5$. This ensures that the performance of unprocessed MALA that we report is not negatively affected by burn-in, in so far as is practical to control.

Full results are presented in Figure S4. These results broadly support the interpretation that SΠIS-MALA usually outperforms SIS-MALA, or otherwise both methods provide a similar level of performance, for the sufficiently large sample sizes $n$ considered. The sample size threshold at which SΠIS-MALA outperforms SIS-MALA appears to be dimension-dependent. A notable exception is panel 29 of Figure S4, a $d = 10$ dimensional task for which SΠIS-MALA provided a substantially worse approximation in KSD for the range of values of $n$ considered.

Figure S4: Benchmarking on `PosteriorDB`. Here we compared raw output from MALA (dotted lines) with the post-processed output provided by the default Stein importance sampling method of Liu and Lee (2017) (SIS-MALA; solid lines) and the proposed Stein Π-Importance Sampling method (SΠIS-MALA; dashed lines). The Langevin (purple) and KGM3–Stein kernels (blue) were used for SIS-MALA and SΠIS-MALA and the associated KSDs are reported as the number $n$ of iterations of MALA is varied. Ten replicates were computed and standard errors were plotted. The name of each model is shown in the title of the corresponding panel, and the dimension $d$ of the parameter vector is given in parentheses. [Langevin–Stein kernel: ········MALA, ——SIS-MALA, - - - - SΠIS-MALA. KGM3–Stein kernel: ········MALA, ——SIS-MALA, - - - - SΠIS-MALA.]



(S4.1)

(S4.2)

(S4.3)

(S4.4)

mesquite − logmesquite_logvolume(3D)

(S4.5)

arma − arma11(4D)

(S4.6)

earnings − logearn_logheight_male(4D)

(S4.7)

garch − garch11(4D)

(S4.8)

kidiq − kidscore_momhsiq(4D)

(S4.9)

earnings − logearn_interaction_z(5D)

(S4.10)

759

760

761

27

kidiq − kidscore_interaction(5D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.11)

kidiq_with_mom_work − kidscore_interaction_c(5D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.12)

kidiq_with_mom_work − kidscore_interaction_c2(5D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.13)

kidiq_with_mom_work − kidscore_interaction_z(5D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.14)

kidiq_with_mom_work − kidscore_mom_work(5D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.15)

low_dim_gauss_mix − low_dim_gauss_mix(5D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.16)

mesquite − logmesquite_logva(5D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.17)

hmm_example − hmm_example(6D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.18)

sblrc − blr(6D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.19)

sblri − blr(6D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.20)

arK − arK(7D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.21)

mesquite − logmesquite_logvash(7D)

$\mathbb{E}[\text{KSD}]$

$n$

(S4.22)

bball_drive_event_0 − hmm_drive_0(8D)

$\mathbb{E}[\mathrm{KSD}]$

$n$

(S4.23)

bball_drive_event_1 − hmm_drive_1(8D)

$\mathbb{E}[\mathrm{KSD}]$

$n$

(S4.24)

hudson_lynx_hare − lotka_volterra(8D)

$\mathbb{E}[\mathrm{KSD}]$

$n$

(S4.25)

mesquite − logmesquite(8D)

$\mathbb{E}[\mathrm{KSD}]$

$n$

(S4.26)

mesquite − logmesquite_logvas(8D)

$\mathbb{E}[\mathrm{KSD}]$

$n$

(S4.27)

mesquite − mesquite(8D)

$\mathbb{E}[\mathrm{KSD}]$

$n$

(S4.28)

eight_schools − eight_schools_centered($10D$)

(S4.29)

eight_schools − eight_schools_noncentered($10D$)

(S4.30)

nes1972 − nes($10D$)

(S4.31)

nes1976 − nes($10D$)

(S4.32)

nes1980 − nes($10D$)

(S4.33)

nes1984 − nes($10D$)

(S4.34)

771

772

773

31

nes1988 − nes(10*D*)

(S4.35)



nes1992 − nes(10*D*)

(S4.36)



nes1996 − nes(10*D*)

(S4.37)



nes2000 − nes(10*D*)

(S4.38)



diamonds − diamonds(26*D*)

(S4.39)



mcycle_gp − accel_gp(66*D*)

(S4.40)

## D.6 Stein $\Pi$-Thinning for `PosteriorDB`

The results presented in the main text concerned $n = 3 \times 10^3$ samples from MALA, which is near the limit at which the optimal weights $w^\star$ can be computed in a few seconds on a laptop PC. For larger values of $n$, sparse approximation methods are likely to required. In the main text we presented *Stein* $\Pi$-*Thinning*, which employs a greedy optimisation perspective to obtain a sparse approximation to the optimal weights at cost $O(m^2 n)$, where $m$ are the number of greedy iterations performed. Explicit and verifiable conditions for the strong consistency of the resulting S$\Pi$T-MALA algorithm were established in Section 3.3. The purpose of this appendix is to empirically explore the convergence of S$\Pi$T-MALA using the `PosteriorDB` test bed.

32

In the experiments we report the number of MALA samples was fixed to $n = 10^3$ and the number of greedy iterations was varied from $m = 1$ to $m = 10^3$. The results, in Figure S5, indicate that for most models in `PosteriorDB` the minimum value of KSD is approximately reached when $m$ is anywhere from $\frac{n}{10}$ to $\frac{n}{2}$, representing a modest but practically significant reduction in computational cost compared to SΠIS-MALA. This agrees with the qualitative findings reported in the original Stein thinning paper of Riabiz et al. (2022).

Figure S5: Benchmarking on `PosteriorDB`. Here we investigate the convergence of the sparse approximation provided by the proposed Stein Π-Thinning method (SΠT-MALA). The Langevin (purple) and KGM3–Stein kernels (blue) were used for SΠT-MALA and the associated KSDs are reported as the number $m$ of iterations of Stein thinning is varied. Ten replicates were computed and standard errors were plotted. The name of each model is shown in the title of the corresponding panel, and the dimension $d$ of the parameter vector is given in parentheses. [Langevin–Stein kernel: ——— SΠT-MALA. KGM3–Stein kernel: ——— SΠT-MALA.]



(S5.1)

(S5.2)

(S5.3)

(S5.4)

33

mesquite − logmesquite_logvolume(3D)

$\mathbb{E}[\text{KSD}]$

$m$

(S5.5)

arma − arma11(4D)

$\mathbb{E}[\text{KSD}]$

$m$

(S5.6)

earnings − logearn_logheight_male(4D)

$\mathbb{E}[\text{KSD}]$

$m$

(S5.7)

garch − garch11(4D)

$\mathbb{E}[\text{KSD}]$

$m$

(S5.8)

kidiq − kidscore_momhsiq(4D)

$\mathbb{E}[\text{KSD}]$

$m$

(S5.9)

earnings − logearn_interaction_z(5D)

$\mathbb{E}[\text{KSD}]$

$m$

(S5.10)

kidiq − kidscore_interaction(5D)

(S5.11)

kidiq_with_mom_work − kidscore_interaction_c(5D)

(S5.12)

kidiq_with_mom_work − kidscore_interaction_c2(5D)

(S5.13)

kidiq_with_mom_work − kidscore_interaction_z(5D)

(S5.14)

kidiq_with_mom_work − kidscore_mom_work(5D)

(S5.15)

low_dim_gauss_mix − low_dim_gauss_mix(5D)

(S5.16)

mesquite − logmesquite_logva($5D$)

(S5.17)

hmm_example − hmm_example($6D$)

(S5.18)

sblrc − blr($6D$)

(S5.19)

sblri − blr($6D$)

(S5.20)

arK − arK($7D$)

(S5.21)

mesquite − logmesquite_logvash($7D$)

(S5.22)

800

801

802

36

bball_drive_event_0 − hmm_drive_0(8D)

$\mathbb{E}[\mathrm{KSD}]$

$m$

(S5.23)

bball_drive_event_1 − hmm_drive_1(8D)

$\mathbb{E}[\mathrm{KSD}]$

$m$

(S5.24)

hudson_lynx_hare − lotka_volterra(8D)

$\mathbb{E}[\mathrm{KSD}]$

$m$

(S5.25)

mesquite − logmesquite(8D)

$\mathbb{E}[\mathrm{KSD}]$

$m$

(S5.26)

mesquite − logmesquite_logvas(8D)

$\mathbb{E}[\mathrm{KSD}]$

$m$

(S5.27)

mesquite − mesquite(8D)

$\mathbb{E}[\mathrm{KSD}]$

$m$

(S5.28)

803

804

805

37

eight_schools − eight_schools_centered(10D)

(S5.29)

eight_schools − eight_schools_noncentered(10D)

(S5.30)

nes1972 − nes(10D)

(S5.31)

nes1976 − nes(10D)

(S5.32)

nes1980 − nes(10D)

(S5.33)

nes1984 − nes(10D)

(S5.34)

806

807

808

nes1988 − nes(10D)

(S5.35)



nes1992 − nes(10D)

(S5.36)



nes1996 − nes(10D)

(S5.37)



nes2000 − nes(10D)

(S5.38)



diamonds − diamonds(26D)

(S5.39)



mcycle_gp − accel_gp(66D)

(S5.40)

## D.7 Performance of Stein Discrepancies

The properties of Stein discrepancies was out of scope for this work. Nonetheless, there is much interest in better understanding the properties of KSDs, and in this appendix the performance of SΠIS-MALA in terms of 1-Wasserstein divergence is reported. This was made possible since `PosteriorDB` supplies a set of posterior samples obtained from a long run of Hamiltonian Monte Carlo (the No-U-Turn sampler in `Stan`) which we treat as a gold standard.

Full results are presented in S6. Broadly speaking, for most models the minimisation of KSD seems to be associated with minimisation of 1-Wasserstein distance, however there are some models for

which minimisation of KSD is loosely, if at all, related to minimisation of 1-Wasserstein divergence. In these cases, we attribute this performance to the *blindness to mixing proportions* phenomena, described in Wenliang and Kanagawa (2021); Koehler et al. (2022); Liu et al. (2023). Convergence in 1-Wasserstein is equivalent to weak convergence plus convergence of the first moment, so the KGM–Stein kernels of order $s \geq 1$ control convergence in 1-Wasserstein. In Section 2.3 we proved that SΠIS-MALA is strongly consistent in KSD for the KGM–Stein kernel in the case $s = 1$, so we can expect strong consistency in 1-Wasserstein divergence for SΠIS-MALA in this case as well. It is interesting to observe that better 1-Wasserstein quantisations tend to be provided by SΠIS-MALA compared to SIS-MALA when either the Langevin–Stein or KGM–Stein kernel are used.

The development of improved Stein discrepancies is an active area of research, and we emphasise that the methodology developed in this work can be applied to *any* KSDs, including potentially KSDs with better or more direct control over standard notions of convergence (such as 1-Wasserstein) that in the future may be developed.

Figure S6: Performance of Stein discrepancies on `PosteriorDB`. Here we compared raw output from MALA (dotted lines) with the post-processed output provided by the default Stein importance sampling method of Liu and Lee (2017) (SIS-MALA; solid lines) and the proposed Stein Π-Importance Sampling method (SΠIS-MALA; dashed lines). The Langevin (purple) and KGM3–Stein kernels (blue) were used for SIS-MALA and SΠIS-MALA, and the 1-Wasserstein divergence is reported as the number $n$ of iterations of MALA is varied. Ten replicates were computed and standard errors were plotted. The name of each model is shown in the title of the corresponding panel, and the dimension $d$ of the parameter vector is given in parentheses. [Legend: ⋯⋯⋯Raw MALA. Langevin–Stein kernel: ——— SIS-MALA, - - - - SΠIS-MALA. KGM3–Stein kernel: ——— SIS-MALA, - - - - SΠIS-MALA.]
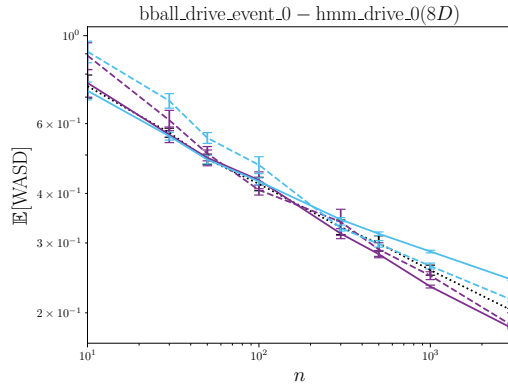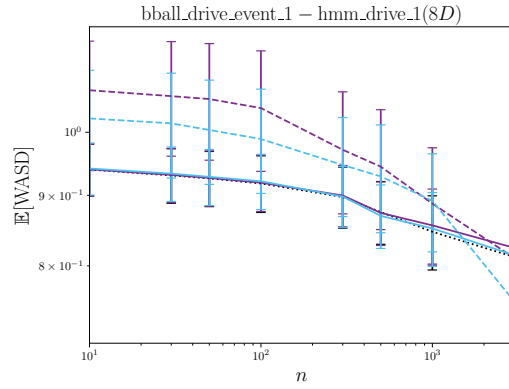


(S6.1)

(S6.2)

(S6.3)

(S6.4)

mesquite − logmesquite_logvolume($3D$)

(S6.5)

arma − arma11($4D$)

(S6.6)

earnings − logearn_logheight_male($4D$)

(S6.7)

garch − garch11($4D$)

(S6.8)

kidiq − kidscore_momhsiq($4D$)

(S6.9)

earnings − logearn_interaction_z($5D$)

(S6.10)

835

836

837

41

kidiq − kidscore_interaction(5D)

(S6.11)

kidiq_with_mom_work − kidscore_interaction_c(5D)

(S6.12)

kidiq_with_mom_work − kidscore_interaction_c2(5D)

(S6.13)

kidiq_with_mom_work − kidscore_interaction_z(5D)

(S6.14)

kidiq_with_mom_work − kidscore_mom_work(5D)

(S6.15)

low_dim_gauss_mix − low_dim_gauss_mix(5D)

(S6.16)

838

839

840

mesquite − logmesquite_logva(5D)

(S6.17)

hmm_example − hmm_example(6D)

(S6.18)

sblrc − blr(6D)

(S6.19)

sblri − blr(6D)

(S6.20)

arK − arK(7D)

(S6.21)

mesquite − logmesquite_logvash(7D)

(S6.22)

## bball_drive_event_0 − hmm_drive_0(8D)



(S6.23)

## bball_drive_event_1 − hmm_drive_1(8D)



(S6.24)

## hudson_lynx_hare − lotka_volterra(8D)



(S6.25)

## mesquite − logmesquite(8D)



(S6.26)

## mesquite − logmesquite_logvas(8D)



(S6.27)

## mesquite − mesquite(8D)



(S6.28)

eight_schools − eight_schools_centered(10D)

(S6.29)

eight_schools − eight_schools_noncentered(10D)

(S6.30)

nes1972 − nes(10D)

(S6.31)

nes1976 − nes(10D)

(S6.32)

nes1980 − nes(10D)

(S6.33)

nes1984 − nes(10D)

(S6.34)

847

848

849

45

nes1988 − nes(10D)

(S6.35)

nes1992 − nes(10D)

(S6.36)

nes1996 − nes(10D)

(S6.37)

nes2000 − nes(10D)

(S6.38)

diamonds − diamonds(26D)

(S6.39)

mcycle_gp − accel_gp(66D)
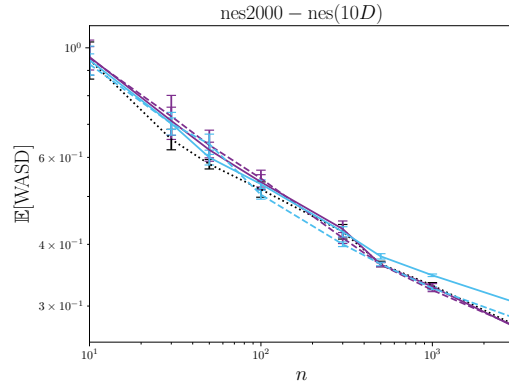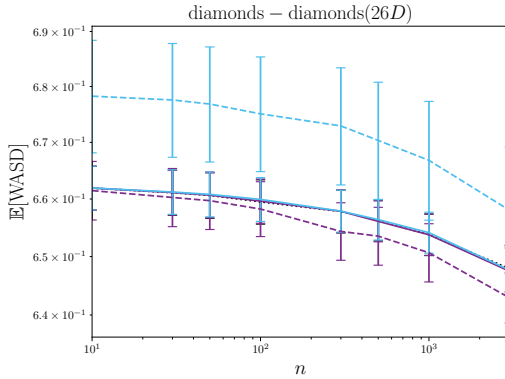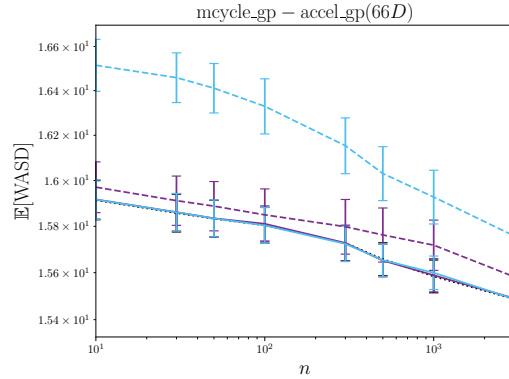
(S6.40)

## D.8 Investigation for a Skewed Target

This final appendix contrasts the 1-Wasserstein optimal sampling distribution $\Pi_1$ (c.f. Section 2.1), with the choice of $\Pi$ that we recommended in (8). In particular, we focus on the KGM3–Stein kernel under a heavily skewed $P$, for which $\Pi_1$ and $\Pi$ can be markedly different.

For this investigation a bivariate skew-normal target was constructed, where the density is given by $p(x_1, x_2) = 4\phi(x_1)\Phi(6x_1)\phi(x_2)\Phi(-3x_2)$, with $\phi$ and $\Phi$ respectively denoting the density and distribution functions of a standard Gaussian. The density $p$ of $P$, together with the marginal densities of $\Pi_1$ and $\Pi$, are plotted in Figure S7. It can be seen that, while both $\Pi_1$ and $\Pi$ are over-dispersed
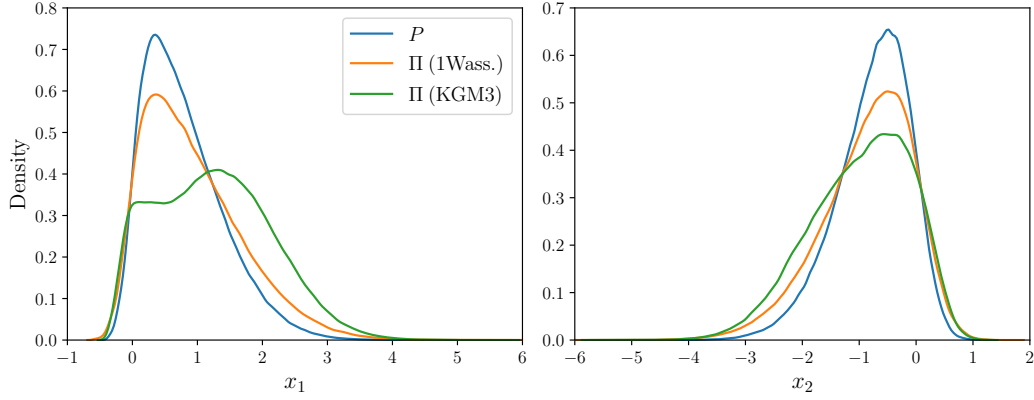
46

Figure S7: Comparing the proposed distribution $\Pi$ (KGM3; based on the KGM3–Stein kernel) to $\Pi_1$ (1Wass.; the optimal choice for 1-Wasserstein quantisation from Section 2.1) for a bivariate skew-normal target ($d = 2$). The marginal density functions of each distribution were approximated using $10^6$ samples from MCMC.
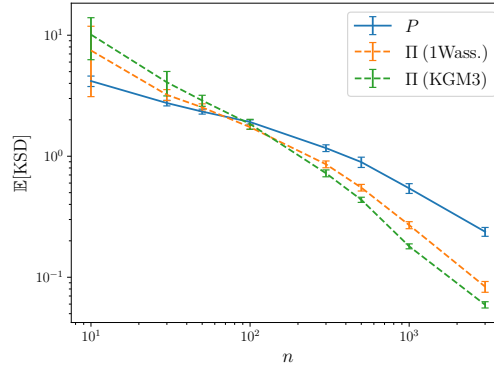


Figure S8: Comparing the performance of using the proposed distribution $\Pi$ (KGM3; based on the KGM3–Stein kernel) to $\Pi_1$ (1Wass.; the optimal choice for 1-Wasserstein quantisation from Section 2.1) for a bivariate skew-normal target ($d = 2$). The mean kernel Stein discrepancy (KSD) for Stein $\Pi$-Importance Sampling was estimated; in each case, the KSD based on the KGM3–Stein kernel was computed. Solid lines indicate the baseline case of sampling from $P$, while dashed lines indicate sampling from $\Pi$. (The experiment was repeated 10 times and standard error bars are plotted.)

with respect to $P$, our recommended $\Pi$ assigns proportionally more mass to the tail that is positively skewed.

The performance of Stein $\Pi$-Importance Sampling based on $\Pi_1$ and $\Pi$ is compared in Figure S8. Though both choices lead to an improvement relative to Stein importance sampling algorithm with $\Pi = P$, the use of $\Pi$ leads to a significant further reduction (on average) in KSD compared to $\Pi_1$. Based on our investigations, this finding seems general; the use of $\Pi_1$ does not realise the full potential of Stein $\Pi$-Imporance sampling when the target is skewed.