
Stein Discrepancies

PART III ESSAY

ABSTRACT

Exact Bayesian inference is rarely feasible in practice, necessitating the use of approximate inference methods. One such method is Markov Chain Monte Carlo (MCMC), which outputs a discrete measure with the aim of approximating the posterior distribution. Increasingly, biased MCMC is being used to speed up mixing times at the expense of asymptotic exactness. It is desirable in this context to be able to quantify how well our discrete measure approximates the target posterior.

In this essay, we introduce the *Stein discrepancy*, which quantifies distance between a discrete measure and a target measure, whose density we need only know up to a normalising constant. We show that the Stein discrepancy can be *kernelised* to yield an analytic solution, and that a specific choice of kernel, namely the *IMQ kernel*, yields a kernelised Stein discrepancy (KSD) which can detect when a discrete measure converges to the target, and when it doesn't.

We also use the IMQ KSD to derive *Stein variational gradient descent* (SVGD), an expressive variational inference algorithm which is asymptotically exact. Finally, we highlight scalability issues present when calculating the IMQ KSD and using SVGD, and propose solutions for these.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Kernelised Stein discrepancies	2
1.3	Outline of chapters	3
2	Assessing Convergence with Kernelised Stein Discrepancies	4
2.1	Detecting convergence: upper bounding the KSD	4
2.2	Detecting non-convergence: lower bounding the KSD	5
2.2.1	Tightness lower bound	6
2.2.2	IMQ KSD detects non-convergence	10
3	Stein Variational Gradient Descent	14
3.1	SVGD algorithm	14
3.2	Convergence of SVGD	15
3.2.1	Convergence to large sample limit	15
3.2.2	Convergence of large sample limit	16
4	Scalable Stein Discrepancies	19
4.1	Random feature Stein discrepancies	19
4.2	Stochastic Stein discrepancies	21
4.2.1	Stochastic SVGD	22
5	Experiments	24
5.1	Measuring sample quality	24
5.1.1	Comparison to the Wasserstein-1 distance in one dimension	24
5.1.2	Hyperparameter selection in Stochastic Gradient Langevin Dynamics	25
5.2	Particle-based variational inference	26
5.2.1	Convergence of SVGD	27
6	Conclusion	28
	Bibliography	29

Chapter 1

Introduction

1.1 Motivation

Many problems in statistics and machine learning require the approximation of a target probability measure P with another Q . Crucial to approximating P accurately is the ability to obtain a measure of “distance” between P and Q . A popular approach to measuring such distance is to consider an **integral probability metric** (IPM):

$$d_{\mathcal{H}}(Q, P) := \sup_{h \in \mathcal{H}} |\mathbb{E}_Q[h(\mathbf{X})] - \mathbb{E}_P[h(\mathbf{X})]|,$$

where \mathcal{H} is some class of test functions.

Often, our ultimate goal is to approximate an integral under P of the form $\mathbb{E}_P[f(\mathbf{X})]$. With this in mind, given a sequence of approximating measures $(Q_n)_{n \geq 1}$, a natural requirement for \mathcal{H} is that, as $n \rightarrow \infty$, $d_{\mathcal{H}}(Q_n, P) \rightarrow 0$ if and only if Q_n converges weakly to P , which we denote by $Q_n \Rightarrow P$. When this is the case, we say that $d_{\mathcal{H}}(\cdot, \cdot)$ **metrises** weak convergence.

One choice of \mathcal{H} satisfying this requirement is $\mathcal{H} = \mathcal{W}_{\|\cdot\|_2} := \{h : \mathbb{R}^d \rightarrow \mathbb{R} \mid \sup_{x \neq y} \frac{|h(x) - h(y)|}{\|x - y\|_2} \leq 1\}$, which corresponds to the **Wasserstein-1 metric**. Another choice is $\mathcal{H} = \text{BL}_{\|\cdot\|_2} := \{h : \mathbb{R}^d \rightarrow \mathbb{R} : \sup_{x \in \mathbb{R}^d} |h(x)| + \sup_{x \neq y} \frac{|h(x) - h(y)|}{\|x - y\|_2} \leq 1\}$, which corresponds to the **bounded Lipschitz metric**.

These two metrics are known to metrise weak convergence. However, they both exhibit a common feature which inhibits calculation in many settings: they require the ability to integrate over both P and Q . In many cases, we are unable to integrate over the target P , and so it is of significant practical interest to find a choice of \mathcal{H} which metrises weak convergence while only requiring integration over the approximating distribution Q .

Although not immediately obvious, such a choice of \mathcal{H} does exist, as outlined by [Gorham and Mackey \[2015\]](#). Their discovery of such an \mathcal{H} was based on **Stein’s method** [[Stein, 1972](#)], and proceeds as follows, for some class of sufficiently well-behaved target measures $P \in \mathcal{P}$ with differentiable densities p :

1. Identify an operator \mathcal{T}_P which maps a set \mathcal{G} of vector-valued functions $\mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ to real-valued functions $\mathcal{T}_P \mathbf{g}$ such that

$$\mathbb{E}_P[(\mathcal{T}_P \mathbf{g})(\mathbf{X})] = 0 \quad \text{for all } \mathbf{g} \in \mathcal{G}. \tag{1.1}$$

Such an operator \mathcal{T}_P is called a **Stein operator**; such a set \mathcal{G} is called a *Stein set*. One can then define the **Stein discrepancy** as

$$\mathcal{S}(Q, \mathcal{T}_P, \mathcal{G}) := \sup_{\mathbf{g} \in \mathcal{G}} |\mathbb{E}_Q[(\mathcal{T}_P \mathbf{g})(\mathbf{X})]| = d_{\mathcal{T}_P \mathcal{G}}(Q, P).$$

2. Lower bound the Stein discrepancy by an IPM which is known to dominate weak convergence, so that $\mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}) \rightarrow 0$ implies $Q_n \Rightarrow P$ as $n \rightarrow \infty$. This allows us to detect when Q_n is not converging to P .
3. Upper bound the Stein discrepancy by an IPM known to be dominated by weak convergence, ensuring that $\mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}) \rightarrow 0$ whenever $Q_n \Rightarrow P$ as $n \rightarrow \infty$, allowing us to detect when Q_n is converging to P .

The most popular choice for the Stein operator is the **Langevin Stein operator**

$$(\mathcal{T}_P \mathbf{g})(\mathbf{x}) := \frac{1}{p(\mathbf{x})} \nabla_{\mathbf{x}} \cdot (p(\mathbf{x}) \mathbf{g}(\mathbf{x})) = \mathbf{g}(\mathbf{x})^\top \mathbf{s}_p(\mathbf{x}) + \nabla_{\mathbf{x}} \cdot \mathbf{g}(\mathbf{x}),$$

where $\mathbf{s}_p(\mathbf{x}) := \nabla_{\mathbf{x}} \log p(\mathbf{x})$ is the score function. To see that this satisfies (1.1), we integrate by parts:

$$\mathbb{E}_P[\mathbf{g}(\mathbf{X})^\top \mathbf{s}_p(\mathbf{X})] = \int_{\mathbb{R}^d} \mathbf{g}(\mathbf{x})^\top \nabla_{\mathbf{x}} p(\mathbf{x}) d\mathbf{x} = - \int_{\mathbb{R}^d} \nabla_{\mathbf{x}} \cdot \mathbf{g}(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = -\mathbb{E}_P[\nabla_{\mathbf{x}} \cdot \mathbf{g}(\mathbf{x})],$$

assuming appropriate boundary conditions which will hold for any reasonable class of test functions \mathcal{G} .

An additional appealing property of the Langevin Stein operator is that it depends on P only through the score \mathbf{s}_p , and so it is only required for us to know the density p up to a multiplicative constant. This is particularly useful in the context of Bayesian inference, where P is a posterior distribution and, for example, Q_n is the discrete measure outputted by some Markov Chain Monte Carlo (MCMC) algorithm.

1.2 Kernelised Stein discrepancies

Computing the Stein discrepancy involves taking the supremum over a complicated set of functions. While [Gorham and Mackey \[2015\]](#) proposed methods for solving such an optimisation problem, an alternative, proposed by [Chwialkowski et al. \[2016\]](#) and [Gorham and Mackey \[2017\]](#), is to *kernelise* the Stein discrepancy, which allows us to calculate the supremum analytically.

Let \mathcal{H}_k be a reproducing kernel Hilbert space (RKHS), with reproducing kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ and inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}_k}$. For $d \in \mathbb{N}$, we can associate with the product space \mathcal{H}_k^d the inner product $\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{H}_k^d} := \sum_{i=1}^d \langle u_i, v_i \rangle_{\mathcal{H}_k}$ for $\mathbf{u}, \mathbf{v} \in \mathcal{H}_k^d$. We then define the **kernel Stein set** to be the unit ball with respect to the norm induced by this inner product:

$$\mathcal{G}_k := \left\{ \mathbf{g} \in \mathcal{H}_k^d : \|\mathbf{g}\|_{\mathcal{H}_k^d} \leq 1 \right\}. \quad (1.2)$$

Since the RKHS \mathcal{H}_k is the completion of the linear span $\{\sum_{i=1}^n \alpha_i k(\cdot, \mathbf{x}_i) : n \in \mathbb{N}, \alpha_i \in \mathbb{R}, \mathbf{x}_i \in \mathbb{R}^d\}$, we have $\mathbf{g} \in \mathcal{G}_k \iff -\mathbf{g} \in \mathcal{G}_k$. Thus, since $\mathcal{T}_P(-\mathbf{g}) = -\mathcal{T}_P(\mathbf{g})$, we can remove the modulus in the Stein discrepancy. Furthermore, the Langevin Stein operator can be decomposed as $\mathcal{T}_P \mathbf{g}(\cdot) = \sum_{i=1}^d \mathcal{T}_P^{(i)} g_i(\cdot)$, where $\mathcal{T}_P^{(i)} h(\mathbf{x}) = h(\mathbf{x}) \nabla_{x_i} \log p(\mathbf{x}) + \nabla_{x_i} h(\mathbf{x})$, so

$$\begin{aligned} \mathcal{S}(Q, \mathcal{T}_P, \mathcal{G}_k) &= \sup_{\mathbf{g} \in \mathcal{G}_k} \sum_{i=1}^d \mathbb{E}_Q[(\mathcal{T}_P^{(i)} g_i)(\mathbf{X})] = \sup_{\mathbf{g} \in \mathcal{G}_k} \sum_{i=1}^d \mathbb{E}_Q[\mathcal{T}_P^{(i)} \langle g_i, k(\mathbf{X}, \cdot) \rangle_{\mathcal{H}_k}] \\ &= \sup_{\mathbf{g} \in \mathcal{G}_k} \sum_{i=1}^d \left\langle g_i, \mathbb{E}_Q \left[\mathcal{T}_P^{(i)} k(\mathbf{X}, \cdot) \right] \right\rangle_{\mathcal{H}_k} = \sup_{\mathbf{g} \in \mathcal{G}_k} \langle \mathbf{g}, \mathbb{E}_Q \xi_p(\mathbf{X}, \cdot) \rangle_{\mathcal{H}_k^d}, \end{aligned}$$

where

$$\xi_p(\mathbf{x}, \mathbf{y}) := \mathbf{s}_p(\mathbf{x}) k(\mathbf{x}, \mathbf{y}) + \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y}). \quad (1.3)$$

By the Cauchy-Schwarz inequality on the Cartesian product Hilbert space \mathcal{H}_k^d , we have

$$\langle \mathbf{g}, \mathbb{E}_Q \xi_p(\mathbf{X}, \cdot) \rangle_{\mathcal{H}_k^d} \leq \|\mathbf{g}\|_{\mathcal{H}_k^d} \|\mathbb{E}_Q \xi_p(\mathbf{X}, \cdot)\|_{\mathcal{H}_k^d},$$

and so for $\mathbf{g} \in \mathcal{G}_k$,

$$\langle \mathbf{g}, \mathbb{E}_Q \boldsymbol{\xi}_p(\mathbf{X}, \cdot) \rangle_{\mathcal{H}_k^d} \leq \|\mathbb{E}_Q \boldsymbol{\xi}_p(\mathbf{X}, \cdot)\|_{\mathcal{H}_k^d}. \quad (1.4)$$

Therefore if $\mathbf{g}(\cdot) := \mathbb{E}_Q \boldsymbol{\xi}_p(\mathbf{X}, \cdot) / \|\mathbb{E}_Q \boldsymbol{\xi}_p(\mathbf{X}, \cdot)\|_{\mathcal{H}_k^d}$ is in the kernel Stein set \mathcal{G}_k , then this will achieve the upper bound in (1.4) and thus attain the supremum defined in the Stein discrepancy. This choice of \mathbf{g} does indeed belong to \mathcal{G}_k : it has unit norm, and $\xi_{p,i}(\mathbf{x}, \cdot) = k(\mathbf{x}, \cdot) \nabla_{x_i} \log p(\mathbf{x}) + \nabla_{x_i} k(\mathbf{x}, \cdot)$ is a linear combination of $k(\mathbf{x}, \cdot)$, which belongs to \mathcal{H}_k by definition, and $\nabla_{x_i} k(\mathbf{x}, \cdot)$, which is also in \mathcal{H}_k , as we now justify. Take $\frac{1}{h_n}(k(\mathbf{x} + h_n \mathbf{e}_i, \cdot) - k(\mathbf{x}, \cdot)) \in \mathcal{H}_k$ for all n , where \mathbf{e}_i is the i^{th} standard basis vector and $(h_n)_{n \geq 1}$ is such that $h_n \rightarrow 0$ as $n \rightarrow \infty$. Taking $n \rightarrow \infty$ gives $\nabla_{x_i} k(\mathbf{x}, \cdot) \in \mathcal{H}_k$. Then we have

$$\mathcal{S}(Q, \mathcal{T}_P, \mathcal{G}_k) = \left\langle \frac{\mathbb{E}_Q \boldsymbol{\xi}_p(\mathbf{X}, \cdot)}{\|\mathbb{E}_Q \boldsymbol{\xi}_p(\mathbf{X}, \cdot)\|_{\mathcal{H}_k^d}}, \mathbb{E}_Q \boldsymbol{\xi}_p(\mathbf{X}, \cdot) \right\rangle_{\mathcal{H}_k^d} = \|\mathbb{E}_Q \boldsymbol{\xi}_p(\mathbf{X}, \cdot)\|_{\mathcal{H}_k^d}.$$

Observe further that

$$\begin{aligned} \langle \xi_{p,i}(\mathbf{x}, \cdot), \xi_{p,i}(\mathbf{y}, \cdot) \rangle_{\mathcal{H}_k} &= \langle s_{p,i}(\mathbf{x})k(\mathbf{x}, \cdot) + \nabla_{x_i} k(\mathbf{x}, \cdot), s_{p,i}(\mathbf{y})k(\mathbf{y}, \cdot) + \nabla_{y_i} k(\mathbf{y}, \cdot) \rangle_{\mathcal{H}_k} \\ &= s_{p,i}(\mathbf{x})s_{p,i}(\mathbf{y})k(\mathbf{x}, \mathbf{y}) + s_{p,i}(\mathbf{x})\nabla_{y_i} k(\mathbf{x}, \mathbf{y}) + s_{p,i}(\mathbf{y})\nabla_{x_i} k(\mathbf{x}, \mathbf{y}) + \nabla_{x_i} \nabla_{y_i} k(\mathbf{x}, \mathbf{y}), \end{aligned}$$

so, defining the **Stein reproducing kernel**

$$\begin{aligned} k_p(\mathbf{x}, \mathbf{y}) &:= \langle \boldsymbol{\xi}_p(\mathbf{x}, \cdot), \boldsymbol{\xi}_p(\mathbf{y}, \cdot) \rangle_{\mathcal{H}_k^d} \\ &= s_p(\mathbf{x})^\top s_p(\mathbf{y})k(\mathbf{x}, \mathbf{y}) + s_p(\mathbf{x})^\top \nabla_{\mathbf{y}} k(\mathbf{x}, \mathbf{y}) + s_p(\mathbf{y})^\top \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y}) + \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{y}} k(\mathbf{x}, \mathbf{y}), \end{aligned}$$

we get the squared **kernelised Stein discrepancy** (KSD)

$$\begin{aligned} \mathcal{S}^2(Q, \mathcal{T}_P, \mathcal{G}_k) &= \langle \mathbb{E}_{\mathbf{X} \sim Q} [\boldsymbol{\xi}_p(\mathbf{X}, \cdot)], \mathbb{E}_{\mathbf{X}' \sim Q} [\boldsymbol{\xi}_p(\mathbf{X}', \cdot)] \rangle_{\mathcal{H}_k^d} \\ &= \mathbb{E}_{\mathbf{X}, \mathbf{X}' \sim Q} [k_p(\mathbf{X}, \mathbf{X}')], \end{aligned}$$

where $\mathbf{X}, \mathbf{X}' \stackrel{\text{i.i.d.}}{\sim} Q$. When $Q(\cdot) = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}(\cdot)$, we have

$$\mathcal{S}^2(Q, \mathcal{T}_P, \mathcal{G}_k) = \frac{1}{n^2} \sum_{i,j=1}^n k_p(\mathbf{x}_i, \mathbf{x}_j).$$

1.3 Outline of chapters

In Chapter 2, we prove that a particular choice of kernel, namely the IMQ kernel, with suitable parameters, yields a KSD that metrises weak convergence for a large class \mathcal{P} of target measures, following the approach of Gorham and Mackey [2017].

In Chapter 3, we demonstrate how the formulation of the KSD can be used to develop a flexible variational inference algorithm known as **Stein variational gradient descent** (SVGD) [Liu and Wang, 2016]. We will prove that it is asymptotically exact using results from Liu [2017] and Gorham et al. [2020], the latter of which receives greater focus in Chapter 4, where we address scalability issues associated with KSDs. In this chapter, we also discuss an approach by Huggins and Mackey [2018] to alleviating the cost of calculating the KSD, which grows quadratically in the number of samples.

In Chapter 5, we put this theory into practice, demonstrating the practicality of the various Stein discrepancies we have introduced, as well as the SVGD algorithm, drawing on experiments carried out in Gorham and Mackey [2017] and Gorham et al. [2020]. We also focus on the scalability issues presented in Chapter 4.

Chapter 2

Assessing Convergence with Kernelised Stein Discrepancies

The goal of this chapter is to find a kernel $k(\cdot, \cdot)$ whose corresponding KSD metrises weak convergence for a large class \mathcal{P} of target measures P with differentiable densities p . We will be interested in so-called *distantly dissipative* target distributions.

2.1 Definition (Distant dissipativity). *Let $r > 0$. We say that a distribution P is **distantly dissipative** if*

$$\kappa_0 := \liminf_{r \rightarrow \infty} \kappa(r) > 0, \quad (2.1)$$

where

$$\kappa(r) := \inf_{\mathbf{x}, \mathbf{y}: \|\mathbf{x} - \mathbf{y}\|_2 = r} \left\{ -2 \frac{\langle s_p(\mathbf{x}) - s_p(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle}{\|\mathbf{x} - \mathbf{y}\|_2^2} \right\}.$$

We let \mathcal{P} be the class of distantly dissipative distributions with Lipschitz score function s_p . As an example, \mathcal{P} includes all distributions which are strongly log-concave outside of a compact set, including Bayesian linear, logistic and Huber regression posteriors with Gaussian priors.

As mentioned above, in order to show that the KSD metrises weak convergence, we aim to both upper and lower bound it by IPMs known to metrise weak convergence. We first consider the easier – and more general – direction, where we obtain an upper bound on the KSD for a variety of kernels $k(\cdot, \cdot)$.

2.1 Detecting convergence: upper bounding the KSD

A KSD detects convergence if $\mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k) \rightarrow 0$ whenever $Q_n \Rightarrow P$. To aid our discussion, we define the following quantities. For any function $\mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$, set

$$M_0(\mathbf{g}) := \sup_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{g}(\mathbf{x})\|_2, \quad M_1(\mathbf{g}) := \sup_{\mathbf{x} \neq \mathbf{y}} \frac{\|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})\|_2}{\|\mathbf{x} - \mathbf{y}\|_2}, \quad M_2(\mathbf{g}) := \sup_{\mathbf{x} \neq \mathbf{y}} \frac{\|\nabla_{\mathbf{x}} \mathbf{g}(\mathbf{x}) - \nabla_{\mathbf{y}} \mathbf{g}(\mathbf{y})\|_{\text{op}}}{\|\mathbf{x} - \mathbf{y}\|_2}. \quad (2.2)$$

We first prove the following lemma.

2.2 Lemma. *Let $\mathbf{Z} \sim P$ and $\mathbf{X} \sim \mu$. Suppose that $\mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is such that $\max\{M_0(\mathbf{g}), M_1(\mathbf{g}), M_2(\mathbf{g})\} < \infty$. If s_p is Lipschitz with $\mathbb{E}_P[\|s_p(\mathbf{Z})\|_2^2] < \infty$, then*

$$|\mathbb{E}_{\mu}[(\mathcal{T}_P \mathbf{g})(\mathbf{X})]| \leq (M_0(\mathbf{g})M_1(s_p) + M_2(\mathbf{g})d)W_2(\mu, P) + \sqrt{2M_0(\mathbf{g})M_1(\mathbf{g})\mathbb{E}_P[\|s_p(\mathbf{Z})\|_2^2]W_2(\mu, P)},$$

where $W_p(\cdot, \cdot)$ denotes the Wasserstein- p distance.

Proof. By Jensen's inequality, we have $\mathbb{E}_P[\|s_p(\mathbf{Z})\|_2] \leq \sqrt{\mathbb{E}_P[\|s_p(\mathbf{Z})\|_2^2]} < \infty$, which is sufficient for us to be able to integrate by parts and obtain $\mathbb{E}_P[(\mathcal{T}_P \mathbf{g})(\mathbf{Z})] = 0$. Then, using the triangle inequality, Jensen's inequality and the Fenchel-Young inequality for dual norms,

$$\begin{aligned} |\mathbb{E}_\mu[(\mathcal{T}_P \mathbf{g})(\mathbf{X})]| &= |\mathbb{E}_P[(\mathcal{T}_P \mathbf{g})(\mathbf{Z})] - \mathbb{E}_\mu[(\mathcal{T}_P \mathbf{g})(\mathbf{X})]| \\ &= |\mathbb{E}[\langle s_p(\mathbf{Z}), \mathbf{g}(\mathbf{Z}) - \mathbf{g}(\mathbf{X}) \rangle + \langle s_p(\mathbf{Z}) - s_p(\mathbf{X}), \mathbf{g}(\mathbf{X}) \rangle + \nabla_{\mathbf{Z}} \cdot \mathbf{g}(\mathbf{Z}) - \nabla_{\mathbf{X}} \cdot \mathbf{g}(\mathbf{X})]| \\ &\leq \mathbb{E}[|\langle s_p(\mathbf{Z}), \mathbf{g}(\mathbf{Z}) - \mathbf{g}(\mathbf{X}) \rangle|] + (M_0(\mathbf{g})M_1(s_p) + M_2(\mathbf{g})d)\mathbb{E}[\|\mathbf{X} - \mathbf{Z}\|_2]. \end{aligned}$$

By the Cauchy-Schwarz inequality and the fact that $\min\{a, b\} \leq \sqrt{ab}$ for $a, b \geq 0$, the first term can be upper bounded as follows:

$$\begin{aligned} \mathbb{E}[|\langle s_p(\mathbf{Z}), \mathbf{g}(\mathbf{Z}) - \mathbf{g}(\mathbf{X}) \rangle|] &\leq \mathbb{E}[\min\{2M_0(\mathbf{g}), M_1(\mathbf{g})\}\|\mathbf{X} - \mathbf{Z}\|_2]\|s_p(\mathbf{Z})\|_2 \\ &\leq (2M_0(\mathbf{g})M_1(\mathbf{g}))^{1/2}\mathbb{E}\left[\|\mathbf{X} - \mathbf{Z}\|_2^{1/2}\|s_p(\mathbf{Z})\|_2\right] \\ &\leq \sqrt{2M_0(\mathbf{g})M_1(\mathbf{g})\mathbb{E}[\|\mathbf{X} - \mathbf{Z}\|_2]\mathbb{E}_P[\|s_p(\mathbf{Z})\|_2^2]}. \end{aligned}$$

Taking the infimum of these bounds over all couplings of $\mathbf{X} \sim \mu$ and $\mathbf{Z} \sim P$ yields the desired result. \square

2.3 Theorem (KSDs detect convergence). *Let k be a kernel such that $\nabla_{\mathbf{x}}^\ell \nabla_{\mathbf{y}}^\ell k(\mathbf{x}, \mathbf{y})$ is continuous and uniformly bounded for $\ell = 0, 1, 2$. If s_p is Lipschitz with $\mathbb{E}_P[\|s_p(\mathbf{Z})\|_2^2] < \infty$, then $\mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k) \rightarrow 0$ whenever $Q_n \Rightarrow P$.*

Proof. For any function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ for which the appropriate derivatives exist, define

$$D^\alpha f(\mathbf{x}) := \frac{\partial^{|\alpha|}}{(dx_1)^{\alpha_1} \dots (dx_d)^{\alpha_d}} f(\mathbf{x}),$$

where $|\alpha| = \sum_{j=1}^d \alpha_j$. Take any $\mathbf{g} \in \mathcal{G}_k$, and choose any $\alpha \in \mathbb{N}^d$ such that $|\alpha| \leq 2$. Then by Cauchy-Schwarz we have

$$\sup_{\mathbf{x} \in \mathbb{R}^d} |D^\alpha g_j(\mathbf{x})| = \sup_{\mathbf{x} \in \mathbb{R}^d} |D^\alpha \langle g_j, k(\mathbf{x}, \cdot) \rangle_{\mathcal{H}_k}| \leq \sup_{\mathbf{x} \in \mathbb{R}^d} \|g_j\|_{\mathcal{H}_k} \|D^\alpha k(\mathbf{x}, \cdot)\|_{\mathcal{H}_k} = \|g_j\|_{\mathcal{H}_k} \sup_{\mathbf{x} \in \mathbb{R}^d} (D_{\mathbf{x}}^\alpha D_{\mathbf{y}}^\alpha k(\mathbf{x}, \mathbf{x}))^{1/2}.$$

Since $\sum_{j=1}^d \|g_j\|_{\mathcal{H}_k}^2 \leq 1$ for all $\mathbf{g} \in \mathcal{G}_k$ and $D_{\mathbf{x}}^\alpha D_{\mathbf{y}}^\alpha k(\mathbf{x}, \mathbf{x})$ is uniformly bounded in \mathbf{x} for all $|\alpha| \leq 2$, the entries of $\mathbf{g}(\mathbf{x})$, $\nabla_{\mathbf{x}} \mathbf{g}(\mathbf{x})$ and $\nabla_{\mathbf{x}}^2 \mathbf{g}(\mathbf{x})$ are uniformly bounded in $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{g} \in \mathcal{G}_k$. So there exists $\lambda > 0$ such that $\sup_{\mathbf{g} \in \mathcal{G}_k} \max\{M_0(\mathbf{g}), M_1(\mathbf{g}), M_2(\mathbf{g})\} \leq \lambda < \infty$. The result then follows from Lemma 2.2 as

$$\mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k) \leq \lambda \left((M_1(s_p) + d)W_2(\mu, P) + \sqrt{2\mathbb{E}_P[\|s_p(\mathbf{Z})\|_2^2]W_2(\mu, P)} \right). \quad \square$$

This more general result shows in particular that the IMQ KSD detects convergence for targets $P \in \mathcal{P}$ with $\mathbb{E}_P[\|s_p(\mathbf{Z})\|_2^2] < \infty$.

2.2 Detecting non-convergence: lower bounding the KSD

A KSD detects non-convergence if $\mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k) \rightarrow 0$ implies $Q_n \Rightarrow P$. The proof of this direction is eminently more involved. While Gorham and Mackey [2017] showed that KSDs based on many different popular kernels detect non-convergence in the univariate case, many of these fail to replicate this property in higher dimensions. This failure lies in their inability to enforce *uniform tightness* (defined below).

Fortunately, the **IMQ kernel** $k(\mathbf{x}, \mathbf{y}) = (c^2 + \|\mathbf{x} - \mathbf{y}\|_2^2)^\beta$ with $c > 0$ and $\beta \in (-1, 0)$ has sufficiently heavy tails for its corresponding KSD to enforce such tightness, as we show in this section. We break the proof down into two parts: we first prove a lower bound for a general KSD in terms of the *tightness rate* (also defined below), and then show that the IMQ KSD automatically enforces tightness.

2.2.1 Tightness lower bound

As promised, we start with two definitions.

2.4 Definition (Uniform tightness). *A sequence of probability measures $(\mu_n)_{n \geq 1}$ is said to be **uniformly tight** if for every $\epsilon > 0$, there exists $R(\epsilon) < \infty$ such that*

$$\limsup_n \mu_n(\|\mathbf{X}\|_2 > R(\epsilon)) \leq \epsilon.$$

The uniform tightness condition essentially states that no mass in the sequence $(\mu_n)_{n \geq 1}$ escapes to infinity. When k decays more rapidly than s_p grows, the KSD ignores excess mass in the tails and hence can be driven to zero by a non-tight sequence of increasingly diffuse probability measures.

2.5 Definition (Tightness rate). *Let μ be a probability measure on \mathbb{R}^d and $\epsilon > 0$. The **tightness rate** is defined as*

$$R(\mu, \epsilon) := \inf\{r \geq 0 : \mu(\|\mathbf{X}\|_2 > r) \leq \epsilon\}. \quad (2.3)$$

The main theorem of this section, which lower bounds the KSD in terms of the tightness rate, is stated as follows.

2.6 Theorem (KSD tightness lower bound). *Suppose $P \in \mathcal{P}$ and let μ be a probability measure with tightness rate $R(\mu, \epsilon)$ defined in (2.3). Suppose further that $k(\mathbf{x}, \mathbf{y}) = \Psi(\mathbf{x} - \mathbf{y})$ with $\Psi \in C^2$ and $F(t) := \sup_{\|\omega\|_\infty \leq t} \hat{\Psi}(\omega)^{-1}$ finite for all $t > 0$, where $\hat{\Psi}$ denotes the generalised Fourier transform of Ψ . Then there exists a constant \mathcal{M}_P such that, for all $\rho, \epsilon, \delta > 0$,*

$$\begin{aligned} d_{BL_{\|\cdot\|_2}}(\mu, P) &\leq \rho\sqrt{d}(1 + M_1(s_p)\mathcal{M}_P) + \epsilon + \min\{\epsilon, 1\} \left(2 + \epsilon + \left(M_1(s_p)\rho\sqrt{d} + \frac{d\delta^{-1}\theta_{d-1}}{\theta_d}\right)\mathcal{M}_P\right) \\ &\quad + (2\pi)^{-d/4}V_d^{1/2}\mathcal{M}_P(R(\mu, \epsilon) + 2\delta)^{d/2}F\left(\frac{12d\log 2}{\pi}(c_{\rho, \delta} + M_1(s_p)\mathcal{M}_P)\epsilon^{-1}\right)^{1/2}\mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k), \end{aligned}$$

where $\theta_d := d \int_0^1 \exp\{-1/(1-r^2)\} r^{d-1} dr$ for $d > 0$, $\theta_0 := e^{-1}$, V_d is the volume of the unit Euclidean ball in dimension d , and

$$c_{\rho, \delta} := 1 + M_1(s_p)\mathcal{M}_P(1+d) + \left\{2 + (M_1(s_p)\rho + 1/\rho)\sqrt{d}\mathcal{M}_P\right\} \frac{d\delta^{-1}\theta_{d-1}}{\theta_d} + \delta^{-2}\frac{22}{\theta_d}\mathcal{M}_P.$$

In order to prove this, we first prove two lemmas.

2.2.1.1 Stein approximations with finite RKHS norm

2.7 Lemma. *Suppose that $\mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is bounded and belongs to $L^2(\mathbb{R}^d) \cap C^1(\mathbb{R}^d)$. Suppose further that $h = \mathcal{T}_P \mathbf{g}$ and s_p are Lipschitz, and that $k(\mathbf{x}, \mathbf{y}) = \Psi(\mathbf{x} - \mathbf{y})$ for some $\Psi \in C^2$ with generalised Fourier transform $\hat{\Psi}$. Then, for every $\epsilon > 0$, we can find $\mathbf{g}_\epsilon \in \mathcal{H}_k^d$ such that $\sup_{\mathbf{x} \in \mathbb{R}^d} |(\mathcal{T}_P \mathbf{g}_\epsilon)(\mathbf{x}) - (\mathcal{T}_P \mathbf{g})(\mathbf{x})| \leq \epsilon$ and*

$$\|\mathbf{g}_\epsilon\|_{\mathcal{H}_k^d} \leq (2\pi)^{-d/4}F\left(\frac{12d\log 2}{\pi}(M_1(h) + M_1(s_p)M_0(\mathbf{g}))\epsilon^{-1}\right)^{1/2} \|\mathbf{g}\|_{L^2},$$

where $F(t) := \sup_{\|\omega\|_\infty \leq t} \hat{\Psi}(\omega)^{-1}$.

Proof. First, define the function $S : \mathbb{R}^d \rightarrow \mathbb{R}$ by $S(\mathbf{x}) := \prod_{j=1}^d \frac{\sin x_j}{x_j}$. Then $S \in L^2$ and $\int_{\mathbb{R}^d} \|\mathbf{x}\|_2 S(\mathbf{x})^4 d\mathbf{x} < \infty$. If we also define the density function $\rho(\mathbf{x}) := Z^{-1}S(\mathbf{x})^4$, where $Z := \int_{\mathbb{R}^d} S(\mathbf{x})^4 d\mathbf{x} = (2\pi/3)^d$ is the normalisation constant, one can then check that $\hat{\rho}(\omega)^2 \leq (2\pi)^{-d} \mathbb{1}_{\{\|\omega\|_\infty \leq 4\}}$.

Let $\mathbf{Y} \sim \rho$. For each $\delta > 0$, define $\rho_\delta(\mathbf{x}) := \delta^{-d}\rho(\mathbf{x}/\delta)$, and for any function \mathbf{f} , set $\mathbf{f}_\delta(\mathbf{x}) := \mathbb{E}[\mathbf{f}(\mathbf{x} - \delta\mathbf{Y})]$. Since $h = \mathcal{T}_P \mathbf{g}$ is assumed to be Lipschitz, this implies that $|h_\delta(\mathbf{x}) - h(\mathbf{x})| \leq \delta M_1(h) \mathbb{E}_\rho[\|\mathbf{Y}\|_2]$ for all $\mathbf{x} \in \mathbb{R}^d$.

Now, since

$$(\mathcal{T}_P \mathbf{g}_\delta)(\mathbf{x}) = \mathbb{E}_\rho[\langle s_p(\mathbf{x}), \mathbf{g}(\mathbf{x} - \delta\mathbf{Y}) \rangle] + \mathbb{E}[\nabla_{\mathbf{x}} \cdot \mathbf{g}(\mathbf{x} - \delta\mathbf{Y})]$$

and

$$h_\delta(\mathbf{x}) = \mathbb{E}_\rho[\langle \mathbf{s}_p(\mathbf{x} - \delta \mathbf{Y}), \mathbf{g}(\mathbf{x} - \delta \mathbf{Y}) \rangle] + \mathbb{E}[\nabla_{\mathbf{x}} \cdot \mathbf{g}(\mathbf{x} - \delta \mathbf{Y})]$$

for all $\mathbf{x} \in \mathbb{R}^d$, we deduce that

$$\begin{aligned} |(\mathcal{T}_P \mathbf{g}_\delta)(\mathbf{x}) - h_\delta(\mathbf{x})| &= |\mathbb{E}_\rho[\langle \mathbf{s}_p(\mathbf{x}) - \mathbf{s}_p(\mathbf{x} - \delta \mathbf{Y}), \mathbf{g}(\mathbf{x} - \delta \mathbf{Y}) \rangle]| \\ &\leq \mathbb{E}_\rho[\|\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_p(\mathbf{x} - \delta \mathbf{Y})\|_2 \|\mathbf{g}(\mathbf{x} - \delta \mathbf{Y})\|_2] \\ &\leq M_0(\mathbf{g}) M_1(\mathbf{s}_p) \delta \mathbb{E}_\rho[\|\mathbf{Y}\|_2], \end{aligned}$$

by Cauchy-Schwarz, Jensen's inequality and the fact that \mathbf{s}_p is Lipschitz. Thus, for any $\delta > 0$,

$$|(\mathcal{T}_P \mathbf{g}_\delta)(\mathbf{x}) - (\mathcal{T}_P \mathbf{g})(\mathbf{x})| \leq |(\mathcal{T}_P \mathbf{g}_\delta)(\mathbf{x}) - h_\delta(\mathbf{x})| + |h_\delta(\mathbf{x}) - h(\mathbf{x})| \leq \delta(M_1(h) + M_1(\mathbf{s}_p)M_0(\mathbf{g}))\mathbb{E}_\rho[\|\mathbf{Y}\|_2], \quad (2.4)$$

by the triangle inequality. Letting $\tilde{\epsilon} = \epsilon / ((M_1(h) + M_1(\mathbf{s}_p)M_0(\mathbf{g}))\mathbb{E}_\rho[\|\mathbf{Y}\|_2])$, we have $M_0(\mathcal{T}_P \mathbf{g}_{\tilde{\epsilon}} - \mathcal{T}_P \mathbf{g}) \leq \epsilon$.

By Theorem 10.12 of Wendland [2004], \mathcal{H}_k^d is precisely the subset of $L^2(\mathbb{R}^d) \cap C^1(\mathbb{R}^d)$ whose elements \mathbf{f} have finite RKHS norm:

$$\|\mathbf{f}\|_{\mathcal{H}_k^d}^2 = \frac{1}{(2\pi)^{d/2}} \sum_{i=1}^d \int_{\mathbb{R}^d} \frac{|\hat{f}_i(\boldsymbol{\omega})|^2}{\hat{\Psi}(\boldsymbol{\omega})} d\boldsymbol{\omega} < \infty.$$

It is therefore enough to prove the bound on the RKHS norm of \mathbf{g}_δ . First note that, by the convolution theorem, $\|\mathbf{g}_\delta\|_{\mathcal{H}_k^d}^2$ is equal to

$$\begin{aligned} \frac{1}{(2\pi)^{d/2}} \sum_{i=1}^d \int_{\mathbb{R}^d} \frac{|\hat{g}_{\delta,i}(\boldsymbol{\omega})|^2}{\hat{\Psi}(\boldsymbol{\omega})} d\boldsymbol{\omega} &= (2\pi)^{d/2} \sum_{i=1}^d \int_{\mathbb{R}^d} \frac{|\hat{g}_i(\boldsymbol{\omega})|^2 \hat{\rho}_\delta(\boldsymbol{\omega})^2}{\hat{\Psi}(\boldsymbol{\omega})} d\boldsymbol{\omega} \\ &\leq (2\pi)^{-d/2} \left\{ \sup_{\|\boldsymbol{\omega}\|_\infty \leq 4\delta^{-1}} \hat{\Psi}(\boldsymbol{\omega})^{-1} \right\} \sum_{i=1}^d \int_{\mathbb{R}^d} |\hat{g}_i(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}, \end{aligned}$$

where the inequality follows from the fact that $\hat{\rho}_\delta(\boldsymbol{\omega}) = \hat{\rho}(\delta\boldsymbol{\omega})$. Therefore, overloading the definition $\|\mathbf{f}\|_{L^2} := \sum_{i=1}^d \|f_i\|_{L^2}$, we have

$$\|\mathbf{g}_\delta\|_{\mathcal{H}_k^d} \leq (2\pi)^{-d/4} F(4\delta^{-1})^{1/2} \|\mathbf{g}\|_{L^2},$$

since $\|\mathbf{g}\|_{L^2} = \|\hat{\mathbf{g}}\|_{L^2}$ by Plancherel's theorem. The result follows from computing $\int_{\mathbb{R}} \sin^4(x)/x^4 dx = 2\pi/3$ and observing that

$$\int_{\mathbb{R}^d} \|\mathbf{x}\|_2 \prod_{j=1}^d \frac{\sin^4 x_j}{x_j^4} d\mathbf{x} \leq \int_{\mathbb{R}^d} \|\mathbf{x}\|_1 \prod_{j=1}^d \frac{\sin^4 x_j}{x_j^4} d\mathbf{x} = \sum_{j=1}^d \int_{\mathbb{R}^d} \frac{\sin^4 x_j}{|x_j|^3} \prod_{k \neq j} \frac{\sin^4 x_k}{x_k^4} d\mathbf{x} = 2d(\log 2) \left(\frac{2\pi}{3}\right)^{d-1},$$

which in turn implies that $\mathbb{E}_\rho[\|\mathbf{Y}\|_2] \leq \frac{3d \log 2}{\pi}$. \square

2.2.1.2 Smoothed indicator function

2.8 Lemma. *For any compact set $K \subseteq \mathbb{R}^d$ and $\delta > 0$, define the **set inflation** $K^{2\delta} := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{y}\|_2 \leq 2\delta \text{ for all } \mathbf{y} \in K\}$. Then there exists a function $v_{K,\delta} : \mathbb{R}^d \rightarrow [0, 1]$ such that*

1. $v_{K,\delta}(\mathbf{x}) = 1$ for all $\mathbf{x} \in K$ and $v_{K,\delta}(\mathbf{x}) = 0$ for all $\mathbf{x} \notin K^{2\delta}$;
2. $\|\nabla_{\mathbf{x}} v_{K,\delta}(\mathbf{x})\|_2 \leq \frac{d\delta^{-1}\theta_{d-1}}{\theta_d} \mathbb{1}_{\{\mathbf{x} \in K^{2\delta} \setminus K\}}$;
3. $\|\nabla_{\mathbf{x}}^2 v_{K,\delta}(\mathbf{x})\|_{op} \leq \delta^{-2} \frac{22}{\theta_d} \mathbb{1}_{\{\mathbf{x} \in K^{2\delta} \setminus K\}}$.

Proof. First, define the standard normalised *bump function*

$$\psi(\mathbf{x}) := J_d^{-1} \exp\{-1/(1 - \|\mathbf{x}\|_2^2)\} \mathbb{1}_{\{\|\mathbf{x}\|_2 < 1\}},$$

where $J_d = \int_{\mathcal{B}(\mathbf{0},1)} \exp\{-1/(1 - \|\mathbf{x}\|_2^2)\} d\mathbf{x} = \theta_d V_d$ is the normalising constant.

Letting $\mathbf{W} \sim \psi$, define $v_{K,\delta}(\mathbf{x}) := \mathbb{E}[\mathbb{1}_{\{\mathbf{x} + \delta \mathbf{W} \in K^\delta\}}]$ as the smoothed approximation of $\mathbf{x} \mapsto \mathbb{1}_{\{\mathbf{x} \in K\}}$. The support $\text{supp}(\mathbf{W})$ of \mathbf{W} is exactly $\mathcal{B}(\mathbf{0},1)$, so we can immediately conclude property 1.

It is also clear that $\text{supp}(\nabla v_{K,\delta}) \subseteq K^{2\delta} \setminus K$. Therefore, to prove property 2, we need only consider $\mathbf{x} \in K^{2\delta} \setminus K$. Since

$$\nabla_{\mathbf{x}} v_{K,\delta}(\mathbf{x}) = \frac{1}{\delta^{d+1}} \int_{\mathcal{B}(\mathbf{x},\delta)} \nabla_{\mathbf{x}} \psi\left(\frac{\mathbf{x}-\mathbf{y}}{\delta}\right) \mathbb{1}_{\{\mathbf{y} \in K^\delta\}} d\mathbf{y}$$

by the Leibniz rule, setting $K_{\mathbf{x}}^\delta := \delta^{-1}(K^\delta - \mathbf{x})$ and using Jensen's inequality, we have

$$\begin{aligned} \|\nabla_{\mathbf{x}} v_{K,\delta}(\mathbf{x})\|_2 &\leq \frac{1}{\delta^{d+1}} \int_{\mathcal{B}(\mathbf{x},\delta) \cap K^\delta} \|\nabla \psi\left(\frac{\mathbf{x}-\mathbf{y}}{\delta}\right)\|_2 d\mathbf{y} \\ &= \delta^{-1} \int_{\mathcal{B}(\mathbf{0},1) \cap K_{\mathbf{x}}^\delta} \|\nabla_{\mathbf{z}} \psi(\mathbf{z})\|_2 d\mathbf{z} \leq \delta^{-1} \int_{\mathcal{B}(\mathbf{0},1)} \|\nabla_{\mathbf{z}} \psi(\mathbf{z})\|_2 d\mathbf{z}, \end{aligned}$$

where we used the substitution $\mathbf{z} := (\mathbf{x} - \mathbf{y})/\delta$. By substituting $r = \|\mathbf{z}\|_2$ and integrating by parts, we have

$$\begin{aligned} \int_{\mathcal{B}(\mathbf{0},1)} \|\nabla_{\mathbf{z}} \psi(\mathbf{z})\|_2 d\mathbf{z} &= J_d^{-1} \int_0^1 \frac{2r}{(1-r^2)^2} \exp\left\{-\frac{1}{1-r^2}\right\} (dV_d r^{d-1}) dr \\ &= \frac{d}{\theta_d} \left(\left[-r^{d-1} \exp\left\{-\frac{1}{1-r^2}\right\} \right] \Big|_{r=0}^{r=1} + \int_0^1 (d-1)r^{d-2} \exp\left\{-\frac{1}{1-r^2}\right\} dr \right) \\ &= \frac{d}{\theta_d} [e^{-1} \mathbb{1}_{\{d=1\}} + \mathbb{1}_{\{d>1\}} \theta_{d-1}] = \frac{d\theta_{d-1}}{\theta_d}, \end{aligned}$$

yielding property 2.

Finally, for property 3, first observe that we again have $\text{supp}(\nabla^2 v_{K,\delta}) \subseteq K^{2\delta} \setminus K$, so we only need to check for $\mathbf{x} \in K^{2\delta} \setminus K$. Analogously to above, we have

$$\|\nabla^2 v_{K,\delta}(\mathbf{x})\|_{op} \leq \frac{1}{\delta^{d+2}} \int_{\mathcal{B}(\mathbf{x},\delta) \cap K^\delta} \|\nabla^2 \psi\left(\frac{\mathbf{x}-\mathbf{y}}{\delta}\right)\|_{op} d\mathbf{y} \leq \delta^{-2} M_1(\nabla \psi) \int_{\mathcal{B}(\mathbf{0},1) \cap K_{\mathbf{x}}^\delta} d\mathbf{z} \leq \delta^{-2} M_1(\nabla \psi) V_d.$$

It is simple to show that

$$\nabla^2 \psi(\mathbf{x}) = J_d^{-1} \frac{\exp\{-1/(1 - \|\mathbf{x}\|_2^2)\}}{(1 - \|\mathbf{x}\|_2^2)^2} \cdot \left[\frac{4}{(1 - \|\mathbf{x}\|_2^2)^2} \mathbf{x} \mathbf{x}^\top - \frac{8}{1 - \|\mathbf{x}\|_2^2} \mathbf{x} \mathbf{x}^\top - 2I_d \right] \mathbb{1}_{\{\|\mathbf{x}\|_2 < 1\}},$$

so by the triangle inequality we have

$$\|\nabla^2 \psi(\mathbf{x})\|_{op} \leq J_d^{-1} e^{-1/(1-r^2)} \left(\frac{4r^2}{(1-r^2)^4} + \frac{8r^2}{(1-r^2)^3} + \frac{2}{(1-r^2)^2} \right) \mathbb{1}_{\{r < 1\}},$$

where $r = \|\mathbf{x}\|_2$. One can then observe by plotting the graph of the expression on the right that $M_1(\nabla \psi) \leq 22J_d^{-1}$ and so $M_1(\nabla v_{K,\delta}) \leq \delta^{-2} \frac{22}{\theta_d}$, as desired. \square

2.2.1.3 Proof of Theorem 2.6

We will now prove the tightness lower bound of Theorem 2.6. To do so, we first note that by Theorem 5 and Section 4.2 of Gorham et al. [2018], for any $h \in \text{BL}_{\|\cdot\|_2}$ there exists $\mathbf{g} \in C^1$ which solves the Stein equation $\mathcal{T}_P \mathbf{g} = h - \mathbb{E}_P[h(\mathbf{Z})]$ and satisfies $M_0(\mathbf{g}) \leq \mathcal{M}_P$, where \mathcal{M}_P is a constant independent of h and \mathbf{g} .

Proof of Theorem 2.6. Let $h \in \text{BL}_{\|\cdot\|_2}$, and choose $\mathbf{g} \in C^1$ which solves the Stein equation and satisfies $M_0(\mathbf{g}) \leq \mathcal{M}_P$. Next, we will show that we can approximate $\mathcal{T}_P \mathbf{g}$ arbitrarily well by a function in a scaled copy of $\mathcal{T}_P \mathcal{G}_k$.

Smoothing \mathbf{g} by convolution. Fix any $\rho > 0$, and define $\mathbf{g}_\rho(\mathbf{x}) := \mathbb{E}[\mathbf{g}(\mathbf{x} - \rho\mathbf{U})]$, where $\mathbf{U} \sim \mathcal{N}(\mathbf{0}, I_d)$. We have $M_0(\mathbf{g}_\rho) \leq M_0(\mathbf{g}) \leq \mathcal{M}_P$ and, integrating by parts, we get

$$M_1(\mathbf{g}_\rho) = \sup_{\mathbf{x} \in \mathbb{R}^d} \|\mathbb{E}[\nabla_{\mathbf{x}} \mathbf{g}(\mathbf{x} - \rho\mathbf{U})]\|_2 = \sup_{\mathbf{x} \in \mathbb{R}^d} \|\mathbb{E}[\mathbf{g}(\mathbf{x} - \rho\mathbf{U})\mathbf{U}]\|_2 / \rho \leq M_0(\mathbf{g})\mathbb{E}[\|\mathbf{U}\|_2] / \rho \leq \mathcal{M}_P\sqrt{d}/\rho.$$

The same argument as in (2.4) gives

$$M_0(\mathcal{T}_P \mathbf{g}_\rho - \mathcal{T}_P \mathbf{g}) \leq \rho(M_1(\mathcal{T}_P \mathbf{g}) + M_1(\mathbf{s}_p)\mathcal{M}_P)\mathbb{E}[\|\mathbf{U}\|_2] \leq \rho\sqrt{d}(1 + M_1(\mathbf{s}_p)\mathcal{M}_P).$$

Finally, we show that $M_0(\mathcal{T}_P \mathbf{g}_\rho)$ and $M_1(\mathcal{T}_P \mathbf{g}_\rho)$ are bounded uniformly in \mathbf{g} . Setting $h_\rho(\mathbf{x}) := \mathbb{E}[h(\mathbf{x} - \rho\mathbf{U})]$, note that $(\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x}) = h_\rho(\mathbf{x}) - \mathbb{E}_{\mathbf{Z} \sim P}[h(\mathbf{Z})] + \mathbb{E}_{\mathbf{U} \sim \mathcal{N}(\mathbf{0}, I_d)}[\langle \mathbf{s}_p(\mathbf{x}) - \mathbf{s}_p(\mathbf{x} - \rho\mathbf{U}), \mathbf{g}(\mathbf{x} - \rho\mathbf{U}) \rangle]$. Thus,

$$M_0(\mathcal{T}_P \mathbf{g}_\rho) \leq M_0(h_\rho - \mathbb{E}_P[h_\rho(\mathbf{Z})]) + M_1(\mathbf{s}_p)\rho\mathbb{E}[\|\mathbf{U}\|_2]M_0(\mathbf{g}_\rho) \leq 2 + M_1(\mathbf{s}_p)\rho\sqrt{d}\mathcal{M}_P,$$

where we have used the fact that $h_\rho \in \text{BL}_{\|\cdot\|_2}$. Furthermore, integrating by parts gives

$$\begin{aligned} M_1(\mathcal{T}_P \mathbf{g}_\rho) &\leq M_1(h_\rho - \mathbb{E}_P[h(\mathbf{Z})]) + M_1(\mathcal{T}_P \mathbf{g}_\rho - h_\rho + \mathbb{E}_P[h(\mathbf{Z})]) \\ &\leq 1 + M_1(\mathbf{s}_p)M_0(\mathbf{g}_\rho) + \sup_{\mathbf{x} \in \mathbb{R}^d} \|\mathbb{E}[\nabla_{\mathbf{U}} \langle \mathbf{s}_p(\mathbf{x} - \rho\mathbf{U}) - \mathbf{s}_p(\mathbf{x}), \mathbf{g}(\mathbf{x} - \rho\mathbf{U}) \rangle] / \rho\|_2 \\ &\leq 1 + M_1(\mathbf{s}_p)\mathcal{M}_P + \sup_{\mathbf{x} \in \mathbb{R}^d} \|\mathbb{E}[\mathbf{U} \langle \mathbf{s}_p(\mathbf{x} - \rho\mathbf{U}) - \mathbf{s}_p(\mathbf{x}), \mathbf{g}(\mathbf{x} - \rho\mathbf{U}) \rangle] / \rho\|_2 \\ &\leq 1 + M_1(\mathbf{s}_p)\mathcal{M}_P + \mathbb{E}[\|\mathbf{U}\|_2^2]M_1(\mathbf{s}_p)M_0(\mathbf{g}_\rho) \leq 1 + (1+d)M_1(\mathbf{s}_p)\mathcal{M}_P. \end{aligned}$$

Truncating \mathbf{g}_ρ . Let $\epsilon, \delta > 0$, and let $K = \mathcal{B}(\mathbf{0}, R(\mu, \epsilon))$ in Lemma 2.8. It then follows that we can set $\mathbf{g}_{K,\delta}(\mathbf{x}) := \mathbf{g}_\rho(\mathbf{x})v_{K,\delta}(\mathbf{x})$ as a smooth, truncated approximation to \mathbf{g}_ρ . We have

$$\begin{aligned} (\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x}) - (\mathcal{T}_P \mathbf{g}_{K,\delta})(\mathbf{x}) &= (1 - v_{K,\delta}(\mathbf{x})) \{ \mathbf{s}_p(\mathbf{x})^\top \mathbf{g}_\rho(\mathbf{x}) + \nabla_{\mathbf{x}} \cdot \mathbf{g}_\rho(\mathbf{x}) \} + \nabla_{\mathbf{x}} v_{K,\delta}(\mathbf{x})^\top \mathbf{g}_\rho(\mathbf{x}) \\ &= (1 - v_{K,\delta}(\mathbf{x}))(\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x}) + \nabla_{\mathbf{x}} v_{K,\delta}(\mathbf{x})^\top \mathbf{g}_\rho(\mathbf{x}), \end{aligned}$$

so properties 1 and 2 of Lemma 2.8 imply that $(\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x}) = (\mathcal{T}_P \mathbf{g}_{K,\delta})(\mathbf{x})$ for $\mathbf{x} \in K$, $(\mathcal{T}_P \mathbf{g}_{K,\delta})(\mathbf{x}) = 0$ for $\mathbf{x} \notin K^{2\delta}$, and

$$\begin{aligned} |(\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x}) - (\mathcal{T}_P \mathbf{g}_{K,\delta})(\mathbf{x})| &\leq |(\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x})| + \|\nabla_{\mathbf{x}} v_{K,\delta}(\mathbf{x})\|_2 \|\mathbf{g}_\rho(\mathbf{x})\|_2 \\ &\leq |(\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x})| + \frac{d\delta^{-1}\theta_{d-1}}{\theta_d} \|\mathbf{g}_\rho(\mathbf{x})\|_2 \leq 2 + (M_1(\mathbf{s}_p)\rho\sqrt{d} + \frac{d\delta^{-1}\theta_{d-1}}{\theta_d})\mathcal{M}_P \end{aligned}$$

for $\mathbf{x} \in K^{2\delta} \setminus K$ by Cauchy-Schwarz. Furthermore,

$$\begin{aligned} M_1(\mathcal{T}_P \mathbf{g}_{K,\delta}) &\leq M_1((\mathcal{T}_P \mathbf{g}_\rho)v_{K,\delta}) + M_1(\nabla v_{K,\delta}^\top \mathbf{g}_\rho) \\ &\leq M_1(\mathcal{T}_P \mathbf{g}_\rho)M_0(v_{K,\delta}) + M_0(\mathcal{T}_P \mathbf{g}_\rho)M_1(v_{K,\delta}) + M_1(\nabla v_{K,\delta})M_0(\mathbf{g}_\rho) + M_0(\nabla v_{K,\delta})M_1(\mathbf{g}_\rho) \\ &\leq 1 + M_1(\mathbf{s}_p)\mathcal{M}_P(1+d) + (2 + M_1(\mathbf{s}_p)\rho\sqrt{d}\mathcal{M}_P)\frac{d\delta^{-1}\theta_{d-1}}{\theta_d} + \delta^{-2}\frac{22}{\theta_d}\mathcal{M}_P + \frac{d\delta^{-1}\theta_{d-1}}{\theta_d}\mathcal{M}_P\sqrt{d}/\rho \\ &=: c_{\rho,\delta}. \end{aligned}$$

By Lemma 2.7, there is a function $\tilde{\mathbf{g}}_\epsilon \in \mathcal{H}_k^d$ such that $|(\mathcal{T}_P \tilde{\mathbf{g}}_\epsilon)(\mathbf{x}) - (\mathcal{T}_P \mathbf{g}_{K,\delta})(\mathbf{x})| \leq \epsilon$ and

$$\|\tilde{\mathbf{g}}_\epsilon\|_{\mathcal{H}_k^d} \leq (2\pi)^{-d/4} F\left(\frac{12d \log 2}{\pi}(c_{\rho,\delta} + M_1(\mathbf{s}_p)\mathcal{M}_P\epsilon^{-1})\right)^{1/2} \|\mathbf{g}_{K,\delta}\|_{L^2}. \quad (2.5)$$

Then, since $\mathcal{T}_P \mathbf{g}_{K,\delta}$ and $\mathcal{T}_P \mathbf{g}_\rho$ are identical on K , we have $|(\mathcal{T}_P \tilde{\mathbf{g}}_\epsilon)(\mathbf{x}) - (\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x})| \leq \epsilon$ for all $\mathbf{x} \in K$. For $\mathbf{x} \notin K$, the triangle inequality gives

$$\begin{aligned} |(\mathcal{T}_P \tilde{\mathbf{g}})(\mathbf{x}) - (\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x})| &\leq |(\mathcal{T}_P \tilde{\mathbf{g}}_\epsilon)(\mathbf{x}) - (\mathcal{T}_P \mathbf{g}_{K,\delta})(\mathbf{x})| + |(\mathcal{T}_P \mathbf{g}_{K,\delta})(\mathbf{x}) - (\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{x})| \\ &\leq 2 + \epsilon + (M_1(\mathbf{s}_p)\rho\sqrt{d} + \frac{d\delta^{-1}\theta_{d-1}}{\theta_d})\mathcal{M}_P. \end{aligned}$$

We can upper bound $\|\mathbf{g}_{K,\delta}\|_{L^2}$ in (2.5) as $\|\mathbf{g}_{K,\delta}\|_{L^2} \leq \text{Vol}(K^{2\delta})^{1/2} M_0(\mathbf{g}_\rho) \leq \text{Vol}(K^{2\delta})^{1/2} \mathcal{M}_P$.

Our choice of K ensures that $\mu(\mathbb{1}_{\{\mathbf{X} \notin K\}}) \leq \min\{\epsilon, 1\}$, so

$$\begin{aligned} & |\mathbb{E}_\mu[h(\mathbf{X})] - \mathbb{E}_P[h(\mathbf{Z})]| = |\mathbb{E}_\mu[(\mathcal{T}_P \mathbf{g})(\mathbf{X})]| \\ & \leq |\mathbb{E}_\mu[(\mathcal{T}_P \mathbf{g})(\mathbf{X}) - (\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{X})]| + |\mathbb{E}_\mu[(\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{X}) - (\mathcal{T}_P \tilde{\mathbf{g}}_\epsilon)(\mathbf{X})]| + |\mathbb{E}_\mu[(\mathcal{T}_P \tilde{\mathbf{g}}_\epsilon)(\mathbf{X})]| \\ & \leq M_0(\mathcal{T}_P \mathbf{g} - \mathcal{T}_P \mathbf{g}_\rho) + |\mathbb{E}_\mu[(\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{X}) - (\mathcal{T}_P \tilde{\mathbf{g}}_\epsilon)(\mathbf{X}) \mathbb{1}_{\{\mathbf{X} \in K\}}]| \\ & \quad + |\mathbb{E}_\mu[(\mathcal{T}_P \mathbf{g}_\rho)(\mathbf{X}) - (\mathcal{T}_P \tilde{\mathbf{g}}_\epsilon)(\mathbf{X}) \mathbb{1}_{\{\mathbf{X} \notin K\}}]| + |\mathbb{E}_\mu[(\mathcal{T}_P \tilde{\mathbf{g}}_\epsilon)(\mathbf{X})]| \\ & \leq \rho \sqrt{d}(1 + M_1(\mathbf{s}_p) \mathcal{M}_P) + \epsilon + \min\{\epsilon, 1\}(2 + \epsilon + (M_1(\mathbf{s}_p) \rho \sqrt{d} + \frac{d\delta^{-1}\theta_{d-1}}{\theta_d}) \mathcal{M}_P) + \|\tilde{\mathbf{g}}_\epsilon\|_{\mathcal{H}_k^d} \mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k) \\ & \leq \rho \sqrt{d}(1 + M_1(\mathbf{s}_p) \mathcal{M}_P) + \epsilon + \min\{\epsilon, 1\}(2 + \epsilon + (M_1(\mathbf{s}_p) \rho \sqrt{d} + \frac{d\delta^{-1}\theta_{d-1}}{\theta_d}) \mathcal{M}_P) \\ & \quad + (2\pi)^{-d/4} F\left(\frac{12d \log 2}{\pi} (c_{\rho,\delta} + M_1(\mathbf{s}_p) \mathcal{M}_P \epsilon^{-1})\right)^{1/2} \text{Vol}(\mathcal{B}(\mathbf{0}, R(\mu, \epsilon) + 2\delta))^{1/2} \mathcal{M}_P \mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k) \end{aligned}$$

The result follows from substituting in $\text{Vol}(\mathcal{B}(\mathbf{0}, r)) = V_d r^d$ and then taking the supremum over all $h \in \text{BL}_{\|\cdot\|_2}$. \square

2.2.2 IMQ KSD detects non-convergence

Our next theorem shows that KSDs based on IMQ kernels automatically enforce tightness, and hence detect non-convergence, whenever $\beta \in (-1, 0)$.

2.9 Theorem (IMQ KSD detects non-convergence). *Suppose $P \in \mathcal{P}$ and $k(\mathbf{x}, \mathbf{y}) = (c^2 + \|\mathbf{x} - \mathbf{y}\|_2^2)^\beta$ for some $c > 0$ and $\beta \in (-1, 0)$. If $\mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k) \rightarrow 0$, then $Q_n \Rightarrow P$ as $n \rightarrow \infty$.*

In order to prove this we first prove two lemmas, and then Theorem 2.12, from which Theorem 2.9 quickly follows.

2.2.2.1 Generalised multiquadric Stein sets yield coercive functions

Set $\Psi_{c,\beta}(\mathbf{x}) := (c^2 + \|\mathbf{x}\|_2^2)^\beta$, where $c > 0$ and $\beta \in (-1, 0)$. In this section, we will show that we can find functions in $\mathcal{T}_P \mathcal{G}_k$ which grow sufficiently rapidly as $\|\mathbf{x}\|_2 \rightarrow \infty$. To do so we use the following results from Wendland [2004]:

- $\Psi_{c,\beta}$ has generalised Fourier transform

$$\hat{\Psi}_{c,\beta}(\boldsymbol{\omega}) = \frac{2^{1+\beta}}{\Gamma(-\beta)} \left(\frac{\|\boldsymbol{\omega}\|_2}{c} \right)^{-\beta-d/2} K_{\beta+d/2}(c\|\boldsymbol{\omega}\|_2), \quad (2.6)$$

where $K_v(z)$ is the *modified Bessel function of the third kind*.

- If we further define

$$\tau_v := \begin{cases} \sqrt{\pi/2} & \text{if } |v| \geq 1/2, \\ \frac{\sqrt{\pi} 3^{|v|-1/2}}{2^{|v|+1} \Gamma(|v|+1/2)} & \text{if } |v| < 1/2, \end{cases}$$

then we have the following bounds on $K_v(z)$ for $v \in \mathbb{R}$ and $z > 0$:

- B1. $K_v(z) \geq \tau_v e^{-z}/\sqrt{z}$ for $z \geq 1$;
- B2. $K_v(z) \geq e^{-1} \tau_v z^{-|v|}$ for $z \leq 1$;
- B3. $K_v(z) \leq \sqrt{2\pi/z} e^{-z+v^2/(2z)}$ for $z > 0$;
- B4. $K_v(z) \leq 2^{|v|-1} \Gamma(|v|) z^{-|v|}$ for $v \neq 0, z > 0$.

2.10 Lemma. Let $P \in \mathcal{P}$ and $k(\mathbf{x}, \mathbf{y}) = \Psi_{c,\beta}(\mathbf{x} - \mathbf{y})$. Then for any $\alpha \in (0, \frac{1}{2}(\beta + 1))$ and $a > c/2$, there exist functions $\dot{\mathbf{g}} \in \mathcal{G}$ and $\mathcal{D}(a, c, \alpha, \beta, d)$ such that $\mathcal{T}_P \dot{\mathbf{g}}$ is bounded below by

$$\zeta(a, c, \alpha, \beta, d) := -\frac{\mathcal{D}(a, c, \alpha, \beta, d)^{1/2}}{2\alpha} \left[\frac{M_1(\mathbf{s}_p)R_0^2 + \|\mathbf{s}_p(\mathbf{0})\|_2 R_0 + d}{a^{2(1-\alpha)}} \right],$$

where $R_0 := \inf\{r > 0 : \omega(r') \geq 0 \forall r' \geq r\}$. Furthermore,

$$\liminf \|\mathbf{x}\|_2^{-2\alpha} (\mathcal{T}_P \dot{\mathbf{g}})(\mathbf{x}) \geq \frac{\alpha}{\mathcal{D}(a, c, \alpha, \beta, d)^{1/2}} \omega_0 \text{ as } \|\mathbf{x}\|_2 \rightarrow \infty. \quad (2.7)$$

Proof. Define the function $\mathbf{g} : \mathcal{X} \rightarrow \mathbb{R}^d$ via

$$g_j(\mathbf{x}) := \nabla_{x_j} \Psi_{a,\alpha}(\mathbf{x}) = 2\alpha x_j (a^2 + \|\mathbf{x}\|_2^2)^{\alpha-1}, \quad j = 1, \dots, d.$$

Recall that $\mathbf{g} \in \mathcal{H}_k^d$ if and only if $\|g_j\|_{\mathcal{H}_k} = \left\| \hat{g}_j / \sqrt{\hat{\Psi}_{c,\beta}} \right\|_{L^2} < \infty$ for each $j = 1, \dots, d$. Noting that $\hat{g}_j(\boldsymbol{\omega}) = (i\omega_j) \hat{\Psi}_{c,\beta}(\boldsymbol{\omega})$, we have

$$\begin{aligned} \|\mathbf{g}\|_{\mathcal{H}_k^d}^2 &= \sum_{j=1}^d \|g_j\|_{\mathcal{H}_k}^2 = \sum_{j=1}^d (2\pi)^{-d/2} \int_{\mathbb{R}^d} \hat{g}_j(\boldsymbol{\omega}) \overline{\hat{g}_j(\boldsymbol{\omega})} / \hat{\Psi}_{c,\beta}(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ &= (2\pi)^{-d/2} \sum_{j=1}^d \int_{\mathbb{R}^d} \frac{2^{2(1+\alpha)} / \Gamma(-\alpha)^2}{2^{1+\beta} / \Gamma(-\beta)} \frac{a^{2\alpha+d}}{c^{\beta+d/2}} \omega_j^2 \|\boldsymbol{\omega}\|_2^{\beta-2\alpha-d/2} \frac{K_{\alpha+d/2}(a\|\boldsymbol{\omega}\|_2)^2}{K_{\beta+d/2}(c\|\boldsymbol{\omega}\|_2)} d\boldsymbol{\omega} \\ &= (2\pi)^{-d/2} \underbrace{\frac{2^{2(1+\alpha)} / \Gamma(-\alpha)^2}{2^{1+\beta} / \Gamma(-\beta)} \frac{a^{2\alpha+d}}{c^{\beta+d/2}}}_{=: \mathcal{C}_0(a,c,\alpha,\beta,d)} \int_{\mathbb{R}^d} \|\boldsymbol{\omega}\|_2^{\beta-2\alpha-d/2+2} \frac{K_{\alpha+d/2}(a\|\boldsymbol{\omega}\|_2)^2}{K_{\beta+d/2}(c\|\boldsymbol{\omega}\|_2)} d\boldsymbol{\omega}. \end{aligned}$$

We decompose the integral into two parts: one over the unit ball $\mathcal{B}(\mathbf{0}, 1)$ and the other over its complement $\mathcal{B}(\mathbf{0}, 1)^c$. Using bounds B2 and B4, and the substitution $r = \|\boldsymbol{\omega}\|_2$, we have

$$\begin{aligned} \int_{\mathcal{B}(\mathbf{0}, 1)} \|\boldsymbol{\omega}\|_2^{\beta-2\alpha-d/2+2} \frac{K_{\alpha+d/2}(a\|\boldsymbol{\omega}\|_2)^2}{K_{\beta+d/2}(c\|\boldsymbol{\omega}\|_2)} d\boldsymbol{\omega} &\leq \int_{\mathcal{B}(\mathbf{0}, 1)} \|\boldsymbol{\omega}\|_2^{\beta-2\alpha-d/2+2} \frac{2^{2\alpha+d-2} \Gamma(\alpha+d/2)^2 (a\|\boldsymbol{\omega}\|_2)^{-2\alpha-d}}{e^{-1} \tau_{\beta+d/2} \|c\boldsymbol{\omega}\|_2^{-\beta-d/2}} d\boldsymbol{\omega} \\ &= 2^{2\alpha+d-2} \Gamma(\alpha+d/2)^2 \frac{e}{\tau_{\beta+d/2}} \frac{c^{\beta+d/2}}{a^{2\alpha+d}} \int_{\mathcal{B}(\mathbf{0}, 1)} \|\boldsymbol{\omega}\|_2^{2\beta-4\alpha-d+2} e^{c\|\boldsymbol{\omega}\|_2} d\boldsymbol{\omega} \\ &= dV_d 2^{2\alpha+d-2} \Gamma(\alpha+d/2)^2 \underbrace{\frac{e}{\tau_{\beta+d/2}} \frac{c^{\beta+d/2}}{a^{2\alpha+d}} \int_0^1 r^{2\beta-4\alpha+1} e^{cr} dr}_{=: \mathcal{C}_1(a,c,\alpha,\beta,d)}, \end{aligned}$$

where V_d is the volume of the d -dimensional unit ball. By our choice of α , we can bound the integral as

$$\int_0^1 r^{2\beta-4\alpha+1} e^{cr} dr \leq e^c \int_0^1 r^{2\beta-4\alpha+1} dr = \frac{e^c}{2\beta - 4\alpha + 2}.$$

Using bounds B1 and B3, we have

$$\begin{aligned} \int_{\mathcal{B}(\mathbf{0}, 1)^c} \|\boldsymbol{\omega}\|_2^{\beta-2\alpha-d/2+2} \frac{K_{\alpha+d/2}(a\|\boldsymbol{\omega}\|_2)^2}{K_{\beta+d/2}(c\|\boldsymbol{\omega}\|_2)} d\boldsymbol{\omega} &\leq \int_{\mathcal{B}(\mathbf{0}, 1)^c} \|\boldsymbol{\omega}\|_2^{\beta-2\alpha-d/2+2} \frac{2\pi/(a\|\boldsymbol{\omega}\|_2) \cdot e^{-2a\|\boldsymbol{\omega}\|_2 + (\alpha+d/2)^2/(a\|\boldsymbol{\omega}\|_2)}}{\tau_{\beta+d/2} e^{-c\|\boldsymbol{\omega}\|_2} / \sqrt{c\|\boldsymbol{\omega}\|_2}} d\boldsymbol{\omega} \\ &\leq \frac{2\pi\sqrt{c}}{a\tau_{\beta+d/2}} \int_{\mathcal{B}(\mathbf{0}, 1)^c} \|\boldsymbol{\omega}\|_2^{\beta-2\alpha-d/2+3/2} e^{(c-2a)\|\boldsymbol{\omega}\|_2 + (\alpha+d/2)^2/(a\|\boldsymbol{\omega}\|_2)} d\boldsymbol{\omega} \\ &= dV_d \frac{2\pi\sqrt{c}}{a\tau_{\beta+d/2}} \int_1^\infty r^{\beta-2\alpha+d/2+1/2} e^{(c-2a)r + (\alpha+d/2)^2/(ar)} dr. \end{aligned}$$

Recalling that $c - 2a < 0$, and using the substitution $t = (2a - c)r$, the integral above can be upper bounded as follows:

$$\begin{aligned} \int_1^\infty r^{\beta-2\alpha+d/2+1/2} e^{(c-2\alpha)r+(\alpha+d/2)^2/(ar)} dr &\leq e^{(c-2a)+(\alpha+d/2)^2/a} \int_1^\infty r^{\beta-2\alpha+d/2+1/2} e^{(c-2a)r} dr \\ &= e^{(c-2a)+(\alpha+d/2)^2/a} (2a - c)^{-\beta+2\alpha-d/2-3/2} \Gamma(\beta - 2\alpha + d/2 + 3/2, 2a - c) =: \mathcal{C}_2(a, c, \alpha, \beta, d), \end{aligned}$$

where $\Gamma(s, x) := \int_x^\infty t^{s-1} e^{-t} dt$ is the *upper incomplete gamma function*. Thus, $\mathbf{g} \in \mathcal{H}_k^d$ with norm upper bounded by $\mathcal{D}(a, b, \alpha, \beta, d)^{1/2}$, where

$$\mathcal{D}(a, b, \alpha, \beta, d) := \mathcal{C}_0(a, b, \alpha, \beta, d) \left(\mathcal{C}_1(a, c, \alpha, \beta, d) \frac{e^c}{2\beta - 4\alpha + 2} + \mathcal{C}_2(a, c, \alpha, \beta, d) \frac{2\pi d V_d \sqrt{c}}{a \tau_{\beta+d/2}} \right).$$

Now set $\dot{\mathbf{g}} := -\mathcal{D}(a, c, \alpha, \beta)^{-1/2} \mathbf{g} \in \mathcal{G}_k$. By our definition of \mathbf{g} , we have

$$\frac{\mathcal{D}(a, c, \alpha, \beta, d)^{1/2}}{2\alpha} (\mathcal{T}_P \dot{\mathbf{g}})(\mathbf{x}) = -\frac{\mathbf{s}_p(\mathbf{x})^\top \mathbf{x}}{(a^2 + \|\mathbf{x}\|_2^2)^{1-\alpha}} - \frac{d}{(a^2 + \|\mathbf{x}\|_2^2)^{1-\alpha}} + \frac{2(1-\alpha)\|\mathbf{x}\|_2^2}{(a^2 + \|\mathbf{x}\|_2^2)^{2-\alpha}} \quad (2.8)$$

$$\geq -\frac{\mathbf{s}_p(\mathbf{x})^\top \mathbf{x}}{(a^2 + \|\mathbf{x}\|_2^2)^{1-\alpha}} - \frac{d}{a^{2(1-\alpha)}}. \quad (2.9)$$

By the distant dissipativity of P , we have $\limsup_{\|\mathbf{x}\|_2 \rightarrow \infty} \frac{1}{\|\mathbf{x}\|_2^2} \mathbf{s}_p(\mathbf{x})^\top \mathbf{x} \leq -\frac{1}{2}\omega_0$, with $\omega_0 > 0$ as defined in (2.1). Thus the first term in (2.8) has a growth rate of at least $\omega_0 \|\mathbf{x}\|_2^{2\alpha}/2$ as $\|\mathbf{x}\|_2 \rightarrow \infty$. From this we deduce (2.7).

Indeed, by our choice of R_0 and the distant dissipativity of P , we have $-\mathbf{s}_p(\mathbf{x})^\top \mathbf{x} \geq 0$ for $\mathbf{x} \notin \mathcal{B}(\mathbf{0}, R_0)$. Finally, using the Cauchy-Schwarz inequality,

$$|\mathbf{s}_p(\mathbf{x})^\top \mathbf{x}| \leq |(\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_p(\mathbf{0}))^\top \mathbf{x}| + |\mathbf{s}_p(\mathbf{0})^\top \mathbf{x}| \leq M_1(\mathbf{s}_p) \|\mathbf{x}\|_2^2 + \|\mathbf{s}_p(\mathbf{0})\|_2 \|\mathbf{x}\|_2.$$

Then for any $\mathbf{x} \in \mathcal{B}(\mathbf{0}, R_0)$, we have $-\mathbf{s}_p(\mathbf{x})^\top \mathbf{x} \geq -M_1(\mathbf{s}_p)R_0^2 - \|\mathbf{s}_p(\mathbf{0})\|_2 R_0$. Applying these lower bounds to (2.9) yields the uniform lower bound $\zeta(a, c, \alpha, \beta, d)$. \square

2.2.2.2 Coercive functions yield tightness

2.11 Lemma. Suppose $\mathbf{g} \in \mathcal{G}_k$ is such that $\mathcal{T}_P \mathbf{g}$ is bounded below by $\zeta \in \mathbb{R}$ and $\liminf_{\|\mathbf{x}\|_2 \rightarrow \infty} \|\mathbf{x}\|_2^{-u} (\mathcal{T}_P \mathbf{g})(\mathbf{x}) > \eta$ for some $\eta, u > 0$. Then for any probability measure μ , and ϵ small enough, we have

$$R(\mu, \epsilon) \leq \left[\frac{1}{\epsilon \eta} (\mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k) - \zeta) \right]^{1/u}. \quad (2.10)$$

Thus, provided $\limsup_{m \rightarrow \infty} \mathcal{S}(\mu_m, \mathcal{T}_P, \mathcal{G}_k)$ is finite, $(\mu_m)_{m \geq 1}$ is uniformly tight.

Proof. Define $\gamma(r) := \inf\{(\mathcal{T}_P \mathbf{g})(\mathbf{x}) - \zeta \mid \|\mathbf{x}\|_2 \geq r\}$. This is non-negative and non-decreasing, so by Markov's inequality we have

$$\mu(\|\mathbf{X}\|_2 \geq r) \leq \frac{\mathbb{E}_\mu[\gamma(\|\mathbf{X}\|_2)]}{\gamma(r)} \leq \frac{\mathbb{E}_\mu[(\mathcal{T}_P \mathbf{g})(\mathbf{X}) - \zeta]}{\gamma(r)}$$

for any measure μ . It is therefore clear that if $\epsilon \geq (\mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k) - \zeta)/\gamma(r_\epsilon)$, then $\mu(\|\mathbf{X}\|_2 \geq r_\epsilon) \leq \epsilon$. The assumption of the lemma implies that $\gamma(r) \geq \eta r^u$ for large enough r . Thus, if we suppose that

$$r_\epsilon \geq \left[\frac{1}{\epsilon \eta} (\mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k) - \zeta) \right]^{1/u},$$

where ϵ is small enough to yield $\gamma(r_\epsilon) \geq \eta r_\epsilon^u$, then

$$\epsilon \geq \frac{\mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k) - \zeta}{\eta r_\epsilon^u} \geq \frac{\mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k) - \zeta}{\gamma(r_\epsilon)},$$

and so we have $\mu(\|\mathbf{X}\|_2 \geq r_\epsilon) \leq \epsilon$, and (2.10) follows from the definition of $R(\mu, \epsilon)$. Furthermore, if $\limsup_{m \rightarrow \infty} \mathcal{S}(\mu_m, \mathcal{T}_P, \mathcal{G}_k)$ is finite, then so is $\limsup_{m \rightarrow \infty} R(\mu_m, \epsilon)$, so $(\mu_m)_{m \geq 1}$ is uniformly tight. \square

2.2.2.3 Proof of Theorem 2.9

We are now ready to prove Theorem 2.9, which follows quickly from the following bound.

2.12 Theorem (IMQ KSD lower bound). *Let $P \in \mathcal{P}$ and let $k(\mathbf{x}, \mathbf{y}) = (c^2 + \|\mathbf{x} - \mathbf{y}\|_2^2)^\beta$ for some $c > 0$ and $\beta \in (-1, 0)$. Let $\alpha \in (0, \frac{1}{2}(\beta + 1))$ and $a > c/2$. Then there exists $\epsilon_0 > 0$ and a constant \mathcal{M}_P such that, for all μ ,*

$$d_{BL\|\cdot\|_2}(\mu, P) \leq \inf_{\epsilon \in [0, \epsilon_0), \delta, \rho > 0} \rho \sqrt{d} (1 + M_1(\mathbf{s}_p) \mathcal{M}_P) + \left(3 + \epsilon + (M_1(\mathbf{s}_p) \rho \sqrt{d} + \frac{d\delta^{-1}\theta_{d-1}}{\theta_d}) \mathcal{M}_P \right) \epsilon + (2\pi)^{-d/4} \mathcal{M}_P V_d^{1/2} \times \\ \left[\left(\frac{\mathcal{D}(a, c, \alpha, \beta, d)^{1/2} (\mathcal{S}(\mu, \mathcal{T}, \mathcal{G}_k) - \zeta(a, c, \alpha, \beta, d))}{\alpha \omega_0 \epsilon} \right)^{1/\alpha} + 2\delta \right]^{d/2} \sqrt{F_{IMQ} \left(\frac{12d \log 2}{\pi} (c_\rho, \delta + M_1(\mathbf{s}_p) \mathcal{M}_P) \epsilon^{-1} \right)} \mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k),$$

where

$$F_{IMQ}(t) := \frac{\Gamma(-\beta)}{2^{1+\beta}} \left(\frac{\sqrt{d}}{c} \right)^{\beta+d/2} \frac{t^{\beta+d/2}}{K_{\beta+d/2}(c\sqrt{dt})}. \quad (2.11)$$

Proof. For any $\alpha \in (0, \frac{1}{2}(\beta+1))$ and $a > c/2$, we can apply Lemma 2.10 to find $\mathring{g} \in \mathcal{G}_k$ such that $\mathcal{T}_P \mathring{g}$ is bounded below by a constant $\zeta(a, c, \alpha, \beta, d)$ and grows at a rate of $\|\mathbf{x}\|_2^{2\alpha}$ as $\|\mathbf{x}\|_2 \rightarrow \infty$. We can plug the upper bound on $R(\mu, \epsilon)$ given by Lemma 2.11 into the bound given by Theorem 2.6. It is clear from (2.6) that $\hat{\Psi}_{c,\beta}$ is monotonically decreasing in $\|\omega\|_2$, so by noting that $\|\omega\|_\infty \leq t$ implies $\|\omega\|_2 \leq \sqrt{dt}$, we can also substitute (2.11) into the bound from Theorem 2.6. Finally, as indicated by (2.7), we can take $\eta \rightarrow \frac{\alpha}{\mathcal{D}(a, c, \alpha, \beta, d)^{1/2}} \omega_0$, which yields the desired result. \square

Finally, we prove Theorem 2.9.

Proof of Theorem 2.9. Recall that $F_{IMQ}(t)$ is finite for all $t > 0$, so if $\mathcal{S}(\mu_m, \mathcal{T}_P, \mathcal{G}_k) \xrightarrow{m \rightarrow \infty} 0$, then by Theorem 2.12 we have

$$\limsup_{m \rightarrow \infty} d_{BL\|\cdot\|_2}(\mu_m, P) \leq \rho \sqrt{d} (1 + M_1(\mathbf{s}_p) \mathcal{M}_P) + \left(3 + \epsilon + (M_1(\mathbf{s}_p) \rho \sqrt{d} + \frac{d\delta^{-1}\theta_{d-1}}{\theta_d}) \mathcal{M}_P \right) \epsilon.$$

Taking $\epsilon, \rho \rightarrow 0$ then yields $d_{BL\|\cdot\|_2}(\mu_m, P) \xrightarrow{m \rightarrow \infty} 0$. The result follows since $d_{BL\|\cdot\|_2}(\cdot, \cdot)$ metrises weak convergence. \square

The result of all of this is that we finally have a KSD that metrises weak convergence over a large class of probability distributions. In the next chapter, we explore how the formulation of the KSD, along with the theory introduced in this chapter, can be used to formulate, and prove the convergence of, an expressive variational inference algorithm.

Chapter 3

Stein Variational Gradient Descent

3.1 SVGD algorithm

Variational inference approximates the target density p by minimising the *Kullback-Leibler (KL) divergence*

$$\text{KL}(q \parallel p) := \mathbb{E}_Q[\log q(\mathbf{X})] - \mathbb{E}_Q[\log p(\mathbf{X})]$$

over $Q \in \mathcal{Q}$ for some family of distributions \mathcal{Q} . Usually, \mathcal{Q} is taken to be some parametric family $\mathcal{Q} = \{q_\theta : \theta \in \Theta\}$. However, such families are often too restrictive to return an accurate approximation to the target, and have to be carefully constructed depending on the specific modelling problem at hand. A more expressive family \mathcal{Q} would lead to a more general-purpose algorithm; the tradeoff, however, is that the corresponding optimisation of $\text{KL}(q \parallel p)$ is more difficult to solve for such \mathcal{Q} .

Fortunately, the KSD allows for an expressive yet tractable choice of \mathcal{Q} . Liu and Wang [2016] proposed considering the following transformation $\mathbf{T} : \mathcal{X} \rightarrow \mathcal{X}$ defined as a small perturbation of the identity map:

$$\mathbf{T}(\mathbf{x}) = \mathbf{x} + \epsilon \phi(\mathbf{x}). \quad (3.1)$$

An intuitive approach to minimising $\text{KL}(q \parallel p)$ is to perform gradient descent by iteratively applying transformations of the form (3.1) to \mathbf{x} , each time choosing ϕ so as to maximise the gradient $\nabla_\epsilon \text{KL}(q_{[\mathbf{T}]} \parallel p)|_{\epsilon=0}$, where $q_{[\mathbf{T}]}$ is as defined below.

3.1 Proposition. Consider the transformation \mathbf{T} defined in (3.1), and let $q_{[\mathbf{T}]}(\mathbf{z})$ be the density of $\mathbf{z} := \mathbf{T}(\mathbf{x})$, where $\mathbf{x} \sim q(\mathbf{x})$. We have

$$\nabla_\epsilon \text{KL}(q_{[\mathbf{T}]} \parallel p)|_{\epsilon=0} = -\mathbb{E}_Q[\mathbf{s}_p(\mathbf{x})^\top \phi(\mathbf{x}) + \nabla_{\mathbf{x}} \cdot \phi(\mathbf{x})].$$

Proof. By the change of variables formula,

$$q_{[\mathbf{T}]}(\mathbf{z}) = q(\mathbf{T}^{-1}(\mathbf{z})) \cdot |\det(\nabla_{\mathbf{z}} \mathbf{T}^{-1}(\mathbf{z}))|.$$

Letting $\tilde{\mathbf{x}} \sim p(\tilde{\mathbf{x}})$ and setting $\tilde{\mathbf{z}} := \mathbf{T}^{-1}(\tilde{\mathbf{x}})$,

$$p_{[\mathbf{T}^{-1}]}(\tilde{\mathbf{z}}) = p(\mathbf{T}(\tilde{\mathbf{z}})) \cdot |\det(\nabla_{\tilde{\mathbf{z}}} \mathbf{T}(\tilde{\mathbf{z}}))|.$$

Then

$$\text{KL}(q_{[\mathbf{T}]} \parallel p) = \mathbb{E}_{\mathbf{z} \sim q_{[\mathbf{T}]}(\mathbf{z})} \left[\log \frac{q_{[\mathbf{T}]}(\mathbf{z})}{p(\mathbf{z})} \right] = \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\log \frac{q(\mathbf{x})}{p_{[\mathbf{T}^{-1}]}(\mathbf{x})} \right] = \text{KL}(q \parallel p_{[\mathbf{T}^{-1}]}), \quad (3.2)$$

and so

$$\begin{aligned} \nabla_\epsilon \text{KL}(q_{[\mathbf{T}]} \parallel p) &= -\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [\nabla_\epsilon \log p_{[\mathbf{T}^{-1}]}(\mathbf{x})] \\ &= -\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [\mathbf{s}_p(\mathbf{T}(\mathbf{x}))^\top \nabla_{\mathbf{x}} \mathbf{T}(\mathbf{x}) + \text{Trace}((\nabla_{\mathbf{x}} \mathbf{T}(\mathbf{x}))^{-1} \cdot \nabla_\epsilon \nabla_{\mathbf{x}} \mathbf{T}(\mathbf{x}))]. \end{aligned}$$

Evaluating at $\epsilon = 0$ yields the desired result. \square

This result allows us to find the steepest descent direction $\tilde{\phi}_{q,p}$ over $\phi \in \mathcal{G}_k$, when \mathcal{G}_k is the unit ball of some product Hilbert space.

3.2 Corollary. *Let \mathcal{G}_k be the kernel Stein set (1.2). Over all $\phi \in \mathcal{G}_k$, the direction of steepest descent is given by*

$$\tilde{\phi}_{q,p}(\cdot) = \frac{\mathbb{E}_Q \xi_p(\mathbf{X}, \cdot)}{\|\mathbb{E}_Q \xi_p(\mathbf{X}, \cdot)\|_{\mathcal{H}_k^d}}. \quad (3.3)$$

Proof. Observe that

$$\sup_{\phi \in \mathcal{G}_k} -\nabla_\epsilon \text{KL}(q_{[\mathbf{T}]} \| p)|_{\epsilon=0} = \sup_{\phi \in \mathcal{G}_k} \mathbb{E}_Q [\mathbf{s}_p(\mathbf{x})^\top \phi(\mathbf{x}) + \nabla_{\mathbf{x}} \cdot \phi(\mathbf{x})] = \mathcal{S}(Q, \mathcal{T}_P, \mathcal{G}_k).$$

We know that the supremum in the KSD is achieved by

$$\tilde{\phi}_{q,p}(\cdot) = \frac{\mathbb{E}_Q \xi_p(\mathbf{X}, \cdot)}{\|\mathbb{E}_Q \xi_p(\mathbf{X}, \cdot)\|_{\mathcal{H}_k^d}},$$

so this is our direction of steepest descent over the unit ball \mathcal{G}_k of \mathcal{H}_k^d . \square

These results motivate Algorithm 3.1, known as **Stein variational gradient descent** (SVGD). The normalising factor in (3.3) can be absorbed into the learning rates $(\epsilon_\ell)_{\ell=0}^L$.

Algorithm 3.1: Stein Variational Gradient Descent (SVGD)

```

1 Input: The score function  $\mathbf{s}_p(\mathbf{x})$ , initial particles  $\{\mathbf{x}_i^{(0)}\}_{i=1}^n$ , kernel  $k$ , number of iterations  $L$  and step sizes  $\{\epsilon_\ell\}_{\ell=0}^{L-1}$ ;
2 for  $\ell = 0, \dots, L-1$  do
3    $\mathbf{x}_i^{(\ell+1)} \leftarrow \mathbf{x}_i^{(\ell)} + \epsilon_\ell \tilde{\phi}_{\mu_n^{(\ell)}, p}(\mathbf{x}_i^{(\ell)})$ ,  $i = 1, \dots, n$ , where
4   
$$\tilde{\phi}_{\mu_n^{(\ell)}, p}(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n \left\{ \mathbf{s}_p(\mathbf{x}_j^{(\ell)}) k(\mathbf{x}_j^{(\ell)}, \mathbf{x}) + \nabla_{\mathbf{x}_j^{(\ell)}} k(\mathbf{x}_j^{(\ell)}, \mathbf{x}) \right\}.$$


```

3.2 Convergence of SVGD

Let $\mu_n^{(\ell)}(\cdot) := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i^{(\ell)}}(\cdot)$ be the discrete probability measure of the n particles at the ℓ^{th} time step of the SVGD algorithm. In this section, we will show that $\mu_n^{(\ell)} \Rightarrow P$ as $n, \ell \rightarrow \infty$ when k is the IMQ kernel and $P \in \mathcal{P}$. Thus, SVGD has the attractive property that, given sufficient computational resources, it can approximate the target arbitrarily well, which is typically not the case with standard parametric variational inference.

First, let's introduce some notation. We denote the optimal transformation between two probability measures ν and π by $\mathbf{T}^{\nu, \pi, \epsilon}(\mathbf{x}) := \mathbf{x} + \epsilon \tilde{\phi}_{\nu, \pi}(\mathbf{x})$. We then define a mapping $\Phi_{\pi, \epsilon}$ on \mathcal{P} by $\Phi_{\pi, \epsilon} : \mu \mapsto \mathbf{T}_\#^{\mu, \pi, \epsilon} \mu$, where $\mathbf{T}_\#^{\mu, \pi, \epsilon}$ denotes the *pushforward measure* of the transformation $\mathbf{T}^{\mu, \pi, \epsilon}$.

We then have $\mu_n^{(\ell+1)} = \Phi_{\pi, \epsilon_\ell}(\mu_n^{(\ell)})$. Taking the infinite particle limit $n \rightarrow \infty$, and given some initial probability measure $\mu_\infty^{(0)}$, we can define another sequence of probability measures $(\mu_\infty^{(\ell)})_{\ell \geq 0}$ through the updates

$$\mu_\infty^{(\ell+1)} := \Phi_{\pi, \epsilon_\ell}(\mu_\infty^{(\ell)}).$$

3.2.1 Convergence to large sample limit

We first show that for each time step ℓ , the finite particle update converges to the infinite particle update as the number of particles goes to infinity.

3.3 Theorem. *Consider the SVGD algorithm as outlined above, and suppose that $W_1(\mu_n^{(0)}, \mu_\infty^{(0)}) \xrightarrow{n \rightarrow \infty} 0$. Suppose further that there exist constants $C_1, C_2 > 0$ such that*

$$M_1(\xi_p(\mathbf{x}, \cdot)) \leq C_1(1 + \|\mathbf{x}\|_2) \quad \text{and} \quad M_1(\xi_p(\cdot, \mathbf{y})) \leq C_2(1 + \|\mathbf{y}\|_2), \quad (3.4)$$

where ξ_p is as defined in (1.3) and M_1 is as defined in (2.2). Then $W_1(\mu_n^{(\ell)}, \mu_\infty^{(\ell)}) \rightarrow 0$ as $n \rightarrow \infty$ for each $\ell \in \{1, 2, \dots\}$.

Proof. We prove this by induction on the time step ℓ . The base case $\ell = 0$ is satisfied by assumption. Now suppose $W_1(\mu_n^{(\ell)}, \mu_\infty^{(\ell)}) \xrightarrow{n \rightarrow \infty} 0$ for some $\ell \geq 0$. Then there exists $C' > 0$ such that

$$\sup_{n \geq 1} \{1 + \epsilon_\ell C_1(1 + \mu_n^{(\ell)}(\|\cdot\|_2)) + \epsilon_\ell C_2(1 + \mu_\infty^{(\ell)}(\|\cdot\|_2))\} \leq C'.$$

The assumptions (3.4) and Lemma 3.4 below imply that

$$\begin{aligned} W_1(\mu_n^{(\ell+1)}, \mu_\infty^{(\ell+1)}) &= W_1(\Phi_{p,\epsilon_\ell}(\mu_n^{(\ell)}), \Phi_{p,\epsilon_\ell}(\mu_\infty^{(\ell)})) \\ &\leq W_1(\mu_n^{(\ell)}, \mu_\infty^{(\ell)})(1 + \epsilon_\ell C_1(1 + \mu_n^{(\ell)}(\|\cdot\|_2)) + \epsilon_\ell C_2(1 + \mu_\infty^{(\ell)}(\|\cdot\|_2))) \\ &\leq C' W_1(\mu_n^{(\ell)}, \mu_\infty^{(\ell)}) \xrightarrow{n \rightarrow \infty} 0, \end{aligned}$$

which completes the proof. \square

It thus remains to prove the lemma used in the proof above.

3.4 Lemma. Suppose that, for some constants $\tilde{C}_1, \tilde{C}_2 > 0$, we have

$$\sup_{y \in \mathcal{X}} \|\nabla_y \xi_p(x, y)\|_{op} \leq \tilde{C}_1(1 + \|x\|_2), \quad \text{and} \quad \sup_{x \in \mathcal{X}} \|\nabla_x \xi_p(x, y)\|_{op} \leq \tilde{C}_2(1 + \|y\|_2).$$

Then, for any $\epsilon > 0$ and probability measures μ, ν ,

$$W_1(\Phi_{p,\epsilon}(\mu), \Phi_{p,\epsilon}(\nu)) \leq W_1(\mu, \nu)(1 + \epsilon \tilde{C}_1(1 + \mu(\|\cdot\|_2)) + \epsilon \tilde{C}_2(1 + \nu(\|\cdot\|_2))).$$

Proof. Let $(\mathbf{X}', \mathbf{Z}')$ be an optimal 1-Wasserstein coupling of (μ, ν) . Then

$$\begin{aligned} \|\mathbf{T}^{\mu,p,\epsilon}(\mathbf{x}) - \mathbf{T}^{\mu,p,\epsilon}(\mathbf{z})\|_2 &\leq \|\mathbf{x} - \mathbf{z}\|_2 + \epsilon \|\mathbb{E}[\xi_p(\mathbf{X}', \mathbf{x}) - \xi_p(\mathbf{X}', \mathbf{z})]\|_2 + \epsilon \|\mathbb{E}[\xi_p(\mathbf{X}', \mathbf{z}) - \xi_p(\mathbf{Z}', \mathbf{z})]\|_2 \\ &\leq \|\mathbf{x} - \mathbf{z}\|_2(1 + \epsilon \tilde{C}_1(1 + \mathbb{E}[\|\mathbf{X}'\|_2])) + \epsilon \tilde{C}_2 \mathbb{E}[\|\mathbf{X}' - \mathbf{Z}'\|_2](1 + \|\mathbf{z}\|_2) \\ &= \|\mathbf{x} - \mathbf{z}\|_2(1 + \epsilon \tilde{C}_1(1 + \mu(\|\cdot\|_2))) + \epsilon \tilde{C}_2 W_1(\mu, \nu)(1 + \|\mathbf{z}\|_2). \end{aligned}$$

By definition of \mathbf{X}' and \mathbf{Z}' , we have $\mathbf{T}^{\mu,\nu,\epsilon}(\mathbf{X}') \sim \Phi_{p,\epsilon}(\mu)$ and $\mathbf{T}^{\mu,\nu,\epsilon}(\mathbf{Z}') \sim \Phi_{p,\epsilon}(\nu)$. Hence,

$$\begin{aligned} W_1(\Phi_{p,\epsilon}(\mu), \Phi_{p,\epsilon}(\nu)) &\leq \mathbb{E}[\|\mathbf{T}^{\mu,p,\epsilon}(\mathbf{X}') - \mathbf{T}^{\nu,p,\epsilon}(\mathbf{Z}')\|_2] \\ &\leq \mathbb{E}[\|\mathbf{X}' - \mathbf{Z}'\|_2](1 + \epsilon \tilde{C}_1(1 + \mu(\|\cdot\|_2))) + \epsilon \tilde{C}_2 W_1(\mu, \nu)(1 + \mathbb{E}[\|\mathbf{Z}'\|_2]) \\ &= W_1(\mu, \nu)(1 + \epsilon \tilde{C}_1(1 + \mu(\|\cdot\|_2)) + \epsilon \tilde{C}_2(1 + \nu(\|\cdot\|_2))), \end{aligned}$$

as desired. \square

Crucially, the growth assumptions on the kernel k and the target P in the statement of Theorem 3.3 are satisfied by the IMQ kernel, as well as other kernels, and by target measures P with Lipschitz score function s_p .

3.2.2 Convergence of large sample limit

We have shown that $\mu_n^{(\ell)} \Rightarrow \mu_\infty^{(\ell)}$ as $n \rightarrow \infty$ for all ℓ . Our goal now is to show that $\mu_\infty^{(\ell)} \Rightarrow P$ as $\ell \rightarrow \infty$, and hence that $\hat{\mu}_n^{(\ell)} \Rightarrow P$ as $n, \ell \rightarrow \infty$. In order to do so, we first prove the following lemma.

3.5 Lemma. Let $B \in \mathbb{R}^{r \times r}$ and let $0 < \epsilon < \frac{1}{\rho(B+B^\top)}$, where $\rho(\cdot)$ denotes the spectral radius. Then $I_r + \epsilon(B + B^\top)$ is positive definite and

$$\log |\det(I_r + \epsilon B)| \geq \epsilon \text{Trace}(B) - \epsilon^2 \frac{\|B\|_F^2}{1 - \epsilon \rho(B + B^\top)},$$

where $\|\cdot\|_F$ is the Frobenius norm. Thus, if we take $\epsilon \leq \frac{1}{2\rho(B+B^\top)}$, then

$$\log |\det(I_r + \epsilon B)| \geq \epsilon \text{Trace}(B) - 2\epsilon^2 \|B\|_F^2.$$

Proof. First note that when $\epsilon < 1/\rho(B + B^\top)$, we have

$$\rho(I_r + \epsilon(B + B^\top)) \geq 1 - \epsilon\rho(B + B^\top) > 0,$$

so $I_r + \epsilon(B + B^\top)$ is positive definite. Now,

$$\begin{aligned} \log |\det(I_r + \epsilon B)| &= \frac{1}{2} \log \det\{(I_r + \epsilon B)(I_r + \epsilon B)^\top\} = \frac{1}{2} \log \det\{I_r + \epsilon(B + B^\top) + \epsilon^2 BB^\top\} \\ &\geq \frac{1}{2} \log \det\{I_r + \epsilon(B + B^\top)\}, \end{aligned}$$

since $I_r + \epsilon(B + B^\top)$ and $\epsilon^2 BB^\top$ are both positive semi-definite.

Let $\{\lambda_i\}$ be the eigenvalues of $B + B^\top$. Then, setting $f_i(s) := \log(1 + s\epsilon\lambda_i)$ for $s \in [0, 1]$, we have

$$\log(1 + \epsilon\lambda_i) = f_i(1) = \int_0^1 \frac{\epsilon\lambda_i}{1 + s\epsilon\lambda_i} ds$$

by Taylor's theorem, so

$$\begin{aligned} \log \det(I_r + \epsilon(B + B^\top)) - \epsilon \text{Trace}(B + B^\top) &= \sum_i \{\log(1 + \epsilon\lambda_i) - \epsilon\lambda_i\} \\ &= \sum_i \left\{ \int_0^1 \frac{\epsilon\lambda_i}{1 + s\epsilon\lambda_i} ds - \epsilon\lambda_i \right\} \\ &= - \sum_i \int_0^1 \frac{s\epsilon^2\lambda_i^2}{1 + s\epsilon\lambda_i} ds \\ &\geq - \sum_i \frac{\epsilon^2\lambda_i^2}{1 - s\epsilon \max|\lambda_i|} \int_0^1 s ds \\ &\geq \frac{\epsilon^2}{2(1 - \epsilon\rho(B + B^\top))} \sum_i \lambda_i^2 = -\frac{\epsilon^2}{2} \frac{\|B + B^\top\|_F^2}{1 - \epsilon\rho(B + B^\top)}. \end{aligned}$$

Hence,

$$\begin{aligned} \log |\det(I_r + \epsilon B)| &\geq \frac{1}{2} \log \det\{I_r + \epsilon(B + B^\top)\} \\ &\geq \frac{\epsilon}{2} \text{Trace}(B + B^\top) - \frac{\epsilon^2}{4} \frac{\|B + B^\top\|_F^2}{1 - \epsilon\rho(B + B^\top)} \\ &\geq \epsilon \text{Trace}(B) - \epsilon^2 \frac{\|B\|_F^2}{1 - \epsilon\rho(B + B^\top)}, \end{aligned}$$

since $\|B + B^\top\|_F \leq \|B\|_F + \|B^\top\|_F = 2\|B\|_F$. □

We now prove that $\mu_\infty^{(\ell)} \Rightarrow P$ as $\ell \rightarrow \infty$.

3.6 Theorem (Convergence of SVGD). *Suppose that $\text{KL}(\mu_\infty^{(0)} \| P) < \infty$ and that the step sizes ϵ_ℓ are at most*

$$\tilde{\epsilon}_\ell := \left(2 \sup_{\mathbf{x} \in \mathbb{R}^d} \rho \left(\nabla \tilde{\phi}_{\mu_\infty^{(\ell)}, p} + \nabla \tilde{\phi}_{\mu_\infty^{(\ell)}, p}^\top \right) \right)^{-1},$$

where $\rho(\cdot)$ denotes the spectral norm. Set $R := c^{2\beta} (M_1(\mathbf{s}_p)/2 + 4d/c^2)$. Then

$$\text{KL}(\mu_\infty^{(\ell+1)} \| P) - \text{KL}(\mu_\infty^{(\ell)} \| P) \leq -\epsilon_\ell (1 - \epsilon_\ell R) \mathcal{S}^2(\mu_\infty^{(\ell)}, \mathcal{T}_P, \mathcal{G}_k).$$

Thus, if we take $\epsilon_\ell \propto \mathcal{S}(\mu_\infty^{(\ell)}, \mathcal{T}_P, \mathcal{G}_k)^K$ for any $K > 0$, then $\hat{\mu}_\infty^{(\ell)} \Rightarrow P$ as $\ell \rightarrow \infty$.

Proof. Using the argument in (3.2) and the fact that the inverse of $\mathbf{T}^{\mu, \nu, \epsilon}$ is $\mathbf{T}^{\mu, \nu, -\epsilon}$, we have

$$\begin{aligned} \text{KL}(\mu_\infty^{(\ell+1)} \| P) - \text{KL}(\mu_\infty^{(\ell)} \| P) &= \text{KL}(\mathbf{T}^{\mu_\infty^{(\ell)}, P, \epsilon_\ell} \mu_\infty^{(\ell)} \| P) - \text{KL}(\mu_\infty^{(\ell)} \| P) \\ &= \text{KL}(\mu_\infty^{(\ell)} \| \mathbf{T}_\#^{\mu_\infty^{(\ell)}, P, -\epsilon_\ell} P) - \text{KL}(\mu_\infty^{(\ell)} \| P) \\ &= -\mathbb{E}_{\mathbf{x} \sim \mu_\infty^{(\ell)}} \left[\underbrace{\log |\det(\nabla_{\mathbf{x}} \mathbf{T}^{\mu_\infty^{(\ell)}, P, \epsilon_\ell}(\mathbf{x}))|}_{(1)} + \underbrace{\log p(\mathbf{T}^{\mu_\infty^{(\ell)}, P, \epsilon_\ell}(\mathbf{x})) - \log p(\mathbf{x})}_{(2)} \right], \quad (3.5) \end{aligned}$$

where the final equality follows from applying the change of variables formula for the transformation $\mathbf{T}^{\mu_\infty^{(\ell)}, P, -\epsilon_\ell}$. For term (1), take $B := \nabla_{\mathbf{x}} \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})$ in Lemma 3.5 and $\epsilon < 1/(2\rho(B + B^\top))$. Then

$$\log |\det(\nabla_{\mathbf{x}} \mathbf{T}^{\mu_\infty^{(\ell)}, P, \epsilon_\ell}(\mathbf{x}))| \geq \epsilon \text{Trace}(\nabla_{\mathbf{x}} \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})) - 2\epsilon^2 \|\nabla_{\mathbf{x}} \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})\|_F^2.$$

For term (2), by defining $\mathbf{x}_s := \mathbf{T}^{\mu_\infty^{(\ell)}, P, s\epsilon_\ell}(\mathbf{x})$ for $s \in [0, 1]$ and applying Taylor's theorem, we get

$$\begin{aligned} \log p(\mathbf{x}) - \log p(\mathbf{T}^{\mu_\infty^{(\ell)}, P, \epsilon_\ell}(\mathbf{x})) &= - \int_0^1 \nabla_s \log p(\mathbf{x}_s) ds \\ &= - \int_0^1 \nabla_{\mathbf{x}} \log p(\mathbf{x}_s)^\top (\epsilon \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})) ds \\ &= -\epsilon \nabla_{\mathbf{x}} \log p(\mathbf{x})^\top \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x}) - \int_0^1 (\nabla_{\mathbf{x}} \log p(\mathbf{x}_s) - \nabla_{\mathbf{x}} \log p(\mathbf{x}))^\top (\epsilon \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})) ds \\ &\leq -\epsilon \mathbf{s}_p(\mathbf{x})^\top \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x}) + \epsilon^2 M_1(\mathbf{s}_p) \cdot \|\tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})\|_2^2 \int_0^1 ds \\ &= -\epsilon \mathbf{s}_p(\mathbf{x})^\top \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x}) + \frac{\epsilon^2}{2} M_1(\mathbf{s}_p) \cdot \|\tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})\|_2^2. \end{aligned}$$

Thus (3.5) can be upper bounded as follows:

$$\begin{aligned} \text{KL}(\mu_\infty^{(\ell+1)} \| P) - \text{KL}(\mu_\infty^{(\ell)} \| P) &\leq -\epsilon \left(\mathbf{s}_p(\mathbf{x})^\top \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x}) + \nabla_{\mathbf{x}} \cdot \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x}) \right) + \Delta \\ &= -\epsilon \mathcal{S}^2(\mu_\infty^{(\ell)}, \mathcal{T}_P, \mathcal{G}_k) + \Delta, \end{aligned}$$

by definition of $\tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})$, where

$$\Delta := \epsilon^2 \mathbb{E}_{\mu_\infty^{(\ell)}} \left[\frac{1}{2} M_1(\mathbf{s}_p) \cdot \|\tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})\|_2^2 + 2 \|\nabla_{\mathbf{x}} \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})\|_F^2 \right].$$

We now want to bound the terms in Δ . Letting $\tilde{\phi}_{\mu_\infty^{(\ell)}, p} = (\phi_1, \dots, \phi_d)^\top$, by the Cauchy-Schwarz inequality,

$$\begin{aligned} \|\tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})\|_2^2 &= \sum_{i=1}^d \phi_i(\mathbf{x})^2 = \sum_{i=1}^n (\langle k(\mathbf{x}, \cdot), \phi_i(\cdot) \rangle_{\mathcal{H}_k})^2 \\ &\leq \sum_{i=1}^d \|k(\mathbf{x}, \cdot)\|_{\mathcal{H}_k}^2 \cdot \|\phi_i\|_{\mathcal{H}_k}^2 \\ &= k(\mathbf{x}, \mathbf{x}) \cdot \|\tilde{\phi}_{\mu_\infty^{(\ell)}, p}\|_{\mathcal{H}_k^d}^2 = c^{2\beta} \mathcal{S}^2(\mu_\infty^{(\ell)}, \mathcal{T}_P, \mathcal{G}_k), \end{aligned}$$

for the IMQ kernel k . Furthermore,

$$\begin{aligned} \|\nabla_{\mathbf{x}} \tilde{\phi}_{\mu_\infty^{(\ell)}, p}(\mathbf{x})\|_F^2 &= \sum_{i,j=1}^d \partial_{x_j} \phi_i(\mathbf{x})^2 = \sum_{i,j=1}^d (\langle \partial_{x_j} k(\mathbf{x}, \cdot), \phi_i(\cdot) \rangle_{\mathcal{H}_k})^2 \\ &\leq \sum_{i,j=1}^d \|\partial_{x_j} k(\mathbf{x}, \cdot)\|_{\mathcal{H}_k}^2 \cdot \|\phi_i\|_{\mathcal{H}_k}^2 \\ &= \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{x}'} k(\mathbf{x}, \mathbf{x}') \Big|_{\mathbf{x}=\mathbf{x}'} \sum_{i=1}^d \|\phi_i\|_{\mathcal{H}_k}^2 \\ &= -2d\beta c^{2(\beta-1)} \|\tilde{\phi}_{\mu_\infty^{(\ell)}, p}\|_{\mathcal{H}_k^d}^2 \leq 2dc^{2(\beta-1)} \mathcal{S}^2(\mu_\infty^{(\ell)}, \mathcal{T}_P, \mathcal{G}_k) \end{aligned}$$

for $\beta \in (-1, 0)$. Substituting these bounds into Δ , we conclude that

$$\text{KL}(\mu_\infty^{(\ell+1)} \| P) - \text{KL}(\mu_\infty^{(\ell)} \| P) \leq -\epsilon(1 - \epsilon R) \mathcal{S}^2(\mu_\infty^{(\ell)}, \mathcal{T}_P, \mathcal{G}_k). \quad \square$$

Combining the results of Theorems 3.3 and 3.6, we have that $\mu_n^{(\ell)} \Rightarrow P$ as $n, \ell \rightarrow \infty$ with sufficiently small step sizes, and hence that SVGD can approximate the target arbitrarily well. This is a very attractive property to have in an approximate inference algorithm which is not true of standard parametric variational inference approaches; however, it is an asymptotic result, and so depends deeply on its scalability. This is our focus in Chapter 4.

Chapter 4

Scalable Stein Discrepancies

Despite its satisfying theoretical results, the IMQ KSD presents issues in terms of its applicability to large-scale problems. For example, consider the case where the target is the posterior

$$p(\mathbf{x}|\mathcal{D}) = p_0(\mathbf{x}) \prod_{r=1}^R p(\mathcal{D}_r|\mathbf{x}),$$

where $\mathbf{x} \in \mathbb{R}^d$ and $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_R\}$ is a (potentially very large) set of training data. Assuming each likelihood $p(\mathcal{D}_r|\mathbf{x})$ can be evaluated cheaply enough, and letting $Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}(\cdot)$, the computational cost of calculating $\mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k)$ is $\mathcal{O}(n^2 d R)$. In the context of SVGD, at each update we must compute the optimal direction, at cost $\mathcal{O}(ndR)$, n times, thus also giving a complexity of $\mathcal{O}(n^2 d R)$ per update. This setting presents us with several limitations:

1. The most immediately obvious scalability issue is that the computational cost is quadratic in n , the number of samples. In Section 4.1, we introduce a *random feature Stein discrepancy* [Huggins and Mackey, 2018] which is inspired by the IMQ KSD and is computable in near-linear time, as a remedy for this problem.
2. Secondly, the ability to perform reliable inference in higher dimensions is dependent on having sufficient data to do so. Consequently, although a computation time linear in d is in itself acceptable, R will typically have to scale with d , which in turn impedes both the evaluation of the KSD and the application of SVGD to large-scale Bayesian inference. We discuss *stochastic Stein discrepancies* [Gorham et al., 2020] as a solution to this issue in Section 4.2.

There is another potential scalability issue, which we will not discuss here. It is not one of computational complexity, but instead one of degrading performance in higher dimensions:

3. Reddi et al. [2014] demonstrated that many methods based on evaluating kernels and distances in high dimensions exhibit rapidly deteriorating performance. In this sense, the IMQ KSD presents a double threat: it involves evaluating the IMQ kernel on high-dimensional inputs, and it involves computing high-dimensional inner products. A recent proposal for remedying this problem has been to *slice* the KSD, i.e. to project both the input and the score function onto one dimensional “slices” [Gong et al., 2021]. It is not clear, however, as to whether or not sliced KSDs retain the convergence-monitoring properties of the regular IMQ KSD.

4.1 Random feature Stein discrepancies

We now let $\Phi : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{C}$ be a feature function such that $\Phi(\mathbf{x}, \cdot) \in L^r$ and $\Phi(\cdot, \mathbf{z}) \in C^1$ for all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$ and some $r \in [1, \infty)$. Rather than taking the Stein set to be the unit ball of some product Hilbert space, we

instead let

$$\mathcal{G}_{\Phi,r} := \left\{ \mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^d \mid g_i(\mathbf{x}) = \int \Phi(\mathbf{x}, \mathbf{z}) \overline{f_i(\mathbf{z})} d\mathbf{z} \text{ with } \sum_{i=1}^d \|f_i\|_{L^s}^2 \leq 1, \text{ where } s = \frac{r}{r-1} \right\}.$$

We can now define the (squared) **feature Stein discrepancy** (FSD) which, analogously to the KSD, has an analytic solution:

$$\begin{aligned} F_{\Phi,r} \mathcal{S}^2(\mu, P) &:= \sup_{\mathbf{g} \in \mathcal{G}_{\Phi,r}} |\mathbb{E}_\mu[(\mathcal{T}_P \mathbf{g})(\mathbf{X})]|^2 = \sup_{\mathbf{g} \in \mathcal{G}_{\Phi,r}} \left| \sum_{i=1}^d \mathbb{E}_\mu[(\mathcal{T}_P^{(i)} g_d)(\mathbf{X})] \right|^2 \\ &= \sup_{\mathbf{f} : v_d = \|f_d\|_{L^s}, \|\mathbf{v}\|_2 \leq 1} \left| \sum_{i=1}^d \int \mathbb{E}_\mu[(\mathcal{T}_P^{(i)} \Phi)(\mathbf{X}, \mathbf{z})] \overline{f_i(\mathbf{z})} d\mathbf{z} \right|^2 \\ &= \sup_{\mathbf{v} : \|\mathbf{v}\|_2 \leq 1} \left| \sum_{i=1}^d \|\mathbb{E}_\mu[(\mathcal{T}_P^{(i)} \Phi)(\mathbf{X}, \cdot)]\|_{L^r} v_d \right|^2 = \sum_{i=1}^d \|\mathbb{E}_\mu[(\mathcal{T}_P^{(i)} \Phi)(\mathbf{X}, \cdot)]\|_{L^r}^2. \end{aligned}$$

To reduce the cost of large n , we can approximate $F_{\Phi,r} \mathcal{S}^2(\mu, P)$ using an importance sample of size M :

$$RF_{\Phi,r,\nu,M} \mathcal{S}^2(\mu, P) := \sum_{i=1}^d \left(\frac{1}{M} \sum_{j=1}^M \nu(\mathbf{Z}_j)^{-1} |\mathbb{E}_{\mathbf{X} \sim \mu}[(\mathcal{T}_P^{(i)} \Phi)(\mathbf{X}, \mathbf{Z}_j)]|^r \right)^{2/r}, \quad \text{where } \mathbf{Z}_1, \dots, \mathbf{Z}_M \stackrel{\text{i.i.d.}}{\sim} \nu.$$

We call this the (squared) **random feature Stein discrepancy** (RFSD). In what follows, we intend to give an outline of the theory which justifies using the RFSD to monitor the convergence of a discrete measure Q_n to a target P .

4.1 Proposition (KSD-FSD inequality). *Let $k(\mathbf{x}, \mathbf{y}) = \int \mathcal{F}(\Phi(\mathbf{x}, \cdot))(\omega) \overline{\mathcal{F}(\Phi(\mathbf{y}, \cdot))(\omega)} \rho(\omega) d\omega$ be a kernel, where \mathcal{F} denotes the generalised Fourier transform, and let $r \in [1, 2]$. Setting $t := r/(2-r)$, if $\rho \in L^t$, then we have*

$$\mathcal{S}^2(Q_n, \mathcal{T}_P, \mathcal{G}_k) \leq \|\rho\|_{L^t} F_{\Psi_{c,\beta},r} \mathcal{S}^2(Q_n, P).$$

Proof. Let $s = r/(r-1)$. Using Hölder's inequality and the Babenko-Beckner inequality, we have

$$\begin{aligned} \mathcal{S}^2(Q_n, \mathcal{T}_P, \mathcal{G}_k) &= \sum_{i=1}^d \int |\mathcal{F}(\mathbb{E}_{\mathbf{X} \sim Q_n}[(\mathcal{T}_P^{(i)} \Phi)(\mathbf{X}, \cdot)])|(\omega)|^2 \rho(\omega) d\omega \leq \|\rho\|_{L^t} \sum_{i=1}^d \|\mathcal{F}(\mathbb{E}_{\mathbf{X} \sim Q_n}[(\mathcal{T}_P^{(i)} \Phi)(\mathbf{X}, \cdot)])\|_{L^s}^2 \\ &\leq (r^{1/r}/s^{1/s})^d \|\rho\|_{L^t} \sum_{i=1}^d \|\mathbb{E}_{\mathbf{X} \sim Q_n}[(\mathcal{T}_P^{(i)} \Phi)(\mathbf{X}, \cdot)]\|_{L^r}^2 = (r^{1/r}/s^{1/s})^d \|\rho\|_{L^t} F_{\Phi,r} \mathcal{S}^2(Q_n, P), \end{aligned}$$

since $(r^{1/r}/s^{1/s})^{d/2} \leq 1$ for $s = r/(r-1)$. \square

As a consequence of this proposition, if we choose k to be the IMQ kernel with $c > 0$ and $\beta \in (-1, 0)$, then any Φ satisfying the condition in Proposition 4.1 will yield a FSD which detects non-convergence. Since $k(\mathbf{x}, \mathbf{y}) = \Psi_{c,\beta}(\mathbf{x} - \mathbf{y})$, it makes sense to take $\Phi(\mathbf{x}, \mathbf{y}) = F(\mathbf{x} - \mathbf{y})$ for some F which satisfies the following regularity conditions.

4.2 Assumption. *The function $F \in C^1$ is positive, and there exists a norm $\|\cdot\|$ and constants $s, C \geq 0$ such that*

$$|\partial_{x_i} \log F(\mathbf{x})| \leq C(1 + \|\mathbf{x}\|^s), \quad \lim_{\|\mathbf{x}\| \rightarrow \infty} (1 + \|\mathbf{x}\|^s) F(\mathbf{x}) = 0, \quad \text{and } F(\mathbf{x} - \mathbf{z}) \leq C F(\mathbf{z}) / F(\mathbf{x}).$$

Furthermore, there is a constant $\underline{c} \in (0, 1]$ and a continuous, non-increasing function f such that $\underline{c} f(\|\mathbf{x}\|) \leq F(\mathbf{x}) \leq f(\|\mathbf{x}\|)$.

How do we choose F ? An intuitive choice would be $F = \mathcal{F}^{-1}(\hat{\Psi}_{c,\beta}^{1/2})$, so that $\Psi_{c,\beta}(\mathbf{x} - \mathbf{y}) = \int F(\mathbf{x} - \mathbf{z}) F(\mathbf{y} - \mathbf{z}) d\mathbf{z}$, and in particular, $F_{\Phi,2} \mathcal{S}(\mu, P) = \mathcal{S}(\mu, \mathcal{T}_P, \mathcal{G}_k)$. However, this choice of F is difficult to compute in practice, so we instead only require F to be a good approximation to $\mathcal{F}^{-1}(\hat{\Psi}_{c,\beta})$.

4.3 Assumption. *Assumption 4.2 holds, and there exists a smoothness parameter $\bar{\lambda} \in (1/2, 1]$ such that $\hat{F}/\hat{\Psi}_{c,\beta}^{\lambda/2} \in L^2$ for all $\lambda \in (1/2, \bar{\lambda})$.*

The following proposition shows that certain RFSDs detect weak convergence, in probability.

4.4 Proposition. *Suppose that Assumption 4.2 holds with $F \in L^r$ and that $\mathbb{E}_{\mathbf{X} \sim P}[\|\mathbf{X}\|_2^2] < \infty$. If $W_1(Q_n, P) \rightarrow 0$ as $n \rightarrow \infty$, then $\text{F}_{\Phi,r}\mathcal{S}(Q_n, P) \rightarrow 0$ and $\text{RF}_{\Phi,r,\nu_n,M_n}\mathcal{S}(Q_n, P) \xrightarrow{P} 0$ for any $r \in [1, 2]$, ν_n and M_n .*

The aim when choosing ν is to ensure that the second moment of each RFSD feature

$$w_i(\mathbf{Z}, Q_n) := |\mathbb{E}_{\mathbf{X} \sim Q_n}[(\mathcal{T}_P^{(i)}\Phi)(\mathbf{X}, \mathbf{Z})]|^r / \nu(\mathbf{Z})$$

is bounded by a power of its mean, as outlined by the following definition.

4.5 Definition $((C, \gamma)$ second moments). *We say that (Φ, r, ν) yields (C, γ) second moments for P and Q_n if $\mathbb{E}[w_i(\mathbf{Z}, Q_n)^2] \leq C\mathbb{E}[w_i(\mathbf{Z}, Q_n)]^{2-\gamma}$ for all $i = 1, \dots, d$ and $n \geq 1$.*

4.6 Proposition. *Suppose that (Φ, r, ν) yields (C, γ) second moments for P and Q_n . Under the assumptions of Proposition 4.1, if the reference KSD $\mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k) \gtrsim n^{-1/2}$, and $M \gtrsim n^{\gamma r/2}C\|\rho\|_{L^t}^{\gamma r/2}\log(d/\delta)/\epsilon^2$, then with probability at least $1 - \delta$,*

$$\|\rho\|_{L^t}^{1/2}RF_{\Phi,r,\nu,M}\mathcal{S}(Q_n, P) \geq (1 - \epsilon)^{1/r}\mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k).$$

The following final assumption ensures in particular that Proposition 4.1, and hence Proposition 4.6, applies.

4.7 Assumption. *Assumption 4.2 holds, $\omega_1^2\hat{\Psi}_{c,\beta}^{1/2}(\boldsymbol{\omega}) \in L^1$, and $\hat{\Psi}_{c,\beta}/\hat{F}^2 \in L^t$, where $t = r/(2-r)$.*

Huggins and Mackey [2018] proceed to show that, if Assumptions 4.2, 4.3 and 4.7 hold, and there exists $C > 0$ such that

$$\mathbb{E}_{\mathbf{X} \sim Q_n}[(1 + \|\mathbf{X}\| + \|\mathbf{X} - \mathbb{E}_{\mathbf{X}' \sim Q_n}[\mathbf{X}']\|^s)/F(\mathbf{X} - \mathbb{E}_{\mathbf{X}' \sim Q_n}[\mathbf{X}'])] \leq C \quad \text{for all } n \geq 1,$$

then there exists $b \in [0, 1)$ such that for any $\xi \in (0, 1-b)$, $c > 0$ and $\alpha > 2(1-\bar{\lambda})$, if $\nu(\mathbf{z}) \geq c\Psi_{c,\beta}(\mathbf{z} - \mathbb{E}_{\mathbf{X} \sim Q_n}[\mathbf{X}])^{\xi r}$, then there exists a constant $C_\alpha > 0$ such that (Φ, r, ν) yields $(C_\alpha, \gamma_\alpha)$ second moments for P and Q_n , where $\gamma_\alpha := \alpha + (2-\alpha)\xi/(2-b-\xi)$. Thus, by increasing the smoothness $\bar{\lambda}$ of F and decreasing the overdispersion parameter ξ of ν , we can make γ arbitrarily close to 0, and hence yield RFSDs which are computable in arbitrarily-close-to-linear time.

A particular case, which they show to satisfy the necessary assumptions, is choosing $\bar{\lambda}$ and $\underline{\xi} \in (0, 1/2)$ – the minimum ξ that we can choose when constructing ν – and then taking $F = \Psi_{c',\beta'}$, where $c' = \bar{\lambda}c/2$, $\beta' \in [-d/(2\underline{\xi}), -\beta/(2\underline{\xi}) - d/(2\underline{\xi})]$, $r = -d/(2\beta'\underline{\xi})$, $\xi \in (\underline{\xi}, 1)$ and $\nu(\mathbf{z}) \propto \Psi_{c',\beta'}(\mathbf{z} - \mathbb{E}_{\mathbf{X} \sim Q_n}[\mathbf{X}])^{\xi r}$. This is known as the L^r IMQ RFSD. In particular, taking $\beta' = -d/(2\underline{\xi})$ yields $r = 1$. The importance sampling distribution ν in this case is easy to sample from since it is a multivariate t -distribution.

4.2 Stochastic Stein discrepancies

Consider the case where our target is the posterior

$$p(\mathbf{x} | \mathcal{D}) \propto p_0(\mathbf{x}) \prod_{r=1}^R p(\mathcal{D}_r | \mathbf{x}).$$

We can decompose the Langevin Stein operator as $\mathcal{T}_P = \sum_{r=1}^R \mathcal{T}_P^{(r)}$, where

$$(\mathcal{T}_P^{(r)} \mathbf{g})(\mathbf{x}) = \mathbf{s}_{p_r}(\mathbf{x})^\top \mathbf{g}(\mathbf{x}) + \frac{1}{R} \nabla_{\mathbf{x}} \cdot \mathbf{g}(\mathbf{x}), \quad \text{and } p_r(\mathbf{x}) := p_0(\mathbf{x})^{1/R} p(\mathcal{D}_r | \mathbf{x}).$$

(Note that this is different to the decomposition that we have used previously.) The idea behind reducing the cost induced by large R is to randomly subsample the base operators $\mathcal{T}_P^{(r)}$. To this end, for the discrete measure $Q_n(\cdot) := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}(\cdot)$, we define the **stochastic Stein discrepancy** as

$$\text{stoch}_m \mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}) := \sup_{\mathbf{g} \in \mathcal{G}} \left| \frac{1}{n} \sum_{i=1}^n \frac{R}{m} (\mathcal{T}_{\sigma_i} \mathbf{g})(\mathbf{x}_i) \right|,$$

where m is our chosen batch size and the $\sigma_i \subseteq \{1, \dots, R\}$ are subsets of size m sampled uniformly at random (without replacement). By using arguments analogous to those used in Section 1.2, one can show that taking the Stein set to be the unit ball of a product Hilbert space, i.e. $\mathcal{G} = \mathcal{G}_k$, again yields an analytic solution to the optimisation over $\mathbf{g} \in \mathcal{G}$:

$$\text{stoch}_m \mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k) = \frac{1}{n^2} \sum_{i,j=1}^n k_p^{\sigma_i, \sigma_j}(\mathbf{x}_i, \mathbf{x}_j),$$

where

$$k_p^{\sigma, \tilde{\sigma}}(\mathbf{x}, \mathbf{y}) := \frac{R^2}{m^2} \mathbf{s}_{p_\sigma}(\mathbf{x})^\top \mathbf{s}_{p_{\tilde{\sigma}}}(\mathbf{y}) k(\mathbf{x}, \mathbf{y}) + \frac{R}{m} \mathbf{s}_{p_\sigma}(\mathbf{x})^\top \nabla_{\mathbf{y}} k(\mathbf{x}, \mathbf{y}) + \frac{R}{m} \mathbf{s}_{p_{\tilde{\sigma}}}(\mathbf{y})^\top \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y}) + \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{y}} k(\mathbf{x}, \mathbf{y}).$$

The stochastic IMQ KSD metrises weak convergence in a general sense, as outlined by the following theorems, whose proofs are given in Gorham et al. [2020].

4.8 Theorem (Stochastic IMQ KSD detects convergence). *If $P \in \mathcal{P}$, then $\text{stoch}_m \mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k) \rightarrow 0$ whenever $Q_n \Rightarrow P$.*

4.9 Theorem (Stochastic IMQ KSD detects non-convergence). *Let $P \in \mathcal{P}$, and suppose that $\sup_{\mathbf{x} \in \mathcal{X}} \frac{\|\log p_\sigma(\mathbf{x})\|_2}{1 + \|\mathbf{x}\|_2} < \infty$ for all $\sigma \in \binom{[R]}{m}$. If $Q_n \not\Rightarrow P$ as $n \rightarrow \infty$, then $\text{stoch}_m \mathcal{S}(Q_n, \mathcal{T}_P, \mathcal{G}_k) \not\rightarrow 0$.*

4.2.1 Stochastic SVGD

As mentioned above, using SVGD to approximate the posterior $p(\mathbf{x}|\mathcal{D})$ requires $\mathcal{O}(n^2 dR)$ operations per update. The stochastic KSD presents us with an alternative algorithm, **stochastic SVGD**, which requires $\mathcal{O}(n^2 dm)$ operations per update, with $m \ll R$. Rather than using the optimal direction

$$\tilde{\phi}_{\mu_n^{(\ell)}, p}(\cdot) = \frac{1}{n} \sum_{j=1}^n \left\{ \mathbf{s}_p(\mathbf{x}_j^{(\ell)}) k(\mathbf{x}_j^{(\ell)}, \cdot) + \nabla_{\mathbf{x}_j^{(\ell)}} k(\mathbf{x}_j^{(\ell)}, \cdot) \right\},$$

as in Algorithm 3.1, we instead use the stochastic direction

$$\tilde{\phi}_{\hat{\mu}_{n,m}^{(\ell)}, p}^{\text{stoch}}(\cdot) = \frac{1}{n} \sum_{j=1}^n \left\{ \frac{R}{m} \mathbf{s}_{p_{\sigma_j}}(\mathbf{x}_j^{(\ell)}) k(\mathbf{x}_j^{(\ell)}, \cdot) + \nabla_{\mathbf{x}_j^{(\ell)}} k(\mathbf{x}_j^{(\ell)}, \cdot) \right\},$$

resulting in Algorithm 4.1. For n particles and batch size m , we denote the discrete measure at time ℓ by $\hat{\mu}_{n,m}^{(\ell)}$, and the stochSVGD update by

$$\hat{\mu}_{n,m}^{(\ell+1)} = \Phi_{\epsilon, p}^m(\hat{\mu}_{n,m}^{(\ell)}).$$

Algorithm 4.1: Stochastic Stein Variational Gradient Descent (stochSVGD)

- 1 **Input:** The score function $\mathbf{s}_p(\mathbf{x})$, initial particles $\{\mathbf{x}_i^{(0)}\}_{i=1}^n$, kernel k , batch size m , number of iterations L and step sizes $\{\epsilon_\ell\}_{\ell=0}^{L-1}$;
 - 2 **for** $\ell = 0, \dots, L-1$ **do**
 - 3 For $i = 1, \dots, n$, sample an independent mini-batch σ_i of size m from $\{1, \dots, R\}$;
 - 4 $\mathbf{x}_i^{(\ell+1)} \leftarrow \mathbf{x}_i^{(\ell)} + \epsilon_\ell \tilde{\phi}_{\hat{\mu}_{n,m}^{(\ell)}, p}^{\text{stoch}}(\mathbf{x}_i^{(\ell)})$, $i = 1, \dots, n$, where
 - 5
$$\tilde{\phi}_{\hat{\mu}_{n,m}^{(\ell)}, p}^{\text{stoch}}(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n \left\{ \frac{R}{m} \mathbf{s}_{p_{\sigma_j}}(\mathbf{x}_j^{(\ell)}) k(\mathbf{x}_j^{(\ell)}, \mathbf{x}) + \nabla_{\mathbf{x}_j^{(\ell)}} k(\mathbf{x}_j^{(\ell)}, \mathbf{x}) \right\}.$$
-

We will now prove that this new algorithm converges almost surely as $n, \ell \rightarrow \infty$. Note that the infinite particle sequence $\{\mu_\infty^{(\ell)}\}_{\ell \geq 0}$ is the same as before, so we only need to show that $\hat{\mu}_{n,m}^{(\ell)} \Rightarrow \mu_\infty^{(\ell)}$ almost surely as $n \rightarrow \infty$, for each ℓ .

4.10 Theorem (Convergence of stochSVGD). *Consider the stochSVGD algorithm as outlined in Algorithm 4.1, and suppose that $W_1(\hat{\mu}_{n,m}^{(0)}, \mu_\infty^{(0)}) \xrightarrow{n \rightarrow \infty} 0$. Suppose that condition (3.4) holds from the original proof, along with the additional conditions that, for some $C_0 > 0$,*

$$\max_{\sigma \in \binom{[R]}{m}} \sup_{\mathbf{y} \in \mathbb{R}^d} \|\mathbf{s}_{p_\sigma}(\mathbf{x}) k(\mathbf{x}, \mathbf{y})\|_2 \leq C_0(1 + \|\mathbf{x}\|_2), \quad (4.1)$$

$$\max_{\sigma \in \binom{[R]}{m}} \sup_{\mathbf{y} \in \mathbb{R}^d} \|\nabla_{\mathbf{x}}(\mathbf{s}_{p_\sigma}(\mathbf{x}) k(\mathbf{x}, \mathbf{y}))\|_{op} \text{ is bounded on compact sets } J, \quad (4.2)$$

where $\binom{[R]}{m}$ is the set of all subsets of R of size m . Then $W_1(\hat{\mu}_{n,m}^{(\ell)}, \mu_\infty^{(\ell)}) \xrightarrow{a.s.} 0$ as $n \rightarrow \infty$.

Proof. We prove this by induction. The base case $\ell = 0$ is satisfied by assumption. Now, for any $\ell \geq 0$, assume that $\mathbb{P}(\mathcal{E}) = 1$, where \mathcal{E} is the event that $W_1(\hat{\mu}_{n,m}^{(\ell)}, \mu_n^{(\ell)}) \rightarrow 0$ as $n \rightarrow \infty$. Then, since $W_1(\mu_n^{(\ell)}, \mu_\infty^{(\ell)}) \rightarrow 0$ as $n \rightarrow \infty$ by Theorem 3.3, we have $W_1(\hat{\mu}_{n,m}^{(\ell)}, \mu_\infty^{(\ell)}) \rightarrow 0$ as $n \rightarrow \infty$ on \mathcal{E} . Therefore, there exists a constant C' such that

$$\sup_{n \geq 1} 1 + \epsilon_\ell C_1(1 + \hat{\mu}_{n,m}^{(\ell)}(\|\cdot\|_2)) + \epsilon_\ell C_2(1 + \mu_\infty^{(\ell)}(\|\cdot\|_2)) \leq C'$$

on \mathcal{E} . By the triangle inequality,

$$\begin{aligned} W_1(\hat{\mu}_{n,m}^{(\ell+1)}, \mu_n^{(\ell+1)}) &= W_1(\Phi_{\epsilon_\ell, p}^m(\hat{\mu}_{n,m}^{(\ell)}), \Phi_{\epsilon_\ell, p}(\mu_n^{(\ell)})) \\ &\leq W_1(\Phi_{\epsilon_\ell, p}^m(\hat{\mu}_{n,m}^{(\ell)}), \Phi_{\epsilon_\ell, p}(\hat{\mu}_{n,m}^{(\ell)})) + W_1(\Phi_{\epsilon_\ell, p}(\hat{\mu}_{n,m}^{(\ell)}), \Phi_{\epsilon_\ell, p}(\mu_n^{(\ell)})). \end{aligned}$$

By the assumptions (3.4) and Lemma 3.4, on \mathcal{E} we have

$$\begin{aligned} W_1(\Phi_{\epsilon_\ell, p}(\hat{\mu}_{n,m}^{(\ell)}), \Phi_{\epsilon_\ell, p}(\mu_n^{(\ell)})) &\leq W_1(\hat{\mu}_{n,m}^{(\ell)}, \mu_n^{(\ell)})(1 + \epsilon_\ell C_1(1 + \hat{\mu}_{n,m}^{(\ell)}(\|\cdot\|_2)) + \epsilon_\ell C_2(1 + \mu_\infty^{(\ell)}(\|\cdot\|_2))) \\ &\leq C' W_1(\hat{\mu}_{n,m}^{(\ell)}, \mu_n^{(\ell)}) \rightarrow 0 \text{ as } n \rightarrow \infty. \end{aligned}$$

Furthermore, using the growth assumptions (4.1) and (4.2), we can apply Lemma 4.11 below, which ensures that $W_1(\Phi_{\epsilon_\ell, p}^m(\hat{\mu}_{n,m}^{(\ell)}), \Phi_{\epsilon_\ell, p}(\hat{\mu}_{n,m}^{(\ell)})) \xrightarrow{a.s.} 0$ as $n \rightarrow \infty$. Hence, on \mathcal{E} , $W_1(\hat{\mu}_{n,m}^{(\ell+1)}, \mu_n^{(\ell+1)}) \rightarrow 0$ as $n \rightarrow \infty$, thus proving our claim. \square

4.11 Lemma. *For any triangular array of points $\{\mathbf{z}_i^n\}_{i \in \{1, \dots, n\}, n \geq 1}$ in \mathbb{R}^d , define the discrete probability measure $\nu_n(\cdot) := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{z}_i^n}(\cdot)$. Suppose that $\mathbf{s}_{p_\sigma}(\cdot) k(\cdot, \mathbf{y})$ is continuous for each $\mathbf{y} \in \mathbb{R}^d$ and $\sigma \in \binom{[R]}{m}$. Set*

$$\begin{aligned} f_0(\mathbf{x}) &:= \sup_{\mathbf{y} \in \mathbb{R}^d, \sigma \in \binom{[R]}{m}} M_0(\mathbf{s}_{p_\sigma}) |k(\mathbf{x}, \mathbf{y})|, \\ f_1(\mathbf{x}) &:= \sup_{\mathbf{y} \in \mathbb{R}^d, \sigma \in \binom{[R]}{m}} \|\nabla_{\mathbf{x}}(\mathbf{s}_p(\mathbf{x}) k(\mathbf{x}, \mathbf{y}))\|_{op}. \end{aligned}$$

If f_0 is ν_n -uniformly integrable and f_0, f_1 are bounded on every compact set, then for any $\epsilon > 0$,

$$W_1(\Phi_{\epsilon, p}^m(\nu_n), \Phi_{\epsilon, p}(\nu_n)) \xrightarrow{a.s.} 0$$

as $n \rightarrow \infty$.

This lemma is proved in Lemma 13 of Gorham et al. [2020]. It is again important to highlight that the additional growth conditions (4.1) and (4.2) are satisfied by the IMQ kernel and P with Lipschitz score function \mathbf{s}_p .

Chapter 5

Experiments

In this chapter, we demonstrate the validity of our results in practice. Section 5.1 compares the various Stein discrepancies that we have introduced and verifies that they behave as we expect, while also introducing an application of these discrepancies in the form of hyperparameter selection in a biased MCMC algorithm. Section 5.2 focuses on approximate inference, and in particular the application of SVGD and stochSVGD as variational inference algorithms. The code is available at https://github.com/alex-hinds/stein_discrepancies.

5.1 Measuring sample quality

In this section, we will compare the different Stein discrepancies we have considered in terms of their ability to monitor convergence in various settings.

5.1.1 Comparison to the Wasserstein-1 distance in one dimension

To check that our Stein discrepancies behave as we expect, we construct a very simple experiment influenced by Section 4.1 of Gorham and Mackey [2017]. We generate synthetic data $\{\mathbf{x}_1, \dots, \mathbf{x}_R\}$ from the $\mathcal{N}(\alpha \mathbf{e}_1, I_D)$ distribution

$$p(\mathbf{x}|\alpha) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2}\|\mathbf{x}-\alpha\mathbf{e}_1\|_2^2},$$

where α is a parameter, taken to be 1 when generating the data. Assuming an improper uninformative prior over α , our posterior is

$$p(\alpha|\mathcal{D}) \propto \prod_{r=1}^R e^{-\frac{1}{2}\|\mathbf{x}_r - \alpha\mathbf{e}_1\|_2^2} \propto \exp\left\{-\frac{1}{2}\left(-2\alpha \sum_{r=1}^R x_{r1} + R\alpha^2\right)\right\} \sim \mathcal{N}\left(\frac{1}{R} \sum_{r=1}^R x_{r1}, \frac{1}{R}\right).$$

The purpose of this example is purely to demonstrate the efficacy of our Stein discrepancies to monitor convergence. Of course, in this example, we can perform exact inference and there is no need for approximate sampling. What's more, since we can sample efficiently from the posterior, we shouldn't expect any computational benefits from using the stochastic IMQ KSD or the L^1 IMQ RFSD over the standard IMQ KSD. Nevertheless, the fact that this posterior is one-dimensional and exact allows us to compare the IMQ KSD and its stochastic variant with the Wasserstein-1 distance.

To make such a comparison, we generate sequences of i.i.d. random variables sampled from

$$\tilde{p}_\gamma(\alpha|\mathcal{D}) \sim \mathcal{N}\left(\gamma \times \frac{1}{R} \sum_{r=1}^R x_{r1}, \frac{1}{R}\right),$$

where $\gamma \in \mathbb{R}$. Clearly, when $\gamma = 1$, the sequence will converge weakly to the target $p(\alpha|\mathcal{D})$, but when $\gamma \neq 1$ it will not. We are interested in how effective our Stein discrepancies are in detecting small deviations of γ .

from 1, particularly in comparison with the Wasserstein-1 distance.

We plot the values of the IMQ KSD, the Stochastic IMQ KSD with $m = 0.1R$ and $m = 0.2R$, the L^1 IMQ RFSD and the Wasserstein distance, for the sequence with three different values of γ and $R = 10$. We took the median over 10 evaluations for the RFSD and took $M = 100$.

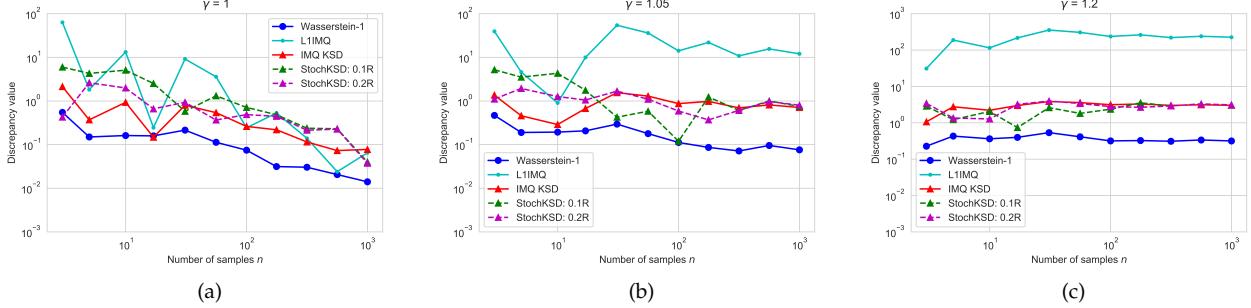


Figure 5.1: Comparing KSDs to the Wasserstein-1 distance.

We see in Figure 5.1 that, as expected, all Stein discrepancies behave similarly to the Wasserstein-1 distance in detecting convergence in the case $\gamma = 1$. For the other values of γ tested, the Stein discrepancies detect non-convergence, and in fact seem to be more discriminative than the Wasserstein-1 distance.

5.1.2 Hyperparameter selection in Stochastic Gradient Langevin Dynamics

An MCMC algorithm which has gained huge popularity in recent years is **Stochastic Gradient Langevin Dynamics** (SGLD) [Welling and Teh, 2011]. Given $p(\mathbf{x}|\mathcal{D}) \propto p_0(\mathbf{x}) \prod_{r=1}^R p(\mathcal{D}_r|\mathbf{x})$, SGLD starts with initial particles \mathbf{x}_0 and makes the updates

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\epsilon_t}{2} \left(\nabla_{\mathbf{x}_t} \log p_0(\mathbf{x}_t) + \frac{R}{m} \sum_{r \in \sigma_t} \nabla_{\mathbf{x}_t} \log p(\mathcal{D}_r|\mathbf{x}_t) \right) + \boldsymbol{\eta}_t,$$

where σ_t is a subset of $\{1, \dots, R\}$ of size m sampled uniformly at random for each t , and $\boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, \epsilon_t I_d)$.

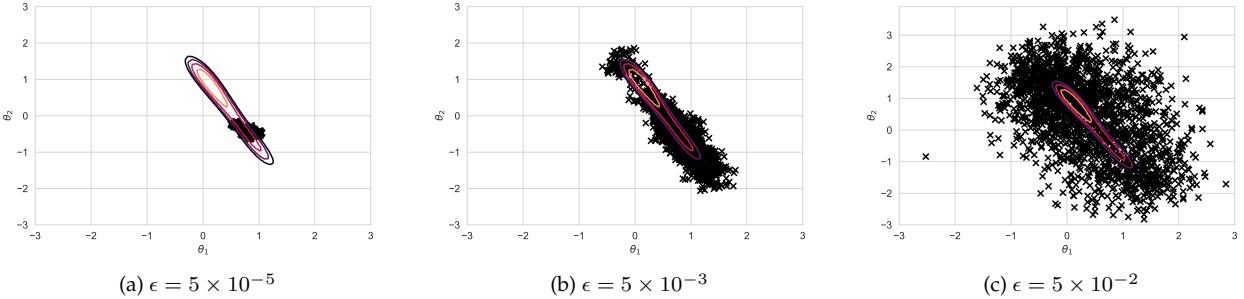


Figure 5.2: Hyperparameter selection in SGLD.

Although convergence guarantees require that the step sizes $(\epsilon_t)_{t \geq 1}$ satisfy

$$\sum_{t=1}^{\infty} \epsilon_t = \infty \quad \text{and} \quad \sum_{t=1}^{\infty} \epsilon_t^2 < \infty,$$

it is common in practice to use a constant step size $\epsilon_t \equiv \epsilon$, which introduces bias. The choice of ϵ is important: if it is too small then the sequence will exhibit long mixing times, while if it is too large then the bias will lead

to a stationary distribution significantly different to the target.

As demonstrated in [Gorham et al. \[2020\]](#), we can calculate KSDs between the generated samples and the target for a range of ϵ . We consider the Gaussian mixture model (GMM) discussed in [Welling and Teh \[2011\]](#):

$$\theta_1 \sim \mathcal{N}(0, \sigma_1^2), \quad \theta_2 \sim \mathcal{N}(0, \sigma_2^2), \\ x_r \sim \frac{1}{2}\mathcal{N}(\theta_1, \sigma_x^2) + \frac{1}{2}\mathcal{N}(\theta_1 + \theta_2, \sigma_x^2), \quad r = 1, \dots, R,$$

where we take $\sigma_1^2 = 10$, $\sigma_2^2 = 1$ and $\sigma_x^2 = 2$. We generate the data x_r using $\theta_1 = 0$ and $\theta_2 = 1$. We see in Figure 5.2 that when ϵ is too small, the samples are underdispersed, and when ϵ is too large, the samples are overdispersed. The most acceptable value of ϵ of those plotted seems to be $\epsilon = 5 \times 10^{-3}$.

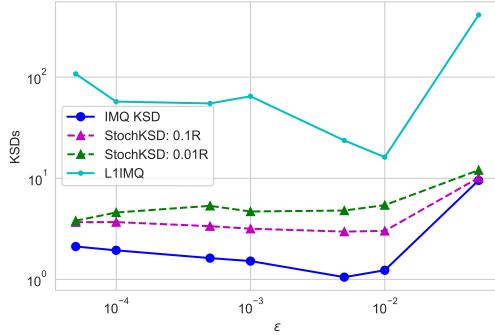


Figure 5.3: Various Stein discrepancies against the hyperparameter ϵ in SGLD.

Figure 5.3 shows the IMQ KSD, stochastic IMQ KSDs with $m = 0.01R$ and $m = 0.1R$, and the L^1 IMQ RFSD between SGLD samples and the target for various ϵ . For the L^1 IMQ RFSD, we take $M = 10$ and take the median over 5 evaluations. For the stochastic KSDs, we take the mean over 5 evaluations.

We see from Figure 5.3 that the KSD, RFSD and stochKSD with $m = 0.1R$ are minimised for ϵ somewhere between 10^{-3} and 10^{-2} , which agrees with Figure 5.2. The result is less clear for the stochastic KSD with $m = 0.01R$, which suggests that more evaluations of this discrepancy may be needed.

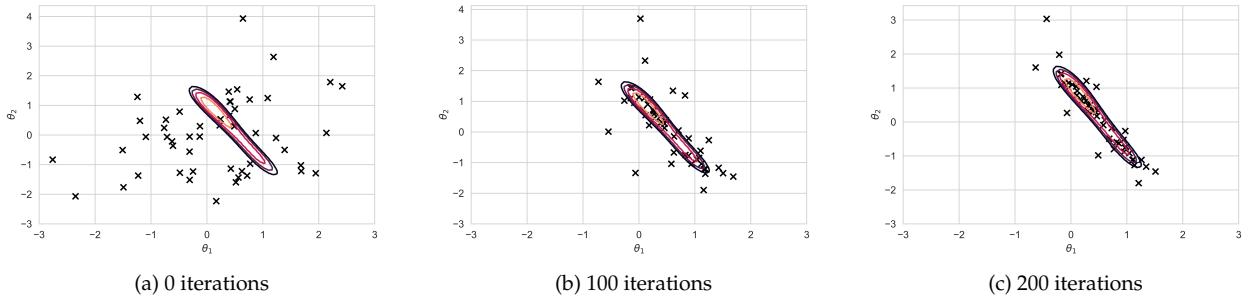


Figure 5.4: SVGD for the GMM posterior at multiple time steps.

5.2 Particle-based variational inference

In this section, we assess the performance of SVGD and stochSVGD as approximate inference algorithms.

5.2.1 Convergence of SVGD

As a simple example, we return to the GMM example from Section 5.1.2. We plot the points outputted by SVGD and stochSVGD each with 50 particles at several stages. We use $R = 100$ likelihood terms and for the stochastic version, we use $m = 10$.

Figure 5.4 exemplifies the performance of SVGD. The particles appear to converge to the target very quickly. Figure 5.5 shows the stochSVGD particles converging, albeit requiring more iterations: this is to be expected since each stochSVGD iteration requires only one tenth of the likelihood evaluations required by full SVGD.

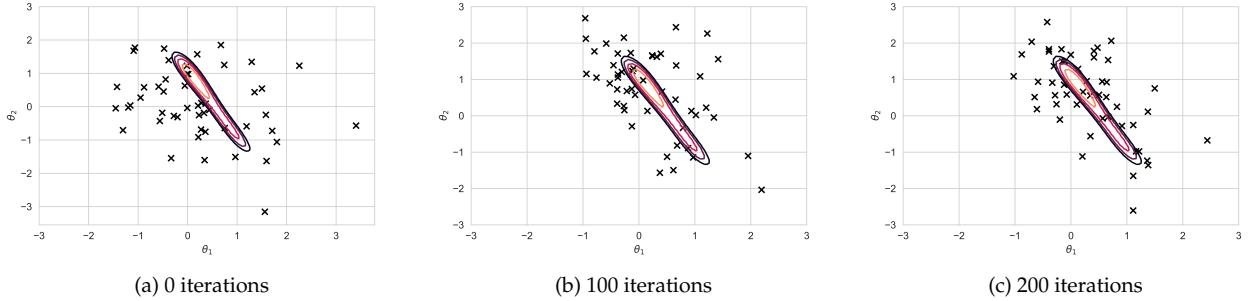


Figure 5.5: stochSVGD for the GMM posterior at multiple time steps.

Since we have samples from the SVGD algorithms and the score of the target, we can compute the IMQ KSD to monitor the convergence of the samples. Figure 5.6 shows the IMQ KSD at early time steps. SVGD repeatedly decreases the IMQ KSD, as to be expected by our theory, while the IMQ KSD for stochSVGD fluctuates due to the randomness of the subsampling procedure.

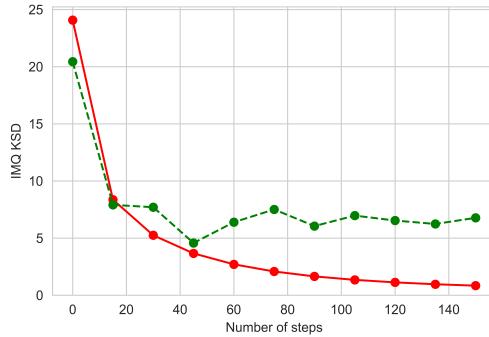


Figure 5.6: The IMQ KSD against the number of steps of SVGD (red) and stochSVGD (green).

We should expect that the IMQ KSD remains bounded away from zero for both algorithms, since we are trying to approximate a continuous distributions with a finite number of particles. The plots of the SVGD particles and the KSD plot both demonstrate that SVGD iteratively improves the particle approximation to the target GMM posterior.

Chapter 6

Conclusion

The recent popularity of biased MCMC algorithms has, among other applications, driven the need for a valid measure of distance between a discrete distribution of samples and a target distribution whose density is only known up to a multiplicative constant. An appropriate characterisation of ‘validity’ is the metrisation of weak convergence. We have shown that the IMQ KSD does indeed metrise weak convergence, while only requiring samples from one distribution and the score function of the other. This makes the IMQ KSD widely applicable, in particular for approximate inference.

Nevertheless, the scalability of the IMQ KSD is key. We have addressed issues concerning approximate inference on large datasets through the L^r IMQ RFSD and the stochastic IMQ KSD. We also introduced SVGD and stochSVGD, approximate inference algorithms in their own right, which are inspired by the KSD and stochastic KSD respectively, and are both (almost surely) asymptotically exact.

While we gave strong theoretical assurances for the IMQ KSD and its stochastic variant, future work could seek to obtain similar theoretical assurances for the sliced KSD introduced by [Gong et al. \[2021\]](#). This would justify its use on high-dimensional data, and would help to combat the curse of dimensionality. It would also be of interest to obtain theoretical assurances for substituting a random feature version of the SVGD update direction in order to reduce the quadratic cost of SVGD in the number of particles.

Bibliography

- K. Chwialkowski, H. Strathmann, and A. Gretton. A kernel test of goodness of fit. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2606–2615, 2016.
- A. Y. K. Foong, D. R. Burt, Y. Li, and R. E. Turner. On the expressiveness of approximate inference in bayesian neural networks, 2020.
- W. Gong, Y. Li, and J. M. Hernez-Lobato. Sliced kernelized stein discrepancy, 2021.
- J. Gorham and L. Mackey. Measuring sample quality with stein’s method. *Advances in Neural Information Processing Systems*, 2015.
- J. Gorham and L. Mackey. Measuring sample quality with kernels. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1292–1301, 2017.
- J. Gorham, A. B. Duncan, S. J. Vollmer, and L. Mackey. Measuring sample quality with diffusions, 2018.
- J. Gorham, A. Raj, and L. Mackey. Stochastic stein discrepancies. In *Advances in Neural Information Processing Systems*, volume 33, pages 17931–17942, 2020.
- J. Huggins and L. Mackey. Random feature stein discrepancies. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Q. Liu. Stein variational gradient descent as gradient flow. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Q. Liu and D. Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- S. J. Reddi, A. Ramdas, B. Ps, A. Singh, and L. Wasserman. On the decreasing power of kernel and distance based nonparametric hypothesis tests in high dimensions, 2014.
- C. Stein. A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 2. University of California Press.* pp. 583–602, 1972.
- M. Welling and Y. W. Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning*, 2011.
- H. Wendland. Scattered data approximation. *Cambridge University Press*, 2004.