

An Introduction to Stein Kernels

Robert Salomone

UNSW Sydney
r.salomone@unsw.edu.au



Introduction

- ▶ Exciting new area of research in Monte Carlo methods.
- ▶ Lots of papers in Machine Learning and Computational Statistics literature in the last few years.
- ▶ Main idea: Combine ideas from **Stein's Method** with **Reproducing Kernel Hilbert Space** (RKHS) theory.

Stein's Method: A History Lesson

- ▶ Stein's (1972) method — general **theoretical tool** for **bounding the distance between two distributions** p and q .

Integral Probability Metrics

- ▶ Let \mathcal{F} is some class of **test functions** $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$.
- ▶ For two random variables, $\mathbf{X} \sim q$, and $\mathbf{Y} \sim p$, an **Integral Probability Metric** (IPM) takes the form

$$d_{\mathcal{F}}(q, p) = \sup_{\phi \in \mathcal{F}} |\mathbb{E}_q \phi(\mathbf{X}) - \mathbb{E}_p \phi(\mathbf{Y})|.$$

- ▶ Essentially, **worst case** error over all functions in \mathcal{F} obtained by integrating using measure associated with q instead of p .
- ▶ Special choice of \mathcal{F} can **recover some familiar friends**: e.g., Total Variation, Maximum Mean Discrepancy, and Wasserstein.
- ▶ For a sequence of distributions $\{q^N\}_{N=1}^{\infty}$, a **sufficiently rich** choice of \mathcal{F} ensures $d_{\mathcal{F}}(q^N, p) \rightarrow 0 \implies q^N \rightarrow p$.

Stein's Insight

- ▶ How do you actually bound (let alone **compute**)

$$\sup_{\phi \in \mathcal{F}} |\mathbb{E}_q \phi(\mathbf{X}) - \mathbb{E}_p \phi(\mathbf{Y})|,$$

when it has **intractable expectations** and **supremum**?

- ▶ **Stein's Insight**: Turn the difference of expectations into a **single computable expectation**.
- ▶ Main idea: Find an Operator \mathcal{A}_p with property

$$\mathbb{E} \mathcal{A}_p \phi(\mathbf{X}) = 0 \text{ for all } \phi \in \tilde{\mathcal{F}} \iff \mathbf{X} \sim p,$$

then, we have the IPM:

$$\mathbb{S}(q, p) = \sup_{\phi \in \tilde{\mathcal{F}}} |\mathbb{E}_q[\mathcal{A}_p \phi(\mathbf{x})]|. \quad (1)$$

- ▶ Stein found such an operator for p being the univariate **Normal** distribution.

Barbour's Generator Approach

- ▶ Big step forward.
- ▶ Let $\{\mathbf{Z}_t\}_{t \geq 0}$ denote a Markov process with stationary distribution p .
- ▶ The **infinitesimal generator** \mathcal{A} of $\{\mathbf{Z}_t\}_{t \geq 0}$ is defined pointwise by

$$\mathcal{A}\phi(\mathbf{x}) = \left. \frac{d}{dt} \mathbb{E}[\phi(\mathbf{Z}_t) \mid \mathbf{Z}_0 = \mathbf{x}_0] \right|_{t=0}.$$

- ▶ **Key Point:** This generator satisfies $\mathbb{E}_p \mathcal{A}\phi(\mathbf{X}) = 0$ under very mild conditions on ϕ and \mathcal{A} .
- ▶ So, let's create a **Stein Operator** for p with support on all of \mathbb{R}^d .

Langevin–Stein Operator

- ▶ **Langevin diffusion** $\{\mathbf{Z}_t\}_{t \geq 0}$ is the process:

$$d\mathbf{Z}_t = \nabla \log p(\mathbf{Z}) + \sqrt{2}d\mathbf{W}_t$$

and has stationary distribution p .

- ▶ For \mathbf{X} with density p , and $f \in \mathcal{C}^2$, the generator of the Langevin SDE with stationary distribution p is

$$\mathcal{A}f(\mathbf{x}) = \Delta f(\mathbf{x}) + \nabla f(\mathbf{x}) \cdot \nabla \log p(\mathbf{x}),$$

- ▶ As $f \in \mathcal{C}^2$, rewriting $\nabla f = \phi$ gives us

$$\mathcal{A}f(\mathbf{x}) = \nabla \phi(\mathbf{x}) + \phi(\mathbf{x}) \cdot \nabla \log p(\mathbf{x}),$$

the RHS of which may seem more familiar.

- ▶ This is often called **the** Stein Operator in ML literature, but it really is just **a** Stein Operator.

So...

- ▶ The general form of **Stein Discrepancy** is

$$\mathbb{S}(q, p) = \sup_{\phi \in \tilde{\mathcal{F}}} |\mathbb{E}_q[\mathcal{A}_p \phi(\mathbf{x})]|. \quad (2)$$

- ▶ Simpler, but to do something **practical** with this would ideally be able to solve the above **exactly**.
- ▶ Seems impossible, but using **kernel methods** get us around it.

Kernel Methods: A Crash Course

- ▶ Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a symmetric positive semi-definite function.
- ▶ Positive semi-definite function: The **Gram Matrix**
 $K = (k(\mathbf{x}_i, \mathbf{x}_j))_{ij}$ is positive definite for any number of points.
- ▶ Concrete example for now: for $\sigma^2 > 0$, **rbf kernel** is
 $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2)$.
- ▶ **For our purposes**, it suffices to know that such kernel functions allow us to work implicitly (e.g., solve optimization problems) in complicated, possibly infinite dimensional, spaces of **functions**.

The Magic

- ▶ Let \mathcal{H} be a **Reproducing Kernel Hilbert Space**¹.
- ▶ **The Magic** — insist that $\tilde{\mathcal{F}}$ be the **unit ball** of \mathcal{H} :

$$\tilde{\mathcal{F}} = \{\phi \in \mathcal{H} : \|\phi\|_{\mathcal{H}} \leq 1\}.$$

- ▶ Then,

$$\mathbb{S}(q, p) = \sup_{\phi \in \tilde{\mathcal{F}}} |\mathbb{E}_q[\mathcal{A}_p \phi(\mathbf{x})]|, \quad (3)$$

has the **closed form** solution

$$\sqrt{\mathbb{E}_{q \times q} k_p(\mathbf{X}, \mathbf{Y})},$$

where k_p is a special type of kernel called a **(Reproducing) Stein Kernel**. More on what precisely k_p is in a little while.

¹intuitively, a space of functions induced by the choice of $k(\cdot, \cdot)$

Properties of Stein Kernels

- 1. We had **Stein's Identity**: $\mathbb{E}_p \mathcal{A}_p \phi(\mathbf{X}) = 0$.
Now, for $\mathbf{X}, \mathbf{Y} \sim p$, we have

$$\underbrace{\mathbb{E}_{p \times p} k_p(\mathbf{X}, \mathbf{Y}) = 0}_{\text{Kernelized Stein Identity}}$$

2. For $q \neq p$: $\mathbb{E}_{q \times q} k_p(\mathbf{X}, \mathbf{Y}) \geq 0$, and thus

$$\underbrace{\sqrt{\mathbb{E}_{q \times q} k_p(\mathbf{X}, \mathbf{Y})}}_{\text{Kernelized Stein Discrepancy}} \geq 0.$$

Empirical Measures

- ▶ Recall, $\mathbb{S}(q, p) = \sup_{\phi \in \tilde{\mathcal{F}}} |\mathbb{E}_q[\mathcal{A}_p \phi(\mathbf{x})]| = \sqrt{\mathbb{E}_{q \times q} k_p(\mathbf{X}, \mathbf{Y})}$ for iid \mathbf{X} and \mathbf{Y} . It is tempting to think we need to **estimate** this.
- ▶ However, often q is the empirical distribution of a **set of N points** $\{\mathbf{x}_k\}_{k=1}^N$, i.e.,

$$\mathbb{E}_{q^N} \phi(\mathbf{X}) = \frac{1}{N} \sum_{k=1}^N \phi(\mathbf{x}_k).$$

- ▶ We do this all the time in Monte Carlo / Quasi Monte Carlo (use an empirical measure).
- ▶ The KSD in such a case is

$$\mathbb{S}(q^N, p) = \sqrt{\mathbb{E}_{q^N \times q^N} k_p(\mathbf{X}, \mathbf{Y})} = \sqrt{\frac{1}{N^2} \sum_{j=1}^N \sum_{k=1}^N k_p(\mathbf{x}_j, \mathbf{x}_k)}.$$

- ▶ This is not a **Monte Carlo estimate** of the KSD, this **is** the KSD of your point set! An **explicitly computable** IPM.

An Explicitly Computable Integral Probability Metric!

How do you construct Stein kernels though?

Step 1: Choose a Base Kernel

- ▶ We begin with a **base** kernel, a symmetric positive definite function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. For example:

kernel	$k(\mathbf{x}, \mathbf{y})$	parameters
Polynomial	$(\mathbf{x}^\top \mathbf{y} + c)^d$	$c \in \mathbb{R}_+, d \in \mathbb{N}$
Gaussian (rbf)	$\exp(-\ \mathbf{x} - \mathbf{y}\ ^2 / 2\sigma^2)$	$\sigma^2 > 0$
Inverse Multiquadric (IMQ)	$(1 + \ \mathbf{x} - \mathbf{y}\ _2^2)^{-\beta}$	$\beta \in (0, 1)$

- ▶ Many others: Matérn, Exponential, you can easily create your own.

Step 2: “Steinalize” the Base Kernel

- ▶ **Main Idea:** Hit the base kernel with a Stein Operator on both arguments, i.e., \mathcal{A}_p^x and \mathcal{A}_p^y .
- ▶ In the case where \mathcal{A} is the **Langevin–Stein** operator, this yields

$$\begin{aligned}k_p(\mathbf{x}, \mathbf{y}) &:= \nabla \log p(\mathbf{x})^\top k(\mathbf{x}, \mathbf{y}) \nabla \log p(\mathbf{y}) \\ &\quad + \nabla \log p(\mathbf{x})^\top \nabla_{\mathbf{y}} k(\mathbf{x}, \mathbf{y}) \\ &\quad + \nabla \log p(\mathbf{y})^\top \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{y}) \\ &\quad + \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{y}} k(\mathbf{x}, \mathbf{y}).\end{aligned}$$

- ▶ Still no need of the normalizing constant!

Properties of Stein Kernels (Continued)

- ▶ Stein Kernels are **actually Kernels**. We have

$$(\text{symmetric } k \succ 0) \implies (\text{symmetric } k_p \succ 0).^2$$

- ▶ KSD with **IMQ kernel** and **Langevin–Stein** operator controls weak convergence³, i.e.,

$$\mathbb{S}(q^N, p) \rightarrow 0 \implies q^N \rightarrow p.$$

- ▶ This property is called **convergence-determining**.
- ▶ Surprisingly, one can construct cases where common kernels **fail** to have this property.

²Oates et al., 2017 for Langevin–Stein kernels. More general results by Hodgkinson et al., 2020.

³Gorham and Mackey, 2017, Theorem 8

Succinctly put...

- ▶ The **Stein Operator** removes the need to compute the expectations $\mathbb{E}_p\phi$.
- ▶ The **Stein Kernel** saves you from solving for the supremum.
- ▶ For IMQ kernel, KSD is still equal to some **convergence determining** IPM.
- ▶ So, what can we actually do with this?

Methods

1. Measuring Sample Quality
2. Construct a Sequence of Quality Samples
3. Improve the Quality of Existing Samples

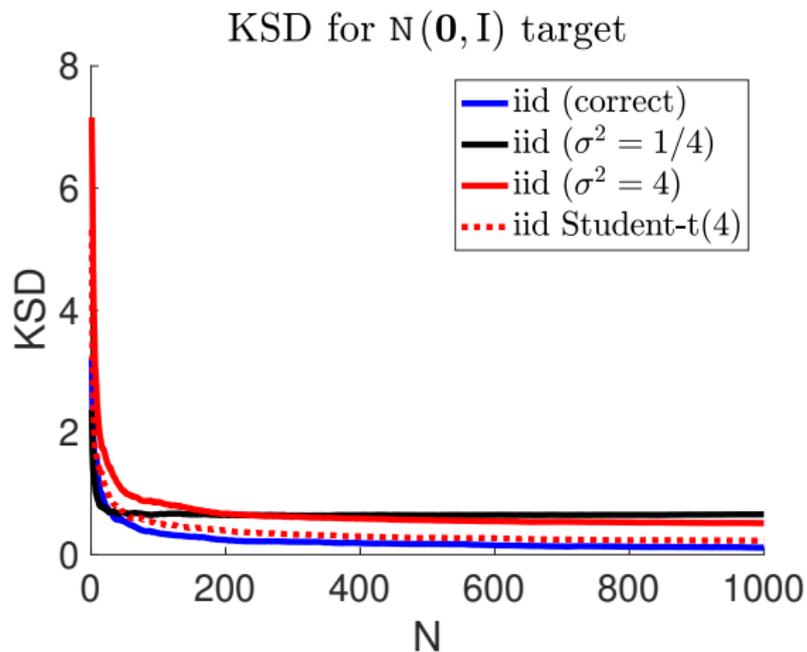
Measuring Sample Quality

- ▶ As mentioned, we have a **general tool** for assessing how well a collection of points approximate expectations with respect to a target distribution.
- ▶ How you got the points **doesn't matter!**
 - ▶ Asymptotically Unbiased MCMC (e.g., Metropolis–Hastings)
 - ▶ Asymptotically Biased MCMC (e.g., SGLD, ULA)
 - ▶ Sequential Monte Carlo / Sequentially Interacting MCMC
 - ▶ Constructed the Point Set Deterministically (e.g., QMC)

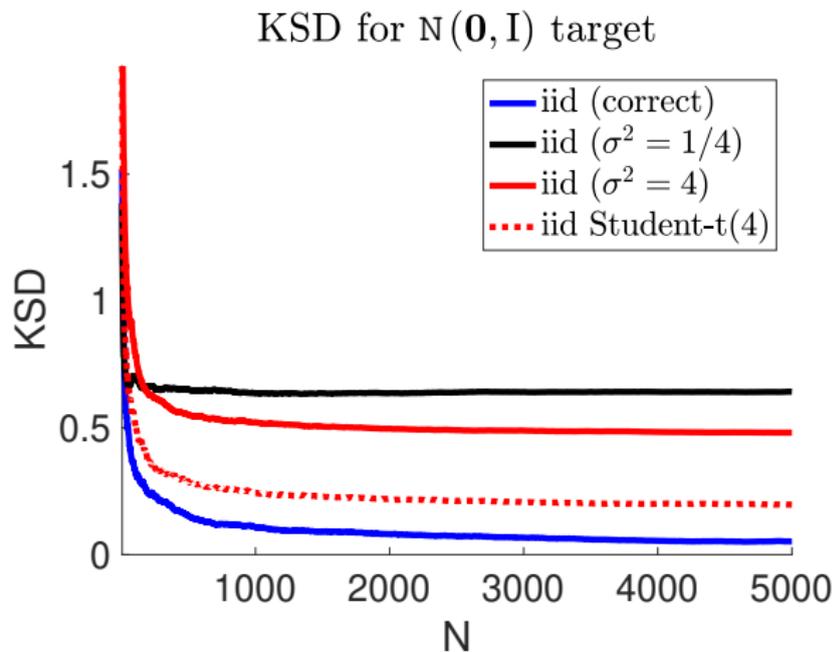
KSD in Action

- ▶ Suppose that our target p is a 10–dimensional standard normal distribution, $\mathcal{N}(\mathbf{0}, \mathbf{I})$.
- ▶ We will evaluate KSD with point sets generated via:
 1. (perfect) iid sampling from p .
 2. iid samples from $\mathcal{N}(\mathbf{0}, 0.25 \times \mathbf{I})$ (underdispersed)
 3. iid samples from $\mathcal{N}(\mathbf{0}, 4 \times \mathbf{I})$ (overdispersed)
 4. iid samples from student- $t(\mathbf{0}, \mathbf{I}, 4)$

KSD in Action



KSD in Action (More Samples)



So...

- ▶ **Practitioners:** You can use KSD to evaluate the quality of points generated by **any** sampling algorithm, and use this to tune hyperparameters (example forthcoming).
- ▶ **Methodologists:** You can use KSD to show your sampling methods are better than other sampling methods.

Hyperparameter Tuning

- ▶ In *Measuring Sample Quality with Kernels* (Gorham and Mackey, 2017), the authors consider tuning an **asymptotically biased** Slice Sampler.
- ▶ Effective Sample Size (ESS) would say set steps as large as possible (minimize autocorrelation), but this is clearly a bad idea (increases bias).
- ▶ However, KSD makes a different recommendation...

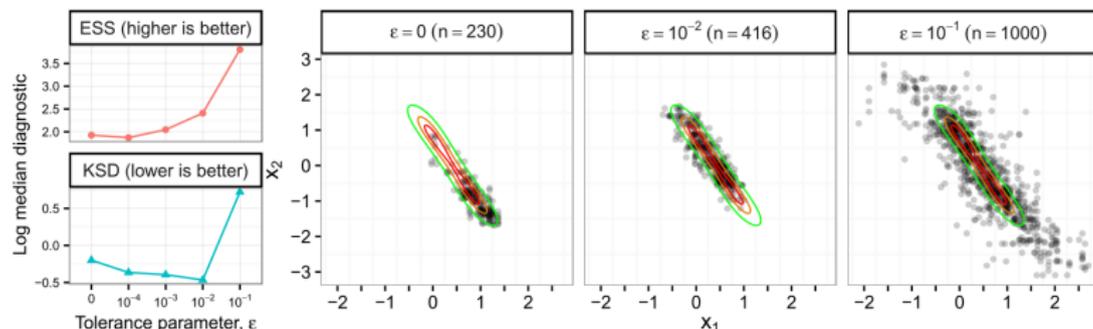


Figure: Figure from *Measuring Sample Quality with Kernels* (Gorham and Mackey, 2017)

Construct a Sequence of Quality Samples

Stein Points (Chen et. al., 2018)

- ▶ Main Idea: Construct point sequences that minimize the KSD!
- ▶ **Extensible Point Sets.**
- ▶ Two methods:
 1. Greedy Algorithm
 2. Herding Algorithm
- ▶ For empirical distribution of N Stein Points q^N obtained by both of the above methods, the authors establish that $\mathbb{S}(q^N, p) \rightarrow 0$.

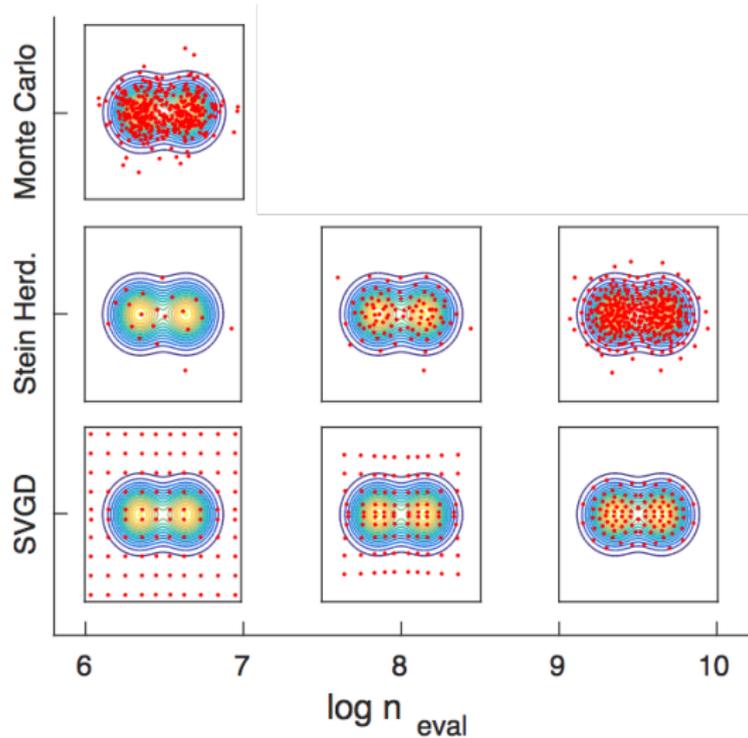


Figure: Modified Figure 1 from *Stein Points* paper.

Stein Importance Sampling (Liu and Lee, 2016)

- ▶ **Classic** Importance sampling. Instead of sampling from p , we sample $\mathbf{X}_1, \dots, \mathbf{X}_N$ according to q and have

$$\sum_{k=1}^N \phi(\mathbf{X}_k) \frac{w_k}{\sum_{k=1}^N w_k}$$

as an estimator of $\mathbb{E}_p \phi(\mathbf{X})$.

- ▶ Calculating the **importance weights** $w_k = p(\mathbf{X}_k)/q(\mathbf{X}_k)$ require evaluation of p and q pointwise, though they need not be normalized.
- ▶ **Consistent**, but has disadvantages. There may be better choices of weight functions.

Improve the Quality of Existing Samples

Stein Importance Sampling

- ▶ Idea is simple, for **fixed points**, consider **weighted** empirical distributions $q_{\mathbf{w}}^N$, i.e.,

$$\mathbb{E}_{q_{\mathbf{w}}^N} \phi(\mathbf{X}) = \sum_{k=1}^N w_k \phi(\mathbf{x}_k), \quad \mathbf{w} \geq 0, \sum_{k=1}^N w_k = 1.$$

- ▶ Define the **Stein Gram Matrix** $K_p := (k_p(\mathbf{x}_i, \mathbf{x}_j))_{ij}$.
- ▶ We solve:

$$\begin{aligned} \mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left\{ \mathbb{S}(q_{\mathbf{w}}^N, p) : \mathbf{w} \geq 0, \sum_{k=1}^N w_k = 1 \right\} \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left\{ \mathbf{w}^\top K_p \mathbf{w} : \mathbf{w} \geq 0, \sum_{k=1}^N w_k = 1 \right\} \end{aligned}$$

- ▶ Interestingly, density q need not be known, and the procedure reduces MSE quite well in practice.

New Result

- ▶ **Research Question:** Is it possible to use the output of a Markov Chain with ergodic distribution $q \neq p$ to construct **consistent** estimators for expectations under p by using **Stein Importance Sampling**?
- ▶ **Answer:** Yes! (under certain conditions) This is our new result. See Theorem 1 of

Hodgkinson, L., Salomone, R., and Roosta, F. (2020). *The reproducing Stein kernel approach for post-hoc corrected sampling*. arXiv:2001.09266.

- ▶ In the above paper, we also provide **general theory** for construction of Stein Kernels on **arbitrary polish spaces** and conditions for such kernels to be convergence-determining.

Conspiring Scholars



Liam Hodgkinson Robert Salomone Fred Roosta

Hodgkinson, L., Salomone, R., and Roosta, F. (2020). *The reproducing Stein kernel approach for post-hoc corrected sampling*. arXiv:2001.09266.

Some Other Methods

- ▶ **Black-Box Variational Inference:** Operator Variational Inference (Ranganath et al., 2016), Stein Variational Gradient Descent (Liu and Chen, 2016), Wild Variational Inference (Liu, 2016)
- ▶ **Reinforcement Learning:** Stein Variational Policy Gradient (Liu, 2017)
- ▶ **Generative Models:** SteinGAN (Wang et al., 2017)
- ▶ **Non-Parametric Control Variates:** Control Functionals (Oates et al., 2017)
- ▶ **Stein Variational Adaptive Importance Sampling** (Han and Liu, 2017).

Thankyou for your attention.