# Kernelized Complete Conditional Stein Discrepancies

**Raghav Singhal**[1], **Xintian Han**[2], **Saad Lahlou**[2], and **Rajesh Ranganath**[1,2]

[1]Courant Institute of Mathematical Sciences, New York University
[2]Center for Data Science, New York University

## Abstract

Much of machine learning relies on comparing distributions with discrepancy measures. Stein's method creates discrepancy measures between two distributions that require only the unnormalized density of one and samples from the other. Stein discrepancies can be combined with kernels to define kernelized Stein discrepancies (KSDs). While kernels make Stein discrepancies tractable, they pose several challenges in high dimensions. We introduce kernelized complete conditional Stein discrepancies (KCC-SDs). Complete conditionals turn a multivariate distribution into multiple univariate distributions. We show that KCC-SDs distinguish distributions. To show the efficacy of KCC-SDs in distinguishing distributions, we introduce a goodness-of-fit test using KCC-SDs. We empirically show that KCC-SDs have higher power over baselines and use KCC-SDs to assess sample quality in Markov chain Monte Carlo.

## 1 Introduction

Discrepancy measures that compare a distribution $p$, known up to normalization, with a distribution $q$, known via samples from it, can be used for finding good variational approximations (Ranganath et al., 2016; Liu and Wang, 2016), checking the quality of MCMC samplers (Gorham and Mackey, 2015, 2017), goodness-of-fit testing (Liu et al., 2016), parameter estimation (Barp et al., 2019) and multiple model comparison (Lim et al., 2019). There are several difficulties with using traditional discrepancies like Wasserstein metrics or total variation distance for these tasks. Mainly, $p$ can be hard to sample so expectations under $p$ cannot be computed. These challenges lead to the following desiderata for a discrepancy $D$ (Gorham and Mackey, 2015).

1. **Tractable** $D$ uses samples from $q$, and evaluations of (unnormalized) $p$.

2. **Distinguishing Distributions** $D(p, q) = 0$ if and only if $p$ is equal in distribution to $q$.

These desiderata ensure that the discrepancy is non zero when $p$ does not equal $q$ and that it can be easily computed. To meet these desiderata, Chwialkowski et al. (2016); Oates et al. (2017); Gorham and Mackey (2017); Liu et al. (2016) developed kernelized Stein discrepancies (KSDs). KSDs measure the expectation of functions under $q$ that have expectation zero under $p$. These functions are constructed by applying Stein's operator to a reproducing kernel Hilbert space (RKHS).

In high dimensions, many popular kernels evaluated on a pair of points are near zero. Thus, KSDs in high dimensions can be near zero, making detecting differences between high dimensional distributions difficult. The median heuristic can be used to address this to some extent, but KSDs with the median heuristic can still have low power in moderately high dimensions (see Figure 1, Jitkrittum et al. (2017)). We develop kernelized complete conditional Stein discrepancies (KCC-SDs). These discrepancies use complete conditionals: the distribution of one variable given the rest. Complete conditionals are univariate distributions. Rather than using multivariate kernels, KCC-SDs use univari-

ate kernels to ensure the complete conditionals match, making it easier to compare distributions in high dimensions.

A given Stein discrepancy relies on a supremum over a class of test functions called the Stein set. KCC-SDs differ from KSDs in that KCC-SDs compute a separate supremum for each complete conditional. An immediate question is whether there is a computable closed form and whether the discrepancy can be used to distinguish distributions. We show that KCC-SDs have a closed form and distinguish between distributions. Computing KCC-SD requires sampling from a complete conditional of $q$, which can be infeasible in some instances. To address this, we introduce approximate KCC-SD that uses a learned sampler for the complete conditional.

To show the efficacy of KCC-SD and approximate KCC-SD in distinguishing distributions we introduce a goodness-of-fit test (Chwialkowski et al., 2016). We show that KCC-SD and approximate KCC-SD have higher power than KSD and other baselines. We empirically show that approximate KCC-SD does not suffer from a loss in power due to an increase in dimension. We also demonstrate that KCC-SD and approximate KCC-SD can be used to select sampler hyperparameters and can be used to assess sample quality in a Gibbs sampler.

**Related Work.**   There have been several lines of work which use factorizations of the distribution $p$ to address the curse of dimensionality. Wang et al. (2017); Zhuo et al. (2017) use the Markov blanket of each node to define a graphical version of KSD to alleviate the curse of dimensionality. Our approach does not presume a graphical structure of $p$ or $q$. Wang et al. (2017) shows that unless the graphical structure for $p, q_n$ match, the graph based KSD converging to zero does not imply that $q_n$ converges in distribution to $p$.

Gong et al. (2020) introduce the maximum sliced kernelized Stein discrepancy (MAXSKSD), which also uses low-dimensional kernels by projecting into a 1-dimensional space. Computing MAXSKSD requires optimizing a projection direction that is specific to both sampling distribution $q$ and the unnormalized distribution $p$. This can be expensive when testing multiple distributions or when changing the parameters of an unnormalized model to fit a collection of samples. Approximate KCC-SD requires learning conditional distributions specific only to the sampling distribution $q$. Similar to approximate KCC-SD with parametric conditional estimates, the closed form for MAXSKSD depends on the optimal direction, therefore the power of their method depends on the quality of the optimization, which can be difficult to guarantee for arbitrary log probabilities.

KSDs suffer from a computational cost that is quadratic in the number of samples. Huggins and Mackey (2018) develop random feature Stein discrepancies R$\Phi$SD, which run in linear time and perform as well as or better than quadratic-time KSDs; these ideas can be applied to KCC-SDs. Chen et al. (2018) introduces the Stein points method which introduces a method to select points to minimize the Stein discrepancy between the empirical distribution supported at the selected points and the posterior.

Chwialkowski et al. (2016) introduced KSD as a test statistic for a goodness-of-fit test, which also suffers from the curse of dimensionality due to the use of kernels in high dimensions, along with a computational cost quadratic in the number of samples. Jitkrittum et al. (2017) introduce a linear-time discrepancy, finite-set Stein discrepancy (FSSD). The authors introduce an optimized version of FSSD which allows one to find features that best indicate the differences between the samples and the target density. FSSD while having a computational cost linear in sample size, also leads to a test with lower power in high dimensions.

## 2   Kernelized Stein Discrepancies

Stein's method provides recipes for constructing expectation zero test functions of distributions known up to a normalization constant. For a distribution $p$ with a integrable score function[1], $\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})$, we can create a *Stein operator*, $\mathcal{A}_p$, that acts on test functions $f : \mathbb{R}^d \to \mathbb{R}^d$ satisfying regularity and boundary conditions (Proposition 1, (Gorham and Mackey, 2015)), such that

---

[1]The score function in general is the gradient of the log-likelihood with respect to the parameter vector. We however refer to the gradient of the log-likelihood with respect to the input (Hyvärinen, 2005).

$$\mathbb{E}_{p(\boldsymbol{x})}\left[\mathcal{A}_{p(\boldsymbol{x})}f(\boldsymbol{x})\right] = 0.$$

This relation called *Stein's identity* is used to create *Stein discrepancies* $\mathcal{S}(q, \mathcal{A}_p, \mathcal{H})$, defined as

$$\mathcal{S}(q, \mathcal{A}_p, \mathcal{H}) = \sup_{f \in \mathcal{H}} \left| \mathbb{E}_{q(\boldsymbol{x})}[\mathcal{A}_{p(\boldsymbol{x})}f(\boldsymbol{x})] - \mathbb{E}_{p(\boldsymbol{x})}[\mathcal{A}_{p(\boldsymbol{x})}f(\boldsymbol{x})] \right|$$
$$= \sup_{f \in \mathcal{H}} \left| \mathbb{E}_{q(\boldsymbol{x})}\left[\mathcal{A}_{p(\boldsymbol{x})}f(\boldsymbol{x})\right] \right| ,$$

where $\mathcal{H}$ is a function space known as the *Stein set*, with its functions satisfying some boundary and regularity conditions. To make the Stein discrepancy simpler to compute, Chwialkowski et al. (2016); Oates et al. (2017); Gorham and Mackey (2017); Liu et al. (2016) used reproducing kernel Hilbert spaces (RKHS) as the Stein set to introduce kernelized Stein discrepancies (KSD). Let $k : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be the kernel of an RKHS $\mathcal{K}_k$, the RKHS consists of functions, $g : \mathbb{R}^d \to \mathbb{R}$, satisfying the reproducing property $g(\boldsymbol{x}) = \langle g, k(\boldsymbol{x}, \cdot) \rangle_{\mathcal{K}_k}$. KSDs are defined by the Stein set

$$\mathcal{G}_k = \left\{ g = (g_1, \ldots, g_d) : g_i \in \mathcal{K}_k, \sum_{i=1}^d \|g_i\|_{\mathcal{K}_k} \leq 1 \right\} .$$

This construction of the Stein set using an RKHS ensures that the Stein discrepancy has a closed form.

**Proposition 1** (Gorham and Mackey, 2017). *Suppose $k \in C^{(1,1)}$ and for each $j \in \{1, \ldots, d\}$, define the Stein kernel as follows:*

$$k_0^j(\boldsymbol{x}, \boldsymbol{y}) = b_j(\boldsymbol{x})b_j(\boldsymbol{y})k(\boldsymbol{x}, \boldsymbol{y}) + \nabla_{x_j}\nabla_{y_j}k(\boldsymbol{x}, \boldsymbol{y}) \tag{1}$$
$$+ b_j(\boldsymbol{x})\nabla_{y_j}k(\boldsymbol{x}, \boldsymbol{y}) + b_j(\boldsymbol{y})\nabla_{x_j}k(\boldsymbol{x}, \boldsymbol{y}) ,$$

*where $b_j(\boldsymbol{x}) = \nabla_{x_j} \log p(\boldsymbol{x})$. If $\sum_{j=1}^d \mathbb{E}_q[k_0^j(\boldsymbol{x}, \boldsymbol{x})^{1/2}] < \infty$, then KSD has a closed form. Given by $\mathcal{S}(q, \mathcal{A}_p, \mathcal{G}_k) = \|\boldsymbol{w}\|_2$, where $w_j^2 \equiv \mathbb{E}_{q(\boldsymbol{x}) \times q(\boldsymbol{y})}\left[k_0^j(\boldsymbol{x}, \boldsymbol{y})\right]$ with $\boldsymbol{x}, \boldsymbol{y} \overset{i.i.d}{\sim} q$.*

When the distribution $p$ lies in the class of distantly dissipative distributions (Eberle, 2016), KSDs provably detect convergence and non-convergence for $d = 1$. That is $\mathcal{S}(q_n, \mathcal{A}_p, \mathcal{G}_k) \to 0$ if and only if $q_n \Rightarrow p$ for sequences $\{q_n\}$, using kernels like the radial basis function or the inverse multi-quadratic (IMQ), (Gorham and Mackey, 2017). In $d > 2$, the KSD with thin tailed kernels like the RBF does not detect non-convergence. But the KSD with the IMQ kernel with $\beta \in (0, 1)$ does detect non-convergence. However, all of these kernels shrink as the $\|\cdot\|_2$ grows, which means their associated KSDs become less sensitive in higher dimensions.

Suppose $\boldsymbol{x}, \boldsymbol{y} \sim N(\boldsymbol{0}, I_d)$ then $\mathbb{E}[\|\boldsymbol{x} - \boldsymbol{y}\|^2] = 2d$, so $k(\boldsymbol{x}, \boldsymbol{y}) = \exp(-\|\boldsymbol{x} - \boldsymbol{y}\|^2/2\sigma^2)$ concentrates around $\exp(-d/\sigma^2)$. The median heuristic, $\sigma = \text{median}(\|\boldsymbol{x}_i - \boldsymbol{x}_j\| ; i < j)$, can be used to deal with this shrinkage. However, (Ramdas et al., 2015) show that even with the median heuristic, kernel based discrepancies can converge to zero as the dimension increases even when the distributions are different.

## 3 Kernelized Complete Conditional Stein Discrepancies.

Complete conditionals are univariate conditional distributions, $p(x_j|\boldsymbol{x}_{-j})$, where $\boldsymbol{x}_{-j} = \{x_1, \ldots x_{j-1}, x_{j+1}, \ldots x_d\}$. Complete conditional distributions are the basis for many inference procedures including the Gibbs sampler (Geman and Geman, 1984), and coordinate ascent variational inference (Ghahramani and Beal, 2001).

Using complete conditionals we construct complete conditional Stein discrepancies (CC-SDs) and their kernelized versions (KCC-SDs). In this work we focus on the Langevin-Stein operator (Barbour, 1990; Gorham and Mackey, 2015), defined for differentiable functions $f : \mathbb{R}^d \to \mathbb{R}^d$ as follows:

$$(\mathcal{A}_{p(\boldsymbol{x})}f)(\boldsymbol{x}) = f(\boldsymbol{x})^T \nabla_{\boldsymbol{x}} \log p(\boldsymbol{x}) + \nabla_{\boldsymbol{x}} \cdot f(\boldsymbol{x}) = \sum_{j=1}^d \mathcal{A}_{p(\boldsymbol{x})}^j f_j(\boldsymbol{x}) .$$

**Definition.** The score function of the complete conditional, $\nabla_{x_j} \log p(x_j \mid \boldsymbol{x}_{-j})$, is the score function of the joint, $\nabla_{x_j} \log p(\boldsymbol{x})$. So for $f_j : \mathbb{R}^d \to \mathbb{R}$,

$$\mathcal{A}^j_{p(x_j \mid \boldsymbol{x}_{-j})} f_j(\boldsymbol{x}) = f_j(\boldsymbol{x}) \nabla_{x_j} \log p(x_j \mid \boldsymbol{x}_{-j}) + \nabla_{x_j} f_j(\boldsymbol{x}) = f_j(\boldsymbol{x}) \nabla_{x_j} \log p(\boldsymbol{x}) + \nabla_{x_j} f_j(\boldsymbol{x})$$
$$= \mathcal{A}^j_{p(\boldsymbol{x})} f_j(\boldsymbol{x})$$

Using this observation, and the fact that the complete conditionals of two distributions $p, q$ match when the distributions match, we define the complete conditional Stein discrepancy (CC-SD), $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C})$ as

$$\sum_{j=1}^d \mathbb{E}_{q(\boldsymbol{x}_{-j})} \left[ \sup_{f_j \in \mathcal{C}^j} \mathbb{E}_{q(x_j \mid \boldsymbol{x}_{-j})} [\mathcal{A}^j_{p(x_j \mid \boldsymbol{x}_{-j})} f_j(\boldsymbol{x})] \right] . \tag{2}$$

The Stein set $\mathcal{C}$ is defined as the set of functions, $f : \mathbb{R}^d \to \mathbb{R}^d$, with each component $f_j(\boldsymbol{x})$ satisfying $\max \left( \|f_j\|_\infty, \|\nabla f_j\|_\infty, Lip(f_j) \right) \leq 1$, where $Lip(f)$ is the Lipschitz constant of $f$. Here, the supremum is taken inside the expectation, so we have to solve optimization problems for each dimension and each conditional. Similar to Stein discrepancies, CC-SDs can be hard to compute. In the next section, we introduce the kernelized version which has a closed form.

### 3.1 Kernelized Complete Conditional Stein Discrepancies.

We now define the Stein set, $\mathcal{C}_k$, for the kernelized version of CC-SD, such that we get a closed form discrepancy.

We use univariate integrally symmetric positive definite (ISPD) kernels, $k : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$, that satisfy the following, for $g : \mathbb{R} \to \mathbb{R}$:

$$\int_{u \in \mathbb{R}} \int_{v \in \mathbb{R}} g(u) k(u, v) g(v) du dv > 0 , \tag{3}$$

with $\|g\|_2 > 0$. Let $\mathcal{K}_k$ denote the reproducing kernel Hilbert space (RKHS) with kernel $k$. Functions $h \in \mathcal{K}_k$ satisfy the reproducing property, $h(x_j) = \langle h, k(x_j, \cdot) \rangle_{\mathcal{K}_k}$ for $x_j \in \mathbb{R}$. The RKHS also satisfies $\Phi_{x_j}(\cdot) = k(x_j, \cdot) \in \mathcal{K}_k$.

We define $\mathcal{C}_k$ with a univariate kernel $k$, as consisting of functions, $f : \mathbb{R}^d \to \mathbb{R}^d$, whose component functions $f_j : \mathbb{R}^d \to \mathbb{R}$ satisfy $f_{j, \boldsymbol{x}_{-j}} \equiv f_j(\cdot, \boldsymbol{x}_{-j}) \in \mathcal{K}_k$ for each $\boldsymbol{x}_{-j}$. So $f_j$ with a fixed $\boldsymbol{x}_{-j}$ is in the RKHS defined by $k$. This means

$$f_{j, \boldsymbol{x}_{-j}}(x_j) = \langle f_{j, \boldsymbol{x}_{-j}}, k(x_j, \cdot) \rangle_{\mathcal{K}_k} . \tag{4}$$

Let $\mathcal{C}_k^j$ denote the set of functions satisfying Equation (4) with norm bounded by

$$\left\| f_{j, \boldsymbol{x}_{-j}} \right\|_{\mathcal{K}_k} \leq \left\| \mathbb{E}_{q(x_j \mid \boldsymbol{x}_{-j})} \left[ \mathcal{A}^j_{p(x_j \mid \boldsymbol{x}_{-j})} \Phi_{x_j} \right] \right\|_{\mathcal{K}_k} , \tag{5}$$

for all $\boldsymbol{x}_{-j} \in \mathbb{R}^{d-1}$.

We define the kernelized complete conditional Stein discrepancy (KCC-SD) $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k)$ as follows,

$$\sum_{j=1}^d \mathbb{E}_{q(\boldsymbol{x}_{-j})} \left[ \left\| \sup_{f_j \in \mathcal{C}_k^j} \mathbb{E}_{q(x_j \mid \boldsymbol{x}_{-j})} \left[ \mathcal{A}^j_{p(x_j \mid \boldsymbol{x}_{-j})} f_j(\boldsymbol{x}) \right] \right\| \right] \tag{6}$$

**KCC-SDs admit a closed form.** In our definition of the Stein set, we can change the kernel or the kernel parameters in each dimension, however for clarity we do not focus on that here. Note that the Stein set depends on both distributions $p$ and $q$. We show that the KCC-SD defined in Eq. (6) has a closed form.

**Theorem 1** (Closed form). *For a kernel $k$ which is differentiable in both arguments, we define the Stein kernel for each $j \in \{1, \ldots, d\}$ as follows:*

$$k_{cc}^j(x_j, y_j; \boldsymbol{x}_{-j}) = \mathcal{A}^j_{p(x_j \mid \boldsymbol{x}_{-j})} \mathcal{A}^j_{p(y_j \mid \boldsymbol{x}_{-j})} k(x_j, y_j) \tag{7}$$
$$= b_j(x_j, \boldsymbol{x}_{-j}) b_j(y_j, \boldsymbol{x}_{-j}) k(x_j, y_j) + b_j(x_j, \boldsymbol{x}_{-j}) \nabla_{y_j} k(x_j, y_j)$$
$$+ b_j(y_j, \boldsymbol{x}_{-j}) \nabla_{x_j} k(x_j, y_j) + \nabla_{x_j} \nabla_{y_j} k(x_j, y_j) ,$$

4

where $b_j(\boldsymbol{x})$ is equal to $\nabla_{x_j} \log p(\boldsymbol{x})$ and if $\mathbb{E}_{q(\boldsymbol{x}_{-j})} \mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})} \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})} \left[ k_{cc}^j (x_j, y_j; \boldsymbol{x}_{-j})^{1/2} \right] < \infty$, then the KCC-SD can be computed in closed form as $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) = \|\boldsymbol{w}\|_2^2$, where the weights, $w_j$ are defined as $w_j^2 = \mathbb{E}_{q(\boldsymbol{x}_{-j})} \mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})} \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})} k_{cc}^j (x_j, y_j; \boldsymbol{x}_{-j})$.

The proof is in Appendix A. Theorem 1 implies that the functions, $f_j^*(x_j; \boldsymbol{x}_{-j})$, which achieve the supremum in Equation (6) are

$$f_j^*(x_j; \boldsymbol{x}_{-j}) = \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})} \left[ \mathcal{A}_{p(y_j|\boldsymbol{x}_{-j})}^j \Phi_{x_j} \right] \tag{8}$$
$$= \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})} [k(x_j, y_j) \nabla_{y_j} \log p(y_j \mid \boldsymbol{x}_{-j}) + \nabla_{y_j} k(x_j, y_j)],$$

where $\nabla_{y_j} \log p(y_j \mid \boldsymbol{x}_{-j}) = \nabla_{y_j} \log p(y_j, \boldsymbol{x}_{-j})$ and $\Phi_{x_j}(\cdot) = k(x_j, \cdot)$ is the feature map.

We can also restrict to functions to the unit ball, $\|f_{j, \boldsymbol{x}_{-j}}\|_{\mathcal{K}_k} \leq 1$, and still get a closed form for the KCC-SD:

$$\sum_j \mathbb{E}_{q(\boldsymbol{x}_{-j})} \sqrt{\mathbb{E}_{x_j, y_j \sim q(\cdot|\boldsymbol{x}_{-j})} k_{cc}^j (x_j, y_j; \boldsymbol{x}_{-j})} \,. \tag{9}$$

However, the closed form cannot be easily manipulated.

**KCC-SDs can distinguish two distributions.** We show that $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) = 0$ if and only if $p = q$. This proof relies on the ISPD property of the kernel and an equivalent form of the Stein operator when the score function of $q$ exists. For $f : \mathbb{R}^d \to \mathbb{R}^d$, note that as $\mathbb{E}_{q(\boldsymbol{x})} \left[ \mathcal{A}_{q(\boldsymbol{x})} f(\boldsymbol{x}) \right] = 0$,

$$\mathbb{E}_{q(\boldsymbol{x})} \left[ \mathcal{A}_{p(\boldsymbol{x})} f(\boldsymbol{x}) \right] = \mathbb{E}_{q(\boldsymbol{x})} \left[ \mathcal{A}_{p(\boldsymbol{x})} f(\boldsymbol{x}) - \mathcal{A}_{q(\boldsymbol{x})} f(\boldsymbol{x}) \right] = \mathbb{E}_{q(\boldsymbol{x})} \left[ f(\boldsymbol{x})^T \nabla_{\boldsymbol{x}} \left( \log p(\boldsymbol{x}) - \log q(\boldsymbol{x}) \right) \right] .$$

Using this representation, we prove that if $p$ is equal to $q$ in distribution, then KCC-SD is zero.

**Theorem 2.** *Suppose $k$ is an ISPD kernel and twice differentiable in both arguments, and $\mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})\|^2], \mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log q(\boldsymbol{x})\|^2] < \infty$ where $p(\boldsymbol{x}), q(\boldsymbol{x}) > 0$ for all $\boldsymbol{x} \in \mathbb{R}^d$. If $p \overset{d}{=} q$, then $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) = 0$.*

This property can be see by noting that when both $p$ and $q$ have score functions, their difference will be zero inside the operator. The proof is available in Appendix C. Similarly if $p$ is not equal to $q$ in distribution, KCC-SD will be able to detect that.

**Theorem 3.** *Let $k$ be integrally strictly positive definite. Suppose if $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) < \infty$, and $\mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})\|^2], \mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log q(\boldsymbol{x})\|^2] < \infty$ with $p(\boldsymbol{x}), q(\boldsymbol{x}) > 0$, then if $p$ is not equal to $q$ in distribution, then $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) > 0$.*

The proof is in Appendix C. Combined with the previous result, this shows that KCC-SDs are non-negative and zero only when the two distributions are equal.

## 4 KCC-SD in practice

Computing the optimal test function in KCC-SDs, $f_j^*(x_j; \boldsymbol{x}_{-j})$, requires sampling from the complete conditionals, $y_j \sim q(\cdot \mid \boldsymbol{x}_{-j})$. In this section, we detail how to compute KCC-SD when the complete conditionals can be sampled. We also present a sampling procedure which can be used to compute a lower bound of KCC-SD when the complete conditionals cannot be exactly sampled.

**Exact KCC-SD.** In Algorithm 1 in Appendix A we describe how to compute KCC-SDs, given a dataset $\{\boldsymbol{x}^i\}$ and complete conditionals $q(\cdot \mid \boldsymbol{x}_{-j})$ which can be sampled. For instance, KCC-SDs can be used to assess the sample quality of samples from a Gibbs sampler. Here the Gibbs sampler can be used to generate multiple auxiliary coordinates $y_j^{(i,k)} \sim p(\cdot \mid \boldsymbol{x}_{-j}^{(i)})$ using the sampling procedure for the complete conditional used in the Gibbs sampler. The auxiliary coordinate variables can be used to compute KCC-SD and can be used to assess the quality of the empirical distribution $q_n$ defined by the samples $\{\boldsymbol{x}^{(i)}\}_{i=1}^n$.

5

**Approximate KCC-SD.** Sampling from the complete conditional can be infeasible in several scenarios. To resolve this, we introduce approximate KCC-SDs, $\mathcal{S}_\lambda(q, \mathcal{A}_p, \mathcal{C}_k)$. Suppose $g_j(\boldsymbol{x}) = \mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}[\mathcal{A}^j_{p(y_j|\boldsymbol{x}_{-j})}\Phi_{x_j}]$, where $r_{\lambda_j}$ is a conditional distribution, then we define approximate KCC-SD as

$$\mathcal{S}_\lambda(q, \mathcal{A}_p, \mathcal{C}_k) = \sum_{j=1}^{d} \mathbb{E}_{q(\boldsymbol{x}_{-j})}\mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathcal{A}^j_{p(x_j|\boldsymbol{x}_{-j})}g_j(\boldsymbol{x}).$$

Algorithm 2 in Appendix B summarizes how to compute approximate KCC-SD. We split the dataset $\{\boldsymbol{x}\}_{i=1}^{n}$ into a training, validation and test set. We train a sampler on the training set and select the model based on the lowest loss on the validation set, and then generate samples $y_j$ from that model. KCC-SD is then computed on the test set.

The reduction to probabilistic regression can make use of powerful models, such as conditional kernel density estimation (Hansen, 2004) or neural network based models. The quality of approximate KCC-SD depends on the performance of the learned sampler on held-out data; this performance can be checked on a validation set. Formally, if the distributions $\{r_{\lambda_j}\}_{j=1}^{d}$ satisfy a $\rho$-transport inequality (Definition 3.58, (Wainwright, 2019)) and satisfy $\sup_{\boldsymbol{x}_{-j}} \text{KL}(q(\cdot|\boldsymbol{x}_{-j}) \| r_{\lambda_j}) < \epsilon_j$, then we can bound the difference between approximate KCC-SD and KCC-SD.

**Lemma 1.** *Suppose the model class $r_{\lambda_j}$ satisfies a $\rho$-transport inequality and $\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})$ is Lipschitz and $\mathbb{E}_q[\|\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})\|], \mathbb{E}_{r_{\lambda_j}}[\|\nabla_{x_j} \log p(x_j \mid \boldsymbol{x}_{-j})\|] < \infty$, and the kernel $k$ is bounded with $\nabla_{x_j} k(x_j, y_j)$ Lipschitz, then*

$$|\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) - \mathcal{S}_\lambda(q, \mathcal{A}_p, C_k)| \leq \sum_{j=1}^{d} K_{1,j}\sqrt{2\rho^2\epsilon_j} + \sqrt{K_{2,j}\sqrt{2\rho^2\epsilon_j}}$$

*where $\sup_{\boldsymbol{x}_{-j}} \text{KL}(q(\cdot|\boldsymbol{x}_{-j}) \| r_{\lambda_j}) < \epsilon_j$ and $K_{1,j}, K_{2,j}$ are positive constants.*

The proof is in Appendix D. This gives us a selection criterion for selecting models, models with a lower validation loss have approximate KCC-SD values closer to KCC-SD.

**Goodness of Fit testing.** To show the efficacy of KCC-SD and approximate KCC-SD in distinguishing distributions, we introduce a goodness-of-fit test to test whether a given set of samples come from a target distribution. Let the null be $H_0 : p = q$, and the alternate be $H_1 : p \neq q$. We do not compute the asymptotic null distribution of the normalized test statistic, instead we use the wild-bootstrap technique (Shao, 2010; Fromont et al., 2012; Chwialkowski et al., 2014, 2016). Define the function $h$ as

$$h(\boldsymbol{x}^{(i)}) = \sum_{j=1}^{d} \frac{1}{m} \sum_{k=1}^{m} k_{cc}(x_j^{(i)}, y_j^{(i,k)}; \boldsymbol{x}_{-j}^{(i)}),$$

where $y_j^{(i,k)} \sim q(\cdot \mid \boldsymbol{x}_{-j}^{(i)})$. The test statistic $T_n$ and the bootstrapped statistic $R_n$ are defined as

$$T_n = \frac{1}{n} \sum_{i=1}^{n} h(\boldsymbol{x}^{(i)}) \text{ and } R_n = \frac{1}{n} \sum_{i=1}^{n} \epsilon_i h(\boldsymbol{x}^{(i)}),$$

where $\epsilon_i$ are independent Rademacher random variables and $\boldsymbol{x}^{(i)}$ are independently and identically distributed from $q$. Sampling from the complete conditional is not always computationally feasible, therefore we propose another test with approximate KCC-SD as the test statistic, which samples $y_j^{(i,k)}$ from the model $r_{\lambda_j}$.

When the null hypothesis is true, the test statistic $T_n$ converges to zero (see Theorem 2 for KCC-SD and Lemma 5 in Appendix E for approximate KCC-SD), while $R_n$ converges to zero under both hypotheses. In Appendix E, we show that $\sqrt{n}R_n$ is a good approximation of $\sqrt{n}T_n$, so we can sample $R_n$ and approximate the quantiles of the null distribution.

When using KCC-SD, under the alternate hypothesis, $T_n$ converges to a positive constant (see Theorem 3) while $R_n$ converges to 0. Therefore, we reject the null hypothesis almost surely. The test can be formulated as
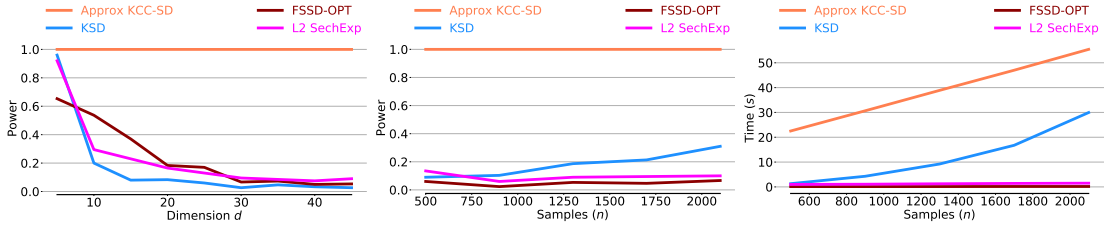
Figure 1: **KCC-SD has more power in high dimensions. Left:** Gaussian vs Laplace, with $n = 1000$ and increasing dimension. Approximate KCC-SD has no loss in power compared to baseline methods. **Middle and Right:** Gaussian vs Laplace, with $d = 30$ and increasing sample size. For all sample sizes studied, approximate KCC-SD has much higher power than the baseline methods, without requiring significantly more compute time.
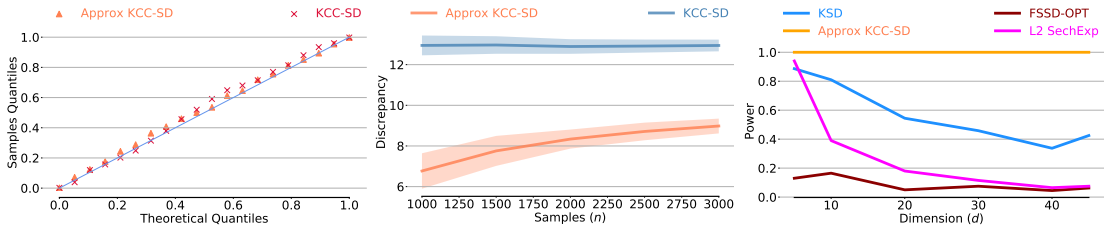


Figure 2: **Left:** Quantile-Quantile plot showing that KCC-SD and approximate KCC-SD have uniform $p$-values under the null, this was computed with $d = 30$ and $n = 3000$. **Middle:** Here we plot the value of KCC-SD and approximate KCC-SD with $p = N(0, I_d)$ and $q = N(0, \Sigma)$. Here the marginals match but $p \neq q$. As the number of samples increase, both discrepancies stay bounded away from zero. **Right:** Correlated Gaussian vs Correlated Gaussian with Laplace noise. As the dimension increases KCC-SD does not see a decrease in performance unlike the baseline methods.

1. Compute the test statistic $T_n$.

2. Compute the estimates $\{R_{n,l}\}_{l=1}^L$.

3. Estimate the $1 - \alpha$ empirical quantile of the samples.

4. Reject the null if $T_n$ exceeds the quantile.

When using approximate KCC-SD, under the null $T_n \to 0$ due to Stein's identity (see Lemma 5 in Appendix E) and the $p$-values are uniform. However, under the alternate the asymptotic behavior of approximate KCC-SD depends on the model class $r_{\lambda_j}$. We show in the experiments that approximate KCC-SD has power 1 in comparison to baselines such as KSD, RΦSD and FSSD-OPT.

## 5 Experiments

We study KCC-SD and approximate KCC-SD on comparing distributions, selecting parameters in samplers for Bayesian neural networks, and assessing the quality of Gibbs samplers for probabilistic matrix factorization on movie ratings.

For computing RΦSD, we use the hyperbolic secant kernel with the median heuristic (Huggins and Mackey, 2018). For the rest, we use the RBF kernel, $k(\boldsymbol{x}, \boldsymbol{y}) = \exp(- \|\boldsymbol{x} - \boldsymbol{y}\|^2 / 2\sigma^2)$. KCC-SD uses $\sigma = 1$, KSD uses the median heuristic, and FSSD-OPT learns the optimal $\sigma$ parameter. For FSSD-OPT we use the code and settings used by the authors in Jitkrittum et al. (2017).

To compute approximate KCC-SD we use a model for $r_{\lambda_j}$ based on histograms. Suppose the samples $x_j$ are in an interval $I$. Divide the interval $I$ into $m$ bins with width $\frac{1}{m}$ and learn a neural network $f_{\theta_j}(\boldsymbol{x}_{-j})$ which predicts the bin of $x_j$ from $\boldsymbol{x}_{-j}$. Sampling proceeds by sampling from the categorical distribution $b_k \sim Cat(f_{\theta_j}(\boldsymbol{x}_{-j}))$, and returning the average of the bin corresponding to $b_k$, the sample from the categorical distribution. See Appendix F for details.
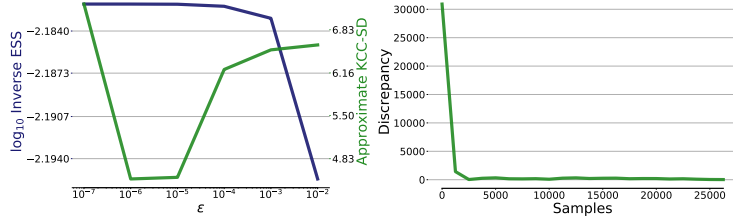
Figure 3: **Left:** Here, we plot the log inverse ESS for comparison to approximate KCC-SD in assessing quality of samples from SGLD. As we can see the inverse ESS is minimized at $10^{-3}$, and KCC-SD is minimized at $10^{-5}$. **Right:** The value of block KCC-SD decreases when the number of iterations goes up in the Gibbs sampler used for Bayesian Probabilistic Matrix Factorization.

**Goodness-of-fit Tests.** In the left panel of Figure 1 we compare samples from $q = \prod_{i=1}^{d} \text{Laplace}(0, 1/\sqrt{2})$ and target density $p = N(\mathbf{0}, I_d)$ with increasing dimension. We generate $n = 1000$ samples to compute the test statistics, and compute the power of the test over 300 repetitions with a significance level $\alpha = 0.05$. We then observe that as the dimension increases, approximate KCC-SD has power 1 while other methods see a substantial decrease in power as dimension increases. We show in Appendix F that similar results hold for the IMQ kernel.

In the middle panel of Figure 1 we plot the power of the test with $q = \prod_{i=1}^{d} \text{Laplace}(0, 1/\sqrt{2})$ and $p = N(\mathbf{0}, I_d)$ with $d = 30$. We then increase the number of samples used to compute the test statistics. And in the right panel of Figure 1 we show the time used to compute approximate KCC-SD, KSD, R$\Phi$SD and FSSD-OPT, the time for approximate KCC-SD also includes the training time for the models. We observe that although approximate KCC-SD requires more time to compute, it has more power than the baselines.

In the left panel of Figure 2 we compare $p = q = N(\mathbf{0}, \Sigma)$ in $d = 30$, with $\Sigma_{i,j} = 0.5$ for all $i \neq j$ otherwise $\Sigma_{i,i} = 1.0$. We show that for $n = 3000$ the distribution of the $p$-values is uniform.

In the middle panel of Figure 2, we have $p = N(\mathbf{0}, I_d)$ and $q = N(\mathbf{0}, \Sigma)$ where $\Sigma_{i,j} = 0.5$ for $i \neq j$ and $\Sigma_{i,i} = 1$. The figure shows that both KCC-SD and approximate KCC-SD detect the differences between these distributions.

In the right panel of Figure 2, we have $p = N(\mathbf{0}, \Sigma)$ with $\Sigma_{i,j} = 0.5$ and $\Sigma_{i,i} = 2$ and samples $\boldsymbol{x}_i = \boldsymbol{z}_i + \boldsymbol{\epsilon}_i$, where $\boldsymbol{\epsilon}_i \sim \prod_{j=1}^{d} \text{Laplace}(0, 1/\sqrt{2})$ and $\boldsymbol{z}_i \sim N(\mathbf{0}, \Sigma_1)$ with $(\Sigma_1)_{i,j} = 0.5$ and $(\Sigma_1)_{i,i} = 1$, and $\boldsymbol{z}_i$ and $\boldsymbol{\epsilon}_i$ are independent. The samples from $q$ have the same mean and variance as $p$. We compute $n = 500$ samples and increase the dimension. As the dimension increases, the power of the test with approximate KCC-SD remains 1, while the baseline methods see a decline in power.

**Selecting Biased Samplers.** In this experiment we do posterior inference for a three-layer neural network, with a sigmoid activation function, for a regression task. The hidden dimensions are 40 and 10. We make use of stochastic gradient Langevin dynamics (SGLD), a biased MCMC sampler (Welling and Teh, 2011). We used the yacht hydrodynamics dataset (Gerritsma et al., 1981) from the UCI dataset repository. Since biased methods trade sampling efficiency for asymptotic exactness, standard MCMC diagnostics like effective sample size are not applicable as they do not account for asymptotic bias. Selecting the stepsize $\epsilon$ is an important task to ensure the samples are approximately from the posterior (Welling and Teh, 2011). For $\epsilon \in [10^{-8}, 10^{-3}]$ we run a chain generating 10,000 samples with a burnin phase of 50,000 samples, with minibatch 256. We compare approximate KCC-SD to effective sample size. The left panel in Figure 3 compares these two metrics. While $\epsilon = 10^{-6}$ has the lowest KCC-SD value, the inverse effective sample size measure is minimized by the value $\epsilon = 10^{-2}$.

**Detecting Convergence of a Gibbs Sampler for Matrix Factorization.** We assess the convergence of a Gibbs sampler for Bayesian probabilistic matrix factorization (Salakhutdinov and Mnih, 2008). We focus on a variant with two mean parameters $\mu_V$ and $\mu_U$ for user and movie feature vectors $U_i \in \mathbb{R}^{10}, V_j \in \mathbb{R}^{10}$ and fixed the covariance matrix to the identity (see Appendix F for details).

8

In this experiment, we chose a subset of the Netflix Prize dataset, with 943 users and 1682 movies. We sampled the posterior $p(\boldsymbol{\mu}_U, \boldsymbol{\mu}_V, \boldsymbol{U}, \boldsymbol{V} \mid \boldsymbol{R})$ in blocks $\{\mu_U, \mu_V, U_1, \ldots, U_N, V_1, \ldots, V_M\}$ by a Gibbs sampler. We ran the sampler for 26K iterations with no burnin. Since the Gibbs sampler samples blocks of variables together, using these blocks of coordinates to compute KCC-SD is more efficient. In Appendix B we describe block KCC-SD. We compute block KCC-SD by taking every $5^{th}$ sample and show the results in the right panel of Figure 3. As the number of samples increases, block KCC-SD goes down. The sample quality of the Gibbs sample increases with the number of iterations.

## 6   Discussion

We developed kernelized complete conditional Stein discrepancies and approximate KCC-SD and corresponding goodness-of-fit tests. We show that these discrepancies can distinguish distributions which have smooth and integrable score functions. We also showed empirically that approximate KCC-SD provides a higher power test than those based on KSD. An interesting avenue of research would be relaxing the score function requirement for $q$ and to compare the relative efficiency of the test based on KCC-SD and approximate KCC-SD with baseline methods.

## Broader Impact

Our work focuses on comparing distributions where one is known in functional form up to a constant. The primary application of this method lies in probabilistic inference. Improvement in inference could help in building models in domains like healthcare and neuroscience especially to propagate uncertainty about the measurements. However, better inference could also mean better predictive models which can have downsides like in surveillance.

## References

Barbour, A. D. (1990). Stein's method for diffusion approximations. *Probability theory and related fields*, 84(3):297–322.

Barp, A., Briol, F.-X., Duncan, A., Girolami, M., and Mackey, L. (2019). Minimum stein discrepancy estimators. In *Advances in Neural Information Processing Systems*, pages 12964–12976.

Carmeli, C., De Vito, E., Toigo, A., and Umanitá, V. (2010). Vector valued reproducing kernel hilbert spaces and universality. *Analysis and Applications*, 8(01):19–61.

Chen, W. Y., Mackey, L., Gorham, J., Briol, F.-X., and Oates, C. J. (2018). Stein points. *arXiv preprint arXiv:1803.10161*.

Chwialkowski, K., Strathmann, H., and Gretton, A. (2016). A kernel test of goodness of fit. JMLR: Workshop and Conference Proceedings.

Chwialkowski, K. P., Sejdinovic, D., and Gretton, A. (2014). A wild bootstrap for degenerate kernel tests. In *Advances in neural information processing systems*, pages 3608–3616.

Eberle, A. (2016). Reflection couplings and contraction rates for diffusions. *Probability theory and related fields*, 166(3-4):851–886.

Fromont, M., Lerasle, M., Reynaud-Bouret, P., et al. (2012). Kernels based tests with non-asymptotic bootstrap approaches for two-sample problems. In *Conference on Learning Theory*, pages 23–1.

Geman, S. and Geman, D. (1984). Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *IEEE Transactions on pattern analysis and machine intelligence*, pages 721–741.

Gerritsma, J., Onnink, R., and Versluis, A. (1981). Geometry, resistance and stability of the delft systematic yacht hull series. *International shipbuilding progress*, 28(328):276–297.

Ghahramani, Z. and Beal, M. J. (2001). Propagation algorithms for variational bayesian learning. In *Advances in neural information processing systems*, pages 507–513.

Gong, W., Li, Y., and Hernández-Lobato, J. M. (2020). Sliced kernelized stein discrepancy. *arXiv preprint arXiv:2006.16531*.

Gorham, J. and Mackey, L. (2015). Measuring sample quality with stein's method. In *Advances in Neural Information Processing Systems*, pages 226–234.

Gorham, J. and Mackey, L. (2017). Measuring sample quality with kernels. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1292–1301. JMLR. org.

Hansen, B. E. (2004). Nonparametric conditional density estimation. *Unpublished manuscript*.

Huggins, J. and Mackey, L. (2018). Random feature stein discrepancies. In *Advances in Neural Information Processing Systems*, pages 1899–1909.

Hyvärinen, A. (2005). Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(Apr):695–709.

Jitkrittum, W., Xu, W., Szabó, Z., Fukumizu, K., and Gretton, A. (2017). A linear-time kernel goodness-of-fit test. In *Advances in Neural Information Processing Systems*, pages 262–271.

Lim, J. N., Yamada, M., Schölkopf, B., and Jitkrittum, W. (2019). Kernel stein tests for multiple model comparison. In *Advances in Neural Information Processing Systems*, pages 2240–2250.

Liu, Q., Lee, J., and Jordan, M. (2016). A kernelized stein discrepancy for goodness-of-fit tests. In *International conference on machine learning*, pages 276–284.

Liu, Q. and Wang, D. (2016). Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances in neural information processing systems*, pages 2378–2386.

Oates, C. J., Girolami, M., and Chopin, N. (2017). Control functionals for monte carlo integration. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3):695–718.

Ramdas, A., Reddi, S. J., Póczos, B., Singh, A., and Wasserman, L. (2015). On the decreasing power of kernel and distance based nonparametric hypothesis tests in high dimensions. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*.

Ranganath, R., Tran, D., Altosaar, J., and Blei, D. (2016). Operator variational inference. In *Advances in Neural Information Processing Systems*, pages 496–504.

Salakhutdinov, R. and Mnih, A. (2008). Bayesian probabilistic matrix factorization using markov chain monte carlo. In *Proceedings of the 25th international conference on Machine learning*, pages 880–887.

Shao, X. (2010). The dependent wild bootstrap. *Journal of the American Statistical Association*, 105(489):218–235.

Steinwart, I. and Christmann, A. (2008). *Support vector machines*. Springer Science & Business Media.

Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press.

Wang, D., Zeng, Z., and Liu, Q. (2017). Stein variational message passing for continuous graphical models. *arXiv preprint arXiv:1711.07168*.

Welling, M. and Teh, Y. W. (2011). Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688.

Zhuo, J., Liu, C., Shi, J., Zhu, J., Chen, N., and Zhang, B. (2017). Message passing stein variational gradient descent. *arXiv preprint arXiv:1711.04425*.

# A Closed Form

*Proof.* Define the Stein operator $\mathcal{A}_{p(\boldsymbol{x})}$ as follows,

$$(\mathcal{A}_{p(\boldsymbol{x})}f)(\boldsymbol{x}) = \sum_{j=1}^{d}(\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}^{j}f_j)(\boldsymbol{x}) = \sum_{j=1}^{d} f_j(\boldsymbol{x})\nabla_{x_j}\log p(\boldsymbol{x}) + \nabla_{x_j}f_j(\boldsymbol{x})$$

then if for all $j$, $f_{j,\boldsymbol{x}_{-j}}$ is in the RKHS of a univariate kernel, $k$, we can use the reproducing property, $f_{j,\boldsymbol{x}_{-j}}(x_j) = \langle f_{j,\boldsymbol{x}_{-j}}, k(x_j,\cdot)\rangle_{\mathcal{K}_k}$ (Steinwart and Christmann (2008)). Now, define the feature map for each kernel $k_j$, $\Phi_{x_j}(\cdot) = k(x_j,\cdot)$, then as

$$\begin{aligned}
\partial_{x_j} f_{j,\boldsymbol{x}_{-j}}(x_j) &= \partial_{x_j}\langle f_{j,\boldsymbol{x}_{-j}}, k(x_j,\cdot)\rangle_{\mathcal{K}_k} \\
&= \langle f_{j,\boldsymbol{x}_{-j}}, \partial_{x_j}k(x_j,\cdot)\rangle_{\mathcal{K}_k} \\
&= \langle f_{j,\boldsymbol{x}_{-j}}, \partial_{x_j}\Phi_{x_j}\rangle_{\mathcal{K}_k}
\end{aligned}$$

then note that we can use the reproducing property for general differential operators, $\mathcal{A}_{p(\boldsymbol{x})}^{j}$, to get

$$\begin{aligned}
(\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}f_j)(\boldsymbol{x}) &= \mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}\langle f_{j,\boldsymbol{x}_{-j}}, k(x_j,\cdot)\rangle_{\mathcal{K}_k} \\
&= \langle f_{j,\boldsymbol{x}_{-j}}, \mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}^{j}\Phi_{x_j}\rangle_{\mathcal{K}_k}
\end{aligned}$$

Then we can define the norm of $\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}\Phi_{x_j}$, as follows:

$$\begin{aligned}
\langle \mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}\Phi_{x_j}, \mathcal{A}_{p(y_j|\boldsymbol{x}_{-j})}\Phi_{y_j}\rangle_{\mathcal{K}_k} &= b_j(x_j,\boldsymbol{x}_{-j})b_j(y_j,\boldsymbol{x}_{-j})k(x_j,y_j) + \nabla_{x_j}\nabla_{y_j}k(x_j,y_j) \\
&\quad + b_j(x_j,\boldsymbol{x}_{-j})\nabla_{y_j}k(x_j,y_j) + b_j(y_j,\boldsymbol{x}_{-j})\nabla k(x_j,y_j) \\
&= k_{cc}^{j}(x_j,y_j;\boldsymbol{x}_{-j}) \quad (10)
\end{aligned}$$

where $b_j(u,\boldsymbol{x}_{-j}) = \nabla_u\log p(u|\boldsymbol{x}_{-j})$. Then we define the following

$$\begin{aligned}
w_j^2 &= \mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}\left[k_j^{cc}(x_j,y_j;\boldsymbol{x}_{-j})\right] \\
&= \mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}\left[\langle\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}\Phi_{x_j}, \mathcal{A}_{p(y_j|\boldsymbol{x}_{-j})}\Phi_{y_j}\rangle_{\mathcal{K}_k}\right] \\
&= \langle\mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}\Phi_{x_j}, \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}\mathcal{A}_{p(y_j|\boldsymbol{x}_{-j})}\Phi_{y_j}\rangle_{\mathcal{K}_k} \quad (11) \\
&= \left\|\mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}\Phi_{x_j}\right\|_{\mathcal{K}_k}^2 \quad (12)
\end{aligned}$$

where $x_j, y_j \overset{i.i.d}{\sim} q(\cdot\mid\boldsymbol{x}_{-j})$ and where we can interchange the inner product and expectation since $\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}\Phi_{x_j}$ is $q$-Bochner integrable, (Steinwart and Christmann (2008), Definition A.5.20).

We can find the closed form for KCC-SD, where KCC-SD is defined as follows:

$$\mathcal{S}(q,\mathcal{A}_p,\mathcal{C}_k) = \sum_{j=1}^{d}\mathbb{E}_{q(\boldsymbol{x}_{-j})}\left[\sup_{f_j\in\mathcal{C}_k}\left|\mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\left[\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}^{j}f_j(\boldsymbol{x})\right]\right|\right]$$

For each $j\in\{1,\ldots,d\}$, and $\boldsymbol{x}_{-j}$

$$\begin{aligned}
\sup_{f_j\in\mathcal{C}_k}\mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\left[\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}^{j}f_j(x)\right] &= \sup_{f_j:\|f_j\|\leq w_j^2}\langle f_j, \mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\left[\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}\Phi_{x_j}\right]\rangle_{\mathcal{K}_k} \\
&= \left\|\mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}\Phi_{x_j}\right\|_{\mathcal{K}_k}^2 \\
&= \mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}\left[k_{cc}^{j}(x_j,y_j;\boldsymbol{x}_{-j})\right]
\end{aligned}$$

hence, KCC-SD can be written in closed form as

$$\mathcal{S}(q,\mathcal{A}_p,\mathcal{C}_k) = \sum_{j=1}^{d}\mathbb{E}_{q(\boldsymbol{x}_{-j})}\mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}\left[k_{cc}^{j}(x_j,y_j;\boldsymbol{x}_{-j})\right]$$

$\square$

Here, we show that KCC-SDs can be expressed as an average of univariate KSDs. We can compute the Stein kernel for KCC-SD as

$$
\begin{aligned}
k_{cc}^j(x_j, y_j; \boldsymbol{x}_{-j}) &= k(x_j, y_j) b_j(x_j, \boldsymbol{x}_{-j}) b_j(y_j, \boldsymbol{x}_{-j}) + \nabla_{x_j} k(x_j, y_j) b_j(y_j, \boldsymbol{x}_{-j}) \\
&\quad + \nabla_{y_j} k(x_j, y_j) b_j(x_j, \boldsymbol{x}_{-j}) + \nabla_{x_j} \nabla_{y_j} k(x_j, y_j), \\
&= \left( \mathcal{A}_{p(x_j | \boldsymbol{x}_{-j})} \mathcal{A}_{p(y_j | \boldsymbol{x}_{-j})} k \right) (x_j, y_j)
\end{aligned}
$$

where $\boldsymbol{x}_{-j} \in \mathbb{R}^{d-1}$ is fixed, $k : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$, and $b_j(x_j, \boldsymbol{x}_{-j}) = \nabla_{x_j} \log p(x_j \mid \boldsymbol{x}_{-j})$. Using the Stein kernel defined above we can compute KSD between $p(\cdot \mid \boldsymbol{x}_{-j})$ and $q(\cdot \mid \boldsymbol{x}_{-j})$ as follows

$$
\begin{aligned}
\mathcal{S}\big(q(\cdot \mid \boldsymbol{x}_{-j}), \mathcal{A}_{p(\cdot | \boldsymbol{x}_{-j})}, \mathcal{G}_k\big)^2 &= \mathbb{E}_{q(x_j | \boldsymbol{x}_{-j})} \mathbb{E}_{q(y_j | \boldsymbol{x}_{-j})} \left[ \left( \mathcal{A}_{p(x_j | \boldsymbol{x}_{-j})} \mathcal{A}_{p(y_j | \boldsymbol{x}_{-j})} k \right) (x_j, y_j) \right] \\
&= \mathbb{E}_{q(x_j | \boldsymbol{x}_{-j})} \mathbb{E}_{q(y_j | \boldsymbol{x}_{-j})} \left[ k_{cc}^j(x_j, y_j; \boldsymbol{x}_{-j}) \right].
\end{aligned}
$$

Therefore, KCC-SD can also be computed as

$$
\begin{aligned}
\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) &= \sum_{j=1}^{d} \mathbb{E}_{q(\boldsymbol{x}_{-j})} \mathbb{E}_{q(x_j | \boldsymbol{x}_{-j})} \mathbb{E}_{q(y_j | \boldsymbol{x}_{-j})} \left[ k_{cc}^j(x_j, y_j; \boldsymbol{x}_{-j}) \right] \\
&= \sum_{j=1}^{d} \mathbb{E}_{q(\boldsymbol{x}_{-j})} \left[ \mathcal{S}\big(q(\cdot \mid \boldsymbol{x}_{-j}), \mathcal{A}_{p(\cdot | \boldsymbol{x}_{-j})}, \mathcal{G}_k\big)^2 \right].
\end{aligned}
$$

In Algorithm 1 we show how to compute KCC-SD exactly when we have samples from the complete conditionals.

---

**Algorithm 1: Computing KCC-SDs with complete conditionals**

**Input:** Dataset $\{\boldsymbol{x}^{(i)}\}_{i=1}^n$, $d$: dimension of $\boldsymbol{x}$, $n_y$: number of $y_j$ samples and complete conditionals $q(\cdot \mid \boldsymbol{x}_{-j})$
**Output:** Estimated KCC-SD $\hat{S}_n(q, \mathcal{A}_p, \mathcal{C}_k)$
**for** $j \in [d]$ **do**
    **for** $i \in [n]$ **do**
        Sample $y_j^{(i,k)} \sim q(\cdot \mid \boldsymbol{x}_{-j}^{(i)})$ for $k \in [n_y]$
    **end**
    Let $\hat{w}_j^2 = \frac{1}{n n_y} \sum_{i=1}^n \sum_{k=1}^{n_y} k_{cc}^j(x_j^{(i)}, y_j^{(i,k)}; \boldsymbol{x}_{-j}^{(i)})$
**end**
Let $\hat{S}_n(q, \mathcal{A}_p, \mathcal{C}_k) = \sum_{j=1}^{d} \hat{w}_j^2$

---

## B   KCC-SD in practice

**Block KCC-SD.** In Gibbs sampling, when variables are sampled together, using blocks of coordinates to compute KCC-SD will be computationally more efficient than using single coordinates. The complete conditional approach still ensures that block KCC-SD distinguishes the distributions $p$ and $q$. For instance, if $\boldsymbol{x} \in \mathbb{R}^d$, then let $I_1, \ldots, I_m$ be disjoint partitions of indices $\{1, \ldots, d\}$ such that $\cup_{j=1}^m I_j = \{1, \ldots, d\}$, then we can define block KCC-SD as

$$
\sum_{j=1}^{m} \mathbb{E}_{q(\boldsymbol{x}_{-I_j})} \sup_{f_{I_j}} \mathbb{E}_{q(\boldsymbol{x}_{I_j} | \boldsymbol{x}_{-I_j})} [\mathcal{A}_{p(\boldsymbol{x}_{I_j} | \boldsymbol{x}_{-I_j})}^j f_{I_j}(\boldsymbol{x})],
$$

here the the dimension of the kernel would depend on the block size, so $k_j : \mathbb{R}^{I_j} \times \mathbb{R}^{I_j} \to \mathbb{R}$. The supremum of the block KCC-SD is

$$
\sum_{j=1}^{m} \mathbb{E}_{q(\boldsymbol{x}_{-I_j})} \mathbb{E}_{\boldsymbol{x}_{I_j}, \boldsymbol{y}_{I_j} \sim q(\cdot | \boldsymbol{x}_{-I_j})} \left[ k_{cc}^j(\boldsymbol{x}_{I_j}, \boldsymbol{y}_{I_j}; \boldsymbol{x}_{-I_j}) \right].
$$

Note that if we take all the coordinates as one block, block KCC-SD is equivalent to KSD.

**Algorithm 2: Computing approximate KCC-SDs**. Given model class $r_{\lambda_j}$, compute approximate KCC-SD.

**Input:** Dataset $\mathcal{D} = \{\boldsymbol{x}^{(i)}\}_{i=1}^{n}$, $d$: dimension of $\boldsymbol{x}$, $n_y$: number of $y_j$ samples, and a model
　　　class $r_{\lambda_j}(\cdot \mid \boldsymbol{x}_{-j})$ for each complete conditional.
**Output:** Approximate KCC-SD
Split the dataset into training, validation and test sets.

**for** $j \in [d]$ **do**
　　Train the sampler $r_{\lambda_j}$ on training set.
　　Select the model $r_{\lambda_j}$ with lowest validation loss.
　　**for** $i \in [n]$ **do**
　　　　Sample $y_j^{(i,l)} \sim r_{\lambda_j}(\cdot \mid \boldsymbol{x}_{-j}^{(i)})$ for $l \in [n_y]$
　　**end**
　　Let $\hat{w}_j^2 = \frac{1}{n}\sum_{i=1}^{n} \frac{1}{n_y}\sum_{l=1}^{n_y} k_{cc}(x_j^{(i)}, y_j^{(i,l)}; \boldsymbol{x}_{-j}^{(i)})$.
**end**
Let $\hat{S}_\lambda(q, \mathcal{A}_p, \mathcal{C}_k) = \sum_{j=1}^{d} \hat{w}_j^2$

## C  Distinguishing Distributions

Here, we rely on the ISPD property of the kernel $k(x_j, y_j)$ so that for any function $f : \mathbb{R} \to \mathbb{R}$, we obtain

$$\int_{u \in \mathbb{R}} \int_{v \in \mathbb{R}} f(u)k(u,v)f(v)dudv > 0$$

for $\|f\| > 0$.

Note that we can write the Stein discrepancy as,

$$
\begin{aligned}
\mathbb{E}_{q(\boldsymbol{x})}\left[\mathcal{A}_{p(\boldsymbol{x})}f(\boldsymbol{x})\right] &= \mathbb{E}_{q(\boldsymbol{x})}\left[\mathcal{A}_{p(\boldsymbol{x})}f(y) - \mathcal{A}_{q(\boldsymbol{x})}f(\boldsymbol{x})\right] \\
&= \mathbb{E}_{q(\boldsymbol{x})}\left[f(\boldsymbol{x})^T\nabla_{\boldsymbol{x}}\log p(\boldsymbol{x}) + \nabla_{\boldsymbol{x}} \cdot f(\boldsymbol{x})\right] - \mathbb{E}_{q(\boldsymbol{x})}\left[f(\boldsymbol{x})^T\nabla_{\boldsymbol{x}}\log q(\boldsymbol{x}) + \nabla_{\boldsymbol{x}} \cdot f(\boldsymbol{x})\right] \\
&= \mathbb{E}_{q(\boldsymbol{x})}\left[f(\boldsymbol{x})^T\left(\nabla_{\boldsymbol{x}}\log p(\boldsymbol{x}) - \log q(\boldsymbol{x}))\right)\right] \\
&= \mathbb{E}_{q(\boldsymbol{x})}\left[f(\boldsymbol{x})^T\nabla_{\boldsymbol{x}}\log \frac{p(\boldsymbol{x})}{q(\boldsymbol{x})}\right] ,
\end{aligned}
\tag{13}
$$

using $\mathbb{E}_{q(\boldsymbol{x})}\left[\mathcal{A}_{q(\boldsymbol{x})}f(\boldsymbol{x})\right] = 0$.

Using this representation for our test function, $f_j^*(\boldsymbol{x}) = \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}[\mathcal{A}_{p(y_j|\boldsymbol{x}_{-j})}^j k(x_j, y_j)]$, where $y_j \sim q(\cdot \mid \boldsymbol{x}_{-j})$, we see that

$$
\begin{aligned}
f_j^*(\boldsymbol{x}) &= \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}[\mathcal{A}_{p(y_j|\boldsymbol{x}_{-j})}^j k(x_j, y_j)] - \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}[\mathcal{A}_{q(y_j|\boldsymbol{x}_{-j})}^j k(x_j, y_j)] \\
&= \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}\left[k(x_j, y_j)\nabla_{y_j}\log \frac{p(y_j \mid \boldsymbol{x}_{-j})}{q(y_j \mid \boldsymbol{x}_{-j})}\right] \\
&= \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})}\left[k(x_j, y_j)\nabla_{y_j}\log \frac{p(y_j, \boldsymbol{x}_{-j})}{q(y_j, \boldsymbol{x}_{-j})}\right] ,
\end{aligned}
\tag{14}
$$

then using the fact that $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) = \sum_{j=1}^{d} \mathbb{E}_{q(\boldsymbol{x})}[\mathcal{A}_{p(\boldsymbol{x})}^j f_j^*(\boldsymbol{x})]$, we obtain using Eq. (13) and Eq. (14)

$$
\begin{aligned}
\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) &= \mathbb{E}_{q(\boldsymbol{x})} \left[ f^*(\boldsymbol{x})^T \nabla_{\boldsymbol{x}} \log \frac{p(\boldsymbol{x})}{q(\boldsymbol{x})} \right] \\
&= \sum_{j=1}^{d} \mathbb{E}_{q(\boldsymbol{x}_{-j})} \left[ \mathbb{E}_{q(x_j | \boldsymbol{x}_{-j})} \left[ f_j^*(\boldsymbol{x}) \nabla_{x_j} \log \frac{p(\boldsymbol{x})}{q(\boldsymbol{x})} \right] \right] \\
&= \sum_{j=1}^{d} \mathbb{E}_{q(\boldsymbol{x}_{-j})} \left[ \mathbb{E}_{q(x_j | \boldsymbol{x}_{-j})} \mathbb{E}_{q(y_j | \boldsymbol{x}_{-j})} \left[ \nabla_{y_j} \log \frac{p(y_j, \boldsymbol{x}_{-j})}{q(y_j, \boldsymbol{x}_{-j})} k(x_j, y_j) \nabla_{x_j} \log \frac{p(\boldsymbol{x})}{q(\boldsymbol{x})} \right] \right] .
\end{aligned}
$$

Now, observe that for each $j \in \{1, \ldots, d\}$, with $r(u, \boldsymbol{x}_{-j}) = \nabla_u \log \frac{p(u, \boldsymbol{x}_{-j})}{q(u, \boldsymbol{x}_{-j})}$, we define a function $h$ over $\boldsymbol{x}_{-j}$

$$
\begin{aligned}
h(\boldsymbol{x}_{-j}) &= \mathbb{E}_{q(x_j | \boldsymbol{x}_{-j})} \mathbb{E}_{q(y_j | \boldsymbol{x}_{-j})} \left[ \nabla_{y_j} \log \frac{p(y_j, \boldsymbol{x}_{-j})}{q(y_j, \boldsymbol{x}_{-j})} k(x_j, y_j) \nabla_{x_j} \log \frac{p(\boldsymbol{x})}{q(\boldsymbol{x})} \right] \\
&= \mathbb{E}_{q(x_j | \boldsymbol{x}_{-j})} \mathbb{E}_{q(y_j | \boldsymbol{x}_{-j})} \left[ r(y_j, \boldsymbol{x}_{-j}) k(x_j, y_j) r(x_j, \boldsymbol{x}_{-j}) \right] \\
&= \int_{x_j} \int_{y_j} q(x_j \mid \boldsymbol{x}_{-j}) r(x_j, \boldsymbol{x}_{-j}) k(x_j, y_j) q(y_j \mid \boldsymbol{x}_{-j}) r(y_j, \boldsymbol{x}_{-j}) dx_j dy_j \\
&= \int_{x_j} \int_{y_j} g_{\boldsymbol{x}_{-j}}(x_j) k(x_j, y_j) g_{\boldsymbol{x}_{-j}}(y_j) dx_j dy_j
\end{aligned}
\tag{15}
$$

where $g_{\boldsymbol{x}_{-j}}(u) = q(u \mid \boldsymbol{x}_{-j}) r(u, \boldsymbol{x}_{-j}) = q(u \mid \boldsymbol{x}_{-j}) \nabla_u \log \frac{p(u, \boldsymbol{x}_{-j})}{q(u, \boldsymbol{x}_{-j})}$.

The proofs in this section rely on the next lemma, which states that if the complete conditionals match, then the distributions also match.

**Lemma 2.** *If $p(\boldsymbol{x})$, $q(\boldsymbol{x}) > 0$ for all $\boldsymbol{x} \in \mathbb{R}^d$ and $p(x_j | \boldsymbol{x}_{-j}) = q(x_j | \boldsymbol{x}_{-j})$ for all $\boldsymbol{x}_{-j}$ and $j$, then $p(\boldsymbol{x}) = q(\boldsymbol{x})$.*

*Proof (Lemma 2).* We prove by induction. If dimension of $x$ is 2, then $p(x_1 | x_2) = q(x_1 | x_2)$ and $p(x_2 | x_1) = q(x_2 | x_1)$. Then we have

$$
\int \frac{p(x_1 | x_2)}{p(x_2 | x_1)} dx_1 = \int \frac{p(x_1)}{p(x_2)} dx_1 = \frac{1}{p(x_2)},
$$

and

$$
\int \frac{q(x_1 | x_2)}{q(x_2 | x_1)} dx_1 = \int \frac{q(x_1)}{q(x_2)} dx_1 = \frac{1}{q(x_2)},
$$

which implies

$$
\frac{1}{p(x_2)} = \int \frac{p(x_1 | x_2)}{p(x_2 | x_1)} dx_1 = \int \frac{q(x_1 | x_2)}{q(x_2 | x_1)} dx_1 = \frac{1}{q(x_2)}.
$$

Therefore, $p(x_2) = q(x_2)$ for all $x_2$. $p(x_1, x_2) = p(x_1 | x_2) p(x_2) = q(x_1 | x_2) q(x_2) = q(x_1, x_2)$.

Assume the dimension of $\boldsymbol{x}$ is $d$. Then we have

$$
\frac{p(\boldsymbol{x}_{-\{i,j\}})}{p(\boldsymbol{x}_{-i})} = \int \frac{p(\boldsymbol{x}_{-j})}{p(\boldsymbol{x}_{-i})} dx_i = \int \frac{p(x_i | \boldsymbol{x}_{-i})}{p(x_j | \boldsymbol{x}_{-j})} dx_i = \int \frac{q(x_i | \boldsymbol{x}_{-i})}{q(x_j | \boldsymbol{x}_{-j})} dx_i = \int \frac{q(\boldsymbol{x}_{-j})}{q(\boldsymbol{x}_{-i})} dx_i = \frac{q(\boldsymbol{x}_{-\{i,j\}})}{q(\boldsymbol{x}_{-i})}
$$

for all $j$. Then $p(\boldsymbol{x}_j | \boldsymbol{x}_{-\{i,j\}}) = q(\boldsymbol{x}_j | \boldsymbol{x}_{-\{i,j\}})$ for all $j$. Since $\boldsymbol{x}_{-i}$ is a $(d-1)$ dimensional distribution, we can use the induction. Since $p(\boldsymbol{x}_j | \boldsymbol{x}_{-\{i,j\}}) = q(\boldsymbol{x}_j | \boldsymbol{x}_{-\{i,j\}})$ for all $j$, by induction, we have $p(\boldsymbol{x}_{-i}) = q(\boldsymbol{x}_{-i})$. Therefore,

$$
p(\boldsymbol{x}) = p(x_i | \boldsymbol{x}_{-i}) p(\boldsymbol{x}_{-i}) = q(x_i | \boldsymbol{x}_{-i}) q(\boldsymbol{x}_{-i}) = q(\boldsymbol{x}).
$$

$\square$

Using Equation (13) we can see that if $p \overset{d}{=} q$, then $\mathbb{E}_q[\mathcal{A}_p f(\boldsymbol{x})] = 0$ for $f$ integrable and smooth. The Stein set for KCC-SD, $\mathcal{C}_k$, consists of such functions. We restate Theorem 2 for clarity.

**Theorem.** *Suppose* $k \in C^{2,2}(\mathbb{R}, \mathbb{R})$ *is an* ISPD *kernel and* $\mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})\|^2], \mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log q(\boldsymbol{x})\|^2] < \infty$ *where* $p(\boldsymbol{x}), q(\boldsymbol{x}) > 0$ *for all* $\boldsymbol{x} \in \mathbb{R}^d$. *If* $p \overset{d}{=} q$, *then* $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) = 0$.

***Proof (Theorem 2).*** If $p \overset{d}{=} q$, then the score functions match and using Equation (13), for all $f$ such that $\mathbb{E}_{q(\boldsymbol{x})}\|f(\boldsymbol{x})\|_2 < \infty$, then

$$\mathbb{E}_{q(\boldsymbol{x})}\left[\mathcal{A}_{p(\boldsymbol{x})} f(\boldsymbol{x})\right] = \mathbb{E}_{q(\boldsymbol{x})}\left[f(\boldsymbol{x})^T \nabla_{\boldsymbol{x}} \log \frac{p(\boldsymbol{x})}{q(\boldsymbol{x})}\right]$$
$$= 0$$

Since all $f \in \mathcal{C}_k$ satisfy $\mathbb{E}_{q(\boldsymbol{x})}\|f(\boldsymbol{x})\|_2 < \infty$, $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) = 0$.

$\square$

Similarly, using Equation (13) we can show that when $p \neq q$, then KCC-SD will be strictly greater than zero. This relies on the fact that if two measures are not equal, then on the set where they are not equal, the complete conditionals will not match. We can exploit this property to show that KCC-SD will not be zero for such distributions. We restate Theorem 3 for clarity.

**Theorem 4.** *Let $k$ be integrally strictly positive definite. Suppose if $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) < \infty$, and* $\mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})\|^2], \mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log q(\boldsymbol{x})\|^2] < \infty$ *with* $p(\boldsymbol{x}), q(\boldsymbol{x}) > 0$, *then if $p$ is not equal to $q$ in distribution, then $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) > 0$.*

***Proof (Theorem 3).*** Suppose $p \neq q$ in distribution, then by Lemma 2 there exists a $j \in \{1, \ldots, d\}$ and a set $B_{-j} \subset \mathbb{R}^{d-1}$, with $m_{d-1}(B_{-j}) > 0$ where $m_{d-1}$ is Lebesgue measure, such that for each $\boldsymbol{x}_{-j} \in B_{-j}$ there exists a set $A_{j,\boldsymbol{x}_{-j}} \subset \mathbb{R}$ with $m_1(A_{j,\boldsymbol{x}_{-j}}) > 0$, where the complete conditional do not match. Then as the complete conditionals, $p(x_j \mid \boldsymbol{x}_{-j}), q(x_j \mid \boldsymbol{x}_{-j})$, do not match on $A_{j,\boldsymbol{x}_{-j}}$, the ratio of the score functions do not match, so for $\boldsymbol{x}_{-j} \in B_{-j}$ and $u \in A_{j,\boldsymbol{x}_{-j}}$,

$$g_{\boldsymbol{x}_{-j}}(u) = q(u \mid \boldsymbol{x}_{-j}) \nabla_{x_j} \log \frac{p(u, \boldsymbol{x}_{-j})}{q(u, \boldsymbol{x}_{-j})} \neq 0 .$$

As $q$ has full support, for all $\boldsymbol{x}_{-j} \in B_{-j}$ we have $g_{\boldsymbol{x}_{-j}}(u) \neq 0$ on $A_{j,\boldsymbol{x}_{-j}}$, this implies that the $L_2$ norm of this function is not zero, $\|g_{\boldsymbol{x}_{-j}}\|_2 \neq 0$. Thus, for $\boldsymbol{x}_{-j} \in B_{-j}$, by the ISPD property of the kernel,

$$h(\boldsymbol{x}_{-j}) = \int_{x_j} \int_{y_j} g_{\boldsymbol{x}_{-j}}(x_j) k(x_j, y_j) g_{\boldsymbol{x}_{-j}}(y_j) dx_j dy_j > 0$$

and since $m_{d-1}(B_{-j})$, $\mathbb{E}_{q(\boldsymbol{x}_{-j})}[h(\boldsymbol{x}_{-j})] > 0$. Thus, $\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) > 0$.

$\square$

# D   Proof of Lemma 4: Bounding the gap in approximate KCC-SD

To prove Lemma 4 we make use of the following lemma (Gorham and Mackey, 2017) to bound the difference between the expectation of the Stein operator on different distributions.

**Lemma 3.** *Suppose $\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})$ is Lipschitz and $L_2(q) \cap L_2(r)$, and $f$ and $\nabla_{\boldsymbol{x}} f$ are uniformly bounded and Lipschitz, then we can show that*

$$\left|\mathbb{E}_{q(\boldsymbol{x})}[\mathcal{A}_{p(\boldsymbol{x})} f(\boldsymbol{x})] - \mathbb{E}_{r(\boldsymbol{y})}[\mathcal{A}_{p(\boldsymbol{y})} f(\boldsymbol{y})]\right| \leq K_1 W_2(q, r) + \sqrt{K_2 W_2(q, r)},$$

*where $K_1, K_2$ are positive constants.*

*Proof.* Suppose the score function, $s_p(\boldsymbol{x}) = \nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})$, is Lipschitz and the function $f$ is bounded with a Lipschitz derivative then we can bound the approximation error as follows

$$\left|\mathbb{E}_{q(\boldsymbol{x})}\left[\mathcal{A}_{p(\boldsymbol{x})} f(\boldsymbol{x})\right] - \mathbb{E}_{r(\boldsymbol{y})}\left[\mathcal{A}_{p(\boldsymbol{y})} f(\boldsymbol{y})\right]\right| \leq \left|\mathbb{E}_{q(\boldsymbol{x})}\left[f(\boldsymbol{x})^T s_p(\boldsymbol{x})\right] - \mathbb{E}_{r(\boldsymbol{y})}\left[f(\boldsymbol{y})^T s_p(\boldsymbol{y})\right]\right|$$
$$+ \left|\mathbb{E}_{q(\boldsymbol{x})} \nabla_{\boldsymbol{x}} \cdot f(\boldsymbol{x}) - \mathbb{E}_{r(\boldsymbol{y})} \nabla_{\boldsymbol{y}} \cdot f(\boldsymbol{y})\right|$$

Now, assume that $f$ is bounded and $\nabla \log p$ is Lipschitz and so is $\nabla_{\boldsymbol{x}} f$. Then, we can bound the second term above as follows

$$\left| \mathbb{E}_{q(\boldsymbol{x})} \nabla_{\boldsymbol{x}} \cdot f(\boldsymbol{x}) - \mathbb{E}_{r(\boldsymbol{y})} \nabla_{\boldsymbol{y}} \cdot f(\boldsymbol{y}) \right| \leq L(\nabla f) \mathbb{E}[\|\boldsymbol{x} - \boldsymbol{y}\|_2],$$

where $L(h)$ is the Lipschitz constant of the function $h$ and $B(h) = \sup_{\boldsymbol{x}} \|h(\boldsymbol{x})\|_2$.

Similarly, we split the first term as follows

$$\left| \mathbb{E}_{q(\boldsymbol{x})} \left[ f(\boldsymbol{x})^T s_p(\boldsymbol{x}) \right] - \mathbb{E}_{r(\boldsymbol{y})} \left[ f(\boldsymbol{y})^T s_p(\boldsymbol{y}) \right] \right| \leq \left| \mathbb{E} \left[ f(\boldsymbol{x})^T \left( s_p(\boldsymbol{x}) - s_p(\boldsymbol{y}) \right) \right] \right| \\ + \left| \mathbb{E}[s_p(\boldsymbol{y})^T (f(\boldsymbol{y}) - f(\boldsymbol{x}))] \right|.$$

We can then bound the first term above using the fact that the function $f$ is bounded and the score function is Lipschitz.

$$\left| \mathbb{E} \left[ f(\boldsymbol{x})^T \left( s_p(\boldsymbol{x}) - s_p(\boldsymbol{y}) \right) \right] \right| \leq B(f) L(\nabla \log p) \mathbb{E}[\|\boldsymbol{x} - \boldsymbol{y}\|_2]$$

and similarly we can bound the second term by using the fact that the function $f$ is bounded and Lipschitz and the the score function is square integrable,

$$\left| \mathbb{E}[s_p(\boldsymbol{y})^T (f(\boldsymbol{y}) - f(\boldsymbol{x}))] \right| \leq \mathbb{E}[\|f(\boldsymbol{y}) - f(\boldsymbol{x})\|_2 \|s_p(\boldsymbol{y})\|_2] \\ \leq \mathbb{E} \left[ \min \left( 2B(f), L(f) \|\boldsymbol{x} - \boldsymbol{y}\|_2 \right) \|s_p(\boldsymbol{y})\| \right]$$

and then using the fact that $\min(a, b) \leq \sqrt{ab}$ for $a, b \geq 0$, and applying Cauchy-Schwarz again we obtain

$$\left| \mathbb{E}[s_p(\boldsymbol{y})^T (f(\boldsymbol{y}) - f(\boldsymbol{x}))] \right| \leq \sqrt{2B(f)L(f)} \mathbb{E} \left[ \|\boldsymbol{x} - \boldsymbol{y}\|_2^{\frac{1}{2}} \|s_p(\boldsymbol{y})\|_2 \right] \\ \leq \sqrt{2B(f)L(f)} \sqrt{\mathbb{E} \left[ \|\boldsymbol{x} - \boldsymbol{y}\|_2 \right] \mathbb{E} \left[ \|s_p(\boldsymbol{y})\|_2^2 \right]}.$$

We can then bound the all the terms by

$$\left| \mathbb{E}_{q(\boldsymbol{x})}[\mathcal{A}_{p(\boldsymbol{x})} f(\boldsymbol{x})] - \mathbb{E}_{r(\boldsymbol{y})}[\mathcal{A}_{p(\boldsymbol{y})} f(\boldsymbol{y})] \right| \leq K_1 \mathbb{E} \left[ \|\boldsymbol{x} - \boldsymbol{y}\|_2 \right] + \sqrt{K_2 \mathbb{E}[\|\boldsymbol{x} - \boldsymbol{y}\|_2]},$$

where $K_1 = L(\nabla f) + B(f) L(\nabla \log p)$ and $K_2 = 2B(f) L(f) \mathbb{E}_r[\|s_p(\boldsymbol{y})\|_2]$ are constants. Now, by taking the infimum over all joint distributions $P$ on $\boldsymbol{x}$ and $\boldsymbol{y}$, where the marginals match with $q$ and $r$, we obtain

$$\left| \mathbb{E}_{q(\boldsymbol{x})}[\mathcal{A}_{p(\boldsymbol{x})} f(\boldsymbol{x})] - \mathbb{E}_{r(\boldsymbol{y})}[\mathcal{A}_{p(\boldsymbol{y})} f(\boldsymbol{y})] \right| \leq K_1 W_2(q, r) + \sqrt{K_2 W_2(q, r)},$$

where the Wasserstein distance is defined as $W_2(p, q) = \inf_{P, \boldsymbol{x} \sim p, \boldsymbol{y} \sim q} \mathbb{E}_P \left[ \|\boldsymbol{x} - \boldsymbol{y}\|_2 \right]$.

$\square$

Suppose a distribution $q$ satisfies a $\rho$-transport inequality (Definition 3.58, (Wainwright, 2019)), then for any distribution $p$ we have the following inequality

$$W_2(q, p) \leq \sqrt{2\rho^2 D(p, q)}, \tag{16}$$

where $D$ is the KL divergence. We make use of the $\rho$-transport inequality to bound the Wasserstein-2 in the bound proved in Lemma 3.

For brevity we refer to the complete conditional $q(x_j \mid \boldsymbol{x}_{-j})$ as $q_{|\boldsymbol{x}_{-j}}(x_j)$. And we restate Lemma 4 below for reference.

**Lemma 4.** *Suppose the model class $r_{\lambda_j}$ satisfies a $\rho$-transport inequality and $\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})$ is Lipschitz and $\mathbb{E}_q[\|\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})\|], \mathbb{E}_{r_{\lambda_j}}[\|\nabla_{x_j} \log p(\boldsymbol{x} \mid \boldsymbol{x}_{-j})\|] < \infty$, and the kernel $k$ is bounded with $\nabla_{x_j} k(x_j, y_j)$ Lipschitz, then*

$$|\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) - \mathcal{S}_\lambda(q, \mathcal{A}_p, C_k)| \leq \sum_{j=1}^d K_{1,j} \sqrt{2\rho^2 \epsilon_j} + \sqrt{K_{2,j} \sqrt{2\rho^2 \epsilon_j}}$$

*where $\sup_{\boldsymbol{x}_{-j}} \text{KL}(q(\cdot | \boldsymbol{x}_{-j}) \| r_{\lambda_j}) < \epsilon_j$ and $K_{1,j}, K_{2,j}$ are positive constants.*

*Proof of Lemma 4.* Suppose the complete conditionals $r_{\lambda_j}$ satisfy a $\rho$-transport inequality Equation (16). Then suppose $f$ is bounded and has a Lipschitz derivative then using Lemma 3 we get the following bound for each $j$,

$$
\begin{aligned}
\left| \mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}[\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}f(x_j)] - \mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}[\mathcal{A}_{p(y_j|\boldsymbol{x}_{-j})}f(y_j)] \right| &\leq K_{1,j}W_2(r_{\lambda_j}, q) + \sqrt{K_{2,j}W_2(r_{\lambda_j}, q)} \\
&\leq K_{1,j}\sqrt{2\rho^2 D(q_{|\boldsymbol{x}_{-j}}, r_{\lambda_j})} \\
&\quad + \sqrt{K_{2,j}\sqrt{2\rho^2 D(q_{|\boldsymbol{x}_{-j}}, r_{\lambda_j})}}.
\end{aligned}
$$

Note Lemma 4 follows from Lemma 3 as the function $h(y_j) = \mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}[\mathcal{A}_{p(x_j|\boldsymbol{x}_{-j})}k(x_j, y_j)]$ satisfies the boundedness and Lipschitz assumption for Lemma 3. Therefore, we can show that if $\epsilon_j = \sup_{\boldsymbol{x}_{-j}} \mathrm{KL}(q_{|\boldsymbol{x}_{-j}}, r_{\lambda_j})$ then

$$
\begin{aligned}
|\mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k) - \mathcal{S}_\lambda(q, \mathcal{A}_p, C_k)| &\leq \sum_{j=1}^d \mathbb{E}_{q(\boldsymbol{x}_{-j})} \left| \mathbb{E}_{r_{\lambda_j}(z_j|\boldsymbol{x}_{-j})} \mathcal{A}_p h(z_j) - \mathbb{E}_{q(y_j|\boldsymbol{x}_{-j})} \mathcal{A}_p h(y_j) \right| \\
&\leq \sum_{j=1}^d K_{1,j}\sqrt{\epsilon_j} + \sqrt{K_{2,j}\sqrt{2\rho^2 \epsilon_j}}.
\end{aligned}
$$

$\square$

# E    Goodness of fit Testing

In this section we show that

1. When the null hypothesis is true, the bootstrapped statistics can be used to approximate quantile of the null distribution, so

$$
\sup_\beta \left| \mathbb{P}\left[\sqrt{n}T_n > \beta\right] - \mathbb{P}\left[\sqrt{n}R_n > \beta \mid \{\boldsymbol{x}^{(i)}\}_{i\leq n}\right] \right| \to 0
$$

as $n \to \infty$. This holds when the test statistic is computed using either KCC-SD or approximate KCC-SD.

2. When the alternate hypothesis is true, the test statistic computed using KCC-SD converges to a positive constant almost surely (Theorem 3), that is $\mathbb{P}[T_n > 0] \to 1$. And this leads to an almost sure rejection of the null asymptotically.

3. When the alternate hypothesis holds, the asymptotic behaviour of the test with approximate KCC-SD depends on the model $r_{\lambda_j}$.

The goodness-of-fit test using approximate KCC-SD makes use of the fact that approximate KCC-SD converges to zero as the number of samples increases, this can be seen immediately using Stein's identity. Stein's identity states that for bounded functions $f$ with a bounded derivative (Proposition 1, Gorham and Mackey (2015)), which vanish at infinity,

$$
\mathbb{E}_{p(\boldsymbol{x})}\left[\mathcal{A}_{p(\boldsymbol{x})}f(\boldsymbol{x})\right] = 0.
$$

Using Stein's identity we show that when $p = q$, approximate KCC-SD, $S_\lambda(q, \mathcal{A}_p, \mathcal{C}_k)$, is zero.

**Lemma 5.** *Suppose $k$ is bounded and twice differentiable in both arguments with bounded derivatives and both $k(x_j, y_j)$ and $\nabla_{x_j}k(x_j, y_j), \nabla_{y_j}k(x_j, y_j)$ vanish at infinity, and the score function $\nabla_{y_j} \log p(y_j \mid \boldsymbol{x}_{-j}) \in L_2(r_{\lambda_j})$ for $q(\boldsymbol{x}_{-j})$ almost surely and $r_{\lambda_j}$ is a density and $\mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})\|_2^2], \mathbb{E}_{q(\boldsymbol{x})}[\|\nabla_{\boldsymbol{x}} \log q(\boldsymbol{x})\|_2^2] < \infty$. Then if $p = q$, we have the following*

$$
\mathcal{S}_\lambda(q, \mathcal{A}_p, \mathcal{C}_k) = \sum_{j=1}^d \mathbb{E}_{q(\boldsymbol{x}_{-j})} \mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})} \mathcal{A}^j_{p(x_j|\boldsymbol{x}_{-j})} g_j(\boldsymbol{x}) = 0,
$$

*where $g_j(\boldsymbol{x}) = \mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}[\mathcal{A}^j_{p(y_j|\boldsymbol{x}_{-j})}k(x_j, y_j)]$.*

*Proof of Lemma 5.* We show that approximate KCC-SD is zero when $p = q$ by using Stein's identity, which states that for bounded functions $f$ with a bounded derivative and which vanish at infinity, we have the following

$$\mathbb{E}_{p(\boldsymbol{x})}\left[\mathcal{A}_{p(\boldsymbol{x})}f(\boldsymbol{x})\right] = 0$$

We show that the function $g_j(x_j; \boldsymbol{x}_{-j}) = \mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}[\mathcal{A}^j_{p(y_j|\boldsymbol{x}_{-j})}k(x_j, y_j)]$ is bounded, with a bounded derivative and vanishes at infinity in $x_j$ with a fixed $\boldsymbol{x}_{-j}$. Now, using Cauchy-Schwarz we have

$$|g_j(x_j; \boldsymbol{x}_{-j})| \leq \left|\mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}k(x_j, y_j)\nabla_{y_j}\log p(y_j, \boldsymbol{x}_{-j})\right| + \mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}\left|\nabla_{y_j}k(x_j, y_j)\right| \quad (17)$$

$$\leq \sqrt{\mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}[k(x_j, y_j)^2]\mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}[\nabla_{y_j}\log p(y_j, \boldsymbol{x}_{-j})^2]}$$

$$+ \mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}\left|\nabla_{y_j}k(x_j, y_j)\right|.$$

As the kernel $k$ and its derivative are bounded and both vanish at infinity and $\nabla_{y_j}\log p(y_j \mid \boldsymbol{x}_{-j}) \in L_2(r_{\lambda_j})$ for $q(\boldsymbol{x}_{-j})$ almost surely, we have that the function $g_j$ is bounded and as $x_j \to \infty$ the function $g_j$ converges to zero.

We also show that the function $g_j$ has a bounded derivative with respect to $x_j$. Using the inequality from Equation (17) we obtain

$$\left|\nabla_{x_j}g_j(x_j; \boldsymbol{x}_{-j})\right| \leq \left|\mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}\nabla_{x_j}k(x_j, y_j)\nabla_{y_j}\log p(y_j, \boldsymbol{x}_{-j})\right| + \mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}\left|\nabla_{y_j}\nabla_{x_j}k(x_j, y_j)\right|$$

$$\leq \sqrt{\mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}[\nabla_{x_j}k(x_j, y_j)^2]\mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}[\nabla_{y_j}\log p(y_j, \boldsymbol{x}_{-j})^2]}$$

$$+ \mathbb{E}_{r_{\lambda_j}(y_j|\boldsymbol{x}_{-j})}\left|\nabla_{x_j}\nabla_{y_j}k(x_j, y_j)\right|.$$

Therefore, as the function is bounded and vanishes at infinity with a bounded derivative, using Stein's identity (Proposition 1, Gorham and Mackey (2015)) for the univariate complete conditionals, we can show that for $q(\boldsymbol{x}_{-j})$ almost surely the following holds

$$\mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathcal{A}^j_{p(x_j|\boldsymbol{x}_{-j})}g_j(x_j; \boldsymbol{x}_{-j}) = 0.$$

as $p = q$.

This implies that $\mathcal{S}_\lambda(q, \mathcal{A}_p, \mathcal{C}_k) = \sum_{j=1}^{d}\mathbb{E}_{q(\boldsymbol{x}_{-j})}\mathbb{E}_{q(x_j|\boldsymbol{x}_{-j})}\mathcal{A}^j_{p(x_j|\boldsymbol{x}_{-j})}g_j(x_j; \boldsymbol{x}_{-j}) = 0.$

$\square$

Now, we show that under the null, the bootstrapped statistics $\sqrt{n}R_n$ can be used to estimate the quantiles of the null distribution. Define the test statistic $T_n$ as follows

$$T_n = \frac{1}{n}\sum_{i=1}^{n}h(\boldsymbol{x}^{(i)})$$

$$h(\boldsymbol{x}^{(i)}) = \sum_{j=1}^{d}\frac{1}{m}\sum_{k=1}^{m}k_{cc}(x_j^{(i)}, y_j^{(i,k)}; \boldsymbol{x}_{-j}^{(i)}),$$

where $y_j^{(i,k)} \sim q(\cdot \mid \boldsymbol{x}_{-j}^{(i)})$ for KCC-SD and $y_j^{(i,k)} \sim r_{\lambda_j}(\cdot \mid \boldsymbol{x}_{-j}^{(i)})$ for approximate KCC-SD. Under the null we have that $T_n \to 0$ (Theorem 2 for KCC-SD and Lemma 5 for approximate KCC-SD) almost surely. Then assuming that $\boldsymbol{x}^{(i)} \overset{i.i.d}{\sim} q$ we have

$$\sqrt{n}T_n = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}h(\boldsymbol{x}^{(i)}) \Rightarrow N(0, \sigma^2_{H_0}).$$

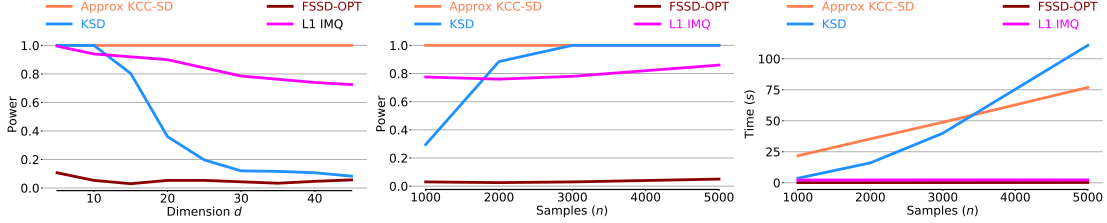as $\mathbb{E}[h(\boldsymbol{x}^{(i)})] = 0$ and $\mathbb{E}\left[h(\boldsymbol{x}^{(i)})^2\right] < \infty.$

Figure 4: **KCC-SD has more power than baseline methods with the IMQ kernel. Left:** We compute $n = 1000$ samples from $q = \prod_{i=1}^{d} \text{Laplace}(0, 1/\sqrt{2})$ with target density $p = N(0, I_d)$, and plot the power with increasing dimension. **Middle and Right:** Here we have samples from the same $p$ and $q$ distributions as before, but for fixed $d = 30$ we increase the number of samples $n$ to show the number of samples required and computation time for baseline methods to achieve similar power as KCC-SD.

In this work, we do not compute the variance and therefore we use the wild bootstrap procedure (Shao, 2010; Fromont et al., 2012; Chwialkowski et al., 2014, 2016). We then define the bootstrapped statistic $R_n$ as

$$R_n = \frac{1}{n} \sum_{i=1}^{n} \epsilon^i h(\boldsymbol{x}^{(i)}),$$

where $\epsilon_i$ are independent Rademacher random variables. Then note that under the null and under the alternate $R_n \to 0$ almost surely. We also observe that under the null

$$\sqrt{n} R_n = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \epsilon^i h(\boldsymbol{x}^{(i)}) \Rightarrow N(0, \sigma_{H_0}^2).$$

Note, the mean and variance of the normalized bootstrapped statistics match that of the normalized test statistic $\sqrt{n} T_n$,

$$\mathbb{E}[\epsilon_i h(\boldsymbol{x})] = 0, \text{ and } \mathbb{E}[\epsilon_i h(\boldsymbol{x})^2] = \mathbb{E}[h(\boldsymbol{x}^{(i)})^2].$$

Therefore, under the null we have $\sup_\beta \left| \mathbb{P}\left[\sqrt{n} T_n > \beta\right] - \mathbb{P}\left[\sqrt{n} R_n > \beta \mid \{\boldsymbol{x}^{(i)}\}_{i=1}^{n}\right] \right| \to 0$.

Under the alternative hypothesis, we note that by Theorem 3, $T_n \to C > 0$, when using KCC-SD. While, $R_n \to 0$ almost surely, therefore as $\mathbb{P}[T_n > 0] \to 1$ we reject the null almost surely.

When using approximate KCC-SD as a test statistic, the probability of rejection asymptotically is controlled by the quality of the model $r_{\lambda_j}$ as can be seen using Lemma 4,

$$|\mathcal{S}_\lambda(q, \mathcal{A}_p, \mathcal{C}_k) - \mathcal{S}(q, \mathcal{A}_p, \mathcal{C}_k)| \leq \sum_{j=1}^{d} K_{1,j} \sqrt{\epsilon_j} + \sqrt{K_{2,j} \sqrt{2\rho^2 \epsilon_j}}.$$

where $\epsilon_j = \sup_{\boldsymbol{x}_{-j}} \text{KL}(q_{|\boldsymbol{x}_{-j}}, r_{\lambda_j})$.

# F   Experiments

For the histogram-based sampler we use in our experiments, we use a two-layer neural network with 15-dimensional hidden-layer with a sigmoid activation function. We train the model with gradient descent for 500 epochs. We select the model with the lowest validation loss.

For FSSD-OPT we use 20% of the samples for training and in approximate KCC-SD we use 20% for training and 10% for validation.

The $p$-value is computed as the proportion of the bootstrapped statistics, $R_n$, greater than the test statistic, $T_n$. And the power is computed as the proportion of $p$-values less than the significance level, in other words the power is the rejection rate of the null when the alternate hypothesis is true.
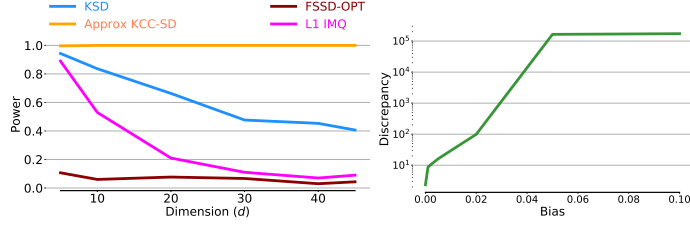
Figure 5: **Left:** Correlated Gaussian vs Correlated Gaussian with Laplace noise. As the dimension increases KCC-SD does not see a decrease in performance unlike the baseline methods. We use the IMQ kernel here. **Right:** As we add larger bias terms to the acceptance probability in the inner Metropolis sampler, samples from Metropolis-within-Gibbs sampler give larger KCC-SD.

**Choice of Kernel and Goodness-of-Fit tests.** All the experiments done with approximate KCC-SD, KSD and FSSD-OPT were done using the RBF kernel. The RBF kernel, a $C_0$ kernel (Definition 4.1, Carmeli et al. (2010)), suffices in defining consistent goodness-of-fit tests when comparing independent samples from a distribution $q$ (Theorem 2.2, Chwialkowski et al. (2016)).

Gorham and Mackey (2017) construct a sequence of empirical distributions, $q_n$, which does not converge to any distribution, a non-tight sequence. However, they prove that when comparing sequences $q_n$ with $p = N(\mathbf{0}, I_d)$ KSD with the RBF kernel still converges to zero. For this purpose they show that if KSD is computed with the IMQ kernel then KSD can enforce uniform tightness (Theorem 6, Gorham and Mackey (2017)).

However, as we have independent samples from a distribution $q$, this situation does not arise and we can use the RBF kernel. In Figure 4 we repeat the experiments from Figure 1 with IMQ kernel. In the left panel, we compare the power of the test using KCC-SD, R$\Phi$SD, KSD and FSSD-OPT with $q = \prod_{i=1}^n \text{Laplace}(0, 1/\sqrt{2})$ and $p = N(\mathbf{0}, I_d)$. We compute $n = 1000$ samples and then increase the dimension.

In the center and right panel of Figure 4, we compare the same distribution as above in $d = 30$ with an increasing number of samples. We observe that KCC-SD requires less samples than the baseline methods to have power 1, and for the baseline methods to have the same amount of power requires a similar amount of computation.

In the left panel of Figure 5, we have $p = N(\mathbf{0}, \Sigma)$ with $\Sigma_{i,j} = 0.5$ and $\Sigma_{i,i} = 2$ and samples $\boldsymbol{x}_i = \boldsymbol{z}_i + \boldsymbol{\epsilon}_i$, where $\boldsymbol{\epsilon}_i \sim \prod_{j=1}^d \text{Laplace}(0, 1/\sqrt{2})$ and $\boldsymbol{z}_i \sim N(\mathbf{0}, \Sigma_1)$ with $(\Sigma_1)_{i,j} = 0.5$ and $(\Sigma_1)_{i,i} = 1$, and $\boldsymbol{z}_i$ and $\boldsymbol{\epsilon}_i$ are independent. The samples from $q$ have the same mean and variance as $p$. We compute $n = 500$ samples and increase the dimension. As the dimension increases, the power of the KCC-SD test with the IMQ kernels remains 1, while the baseline methods with the IMQ kernels see a decline in power.

**Detecting Convergence of a Gibbs Sampler for Matrix Factorization.** Here we provide details of the probabilistic model considered in the experiments section for assessing the convergence of a Gibbs sampler for Bayesian probabilistic matrix factorization (Salakhutdinov and Mnih, 2008). We focus on a variant with two mean parameters $\mu_V$ and $\mu_U$ for user and movie feature vectors $U_i \in \mathbb{R}^{10}, V_j \in \mathbb{R}^{10}$ and fixed the covariance matrix to the identity.

$$p(\boldsymbol{U}|\boldsymbol{\mu}_U) = \prod_{i=1}^{N} N(U_i|\boldsymbol{\mu}_U, I), \quad p(\boldsymbol{\mu}_U) = N(0, I)$$

$$p(\boldsymbol{V}|\boldsymbol{\mu}_V) = \prod_{j=1}^{M} N(V_j|\boldsymbol{\mu}_V, I), \quad p(\boldsymbol{\mu}_V) = N(0, I)$$

$$p(\boldsymbol{R} \mid \boldsymbol{U}, \boldsymbol{V}) = \prod_{i=1}^{N} \prod_{j=1}^{M} \left[ N(R_{ij} \mid U_i^T V_j, I) \right]^{I_{ij}}$$

where $U_i, V_j$ have normal priors and $I_{ij}$ is the indicator variable that is one if user $i$ rated movie $j$ and 0 otherwise (see Appendix F).

**Selecting Biased Samplers.** We use a simple bimodal Gaussian mixture model to demonstrate the power of KCC-SD in distinguishing biased samplers,

$$x_i \sim \frac{1}{2} N(\theta_1, 2) + \frac{1}{2} (\theta_2, 2) \,,$$

where $\theta_1, \theta_2$ have standard normal priors. We draw 100 samples of $x_i$ from the model with $(\theta_1, \theta_2) = (1, -1)$. We choose Metropolis-within-Gibbs to sample from the posterior over $\boldsymbol{\theta}$. This sampler uses a Metropolis sampler to sample each complete conditional inside the Gibbs sampler. We also use the Metropolis step to generate auxiliary variables used to calculate KCC-SD. Denote $q(\boldsymbol{\theta})$ to be the target distribution. The inner Metropolis step accepts the candidate $\boldsymbol{\theta}_{new}$ with probability $\min\left(1, q(\boldsymbol{\theta}_{new})/q(\boldsymbol{\theta}_{old})\right)$. Then we add a bias term to the acceptance probability, $\min\left(1, q(\boldsymbol{\theta}_{new})/q(\boldsymbol{\theta}_{old}) + bias\right)$, thus the sampler is not unbiased anymore. We run for 60,000 iterations in total and drop the first 50,000 for burn-in. We show KCC-SDs versus size of the bias terms in the right panel of Figure 5. KCC-SD increases with the size of the bias.