

A. Terms and Notations

For the clarity of the paper, we give a summary of the commonly used notations in the main text and proof.

Symbols:

$\mathbf{s}_p(\mathbf{x})$	$\nabla_{\mathbf{x}} \log p(\mathbf{x})$
$s_p^r(\mathbf{x})$	Projected score function $\nabla_{\mathbf{x}} \log p(\mathbf{x})^T \mathbf{r}$
\mathcal{X}	A subset of \mathbb{R}^D
\mathcal{K}	A subset of \mathbb{R} .
k_{r,g_r}	kernel function $k : \mathcal{K} \times \mathcal{K} \rightarrow \mathcal{R}$
\mathcal{H}_{r,g_r}	Induced RKHS by the kernel k_{r,g_r} .
$\ \cdot\ _{\mathcal{H}_{r,g_r}}$	RKHS norm of \mathcal{H}_{r,g_r}
\mathbf{g}_r	Input projection direction (e.g. $\mathbf{x}^T \mathbf{g}_r$) for corresponding r .
\mathbf{r}	Score projection direction (e.g. $s_p^r(\mathbf{x}) = \mathbf{s}_p(\mathbf{x})^T \mathbf{r}$)
$S_{\max_{g_r}}$	$maxSSD-g$ (Eq.4).
S_{g_r}	$SSD-g$, i.e. $S_{\max_{g_r}}$ (Eq.4) without \sup_{g_r} . But with summation of O_r .
S_{r,g_r}	$SSD-rg$, i.e. $S_{\max_{g_r}}$ (Eq.4) without \sup_{g_r} and summation of O_r . Instead, we use specific r .
$SK_{\max_{g_r}}$	$maxSKSD-g$. The kernelized version of $S_{\max_{g_r}}$
SK_{g_r}	$SKSD-g$. The kernelized version of S_{g_r}
SK_{r,g_r}	$SKSD-rg$. The kernelized version of S_{r,g_r}
PSD	Projected Stein discrepancy (Eq.9)
PSD $_r$	Projected Stein discrepancy (Eq.9) without summation O_r and use specific r instead.
f_r^*	Optimal test function for PSD. $f_r^*(\mathbf{x}) \propto s_p^r(\mathbf{x}) - s_q^r(\mathbf{x})$
h_{r,g_r}^*	Optimal test function for S_{g_r} with specific r and \mathbf{g}_r , defined in Eq.8.
*	This indicates the optimal test function (e.g. f_r^*)
C_{\sup}	Supremum of Poincaré constant defined in assumption 6.

A.1. "Sub-optimal" variants of SSD

For the ease of the analysis, we want to define the notations without the sup operator over the slice directions \mathbf{r} , \mathbf{g}_r . Here, we define $SSD-g$ (S_{g_r}) as the $maxSSD-g$ ($S_{\max_{g_r}}$ in Eq.4) without the \sup_{g_r} .

$$S_{g_r} = \sum_{\mathbf{r} \in O_r} \sup_{h_{r,g_r} \in \mathcal{F}_q} \mathbb{E}_q[s_p^r(\mathbf{x})h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r) + \mathbf{r}^T \mathbf{g}_r \nabla_{\mathbf{x}^T \mathbf{g}_r} h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)] \quad (16)$$

Similarly, we define $SSD-rg$ (S_{r,g_r}) as $maxSSD-rg$ ($S_{\max_{r,g_r}}$ in Eq.37) without \sup_{r,g_r} :

$$S_{r,g_r} = \sup_{h_{r,g_r} \in \mathcal{F}_q} \mathbb{E}_q[s_p^r(\mathbf{x})h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r) + \mathbf{r}^T \mathbf{g}_r \nabla_{\mathbf{x}^T \mathbf{g}_r} h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)] \quad (17)$$

As for each of the above "optimal" discrepancies, it has the corresponding kernelized version. Therefore, we need to define their "un-optimal" version as well. We define $SKSD-g$ (SK_{g_r}) as $maxSKSD-g$ ($SK_{\max_{g_r}}$ in Eq.6) as

$$SK_{g_r} = \sum_{\mathbf{r} \in O_r} \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2, \quad (18)$$

Similarly, we define $SKSD-rg$ (SK_{r,g_r}) as $maxSKSD-rg$ ($SK_{\max_{r,g_r}}$ in Eq.41) as

$$SK_{r,g_r} = \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 \quad (19)$$

B. Assumptions and Definitions

Definition B.1 (Inner product in Hilbert space). *We denote the algebraic space \mathbb{R}^D refers to a parameter space of dimension D . The Borel sets of \mathbb{R}^D is denoted as $\mathcal{B}(\mathbb{R}^D)$, and we let $\mu(x)$ be a probability measure on \mathbf{x} . We define*

$$\mathcal{H}_\mu = L^2(\mathbb{R}^D, \mathcal{B}(\mathbb{R}^D), \mu) \quad (20)$$

as the Hilbert space which contains all the measurable functions $f : \mathbb{R}^D \rightarrow \mathbb{R}$, such that $\|f\|_{\mathcal{H}_\mu} \leq \infty$, where we define inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}_\mu}$ to be

$$\langle f, g \rangle_{\mathcal{H}_\mu} = \int f(\mathbf{x})g(\mathbf{x})d\mu(\mathbf{x}) \quad (21)$$

for all $f, g \in \mathcal{H}_\mu$

Definition B.2. (*Stein Class (Liu et al., 2016)*) *Assume distribution q has continuous and differentiable density $q(\mathbf{x})$. A function f defined on the domain $\mathcal{X} \subseteq \mathbb{R}^D$, $f : \mathcal{X} \rightarrow \mathbb{R}$ is in the **Stein class of q** if f is smooth and satisfies*

$$\int_{\mathcal{X}} \nabla_{\mathbf{x}}(f(\mathbf{x})q(\mathbf{x}))d\mathbf{x} = 0 \quad (22)$$

We call a function $f(\mathbf{x}) \in \mathcal{F}_q$ if f belongs to the Stein class of q . We say vector-valued function $\mathbf{f}(\mathbf{x}) : \mathcal{X} \subseteq \mathbb{R}^D \rightarrow \mathbb{R}^m \in \mathcal{F}_q$ if each component of \mathbf{f} belongs to the Stein class of q .

Definition B.3 (Stein Identity). *Assume q is a smooth density satisfied assumption 1, then we have*

$$\mathbb{E}_q[s_q(\mathbf{x})f(\mathbf{x})^T + \nabla f(\mathbf{x})] = 0 \quad (23)$$

for any functions $f : \mathcal{X} \subseteq \mathbb{R}^D \rightarrow \mathbb{R}^D$ in Stein class of q .

We can easily see that the above holds true for $\mathcal{X} = \mathbb{R}^D$ if

$$\lim_{\|\mathbf{x}\| \rightarrow \infty} q(\mathbf{x})f(\mathbf{x}) = 0 \quad (24)$$

Assumption 1 (Properties of densities) Assume the two probability distributions p, q has continuous differentiable density $p(\mathbf{x}), q(\mathbf{x})$ supported on $\mathcal{X} \subseteq \mathbb{R}^D$, such that the induced set $\mathcal{K} = \{y \in \mathbb{R} | y = \mathbf{x}^T \mathbf{g}, \|\mathbf{g}\|^2 = 1, \mathbf{x} \in \mathcal{X}\}$ is *locally compact Hausdorff* (LCH) for all possible $\mathbf{g} \in \mathbb{S}^{D-1}$. If $\mathcal{X} = \mathbb{R}^D$, then the density q satisfies: $\lim_{\|\mathbf{x}\| \rightarrow \infty} q(\mathbf{x}) = 0$. If $\mathcal{X} \subset \mathbb{R}^D$ is compact, then $q(\mathbf{x}) = 0$ at boundary $\partial\mathcal{X}$.

Assumption 2 (Regularity of score functions) Denote the score function of $p(\mathbf{x})$ as $\mathbf{s}_p(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) \in \mathbb{R}^D$ and score function of $q(\mathbf{x})$ accordingly. Assume the score functions are bounded continuous differentiable functions and satisfying

$$\begin{aligned} \int_{\mathcal{X}} q(\mathbf{x}) |(s_p(\mathbf{x}) - s_q(\mathbf{x}))^T \mathbf{r}| d\mathbf{x} < \infty \\ \int_{\mathcal{X}} q(\mathbf{x}) \|(s_p(\mathbf{x}) - s_q(\mathbf{x}))^T \mathbf{r}\|^2 d\mathbf{x} < \infty \end{aligned} \quad (25)$$

for all \mathbf{r} where $\mathbf{r} \in \mathbb{S}^{D-1}$.

Assumption 3 (Test functions) Assume the test function $h_{r, g_r} : \mathcal{K} \subseteq \mathbb{R} \rightarrow \mathbb{R}$ is smooth and belongs to the Stein class of q . Specifically, if with assumption 1, we only requires h_{r, g_r} to be a bounded continuous function. Similarly, we assume this also holds for PSD (eq.9) test function $f_r(\mathbf{x})$.

Assumption 4 (Bounded Conditional Expectation) Define

$$h_{r, g_r}^*(y_d) = \mathbb{E}_{q_{G_r}(\mathbf{y}_{-d}|y_d)}[(s_p^T(\mathbf{G}_r^{-1} \mathbf{y}) - s_q^T(\mathbf{G}_r^{-1} \mathbf{y}))] \quad (26)$$

as in proposition 1. We assume h_{r, g_r}^* is uniformly bounded for all possible $\mathbf{g}_r \in \mathbb{S}^{D-1}$.

Assumption 5 (universal kernel): We assume the kernel $k_{r, g} : \mathcal{K} \times \mathcal{K} \rightarrow \mathbb{R}$ is bounded and c_0 -universal.

Assumption 6 (Real analytic translation invariant kernel): We assume the kernel is translation invariant $k(x, y) = \phi(x - y) : \mathcal{K} \rightarrow \mathbb{R}$ and ϕ is a real analytic function. Additionally, we assume if $k(cx, cy) = k'(x, y)$ for a constant $c > 0$ where k' is also a c_0 -universal kernel. For example, *radial basis kernel function* (RBF) and *inverse multiquadric* (IMQ) kernel satisfy these assumptions.

Assumption 7 (Log-concave probabilities) Assume a probability distribution q with density function such that $q(\mathbf{x}) = \exp(-V(\mathbf{x}))$, where $V(\mathbf{x})$ is a convex function.

Assumption 8 (Existence of supremum of Poincaré constant). For the Poincaré constant defined in lemma 5, the essential supremum exists $C_{ess, \mathbf{G}} = \text{ess sup}_{y_d} C_{y_d} < \infty$ and also the $C_{sup} = \sup_{\mathbf{G}} C_{ess, \mathbf{G}} < \infty$ exists over all possible orthogonal matrix \mathbf{G} .

C. Detailed Background

C.1. Stein Discrepancy

Assume we have two differentiable probability density functions $q(\mathbf{x})$ and $p(\mathbf{x})$ where $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^D$. We further define a test function $\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^D$ and a suitable test function family \mathcal{F}_q called *Stein's class of q* . Recall the Stein operator (Eq.1) is defined as

$$\mathcal{A}_p \mathbf{f}(\mathbf{x}) = \mathbf{s}_p(\mathbf{x})^T \mathbf{f}(\mathbf{x}) + \nabla_{\mathbf{x}}^T \mathbf{f}(\mathbf{x}) \quad (27)$$

The function family \mathcal{F}_q is defined as

$$\mathcal{F}_q = \{\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^D \mid \mathbb{E}_q[\mathcal{A}_q \mathbf{f}] = 0\} \quad (28)$$

This function space can be quite general. For example, if $\mathcal{X} = \mathbb{R}^D$, we only require \mathbf{f} to be differentiable and vanishing at infinity. With all the notations, *Stein discrepancy* is defined as follows:

$$D_{SD}(q, p) = \sup_{\mathbf{f} \in \mathcal{F}_q} \mathbb{E}_q[\mathcal{A}_p \mathbf{f}(\mathbf{x})] \quad (29)$$

which can be proved to be a valid discrepancy (Gorham & Mackey, 2017). Stein discrepancy has been shown to be closely related to *Fisher discrepancy* defined as

$$D_F(q, p) = \mathbb{E}_q \|\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x})\|_2^2 \quad (30)$$

Indeed, Hu et al. (2018) shows that the optimal test function for *Stein discrepancy* has the form $\mathbf{f}^*(\mathbf{x}) \propto \mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x})$. By substitution, we can show *Stein discrepancy* is equivalent to *Fisher divergence* up to a multiplicative constant.

Unfortunately, the score difference $\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x})$ may be intractable in practice, making SD intractable as a consequence. Thus, Liu et al. (2016); Chwialkowski et al. (2016) propose an variant of SD by restricting \mathcal{F}_q to be a unit ball inside an RKHS \mathcal{H}_k induced by a c_0 -universal kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. By using the reproducing properties, they propose *kernelized Stein discrepancy* as

$$\begin{aligned} D^2(q, p) &= \|\mathbb{E}_q[\mathbf{s}_p(\mathbf{x})k(\mathbf{x}, \cdot) + \nabla_{\mathbf{x}} k(\mathbf{x}, \cdot)]\|_{\mathcal{H}_k}^2 \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim q}[u_p(\mathbf{x}, \mathbf{x}')] \end{aligned} \quad (31)$$

where $u_p(\mathbf{x}, \mathbf{x}')$ is

$$\begin{aligned} u_p(\mathbf{x}, \mathbf{x}') &= \mathbf{s}_p(\mathbf{x})^T k(\mathbf{x}, \mathbf{x}') \mathbf{s}_p(\mathbf{x}') + \mathbf{s}_p(\mathbf{x})^T \nabla_{\mathbf{x}'} k(\mathbf{x}, \mathbf{x}') \\ &+ \mathbf{s}_p(\mathbf{x}')^T \nabla_{\mathbf{x}} k(\mathbf{x}, \mathbf{x}') + \nabla_{\mathbf{x}, \mathbf{x}'}^2 k(\mathbf{x}, \mathbf{x}') \end{aligned} \quad (32)$$

and \mathbf{x}, \mathbf{x}' are i.i.d. samples from q .

Due to its tractability, it has been extensively used in statistical test e.g. GOF test Liu et al. (2016); Chwialkowski et al. (2016); Huggins & Mackey (2018); Jitkrittum et al. (2017). However, recent work demonstrate KSD suffers from the curse-of-dimensionality problem Gong et al. (2021); Huggins & Mackey (2018); Chwialkowski et al. (2016). One potential fix is to use another variant called *sliced kernelized Stein discrepancy*.

C.2. Sliced Kernelized Stein Discrepancy

In this section, we give a more detailed introduction to sliced kernelized Stein discrepancy (SKSD). Recall the definition of Stein discrepancy:

$$D_{SD}(q, p) = \sup_{\mathbf{f} \in \mathcal{F}_q} \mathbb{E}_q[s_p^T(\mathbf{x})\mathbf{f}(\mathbf{x}) + \nabla_{\mathbf{x}}^T \mathbf{f}(\mathbf{x})] \quad (33)$$

In the original paper of (Gong et al., 2021), they argue that the curse of dimensionality comes from two sources: (i) the high dimensionality of the score function $s_p : \mathcal{X} \subseteq \mathbb{R}^D \rightarrow \mathbb{R}^D$ and (ii) the test function input $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D$. Therefore, authors proposed two slice directions \mathbf{r}, \mathbf{g} to project s_p and \mathbf{x} respectively. However, this projection is equivalent to throwing away most of the information possessed by s_p and \mathbf{x} . To tackle this problem, authors proposed the first member of the SSD family by considering over all possible directions of \mathbf{r} and \mathbf{g} (a distribution over $\mathbf{r} \sim p_r, \mathbf{g} \sim p_g$), called *integrated sliced Stein discrepancy*:

$$S(q, p) = \mathbb{E}_{p_r, p_g} \left[\sup_{h_{r,g} \in \mathcal{F}_q} \mathbb{E}_q[s_p^T(\mathbf{x})h_{r,g}(\mathbf{x}^T \mathbf{g}) + \mathbf{r}^T \mathbf{g} \nabla_{\mathbf{x}^T \mathbf{g}} h_{r,g}(\mathbf{x}^T \mathbf{g})] \right] \quad (34)$$

where $h_{r,g}$ is the test function. Although it is theoretically valid (Theorem 1 in (Gong et al., 2021)), its practical usage is limited by the intractability of the integral over p_r, p_g and the optimal test function $h_{r,g}$. Surprisingly, authors show that the integral over \mathbf{r}, \mathbf{g} is not necessary for discrepancy validity. They achieved this in two steps.

The first step is to replace the expectation w.r.t. \mathbf{r} by a finite summation over orthogonal basis. The author showed that this is a valid discrepancy, called *orthogonal sliced Stein discrepancy* defined as

$$S_O(q, p) = \sum_{\mathbf{r} \in O_r} \mathbb{E}_{p_g} \left[\sup_{h_{r,g} \in \mathcal{F}_q} \mathbb{E}_q[s_p^T(\mathbf{x})h_{r,g}(\mathbf{x}^T \mathbf{g}) + \mathbf{r}^T \mathbf{g} \nabla_{\mathbf{x}^T \mathbf{g}} h_{r,g}(\mathbf{x}^T \mathbf{g})] \right] \quad (35)$$

where O_r is an orthogonal basis (e.g. one-hot vectors). The next step is to get rid of the expectation w.r.t. \mathbf{g} by a supremum operator. This is called *maxSSD-g*, which is defined as Eq.4 in the main text. For a quick recall, we include *maxSSD-g* in here:

$$S_{\max_{g_r}}(q, p) = \sum_{\mathbf{r} \in O_r} \sup_{\substack{h_{r,g_r} \in \mathcal{F}_q \\ \mathbf{g}_r \in \mathbb{S}^{D-1}}} \mathbb{E}_q[s_p^T(\mathbf{x})h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r) + \mathbf{r}^T \mathbf{g}_r \nabla_{\mathbf{x}^T \mathbf{g}_r} h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)] \quad (36)$$

Further, one can also use single optimal direction \mathbf{r} to replace the summation over the orthogonal basis O_r , resulting in *maxSSD-rg* ($S_{\max_{r,g_r}}$):

$$S_{\max_{r,g_r}}(q, p) = \sup_{h_{r,g} \in \mathcal{F}_q, \mathbf{g}_r, \mathbf{r} \in \mathbb{S}^{D-1}} \mathbb{E}_q [s_p^r(\mathbf{x})h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r) + \mathbf{r}^T \mathbf{g}_r \nabla_{\mathbf{x}^T \mathbf{g}_r} h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)] \quad (37)$$

Similar to KSD, authors addressed tractability issue of the optimal h_{r,g_r} by restricting the \mathcal{F}_q to be a one-dimensional RKHS induced by a c_0 -universal kernel $k_{r,g} : \mathcal{K} \times \mathcal{K} \rightarrow \mathbb{R}$ where $\mathcal{K} \subseteq \mathbb{R}$. Thus, for each member of the above SSD family, we have a corresponding kernelized version. They are called *integrated sliced kernelized Stein discrepancy*, *orthogonal SKSD*, and *max sliced kernelized Stein discrepancy* (including *maxSKSD-g* and *maxSKSD-rg*). In practice, *maxSKSD-g* or *maxSKSD-rg* is often preferred over the others due to its computational tractability, where their optimal slices for \mathbf{r} and \mathbf{g}_r are obtained by gradient-based optimization.

By reproducing properties of RKHS, one can define $\xi_{p,r,g_r}(\mathbf{x}, \cdot)$ as in Eq.5, and further define $\mu_{p,r,g_r} = \langle \xi_{p,r,g_r}(\mathbf{x}, \cdot), \xi_{p,r,g_r}(\mathbf{y}, \cdot) \rangle_{\mathcal{H}_{r,g_r}}$

$$\begin{aligned} \mu_{p,r,g_r}(\mathbf{x}, \mathbf{y}) &= s_p^r(\mathbf{x})k_{r,g_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{y}^T \mathbf{g}_r)s_p^r(\mathbf{y}) \\ &\quad + \mathbf{r}^T \mathbf{g}_r s_p^r(\mathbf{y}) \nabla_{\mathbf{x}^T \mathbf{g}_r} k_{r,g_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{y}^T \mathbf{g}_r) \\ &\quad + \mathbf{r}^T \mathbf{g}_r s_p^r(\mathbf{x}) \nabla_{\mathbf{y}^T \mathbf{g}_r} k_{r,g_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{y}^T \mathbf{g}_r) \\ &\quad + (\mathbf{r}^T \mathbf{g}_r)^2 \nabla_{\mathbf{x}^T \mathbf{g}_r, \mathbf{y}^T \mathbf{g}_r}^2 k_{r,g_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{y}^T \mathbf{g}_r). \end{aligned} \quad (38)$$

Then, by simple algebra, one can show that given \mathbf{r}, \mathbf{g}_r , the optimality w.r.t. test functions can be computed analytically:

$$\begin{aligned} D_{r,g_r}^2(q, p) &= \left(\sup_{\substack{h_{r,g_r} \in \mathcal{H}_{r,g_r} \\ \|h_{r,g_r}\|_{\mathcal{H}_{r,g_r}} \leq 1}} \mathbb{E}_q[s_p^r(\mathbf{x})h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r) + \mathbf{r}^T \mathbf{g}_r \nabla_{\mathbf{x}^T \mathbf{g}_r} h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)] \right)^2 \\ &= \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 = \mathbb{E}_q(\mathbf{x} | \mathbf{x}') [\mu_{p,r,g_r}(\mathbf{x}, \mathbf{x}')]. \end{aligned} \quad (39)$$

where \mathcal{H}_{r,g_r} is the RKHS induced by the kernel k_{r,g_r} . Therefore, the *maxSSD-g* and *maxSSD-rg* can be computed as

$$SK_{\max_{g_r}}(q, p) = \sum_{\mathbf{r} \in O_r} \sup_{\mathbf{g}_r \in \mathbb{S}^{D-1}} D_{r,g_r}^2(q, p) \quad (40)$$

and

$$SK_{\max_{r,g_r}}(q, p) = \sup_{\substack{\mathbf{g}_r \in \mathbb{S}^{D-1} \\ \mathbf{r} \in \mathbb{S}^{D-1}}} D_{r,g_r}^2(q, p) \quad (41)$$

D. Goodness-of-fit test

In this section, we give an introduction to the GOF test. To be general, we focus on the $SKSD$ -rg ($SK_{rg_r} = \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{rg_r}}^2$) as other related discrepancy can be easily derived from it. Assuming we have active slices \mathbf{r} and \mathbf{g}_r from algorithm 1. Thus, we can estimate SK_{rg_r} using the minimum variance U-statistics (Hoeffding, 1992; Serfling, 2009):

$$\widehat{SK}_{rg_r}(q, p) = \frac{1}{N(N-1)} \sum_{1 \leq i \neq j \leq N} \mu_{p,r,g_r}(\mathbf{x}_i, \mathbf{x}_j). \quad (42)$$

where $\mu_{\mathbf{x}, \mathbf{y}}$ is defined in Eq.38 which satisfies $\mathbb{E}_{q(\mathbf{x})q(\mathbf{x}')}[\mu_{p,r,g_r}(\mathbf{x}, \mathbf{x}')] = \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{rg_r}}^2$, and \mathbf{x}, \mathbf{x}' are i.i.d. samples from q . With the help of the U-statistics, we characterize its asymptotic distribution.

Theorem 6. *Assume the conditions in theorem 1 are satisfied, we have the following:*

1. If $q \neq p$, then \widehat{SK}_{rg_r} is asymptotically normal. Particularly,

$$\sqrt{N}(\widehat{SK}_{rg_r} - SK_{rg_r}) \xrightarrow{d} \mathcal{N}(0, \sigma_h^2) \quad (43)$$

where $\sigma_h^2 = \text{var}_{\mathbf{x} \sim q}(\mathbb{E}_{\mathbf{x}' \sim q}[\mu_{p,r,g_r}(\mathbf{x}, \mathbf{x}')])$ and $\sigma_h \neq 0$

2. If $q = p$, we have a degenerated U-statistics with $\sigma_h = 0$ and

$$N\widehat{SK}_{rg_r} \xrightarrow{d} \sum_{j=1}^{\infty} c_j(Z_j^2 - 1) \quad (44)$$

where $\{Z_j\}$ are i.i.d standard Gaussian variables, and $\{c_j\}$ are the eigenvalues of the kernel $\mu_{p,r,g_r}(\mathbf{x}, \mathbf{x}')$ under $q(\mathbf{x})$. In other words, they are the solutions of $c_j \phi_j(\mathbf{x}) = \int_{\mathbf{x}'} \mu_{p,r,g_r}(\mathbf{x}, \mathbf{x}') \phi_j(\mathbf{x}') q(\mathbf{x}') d\mathbf{x}'$.

Proof. As the \widehat{SK}_{rg_r} is the second order U-statistic of SK_{rg_r} , thus, we can directly use the results from section 5.5.1 and 5.5.2 in (Serfling, 2009). \square

The above theorem indicates a well-defined asymptotic distribution for SK_{rg_r} , which allows us to use the following bootstrap method to estimate the rejection threshold (Huskova & Janssen, 1993; Arcones & Gine, 1992; Liu et al., 2016). The bootstrap samples can be computed as

$$\widehat{SK}_m^* = \sum_{1 \leq i \neq j \leq N} (w_i^m - \frac{1}{N})(w_j^m - \frac{1}{N}) \mu_{p,r,g_r}(\mathbf{x}_i, \mathbf{x}_j) \quad (45)$$

where $(w_1^m, \dots, w_N^m)_{m=1}^M$ are random weights drawn from multinomial distributions $\text{Multi}(N, \frac{1}{N}, \dots, \frac{1}{N})$. Now, we give the detailed algorithm for GOF test.

Algorithm 2 GOF test with active slices

Input: Samples $\mathbf{x} \sim q$, density p , kernel k_{rg_r} , active slices \mathbf{r}, \mathbf{g}_r , significance level α , and bootstrap sample size M .

Hypothesis: $H_0: p = q$ v.s. $H_1: p \neq q$

Computing U-statistics \widehat{SK}_{rg_r} using Eq.42

Generate M bootstrap samples $\{\widehat{SK}_m^*\}_{m=1}^M$ using Eq.45

Reject null hypothesis H_0 if the proportion of $\widehat{SK}_m^* > \widehat{SK}_{rg_r}$ is less than α

E. Relaxing constraints for kernelized SSD family

E.1. Validity w.r.t \mathbf{r}, \mathbf{g}_r

The key to this proof is to prove the real analyticity of SK_{g_r} (or S_{rg_r}) to slices \mathbf{r} and \mathbf{g}_r . Therefore, let's first give a definition of multivariate real analytic function.

Definition E.1 (Real analytic function). *A function $f: \mathcal{U} \rightarrow \mathbb{R}$ is real analytic if for each $\mathbf{c} \in \mathcal{U}$, there is a power series as in the form*

$$f(\mathbf{x}) = \sum_{\kappa \in \mathbb{N}_0^n} \alpha_{\kappa} (\mathbf{x} - \mathbf{c})^{\kappa}$$

for some choice of $(\alpha_{\kappa})_{\kappa \in \mathbb{N}_0^n} \subset \mathbb{R}$ and all \mathbf{x} in a neighbourhood of \mathbf{c} , and this power series converges absolutely. Namely,

$$\sum_{\kappa \in \mathbb{N}_0^n} |\alpha_{\kappa}| |\mathbf{x} - \mathbf{c}|^{\kappa} < \infty$$

where $\mathbb{N}_0 = \{0, 1, \dots\}$ denotes non-negative integers, $\kappa = (\kappa_1, \dots, \kappa_n)$ are called multiindex, and we define $\mathbf{x}^{\kappa} = x_1^{\kappa_1} \dots x_n^{\kappa_n}$.

Now, we introduce a useful lemma showing that composition of real analytic function is also real analytic.

Lemma 1 (Composition of real analytic function). *Let $\mathcal{U} \subset \mathbb{R}^n$ and $\mathcal{V} \subset \mathbb{R}^m$ be open, and let $\mathbf{f}: \mathcal{U} \rightarrow \mathcal{V}$ and $\mathbf{g}: \mathcal{V} \rightarrow \mathbb{R}^p$ be real analytic. Then $\mathbf{g} \circ \mathbf{f}: \mathcal{U} \rightarrow \mathbb{R}^p$ is real analytic.*

Especially, the real analyticity is not only preserved by function composition, it is also closed under most of the simple operations: addition, multiplication, division (assuming denominator is non-zero), etc. Now we can prove the main proposition to show that the $SKSD$ -rg (SK_{rg_r}) is real analytic w.r.t both \mathbf{g}_r and \mathbf{r} . In the following, we assume the $\mathbf{r}, \mathbf{g}_r \in \mathbb{R}^D$.

Proposition 3 (*SKSD-g is real analytic*). *Assume assumption 1-4 (density regularity), 5-6 (kernel richness and real analyticity) are satisfied, further we let $\mathbf{g}_r \in \mathbb{R}^D$, then SKSD-g (SK_{g_r}) is real analytic w.r.t \mathbf{g}_r and $SK_{r\mathbf{g}_r}$ is real analytic to both $\mathbf{r} \in \mathbb{R}^D$ and \mathbf{g}_r .*

Proof. First, let's focus on the real analyticity w.r.t. \mathbf{g}_r . We re-write the SKSD-g as the following:

$$\begin{aligned} SK_{g_r} &= \sum_{\mathbf{r} \in O_r} \|\xi_{p,r,g_r}(\mathbf{x})\|_{\mathcal{H}_{r\mathbf{g}_r}}^2 \\ &= \sum_{\mathbf{r} \in O_r} \langle \mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x})) \underbrace{k_{r\mathbf{g}_r}(\mathbf{x}^T \mathbf{g}_r, \cdot)}_{f_r^*(\mathbf{x})}], \\ &\quad \mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x})) k_{r\mathbf{g}_r}(\mathbf{x}^T \mathbf{g}_r, \cdot)] \rangle_{\mathcal{H}_{r\mathbf{g}_r}} \\ &= \sum_{\mathbf{r} \in O_r} \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [f_r^*(\mathbf{x}) k_{r\mathbf{g}_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{x}'^T \mathbf{g}_r) f_r^*(\mathbf{x}')] \end{aligned}$$

The second equality is from the definition of RKHS norm $\|\cdot\|_{\mathcal{H}_{r\mathbf{g}_r}}$ and Stein identity. We can observe that \mathbf{g}_r appears inside the kernel $k_{r\mathbf{g}_r}$ in the form of $\mathbf{x}^T \mathbf{g}_r$. So in order to use the function composition lemma (lemma 1), we need to first show that for any given \mathbf{x} , $\mathbf{x}^T \mathbf{g}_r$ is real analytic. By definition of real analytic function, we need a center point $\mathbf{c} \in \mathbb{R}^D$, and \mathbf{g}_r in the neighborhood of \mathbf{c} (i.e. $|\mathbf{g}_r - \mathbf{c}| < R_c$). Then, we define the power series as

$$h_x(\mathbf{g}_r) = \sum_{\kappa_1=0}^{\infty} \dots \sum_{\kappa_D=0}^{\infty} \frac{(g_{r1} - c_1)^{\kappa_1} \dots (g_{rD} - c_D)^{\kappa_D}}{\kappa_1! \dots \kappa_D!} \alpha_{\{\kappa_i\}_i^D} \mathbf{g}_r.$$

with the following coefficient

$$\begin{cases} \alpha_{\{\kappa_i\}_i^D} = 0 & \text{if } \sum_i \kappa_i > 1 \\ \alpha_{\{\kappa_i\}_i^D} = x_i & \text{if } \kappa_i = 1, \sum_i \kappa_i = 1 \\ \alpha_{\{\kappa_i\}_i^D} = \mathbf{c}^T \mathbf{x} & \text{if } \sum_i \kappa_i = 0 \end{cases}$$

Then, by substitution, we have

$$h_x(\mathbf{g}) = \sum_{d=1}^D (g_d - c_d) x_d + \mathbf{c}^T \mathbf{x} \quad (46)$$

$$= \mathbf{x}^T \mathbf{g} \quad (47)$$

which converges with radius of convergence $R_c = \infty$. From assumption 6, we know the kernel $k_{r\mathbf{g}_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{x}'^T \mathbf{g}_r) = \phi((\mathbf{x} - \mathbf{x}')^T \mathbf{g}_r)$ is translation invariant and real analytic. Thus, from lemma 1, we know $k_{r\mathbf{g}_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{x}'^T \mathbf{g}_r)$ is real analytic to \mathbf{g}_r with radius of convergence R_k (R_k is determined by the form of the kernel function). This means we can use a power series to represents this kernel w.r.t. \mathbf{g}_r inside some neighborhood define around center point. Specifically, for a central point $\mathbf{c} \in \mathbb{R}^D$ and any \mathbf{g}_r satisfying $|\mathbf{g}_r - \mathbf{c}| < R_k$, we have

$$k_{r\mathbf{g}_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{x}'^T \mathbf{g}_r) = \sum_{\kappa \in \mathbb{N}_0^D} \alpha_{\kappa}(\mathbf{x}, \mathbf{x}') (\mathbf{g}_r - \mathbf{c})^{\kappa}$$

where this series converges absolutely. We substitute it into SK_{g_r}

$$\begin{aligned} SK_{g_r} &= \sum_{\mathbf{r} \in O_r} \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [f_r^*(\mathbf{x}) k_{r\mathbf{g}_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{x}'^T \mathbf{g}_r) f_r^*(\mathbf{x}')] \\ &= \sum_{\mathbf{r} \in O_r} \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [f_r^*(\mathbf{x}) \sum_{\kappa \in \mathbb{N}_0^D} \alpha_{\kappa}(\mathbf{x}, \mathbf{x}') (\mathbf{g}_r - \mathbf{c})^{\kappa} f_r^*(\mathbf{x}')] \\ &= \sum_{\mathbf{r} \in O_r} \sum_{\kappa \in \mathbb{N}_0^D} \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [\alpha_{\kappa}(\mathbf{x}, \mathbf{x}') f_r^*(\mathbf{x}) f_r^*(\mathbf{x}')] (\mathbf{g}_r - \mathbf{c})^{\kappa} \end{aligned}$$

which also converges absolutely with radius of convergence R_k . The third equality is from the Fubini's theorem. The conditions of Fubini's theorem can be verified by fact that f_r^* is square integrable (assumption 2), and the power series of $k_{r\mathbf{g}_r}$ converges absolutely. Thus, by definition of real analytic function, SKSD-g is real analytic w.r.t each \mathbf{g}_r . This also implies SKSD-rg ($SK_{r\mathbf{g}_r}$) is real analytic w.r.t. \mathbf{g}_r (because $SK_{r\mathbf{g}_r}$ is just SK_{g_r} without summation over O_r).

For the real analyticity w.r.t \mathbf{r} , the proof is almost the same. The inner product $s_p^r(\mathbf{x}) - s_q^r(\mathbf{x})$ is real analytic w.r.t \mathbf{r} obviously for given \mathbf{x} . We also use the fact that real analyticity is preserved under multiplication of two real analytic functions. In addition, note that $k_{r\mathbf{g}_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{x}'^T \mathbf{g}_r)$ act as a constant w.r.t. \mathbf{r} , we can directly apply the Fubini's theorem again to form a power series w.r.t. \mathbf{r} with absolute convergence. Thus, $SK_{r\mathbf{g}_r}$ is real analytic w.r.t. \mathbf{r} for any \mathbf{g}_r . Thus, $SK_{r\mathbf{g}_r}$ is real analytic to both \mathbf{r} and \mathbf{g}_r . \square

Next, we introduce an important property of real analytic function:

Lemma 2 (Zero Set Theorem (Mityagin, 2015)). *Let $f(x)$ be a real analytic function on (a connected open domain \mathcal{U} of) \mathbb{R}^d . If f is not identically 0, then its zero set*

$$S(f) := \{\mathbf{x} \in \mathcal{U} | f(\mathbf{x}) = 0\}$$

has a measure 0, i.e. $mes_d S(f) = 0$

With the help from the zero-set theorem, we can prove the validity of SK_{g_r} (or $SK_{r\mathbf{g}_r}$) with finite random slices \mathbf{g}_r (and \mathbf{r}).

Proof of theorem 1

Proof. We first deal with the validity of \mathbf{g}_r with fixed orthogonal basis O_r . It is trivial that when $p = q$, $SK_{g_r} = 0$ identically. Now, assume $p \neq q$, then, from the theorem 3 in (Gong et al., 2021), the orthogonal SKSD (Eq.48) is a valid discrepancy. Namely, we have

$$\sum_{\mathbf{r} \in O_r} \int q_{g_r}(\mathbf{g}_r) \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r\mathbf{g}_r}}^2 > 0 \quad (48)$$

We should note that the distribution q_{g_r} is originally defined on \mathbb{S}^{D-1} . But, we can easily generalize it to larger spaces.

As for $\mathbf{g}_r \in \mathbb{R}^D$, we can always write $\mathbf{g}_r = c\mathbf{g}'_r$, where $\mathbf{g}'_r \in \mathbb{S}^{D-1}$, and $c \geq 0$. As the domain for \mathbf{g}_r is \mathbb{R}^D , the \mathbf{g}_r can represent all possible directions. Thus, we can follow the same proof logic as theorem 3 in (Gong et al., 2021) to show the corresponding discrepancy is greater than 0 when $p \neq q$.

Therefore, Eq.48 represents there exists a $\mathbf{r} \in O_r$ such that $\|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 > 0$ for a set of \mathbf{g}_r with non-zero measure. Namely, $\|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2$ is not 0 identically. Thus, from the proposition 3 and lemma 2, the set of \mathbf{g}_r that make $\|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 = 0$ has a 0 measure. Then, if \mathbf{g}_r is sampled from some distribution η_g with density supported on \mathbb{R}^D (e.g. Gaussian distribution), we have

$$SK_{g_r} = \sum_{\mathbf{r} \in O_r} \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 > 0$$

almost surely.

Now, we show that SK_{r,g_r} is also a valid discrepancy with $\mathbf{r} \sim \eta_r$. First, due to the validity of *integrated SKSD*, we have

$$\int q_r(\mathbf{r}) \int q_{g_r}(\mathbf{g}_r) \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 d\mathbf{g}_r d\mathbf{r} > 0 \quad (49)$$

Due to the real analyticity of SK_{r,g_r} ($\|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2$) w.r.t \mathbf{r} , we can easily show that

$$\int q_{g_r}(\mathbf{g}_r) \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 d\mathbf{g}_r$$

is real analytic to \mathbf{r} and it is not 0 identically. Thus, by lemma 2, for $\mathbf{r} \sim \eta_r$, we have

$$\int q_{g_r}(\mathbf{g}_r) \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 d\mathbf{g}_r > 0$$

Namely, $\|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 > 0$ for a set of \mathbf{g}_r with non-zero measure. In the beginning of the proof, we show that this set of \mathbf{g}_r is almost everywhere in \mathbb{R}^D due to its real analyticity. Namely, $\|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 > 0$ for $\mathbf{r} \sim \eta_r$ and $\mathbf{g}_r \sim \eta_g$ if $p \neq q$. Thus, we can conclude that for $SK_{r,g_r} = 0$ if and only if $p = q$ almost surely for $\mathbf{r} \sim \eta_r$ and $\mathbf{g}_r \sim \eta_g$. \square

Corollary 6.1 (Normalizing \mathbf{g}_r). *Assume the conditions in theorem 1 are satisfied, then the following operations do not violate the validity of SKSD-rg SK_{r,g_r} . (1) For $\mathbf{g}'_r, \mathbf{r}' \in \mathbb{S}^{D-1}$, we define $\mathbf{g}_r = \mathbf{g}'_r + \gamma_g$ and $\mathbf{r} = \mathbf{r}' + \gamma_r$, where γ_r, γ_g are the noise from Gaussian distribution. (2) Define $\tilde{\mathbf{g}}_r = c_g \times \mathbf{g}_r$ and $\tilde{\mathbf{r}} = c_r \times \mathbf{r}$, where $\tilde{\mathbf{g}}_r, \tilde{\mathbf{r}}$ are unit vectors and $c_r, c_g > 0$. The resulting active slices $\tilde{\mathbf{r}}$ and $\tilde{\mathbf{g}}_r$ do not violate the validity of SK_{r,g_r} .*

Proof. From the theorem 1 with \mathbf{g}_r, \mathbf{r} , when $p \neq q$, we

have

$$\begin{aligned} SK_{r,g_r} &= \|\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})]\|_{\mathcal{H}_{r,g_r}}^2 \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [f_r^*(\mathbf{x}) k_{r,g_r}(\mathbf{x}^T \mathbf{g}_r, \mathbf{x}'^T \mathbf{g}_r) f_r^*(\mathbf{x}')] \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{x}'} [c_r^2 f_{\tilde{r}}^*(\mathbf{x}) k_{r,g_r}(c\mathbf{x}^T \tilde{\mathbf{g}}_r, c\mathbf{x}'^T \tilde{\mathbf{g}}_r) f_{\tilde{r}}^*(\mathbf{x}')] > 0 \end{aligned}$$

From the assumption 6 that $k_{r,g_r}(c\mathbf{x}^T \tilde{\mathbf{g}}_r, c\mathbf{x}'^T \tilde{\mathbf{g}}_r) = k'_{r,g_r}(\mathbf{x}^T \tilde{\mathbf{g}}_r, \mathbf{x}'^T \tilde{\mathbf{g}}_r)$. So this is equivalent to the *SKSD-rg* defined with a new c_0 -universal kernel k'_{r,g_r} and $\tilde{\mathbf{g}}_r, \tilde{\mathbf{r}} \in \mathbb{S}^{D-1}$. Thus, the corresponding *maxSKSD-rg* with $\tilde{\mathbf{g}}_r, \tilde{\mathbf{r}} \in \mathbb{S}^{D-1}$ is a valid discrepancy almost surely. \square

E.2. Relationship between SSD and SKSD

Proof of proposition 1

Proof. We consider the *SSD-rg* (S_{r,g_r}) without the optimal test function:

$$\mathbb{E}_q[s_p^T(\mathbf{x}) h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r) + \mathbf{r}^T \mathbf{g}_r \nabla_{\mathbf{x}^T \mathbf{g}_r} h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)] \quad (50)$$

From the Stein identity (Eq.23), we can let $\mathbf{f}(\mathbf{x}) = [r_1 h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r), r_2 h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r), \dots, r_D h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)]^T$ and then take the trace. Thus, we have

$$\mathbb{E}_q[s_p^T(\mathbf{x}) h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)] = \mathbb{E}_q[\mathbf{r}^T \mathbf{g}_r \nabla_{\mathbf{x}^T \mathbf{g}_r} h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)]$$

Substitute it into Eq.50 and change the variable to $\mathbf{y} = \mathbf{G}_r \mathbf{x}$, we have

$$\begin{aligned} &\mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x})) h_{r,g_r}(\mathbf{x}^T \mathbf{g}_r)] \\ &= \int q_{G_r}(y_d, \mathbf{y}_{-d}) \underbrace{(s_p^r(\mathbf{G}_r^{-1} \mathbf{y}) - s_q^r(\mathbf{G}_r^{-1} \mathbf{y}))}_{f_r^*(\mathbf{G}_r^{-1} \mathbf{y})} h_{r,g_r}(y_d) d\mathbf{y} \\ &= \int q_{G_r}(y_d) \underbrace{\int q_{G_r}(\mathbf{y}_{-d} | y_d) f_r^*(\mathbf{G}_r^{-1} \mathbf{y}) d\mathbf{y}_{-d}}_{h_{r,g_r}^*(y_d)} h_{r,g_r}(y_d) dy_d \\ &\leq \sqrt{\mathbb{E}_{q_{G_r}(y_d)} [h_{r,g_r}^*(y_d)^2]} \sqrt{\mathbb{E}_{q_{G_r}(y_d)} [h_{r,g_r}(y_d)^2]} \end{aligned}$$

where the last inequality is from Cauchy-Schwarz inequality, where the equality holds when

$$\begin{aligned} h_{r,g_r}(y_d) &\propto h_{r,g_r}^*(y_d) \\ &= \mathbb{E}_{q_{G_r}(\mathbf{y}_{-d} | y_d)} [(s_p^r(\mathbf{G}_r^{-1} \mathbf{y}) - s_q^r(\mathbf{G}_r^{-1} \mathbf{y}))] \end{aligned}$$

where $y_d = \mathbf{x}^T \mathbf{g}_r$. \square

Proof of theorem 2

Proof. Let's first re-write of $S_{r_{g_r}}$ and $SK_{r_{g_r}}$.

$$\begin{aligned} S_{r_{g_r}} &= \sup_{h_{r_{g_r}} \in \mathcal{F}_q} \mathbb{E}_q[(s_p^r(\mathbf{y}) - s_q^r(\mathbf{x}))h_{r_{g_r}}(\mathbf{x}^T \mathbf{g}_r)] \\ &= \mathbb{E}_{q_{G_r}(y_d)} \left[\underbrace{\int q_{G_r}(\mathbf{y}-d|y_d) (s_p^r(\mathbf{G}_r^{-1} \mathbf{y}) - s_q^r(\mathbf{G}_r^{-1} \mathbf{y})) d\mathbf{y}_{-d}}_{h_{r_{g_r}}^*(y_d)} \right] \\ &\quad \times h_{r_{g_r}}^*(y_d) \\ &= \mathbb{E}_{q_{G_r}(y_d)} [h_{r_{g_r}}^*(y_d)^2] \end{aligned}$$

where the second equality is from proposition 1.

$$SK_{r_{g_r}} = \langle \mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})], \mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x})] \rangle_{\mathcal{H}_k}$$

where $\xi_{p,r,g_r}(\mathbf{x}, \cdot)$ is defined in Eq.5, and $\langle \cdot, \cdot \rangle_{\mathcal{H}_{r_{g_r}}}$ is the RKHS inner product induced by kernel $k_{r_{g_r}}$. By simple algebraic manipulation and Stein identity (Eq.23), we have

$$\begin{aligned} &\mathbb{E}_q[\xi_{p,r,g_r}(\mathbf{x}, \cdot)] \\ &= \mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x}))k_{r_{g_r}}(\mathbf{x}^T \mathbf{g}_r, \cdot)] \\ &= \mathbb{E}_{q_{G_r}(y_d)} \left[\underbrace{\int q_{G_r}(\mathbf{y}-d|y_d) (s_p^r(\mathbf{G}_r^{-1} \mathbf{y}) - s_q^r(\mathbf{G}_r^{-1} \mathbf{y})) d\mathbf{y}_{-d}}_{h_{r_{g_r}}^*(y_d)} \right] \\ &\quad \times k_{r_{g_r}}(y_d, \cdot) \\ &= \mathbb{E}_{q_{G_r}(y_d)} [h_{r_{g_r}}^*(y_d)k_{r_{g_r}}(y_d, \cdot)] \end{aligned}$$

Thus, we have

$$\begin{aligned} &SK_{r_{g_r}} \\ &= \mathbb{E}_{y_d, y'_d \sim q_{G_r}(y_d)} [h_{r_{g_r}}^*(y_d)k_{r_{g_r}}(y_d, y'_d)h_{r_{g_r}}^*(y'_d)] \\ &\leq \sqrt{\mathbb{E}_{y_d, y'_d} [k_{r_{g_r}}(y_d, y'_d)^2]} \sqrt{\mathbb{E}_{y_d} [h_{r_{g_r}}^*(y_d)^2]} \sqrt{\mathbb{E}_{y'_d} [h_{r_{g_r}}^*(y'_d)^2]} \\ &= M S_{r_{g_r}}^* \end{aligned}$$

where constant M is from the bounded kernel assumption, and the inequality is from Cauchy-Schwarz inequality. Without the loss of generality, we can set $M = 1$. For other value of $M > 0$, one can always set the optimal test function ($h_{r_{g_r}}^*$) for SSD -rg with coefficient M . The the new SSD -g will be M multiplied by the original SSD -rg with $M = 1$.

Thus, SSD -rg is an upper bound for $SKSD$ -rg. From the assumption 1, we know that the induced set $\mathcal{K} = \{y \in \mathbb{R} | y = \mathbf{x}^T \mathbf{g}, ||g|| = 1, \mathbf{x} \in \mathcal{X}\}$ is LCH, and the kernel $k_{r_{g_r}} : \mathcal{K} \times \mathcal{K} \rightarrow \mathbb{R}$ is c_0 -universal. Then, from (Sriperumbudur et al., 2011), c_0 -universal implies L_p -universal. Namely, the induced RKHS $\mathcal{H}_{r_{g_r}}$ is dense in $L^p(\mathcal{K}; \mu)$ with all Borel probability measure μ w.r.t. p -norm, defined as

$$||f||_p = \left(\int |f(\mathbf{x})|^p d\mu(\mathbf{x}) \right)^{\frac{1}{p}}$$

Now, from the assumption 4, we know $h_{r_{g_r}}^*(y_d)$ is bounded for all possible \mathbf{g}_r , we have

$$\int q_{G_r}(y_d) |h_{r_{g_r}}^*(y_d)|^2 dy_d < \infty$$

This means $h_{r_{g_r}}^* \in L^2(\mathcal{K}, \mu_{G_r})$, where μ_{G_r} is the probability measure with density $q_{G_r}(y_d)$

From the L_p -universality, there exists a function $\widetilde{h_{r_{g_r}}^*} \in \mathcal{H}_{r_{g_r}}$, such that for any given $\epsilon > 0$,

$$||h_{r_{g_r}}^* - \widetilde{h_{r_{g_r}}^*}||_2 < \epsilon$$

Let's define $\widetilde{SK}_{r_{g_r}}$ is the $SKSD$ -rg with the specific kernelized test function $\widetilde{h_{r_{g_r}}^*}$, and from the optimality of $SKSD$ -rg, we have

$$SK_{r_{g_r}} \geq \widetilde{SK}_{r_{g_r}}$$

Therefore, we have

$$\begin{aligned} 0 &\leq S_{r_{g_r}} - SK_{r_{g_r}} \\ &\leq S_{r_{g_r}} - \widetilde{SK}_{r_{g_r}} \\ &= \mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x})) (h_{r_{g_r}}^*(\mathbf{x}^T \mathbf{g}_r) - \widetilde{h_{r_{g_r}}^*}(\mathbf{x}^T \mathbf{g}_r))] \\ &\leq \underbrace{\sqrt{\mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x}))^2]}}_{C_r} \\ &\quad \times \sqrt{\mathbb{E}_q[(h_{r_{g_r}}^*(\mathbf{x}^T \mathbf{g}_r) - \widetilde{h_{r_{g_r}}^*}(\mathbf{x}^T \mathbf{g}_r))^2]} \\ &= C_r \sqrt{\int q_{G_r}(y_d, \mathbf{y}_{-d}) (h_{r_{g_r}}^*(y_d) - \widetilde{h_{r_{g_r}}^*}(y_d))^2 d\mathbf{y}} \\ &= C_r ||h_{r_{g_r}}^* - \widetilde{h_{r_{g_r}}^*}||_2 < C_r \epsilon \end{aligned}$$

From assumption 2, we know $s_p^r(\mathbf{x}) - s_q^r(\mathbf{x})$ is square integrable for all possible \mathbf{r} . Therefore, let's define $C = \max_{\mathbf{r} \in \mathbb{S}^{D-1}} C_r$, then,

$$0 \leq S_{r_{g_r}} - SK_{r_{g_r}} < C \epsilon$$

□

F. Theory related to active slice g

F.1. Optimal test function for PSD

Proposition 4 (Optimality of PSD). *Assume the assumption 1 – 3 (density regularity) are satisfied, then the optimal test function for PSD given O_r is proportional to the projected score difference, i.e.*

$$f_r^*(\mathbf{x}) \propto (s_p^r(\mathbf{x}) - s_q^r(\mathbf{x})) \quad (51)$$

Thus,

$$PSD(q, p; O_r) = \sum_{\mathbf{r} \in O_r} \mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x}))^2] \quad (52)$$

if the coefficient of f_r^* to be 1.

Proof. From the Stein identity (Eq.23), we can re-write the inner part of the supremum of Eq.9 as

$$\begin{aligned} &\mathbb{E}_q[s_p^r(\mathbf{x})f_r(\mathbf{x}) + \mathbf{r}^T \nabla_{\mathbf{x}} f_r(\mathbf{x})] \\ &= \mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x}))f_r(\mathbf{x})] \end{aligned}$$

Then, we can upper bound the PSD (Eq.9) as the following

$$\begin{aligned} & \sum_{r \in O_r} \mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x}))f_r(\mathbf{x})] \\ & \leq \sum_{r \in O_r} \sqrt{\mathbb{E}_q[(s_p^r(\mathbf{x}) - s_q^r(\mathbf{x}))^2]} \sqrt{\mathbb{E}_q[(f_r(\mathbf{x}))^2]} \end{aligned}$$

by Cauchy-Schwarz inequality. It is well-known that the equality holds when $f_r(\mathbf{x}) \propto (s_p^r(\mathbf{x}) - s_q^r(\mathbf{x}))$ \square

F.2. Proof of Theorem 3

Proof. The key to this proof is to notice that $h_{r g_r}^*$ is the conditional mean of f_r^* w.r.t. the transformed distribution q_{G_r} . By using the similar terminology of proposition 1, and let $s_p^r = s_p^r(\mathbf{x})$ for abbreviation. Then,

$$\begin{aligned} & \mathbb{E}_q[(s_p^r - s_q^r)f_r^*(\mathbf{x})] - \mathbb{E}_q[(s_p^r - s_q^r)h_{r g_r}^*(\mathbf{x}^T \mathbf{g}_r)] \\ & = \int q(\mathbf{x})[(s_p^r - s_q^r)^2 - (s_p^r - s_q^r)h_{r g_r}^*(\mathbf{x}^T \mathbf{g}_r)]d\mathbf{x} \\ & = \int q_{G_r}(y_d) \left[\int q_{G_r}(\mathbf{y}_{-d}|y_d)(s_p^r(\mathbf{G}_r^{-1}\mathbf{y}) - s_q^r(\mathbf{G}_r^{-1}\mathbf{y}))^2 d\mathbf{y}_{-d} \right. \\ & \quad \left. - \int q_{G_r}(s_p^r(\mathbf{G}_r^{-1}\mathbf{y}) - s_q^r(\mathbf{G}_r^{-1}\mathbf{y}))d\mathbf{y}_{-d} h_{r g_r}^*(y_d) \right] dy_d \\ & = \int q_{G_r}(y_d) \left[\int q_{G_r}(\mathbf{y}_{-d}|y_d)(f_r^*(\mathbf{G}_r^{-1}\mathbf{y}) - h_{r g_r}^*(y_d))^2 d\mathbf{y} \right. \\ & \quad \left. - \mathbb{E}_q[(f_r^*(\mathbf{x}) - h_{r g_r}^*(\mathbf{x}^T \mathbf{g}_r))^2] \right] \geq 0 \end{aligned}$$

where the 3rd equality is due to the fact that $h_{r g_r}^*$ is the conditional mean of f_r^* . Thus,

$$\begin{aligned} & PSD - S_{g_r} \\ & = \sum_{r \in O_r} \mathbb{E}_q[(s_p^r - s_q^r)f_r^*(\mathbf{x})] - \mathbb{E}_q[(s_p^r - s_q^r)h_{r g_r}^*(\mathbf{x}^T \mathbf{g}_r)] \\ & = \sum_{r \in O_r} \mathbb{E}_q[(f_r^*(\mathbf{x}) - h_{r g_r}^*(\mathbf{x}^T \mathbf{g}_r))] \geq 0 \end{aligned}$$

\square

F.3. Proof of Theorem 4

Before we give the details, we introduce the main inequality and its variant for the proof.

Lemma 3 (Poincaré Inequality). *For a probabilistic distribution p that satisfies assumption 7, for all locally Lipschitz function $f(\mathbf{x}) : \mathcal{X} \subseteq \mathbb{R}^D \rightarrow \mathbb{R}$, we have the following inequality*

$$\text{Var}_p(f(\mathbf{x})) \leq C_p \int p(\mathbf{x}) \|\nabla_{\mathbf{x}} f(\mathbf{x})\|^2 d\mathbf{x}$$

where C_p is called Poincaré constant that is only related to p .

One should note that the assumption of log concavity of p is a sufficient condition for Poincaré inequality, which means it may be applied to a broader class of distributions. But it is beyond the scope of this work.

Due to the form of optimal test functions of *SSD-g*, we need to deal with the transformed distribution q_{G_r} and its conditional expectations (see Eq.8). Unfortunately, the original form of Poincaré inequality cannot be applied. In the following, we introduce its variant called *subspace* Poincaré inequality (Constantine et al., 2014; Zahm et al., 2020; Parante et al., 2020) to deal with the conditional expectation. But before that, we need to make sure the transformed distribution and its conditional density still satisfy the conditions of Poincaré inequality, i.e. log concavity.

Lemma 4 (Preservation of log concavity). *Assume distribution $q(\mathbf{x}) = \exp(-V(\mathbf{x}))$ is log-concave. With arbitrary orthogonal matrix \mathbf{G} and corresponding transformed distribution q_G , the conditional distribution $q_G(\mathbf{y}_{-d} | y_d)$ is also log-concave for all $d = 1, \dots, D$.*

Proof. Assume we have $\mathbf{y} = \mathbf{G}\mathbf{x}$. Thus, by change of variable formula, $q_G(\mathbf{y}) = q(\mathbf{x}) = q(\mathbf{G}^{-1}\mathbf{y}) = \exp(-V(\mathbf{G}^{-1}\mathbf{y}))$. Thus, the log conditional distribution

$$\log q_G(\mathbf{y}_{-d} | y_d) = -V(\mathbf{G}^{-1}\mathbf{y}) - \log q_G(y_d)$$

We inspect its Hessian w.r.t \mathbf{y}_{-d}

$$\begin{aligned} & \nabla_{\mathbf{y}_{-d}}^2 (V(\mathbf{G}^{-1}\mathbf{y}) + \log q_G(y_d)) \\ & = \nabla_{\mathbf{y}_{-d}}^2 (V(\mathbf{G}^{-1}\mathbf{y})) \\ & = \nabla_{\mathbf{y}_{-d}} (\mathbf{G}_{\setminus d} V'(\mathbf{G}^{-1}\mathbf{y})) \\ & = \mathbf{G}_{\setminus d} V''(\mathbf{G}^{-1}\mathbf{y}) \mathbf{G}_{\setminus d}^T \end{aligned}$$

where $\mathbf{G}_{\setminus d} = [\mathbf{g}_1, \dots, \mathbf{g}_{d-1}, \mathbf{g}_{d+1}, \dots, \mathbf{g}_D]^T$ and $V'(\mathbf{G}^{-1}\mathbf{y}) = \nabla_{\mathbf{G}^{-1}\mathbf{y}} V(\mathbf{G}^{-1}\mathbf{y})$. We already know that $V(\cdot)$ is a convex function. Thus, for all $\mathbf{u} \in \mathbb{R}^D$, $\mathbf{u}^T V''(\mathbf{x}) \mathbf{u} \geq 0$, therefore,

$$\mathbf{u}^T \mathbf{G}_{\setminus d} V''(\mathbf{G}^{-1}\mathbf{y}) \mathbf{G}_{\setminus d}^T \mathbf{u} = \mathbf{l}^T V''(\mathbf{G}^{-1}\mathbf{y}) \mathbf{l} \geq 0$$

where $\mathbf{l} = \mathbf{G}_{\setminus d}^T \mathbf{u}$. \square

Now, we can introduce the subspace Poincaré inequality

Lemma 5 (Poincaré inequality for conditional expectation). *Assume the assumption 2,4 (density regularity), 7 (Poincaré inequality condition) are satisfied, with arbitrary orthogonal matrix \mathbf{G} , $\mathbf{y} = \mathbf{G}\mathbf{x}$ and $y_d = \mathbf{x}^T \mathbf{g}_d$, we have the following inequality*

$$\begin{aligned} & \int q_G(\mathbf{y}_{-d} | y_d) [f_r^*(\mathbf{G}^{-1}\mathbf{y}) - h_{r g_r}^*(y_d)]^2 d\mathbf{y}_{-d} \\ & \leq C_{y_d} \mathbb{E}_{q_G(\mathbf{y}_{-d}|y_d)} \left[\|\mathbf{G}_{\setminus d} \nabla f_r^*\|^2 \right] \end{aligned}$$

where C_{y_d} is the Poincaré constant, $\mathbf{G}_{\setminus d} = [\mathbf{a}_1, \dots, \mathbf{a}_{d-1}, \mathbf{a}_{d+1}, \dots, \mathbf{a}_D]^T$ is the orthogonal matrix \mathbf{G} excluding $\mathbf{a}_d = \mathbf{g}$ and f_r^*, h_{r, g_r}^* are the optimal test functions defined in proposition 4, 1 respectively with coefficient 1.

Proof. From lemma 4, we know $q_G(\mathbf{y}_{-d} | y_d)$ is a log-concave distribution. Therefore, it satisfies the Poincaré inequality (lemma.3). We have

$$\begin{aligned} & \int q_G(\mathbf{y}_{-d} | y_d) [f_r^*(\mathbf{G}^{-1}\mathbf{y}) - h_{r, g_r}^*(y_d)]^2 d\mathbf{y}_{-d} \\ &= \text{Var}_{q_G(\mathbf{y}_{-d}|y_d)}(f_r^*(\mathbf{G}^{-1}\mathbf{y})) \\ &\leq C_{y_d} \int q_G(\mathbf{y}_{-d} | y_d) \|\nabla_{\mathbf{y}_{-d}} f_r^*(\mathbf{G}^{-1}\mathbf{y})\|^2 d\mathbf{y}_{-d} \\ &= C_{y_d} \int q_G(\mathbf{y}_{-d} | y_d) \|\mathbf{G}_{\setminus d} \nabla_{\mathbf{G}^{-1}\mathbf{y}} f_r^*(\mathbf{G}^{-1}\mathbf{y})\|^2 d\mathbf{y}_{-d} \end{aligned}$$

The first equality comes from the fact that $h_{r, g_r}^*(y_d)$ is actually a conditional mean of $f_r^*(\mathbf{G}^{-1}\mathbf{y})$, and the inequality comes from the direct application of Poincaré inequality on $q_G(\mathbf{y}_{-d}|y_d)$ and $f_r^*(\mathbf{G}^{-1}\mathbf{y})$. \square

With the above tools, it is now easy to prove theorem 4.

Theorem 4

Proof. We can re-write the inner part of controlled approximation (Eq.11) in the following:

$$\begin{aligned} & \int q_{G_r}(y_d, \mathbf{y}_{-d}) [f_r^*(\mathbf{G}_r^{-1}\mathbf{y}) - h_{r, g_r}^*(y_d)]^2 d\mathbf{y} \\ &= \int q_{G_r}(y_d) \mathbb{E}_{q_{G_r}(\mathbf{y}_{-d}|y_d)} [(f_r^*(\mathbf{G}_r^{-1}\mathbf{y}) - h_{r, g_r}^*(y_d))^2] d\mathbf{y} \\ &\leq \int q_{G_r}(y_d, \mathbf{y}_{-d}) C_{y_d} \|\mathbf{G}_{r \setminus d} \nabla f_r^*\|^2 d\mathbf{y} \\ &\leq C_{sup} \int q_{G_r}(y_d, \mathbf{y}_{-d}) \|\mathbf{G}_{r \setminus d} \nabla f_r^*\|^2 d\mathbf{y} \\ &= C_{sup} \int q(\mathbf{x}) \text{tr} \left[(\mathbf{G}_{r \setminus d} \nabla f_r^*) (\mathbf{G}_{r \setminus d} \nabla f_r^*)^T \right] d\mathbf{x} \\ &= C_{sup} \text{tr} \left[\mathbf{G}_{r \setminus d} \mathbf{H}_r \mathbf{G}_{r \setminus d}^T \right] \end{aligned}$$

where the first inequality is directly from lemma 5 and the second inequality is from the definition of C_{sup} .

To minimize this upper bound, we can directly use the theorem 2.1 (Sameh & Tong, 2000) by setting $B = I$ and $\mathbf{X} = \mathbf{G}_{r \setminus d}^T$. Therefore, we only need to check if $\mathbf{G}_{r \setminus d} \mathbf{G}_{r \setminus d}^T = I$. This is trivial as \mathbf{G}_r is an orthogonal matrix. Thus, the proof is complete. \square

G. Theory related to active slice r

G.1. Proof of proposition 2

First, from the theorem 3, we have

$$\text{PSD}_r \geq S_{r, g}$$

Thus, we can establish the following lower bound

$$\begin{aligned} & S_{r_1, g_{r_1}} - S_{r_2, g_{r_2}} \\ &\geq S_{r_1, g_{r_1}} - \text{PSD}_{r_2} \\ &= \underbrace{S_{r_1, g_{r_1}} - \text{PSD}_{r_1}}_{\text{controlled approximation}} + \text{PSD}_{r_1} - \text{PSD}_{r_2} \end{aligned}$$

Thus, from theorem 4, we can obtain

$$\begin{aligned} & S_{r_1, g_{r_1}} - \text{PSD}_{r_1} \\ &= -\mathbb{E}_q \left[(f_{r_1}^*(\mathbf{x}) - h_{r_1, g_{r_1}}^*(\mathbf{x}^T \mathbf{g}_{r_1}))^2 \right] \\ &\geq -C_{\text{sup}} \text{tr}(\mathbf{G}_{r_1 \setminus d} \mathbf{H}_{r_1} \mathbf{G}_{r_1 \setminus d}^T) \\ &= -C_{\text{sup}} \text{tr}(\mathbf{H}) + \underbrace{\mathbf{g}_{r_1}^T \mathbf{H}_{r_1} \mathbf{g}_{r_1}}_{\geq 0} \\ &\geq -C_{\text{sup}} \text{tr}(\mathbf{H}_{r_1}) \end{aligned}$$

where the first inequality is from the upper bound of controlled approximation (theorem 4) and $\mathbf{g}_{r_1}^T \mathbf{H}_{r_1} \mathbf{g}_{r_1} \geq 0$ is due to the positive semi-definiteness of \mathbf{H}_{r_1} . Assume we have an orthogonal basis O_{r_1} that contains \mathbf{r}_1 , thus, for each $\mathbf{r} \in O_{r_1}$, we have $\text{tr}(\mathbf{H}_r) \geq 0$. Then, we can show

$$\begin{aligned} \text{tr}(\mathbf{H}_{r_1}) &\leq \sum_{\mathbf{r} \in O_{r_1}} \text{tr}(\mathbf{H}_r) \\ &= \sum_{\mathbf{r} \in O_{r_1}} \text{tr}(\mathbb{E}_q[\nabla_{\mathbf{x}} \mathbf{f}^*(\mathbf{x}) \mathbf{r} \mathbf{r}^T \nabla_{\mathbf{x}} \mathbf{f}^*(\mathbf{x})^T]) \\ &= \text{tr}(\mathbb{E}_q[\nabla_{\mathbf{x}} \mathbf{f}^*(\mathbf{x}) \sum_{\mathbf{r} \in O_{r_1}} \mathbf{r} \mathbf{r}^T \nabla_{\mathbf{x}} \mathbf{f}^*(\mathbf{x})^T]) \\ &= \text{tr}(\mathbb{E}_q[\nabla_{\mathbf{x}} \mathbf{f}^*(\mathbf{x}) \nabla_{\mathbf{x}} \mathbf{f}^*(\mathbf{x})^T]) \\ &= \sum_{i=1}^D \omega_i = \Omega \end{aligned}$$

where $\{\omega_i\}_i^D$ are the eigenvalues of $\mathbb{E}_q[\nabla_{\mathbf{x}} \mathbf{f}^*(\mathbf{x}) \nabla_{\mathbf{x}} \mathbf{f}^*(\mathbf{x})^T]$, $\mathbf{f}^*(\mathbf{x}) = \mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x})$ and $\sum_{\mathbf{r} \in O_{r_1}} \mathbf{r} \mathbf{r}^T = I$ since $\mathbf{r} \in O_{r_1}$ are orthogonal to each other.

Thus, we can substitute it back, we have

$$S_{r_1, g_{r_1}} - S_{r_2, g_{r_2}} \geq \text{PSD}_{r_1} - \text{PSD}_{r_2} - C_{\text{sup}} \Omega$$

G.2. Proof of theorem 5

Proof. From proposition 4 we know $f_r^*(\mathbf{x}) = (s_p^r(\mathbf{x}) - s_q^r(\mathbf{x}))$, thus, we can substitute into PSD

(Eq.9), we get

$$\text{PSD}_r = \max_{\mathbf{r} \in \mathbb{S}^{D-1}} \mathbb{E}_q \left[\left((\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x}))^T \mathbf{r} \right)^2 \right]$$

To maximize it, we consider the following constraint optimization problem.

$$\max_{\mathbf{r}} \mathbb{E}_q \left[\left((\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x}))^T \mathbf{r} \right)^2 \right] \quad \text{s.t.} \quad \|\mathbf{r}\|^2 = 1$$

We take the derivative of the corresponding Lagrange multiplier w.r.t. \mathbf{r} ,

$$\begin{aligned} & \mathbb{E}_q \left[\nabla_{\mathbf{r}} \left((\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x}))^T \mathbf{r} \right)^2 \right] - 2\lambda \mathbf{r} = 0 \\ \Rightarrow & \mathbb{E}_q \left[(\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x}))^T \mathbf{r} (\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x})) \right] = \lambda \mathbf{r} \\ \Rightarrow & \underbrace{\mathbb{E}_q \left[(\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x})) (\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x}))^T \right]}_{\mathbf{S} = \mathbb{E}_q[f^*(\mathbf{x})f^{*T}(\mathbf{x})]} \mathbf{r} = \lambda \mathbf{r} \\ \Rightarrow & \mathbf{S} \mathbf{r} = \lambda \mathbf{r} \end{aligned}$$

This exactly the problem of finding eigenpair for matrix \mathbf{S} . Let's assume $\mathbf{r} = \mathbf{v}$ which is the eigenvector of \mathbf{S} with corresponding eigenvalue λ . Substituting it back to PSD , we have

$$\begin{aligned} & \mathbb{E}_q \left[\left((\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x}))^T \mathbf{r} \right)^2 \right] \\ = & \mathbb{E}_q \left[(\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x}))^T \mathbf{r} (\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x})) \right]^T \mathbf{r} \\ = & \mathbf{r}^T \mathbf{S} \mathbf{r} \\ = & \lambda \mathbf{v}^T \mathbf{v} = \lambda \end{aligned}$$

Thus, to obtain the active slice \mathbf{r} , we only need to find the eigenvector of \mathbf{S} with the largest eigenvalue. \square

G.3. Greedy algorithm is eigen-decomposition

Corollary 6.2 (Greedy algorithm is eigen-decomposition). *Assume the conditions in theorem 5 are satisfied, then finding the orthogonal basis O_r from the greedy algorithm is equivalent to the eigen-decomposition of \mathbf{S} .*

Proof. Assume we have obtained the active slice \mathbf{r} from theorem 5, thus, we have $\mathbf{S} \mathbf{r} = \lambda \mathbf{r}$. The greedy algorithm for \mathbf{r}' can be translated into the following constrained optimization

$$\begin{aligned} & \max_{\mathbf{r}'} \mathbb{E}_q \left[\left((\mathbf{s}_p(\mathbf{x}) - \mathbf{s}_q(\mathbf{x}))^T \mathbf{r}' \right)^2 \right] \\ \text{s.t.} & \quad \|\mathbf{r}'\|^2 = 1 \\ & \quad \mathbf{r}'^T \mathbf{r}' = 0 \end{aligned}$$

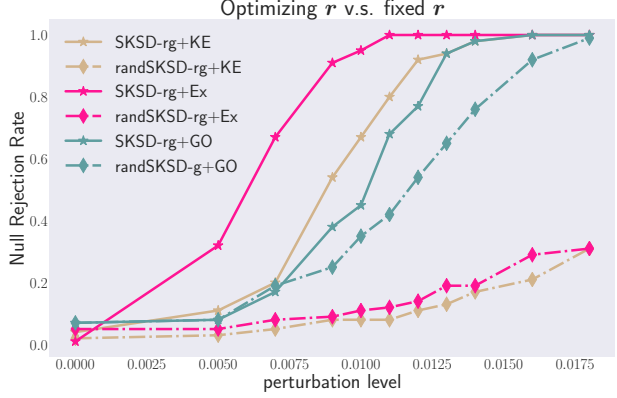


Figure 5. The test power difference with good \mathbf{r} and fixed \mathbf{r} .

By using Lagrange multipliers (μ, γ) , and then take derivative w.r.t. \mathbf{r}' ,

$$\mathbf{S} \mathbf{r}' = \mu \mathbf{r}' + \gamma \mathbf{r}$$

Then taking the inner product with \mathbf{r} in both side, and notice \mathbf{S} is a symmetric matrix, we obtain

$$\begin{aligned} \gamma &= \langle \mathbf{S} \mathbf{r}', \mathbf{r} \rangle \\ &= \langle \mathbf{r}', \mathbf{S}^T \mathbf{r} \rangle \\ &= \langle \mathbf{r}', \lambda \mathbf{r} \rangle = 0 \end{aligned}$$

Therefore, the constrained optimization is the same as the one in theorem 5, which is to find a eigenvector of \mathbf{S} that is different from \mathbf{r} . Repeat the above procedure, the final resulting O_r is a group of eigenvectors of \mathbf{S} . \square

H. Experiment Details

For all experiments in this paper, we use RBF kernel with median heuristics.

H.1. Benchmark GOF test

For gradient based optimization, we use Adam (Kingma & Ba, 2014) with learning rate 0.001 and $\beta = (0.9, 0.99)$. We use random initialization for SKSD-g+GO by drawing \mathbf{g}_r from a Gaussian distribution before normalizing them to unit vectors. For kernel smooth and gradient estimator, we use RBF kernel with median heuristics. Although the algorithm 1 states that small Gaussian noise are needed for active slices, in practice, we found that active slices still have the satisfactory performance without the noise.

The significance level for GOF test $\alpha = 0.05$, and the dimensions of the benchmark problems grow from 2 to 100. We use 1000 bootstrap samples to estimate the threshold and run 100 trials for each benchmark problems.

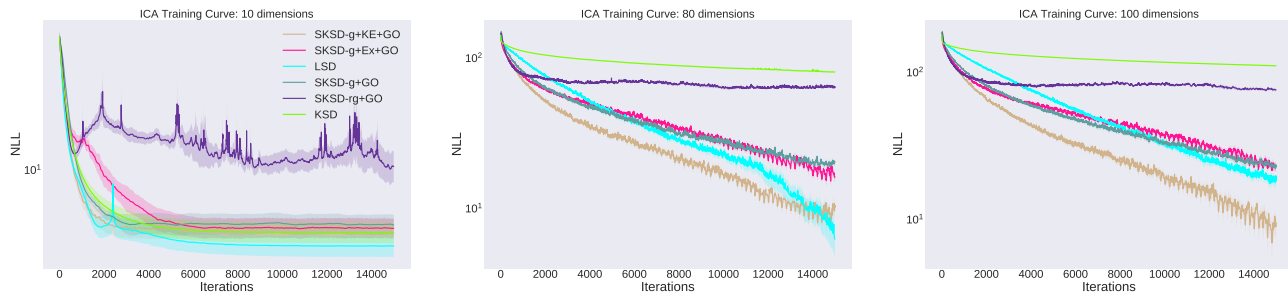


Figure 6. Training curve of ICA model with test NLL for different dimensions.

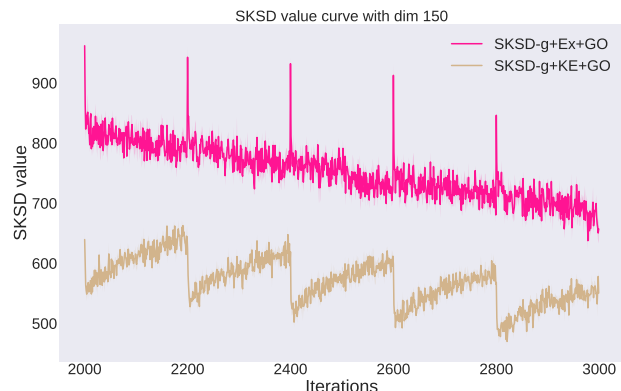


Figure 7. SKSD value curve with 150 dimensional ICA during iteration 5000 to 6000

H.2. RBM GOF test

We set significance level $\alpha = 0.05$ and use 1000 bootstrap samples to compute the threshold. For methods that require training (*SKSD* based method), we need to collect some training samples. Following the same settings as (Gong et al., 2021), to avoid over-fitting to small training set, we collect the pseudo-samples during the early burn-in stage. Note that these pseudo-samples should not be used for testing, as they are not drawn from the q . We collect 2000 samples. For gradient based optimization, we use the same optimizer as benchmark GOF test with the same hyper-parameters. The batch size is 100. For initialization of *SKSD+GO*, we found that if the slices are initialized randomly, the gradient optimization fails to find meaningful slices within a reasonable amount of time, therefore, we have initialize the \mathbf{r} and \mathbf{g}_r as one-hot vectors and set $\mathbf{r} = \mathbf{g}_r$. For pruning ablation study, if the pruning level is set to 50, we initialize \mathbf{r} and \mathbf{g}_r to be the identity matrix. The default number of gradient optimization for *SKSD+GO* is 50. For active slice method, we directly use the active slices without any further optimizations. We run 100 trials for GOF test with 1000 test samples per trial.

(Gong et al., 2021) reports *SKSD-rg+GO* has near optimal test power at perturbation level 0.01. The performance difference is because they train the *SKSD-rg* with 200 batch sizes per burn-in step. Namely, the training set size are $200 \times 2000 = 400000$, which is 200 times larger than ours. They also run 2000 iterations, which is equivalent to 100 epochs in our settings.

Figure 5 shows the test power difference with optimized \mathbf{r} and fixed \mathbf{r} . The legend with *rand* annotation implies we randomly initialized \mathbf{r} as one-hot vectors and fix them while updating \mathbf{g}_r using *GO* or active slice. Without *rand*, it means both \mathbf{r} and \mathbf{g}_r are optimized. We only use 3 \mathbf{r} for active slice method and 50 for gradient-based counterpart. For active slice method with pruning (*randSKSD-g+Ex* or *randSKSD-g+KE*), despite we show that any finite random slices define a valid discrepancy, it is clear that the performance is quite poor with random initialized \mathbf{r} 's. It indicates that using active slices of \mathbf{g}_r alone cannot compensate the poor discriminating power of the random \mathbf{r} 's. Although *SKSD-rg+GO* demonstrates an advantage compared to *randSKSD-g+GO*, the performance boost is less clear compared to active slices method. This is because we do not use any pruning for *randSKSD-g+GO*, and adopt orthogonal basis $O_r = \mathbf{I}$. Despite the orthogonal basis may not capture the important directions, they can provide reasonable discriminating power due to their orthogonality from each other. In summary, using good directions for \mathbf{r} is advantageous compared to fixed \mathbf{r} .

H.3. Model learning: Training ICA

We use Adam optimizer for the model and slice directions with learning rate 0.001 and $\beta = (0.9, 0.99)$. We totally run 15000 iterations. The batch size is 100. We evaluate our method in dimension 10, 80, 100 and 150. For more stable comparisons, we initialized the weight matrix \mathbf{W} until its conditional number is smaller than its dimensions. For active slice method, we use randomly sampled 3000 data from training set to estimate the score difference and the matrices used for eigen-decomposition.

For *SKSD-rg+GO*, we initialize the \mathbf{r} to be a group of one-

hot vectors to form identity matrix and $\mathbf{g}_r = \mathbf{r}$. We use an adversarial training procedure that updates both \mathbf{r} and \mathbf{g}_r using Adam once per iteration before we update the model. For *SKSD-g+GO*, we fix the orthogonal basis O_r to be the identity matrix and only update \mathbf{g}_r . Each results are the average of 5 runs of training.

As for the reason why *SKSD-g+Ex+GO* performs worse than *+KE+GO*, we suspect that *+Ex* only focus on directions with high discriminating power. However, high discriminating power is not necessarily good for model learning. It may focus on very small area that is different from the target but ignore the larger area with small difference. Because our algorithm for finding basis is greedy, this means it can ignore the generally good directions if they are not orthogonal to the directions with high discriminating power.

From figure 7, we can observe there is a spike of *SKSD-g+Ex+GO* value at every 200 iterations due to the new active slices found at the beginning of each training epoch. However, the value drops significantly fast to the one before new active slices. This indicates the *Ex* indeed finds directions with large discriminating power but they do not represents good directions for learning due to the fast drop of SKSD values. On the other hand, the directions provided by KE does not give the highest discriminating power, but it can find generally good directions of \mathbf{g}_r using *GO* refinement steps within a few iterations. This means the directions found by *KE* indeed represents good directions for learning as the model cannot decrease this value quickly. We guess this is due to the smooth estimation of KE, where very small areas with high discriminating power are smoothed out.

Figure 6 shows the ICA training curve of other dimensions. We can observe the convergence speed of LSD deteriorates as the dimension increases due to the poor test function in early training stage, whereas *SKSD-g+KE+GO* maintains the fastest convergence in high dimensions.

I. Perturbation of eigenvectors

The active slice method (algorithm 1) is mainly based on the eigenvalue-decomposition of matrix \mathbf{H} , where

$$\mathbf{H} = \int q(\mathbf{x}) \nabla_{\mathbf{x}} f_r^*(\mathbf{x}) \nabla_{\mathbf{x}} f_r^*(\mathbf{x})^T d\mathbf{x}$$

Obtaining the analytic form of \mathbf{H} involves complicated integration, so Monte Carlo estimation is often used for approximation. We denote it as $\hat{\mathbf{H}}$, with M being the number of samples:

$$\hat{\mathbf{H}} = \frac{1}{M} \sum_{i=1}^M [\nabla_{\mathbf{x}_i} f_r^*(\mathbf{x}_i) \nabla_{\mathbf{x}_i} f_r^*(\mathbf{x}_i)^T] \quad (53)$$

Let \mathbf{g} be the top eigenvector of \mathbf{H} and $\hat{\mathbf{g}}$ be the top eigen-

vector of $\hat{\mathbf{H}}$. Let λ_1, λ_2 be the top two eigenvalues of \mathbf{H} . Assuming the error matrix $\mathbf{E} = \hat{\mathbf{H}} - \mathbf{H}$ is deterministic, (Yu et al., 2015) proved that

$$\|\mathbf{g}\mathbf{g}^T(\mathbf{I} - \hat{\mathbf{g}}\hat{\mathbf{g}}^T)\|_F \leq \frac{2\|\mathbf{E}\|_{op}}{\lambda_1 - \lambda_2} \quad (54)$$

where we define the operator norm for a given $n \times n$ matrix \mathbf{A} as

$$\|\mathbf{A}\|_{op} = \sup\{\|\mathbf{A}\mathbf{x}\| : \mathbf{x} \in \mathbb{R}^n \text{ with } \|\mathbf{x}\| = 1\}$$

We also have (with proof below)

$$\min_{\epsilon \in \{-1, 1\}} \|\mathbf{g} - \epsilon\hat{\mathbf{g}}\|_2 \leq \sqrt{2} \|\mathbf{g}\mathbf{g}^T(\mathbf{I} - \hat{\mathbf{g}}\hat{\mathbf{g}}^T)\|_F \quad (55)$$

Inequality 54 and 55 imply that,

$$\min_{\epsilon \in \{-1, 1\}} \|\mathbf{g} - \epsilon\hat{\mathbf{g}}\|_2 \leq 2^{3/2} \frac{\|\hat{\mathbf{H}} - \mathbf{H}\|_{op}}{\lambda_1 - \lambda_2} \quad (56)$$

I.1. Proof of inequality 55

Proposition 5. *Let \mathbf{S} and \mathbf{U} be two matrices with orthonormal columns and equal rank r . Let $\mathbf{\Pi}_S$ (resp. $\mathbf{\Pi}_U$) indicates the projection matrix to the column space of \mathbf{S} (resp. \mathbf{U}). Then*

$$\min_{O \in \mathbb{R}^{r \times r} \text{ orthogonal}} \|\mathbf{S} - \mathbf{UO}\|_F \leq \sqrt{2} \|\mathbf{\Pi}_S(\mathbf{I} - \mathbf{\Pi}_U)\|_F \quad (57)$$

When $r = 1$, we denote \mathbf{O} as ϵ . Following the definition of orthogonal matrix, we have $\epsilon^T \epsilon = \epsilon^2 = 1$, hence $\epsilon \in \{-1, 1\}$. Substituting $\mathbf{S} = \mathbf{g}$ and $\mathbf{U} = \hat{\mathbf{g}}$, we get inequality 55.

Proof. Let $\mathbf{W}\mathbf{\Sigma}\mathbf{V}^T$ be a singular value decomposition of $\mathbf{S}^T\mathbf{U}$, and use $\mathbf{O} = \mathbf{V}\mathbf{W}^T$. Now,

$$\begin{aligned} \|\mathbf{S} - \mathbf{UO}\|_F^2 &= \text{Tr}((\mathbf{S} - \mathbf{UO})^T(\mathbf{S} - \mathbf{UO})) \\ &= \|\mathbf{S}\|_F^2 + \|\mathbf{U}\|_F^2 - 2\text{Tr}(\mathbf{OS}^T\mathbf{U}) \\ &= 2r - 2\text{Tr}(\mathbf{\Sigma}) \end{aligned}$$

where r is the rank of \mathbf{S} and \mathbf{U} . On the other hand, by Pythagora's theorem

$$\begin{aligned} \|\mathbf{\Pi}_S(\mathbf{I} - \mathbf{\Pi}_U)\|_F^2 &= \|\mathbf{\Pi}_S\|_F^2 - \|\mathbf{\Pi}_S\mathbf{\Pi}_U\|_F^2 \\ &= r - \|\mathbf{\Pi}_S\mathbf{\Pi}_U\|_F^2 \\ &= r - \|\mathbf{S}\mathbf{S}^T\mathbf{U}\mathbf{U}^T\|_F^2 \\ &= r - \text{Tr}(\mathbf{\Sigma}^2) \end{aligned}$$

We claim that the entries of $\mathbf{\Sigma}$ are bounded above by 1, such that $\text{Tr}(\mathbf{\Sigma}) \leq \text{Tr}(\mathbf{\Sigma}^2)$, then

$$\begin{aligned} \min_{O \in \mathbb{R}^{r \times r} \text{ orthogonal}} \|\mathbf{S} - \mathbf{UO}\|_F^2 &\leq 2r - 2\text{Tr}(\mathbf{\Sigma}) \\ &\leq 2r - 2\text{Tr}(\mathbf{\Sigma}^2) \\ &= 2\|\mathbf{\Pi}_S(\mathbf{I} - \mathbf{\Pi}_U)\|_F^2 \end{aligned}$$

Taking the square root of both sides yields the desired inequality. To prove the claim, let $\omega = [\mathbf{S}, \mathbf{S}']$ and $\tilde{\mathbf{U}} = [\mathbf{U}, \mathbf{U}']$ be orthogonal matrices. Then $\mathbf{S}^T \mathbf{U}$ is a diagonal block in $\omega^T \tilde{\mathbf{U}}$. It follows that $\max_i \Sigma_{i,i} = \|\mathbf{S}^T \mathbf{U}\|_{op} \leq \|\omega^T \tilde{\mathbf{U}}\|_{op} = 1$ \square

From Eq.56, we can see if the top two eigenvalues are similar, then their corresponding eigenvectors can be arbitrary different. In terms of our active slice algorithm, it means if the most discriminating directions for two distributions q, p have similar "magnitude of difference", our algorithm may fail under Monte-carlo approximation. On the other hand, if the eigenvalues are different, Eq.56 guarantees that eigenvectors from $\hat{\mathbf{H}}$ are not far-away from the truth.