



## KSD Aggregated Goodness-of-fit Test

Antonin Schrab, Benjamin Guedj, Arthur Gretton

### ► To cite this version:

Antonin Schrab, Benjamin Guedj, Arthur Gretton. KSD Aggregated Goodness-of-fit Test. 2022. hal-03554423

**HAL Id: hal-03554423**

**<https://inria.hal.science/hal-03554423>**

Preprint submitted on 3 Feb 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# KSD Aggregated Goodness-of-fit Test

---

Antonin Schrab<sup>1 2 3</sup> Benjamin Guedj<sup>1 3</sup> Arthur Gretton<sup>2</sup>

## Abstract

We investigate properties of goodness-of-fit tests based on the Kernel Stein Discrepancy (KSD). We introduce a strategy to construct a test, called KSDAGG, which aggregates multiple tests with different kernels. KSDAGG avoids splitting the data to perform kernel selection (which leads to a loss in test power), and rather maximises the test power over a collection of kernels. We provide theoretical guarantees on the power of KSDAGG: we show it achieves the smallest uniform separation rate of the collection, up to a logarithmic term. KSDAGG can be computed exactly in practice as it relies either on a parametric bootstrap or on a wild bootstrap to estimate the quantiles and the level corrections. In particular, for the crucial choice of bandwidth of a fixed kernel, it avoids resorting to arbitrary heuristics (such as median or standard deviation) or to data splitting. We find on both synthetic and real-world data that KSDAGG outperforms other state-of-the-art adaptive KSD-based goodness-of-fit testing procedures.

## 1. Introduction

Kernel selection remains a fundamental question in kernel-based nonparametric hypothesis testing, as it significantly impacts the test power. Kernel selection has attracted a significant interest in the literature, and a number of methods have been proposed in different settings, such as in the two-sample, independence and goodness-of-fit testing frameworks. Those methods include using heuristics (Gretton et al., 2012a), relying on data splitting (Gretton et al., 2012b; Sutherland et al., 2017), learning deep kernels (Grathwohl et al., 2020; Liu et al., 2020), working in the post-selection inference framework (Yamada et al., 2019; Lim et al., 2019; 2020; Freidling et al., 2021), to name but a few.

In this work, we focus on aggregated tests, which have been investigated for the two-sample problem by Fromont et al. (2013), Kim et al. (2020) and Schrab et al. (2021) using the Maximum Mean Discrepancy (MMD, Gretton et al., 2012a) and for the independence problem by Albert et al. (2019) and Kim et al. (2020) using the Hilbert Schmidt Independence Criterion (HSIC, Gretton et al., 2005). We extend the use of aggregated tests to the goodness-of-fit setting, where we are given a model and some samples, and we are interested in deciding whether the samples have been drawn from the model. We employ the Kernel Stein Discrepancy (KSD, Chwialkowski et al., 2016; Liu et al., 2016) as our test statistic, which is an ideal measure of distance for this setting: it admits an estimator which can be computed without requiring samples from the model, and does not require the model to be normalised. To the best of our knowledge, ours represents the first aggregation procedure for the KSD test in the literature.

**Related work.** Fromont et al. (2012; 2013) introduced aggregated tests for the two-sample problem with the equal sample sizes following a Poisson process, using as a test statistic an unscaled version of the MMD. They also provided theoretical results in terms of uniform separation rates using a wild bootstrap. Albert et al. (2019) then proposed an aggregated test for the independence problem using the HSIC, with guarantees using a theoretical quantile, but relying on permutations to obtain the test threshold. Kim et al. (2020) then extended those results to also hold for the estimated quantile, and generalised the two-sample results to hold for the MMD estimator with different sample sizes using a wild bootstrap. All those aforementioned results were proved for the Gaussian kernel only. Schrab et al. (2021) generalised the two-sample results to hold for a wide range of kernels using either a wild bootstrap or a permutation-based procedure, and provided optimality results which hold with fewer restrictions on the class of functions. Our work builds and extends on the above results: we consider the goodness-of-fit framework, where we have samples from only one of the two densities. The main challenges arise from working with the Stein kernel that defines the KSD test statistic: for example, we lose the transition invariant property of the kernel which is crucial to work in the Fourier domain. Finally, Key et al. (2021) addresses the complementary task of Stein test design for a family of models, rather than for a single model.

<sup>1</sup>Centre for Artificial Intelligence, Department of Computer Science, University College London <sup>2</sup>Gatsby Computational Neuroscience Unit, University College London <sup>3</sup>Inria, Lille – Nord Europe research centre and Inria London Programme. Correspondence to: Antonin Schrab <a.schrab@ucl.ac.uk>, Benjamin Guedj <b.guedj@ucl.ac.uk>, Arthur Gretton <arthur.gretton@gmail.com>.

**Contributions.** We introduce our adaptive test KSDAGG which aggregates multiple tests with different kernels. This is a general setting which has many applications, such as kernel bandwidth selection. Our aggregated test allows for two numerical methods for estimating the test thresholds: the parametric bootstrap and the wild bootstrap. We conduct a theoretical analysis: we provide a lower bound on the uniform separation rate (Baraud, 2002) of KSDAGG, a condition which guarantees test power. We discuss the implementation of KSDAGG and its use in practice through experiments on synthetic and real-world data, on which we observe that KSDAGG obtains higher power than other KSD-based adaptive state-of-the-art tests. The code is publicly available at <https://github.com/antoninschrab/ksdagg-paper/>.

**Outline.** Section 2 presents our framework and our notation. Section 3 introduces our algorithm KSDAGG (in Algorithm 1) and contains our main theoretical results, and Section 4 presents numerical experiments to support KSDAGG. We close the paper with avenues for future research in Section 5.

## 2. Notation

We consider the goodness-of-fit problem where given access to a known probability density  $p$  (model) and to some i.i.d.  $d$ -dimensional samples  $\mathbb{X}_N := (X_i)_{i=1}^N$  drawn from an unknown density  $q$ , we want to decide whether  $p \neq q$  holds. This can be expressed as a statistical hypothesis testing problem with null hypothesis  $\mathcal{H}_0 : p = q$  and alternative  $\mathcal{H}_a : p \neq q$ .

As a measure of distance between  $p$  and  $q$ , we use the *Kernel Stein Discrepancy (KSD)* introduced by Chwialkowski et al. (2016) and Liu et al. (2016). For a kernel  $k$ , the KSD is defined as the Maximum Mean Discrepancy (MMD, Gretton et al., 2012a) between  $p$  and  $q$  using the Stein kernel associated to  $k$ , that is

$$\begin{aligned} \text{KSD}_{p,k}^2(q) &:= \text{MMD}_{h_{p,k}}^2(p, q) \\ &:= \mathbb{E}_{q,q}[h_{p,k}(X, Y)] - 2\mathbb{E}_{p,q}[h_{p,k}(X, Y)] \\ &\quad + \mathbb{E}_{p,p}[h_{p,k}(X, Y)] \\ &= \mathbb{E}_{q,q}[h_{p,k}(X, Y)] \end{aligned}$$

where the *Stein kernel*  $h_{p,k} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as

$$\begin{aligned} h_{p,k}(x, y) &:= (\nabla \log p(x)^\top \nabla \log p(y)) k(x, y) \\ &\quad + \nabla \log p(y)^\top \nabla_1 k(x, y) \\ &\quad + \nabla \log p(x)^\top \nabla_2 k(x, y) \\ &\quad + \sum_{i=1}^d \frac{\partial}{\partial x_i \partial y_i} k(x, y) \end{aligned}$$

where  $\nabla_i$  denotes the gradient with respect to the  $i$ -th component. The Stein kernel satisfies the *Stein identity*

$$\mathbb{E}_p[h_{p,k}(X, \cdot)] = 0.$$

A quadratic-time *KSD estimator* can be computed as the  $U$ -statistic (Hoeffding, 1992)

$$\widehat{\text{KSD}}_{p,k}^2(\mathbb{X}_N) := \frac{1}{N(N-1)} \sum_{1 \leq i \neq j \leq N} h_{p,k}(X_i, X_j). \quad (1)$$

In this work, the model density  $p$  is always known, we do not always explicitly write the dependence of  $p$  for all variables as we do for  $\text{KSD}_{p,k}^2$  and  $h_{p,k}$ .

We assume that the kernel  $k$  is such that

$$\text{KSD}_{p,k}^2(q) = \mathbb{E}_{q,q}[h_{p,k}(X, Y)] < \infty \quad (2)$$

and

$$C_k := \mathbb{E}_{q,q}[h_{p,k}(X, Y)^2] < \infty. \quad (3)$$

We now address the requirements for consistency of the Stein test (Chwialkowski et al., 2016, Theorem 2.2): we assume that the kernel  $k$  is  $C_0$ -universal (Carmeli et al., 2010, Definition 4.1) and that  $\mathbb{E}_q \left\| \nabla \left( \log \frac{p(X)}{q(X)} \right) \right\|_2^2 < \infty$ .

We use the notations  $\mathbb{P}_p$  and  $\mathbb{P}_q$  to denote the probability under the model  $p$  and under  $q$ , respectively. Given a kernel  $\kappa : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  and a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , we consider the *integral transform*

$$(\kappa \diamond f)(y) := \int_{\mathbb{R}^d} \kappa(x, y) f(x) dx$$

following the notation of Fromont et al. (2012). When the kernel  $\kappa$  is translation invariant, the operation  $\diamond$  corresponds to a convolution, however this is not true of the Stein kernel.

## 3. Construction of tests and bounds

We now introduce the single and aggregated KSD tests. We show that these control the probability of type I error as desired, and provide conditions for the control of type II error probability.

### 3.1. Single test

We first construct a KSD test for a fixed kernel  $k$  as proposed by Chwialkowski et al. (2016) and Liu et al. (2016). To estimate the test threshold, we can either use a wild bootstrap (Shao, 2010; Leucht & Neumann, 2013; Fromont et al., 2012) or a parametric bootstrap (Stute et al., 1993). Both methods work by simulating sampling values  $(\bar{K}_k^1, \dots, \bar{K}_k^{B_1})$  from the (asymptotic) distribution of

$\widehat{\text{KSD}}_{p,k}^2$  under the null hypothesis and estimating the  $(1-\alpha)$ -quantile using a Monte Carlo approximation<sup>1</sup>

$$\begin{aligned}\hat{q}_{1-\alpha}^k &:= \inf \left\{ u \in \mathbb{R} : 1 - \alpha \leq \frac{1}{B_1} \sum_{b=1}^{B_1} \mathbb{1}(\bar{K}_k^b \leq u) \right\} \\ &= \bar{K}_k^{\bullet \lceil B_1(1-\alpha) \rceil}\end{aligned}$$

where  $\bar{K}_k^{\bullet 1} \leq \dots \leq \bar{K}_k^{\bullet B_1}$  are the sorted elements  $(\bar{K}_k^1, \dots, \bar{K}_k^{B_1})$ . The single test is then defined as

$$\Delta_\alpha^k(\mathbb{X}_N) := \mathbb{1} \left( \widehat{\text{KSD}}_{p,k}^2(\mathbb{X}_N) > \hat{q}_{1-\alpha}^k \right).$$

For the *parametric bootstrap*, we directly draw new samples  $(X'_i)_{i=1}^{N'}$  from the model  $p$  and compute the KSD

$$\bar{K}_k := \frac{1}{N'(N'-1)} \sum_{1 \leq i \neq j \leq N'} h_{p,k}(X'_i, X'_j). \quad (4)$$

For the *wild bootstrap*, we first generate  $n$  i.i.d. Rademacher random variables  $\epsilon_1, \dots, \epsilon_n$  taking values in  $\{-1, 1\}^n$ , and then compute

$$\bar{K}_k := \frac{1}{N(N-1)} \sum_{1 \leq i \neq j \leq N} \epsilon_i \epsilon_j h_{p,k}(X_i, X_j). \quad (5)$$

Both these processes are then repeated  $B_1$  times.

Since it uses samples from the model  $p$ , the parametric bootstrap results in a test with non-asymptotic level  $\alpha$ . This comes at the cost of being computationally more expensive and assuming that we are able to sample from  $p$  (which may be out of reach in some settings). Conversely, the wild bootstrap has the advantage of not requiring to sample from  $p$ , which makes it computationally more efficient as only one kernel matrix needs to be computed, but it only achieves the desired level  $\alpha$  asymptotically (Shao, 2010; Leucht & Neumann, 2013; Chwialkowski et al., 2016; 2014). Note that we cannot obtain a non-asymptotic level for the wild bootstrap by relying on the result of Romano & Wolf (2005, Lemma 1) as done in the two-sample framework by Fromont et al. (2013) and Schrab et al. (2021). This is because in our case  $\bar{K}_k$  and  $\widehat{\text{KSD}}_{p,k}^2(\mathbb{X}_N)$  are not exchangeable variables under the null hypothesis, due to the asymmetry of the KSD statistic with respect to  $p$  and  $q$ .

Having discussed control of the probability of type I error of the single test  $\Delta_\alpha^k$ , we now provide a condition on  $\|p - q\|_2$  which ensures that the probability of type II error is controlled by some  $\beta \in (0, 1)$ . The smallest such value of  $\|p - q\|_2$ , provided that  $p - q$  lies in some given class of functions, is called the *uniform separation rate* (Baraud, 2002).

<sup>1</sup>We do not write explicitly the dependence of  $\hat{q}_{1-\alpha}^k$  on other variables, but those are implicitly considered when writing probabilistic statements.

**Theorem 3.1.** *Let  $\psi := p - q$  and assume that  $\max(\|p\|_\infty, \|q\|_\infty) \leq M$ . Consider  $C_k$  as defined in Equation (3),  $\alpha \in (0, e^{-1})$ ,  $\beta \in (0, 1)$  and  $B_1 \in \mathbb{N}$  satisfying  $B_1 \geq \frac{3}{\alpha^2} (\ln(\frac{8}{\beta}) + \alpha(1 - \alpha))$ . There exists a positive constant  $C$  such that the condition*

$$\|\psi\|_2^2 \geq \|\psi - h_{p,k} \diamond \psi\|_2^2 + C \log\left(\frac{1}{\alpha}\right) \frac{\sqrt{C_k}}{\beta N}$$

*guarantees control over the probability of type II error*

$$\mathbb{P}_q(\Delta_\alpha^k(\mathbb{X}_N) = 0) \leq \beta.$$

*Proof.* We can obtain this result by following the reasoning of Schrab et al. (2021, Theorem 5) for the two-sample case. We highlight the main differences in the proof. Similarly to the result of Schrab et al. (2021, Lemma 2), using Chebyshev's inequality, we obtain the sufficient condition for power

$$\text{KSD}_{p,k}^2(q) \geq \sqrt{\frac{2}{\beta} \text{var}(\widehat{\text{KSD}}_{p,k}^2(\mathbb{X}_N))} + \hat{q}_{1-\alpha}^k.$$

We then need to upper bound the variance and quantile terms. One can obtain using the explicit variance formula (Lee, 1990) that

$$\text{var}(\widehat{\text{KSD}}_{p,k}^2(\mathbb{X}_N)) \leq C \left( \frac{\sigma_1^2}{N} + \frac{\sigma_2^2}{N^2} \right)$$

where

$$\sigma_1^2 := \mathbb{E}_{X'}[(\mathbb{E}_{X'}[h_{p,k}(X, X')])^2] \leq M \|h_{p,k} \diamond \psi\|_2^2$$

by Stein's identity and

$$\sigma_2^2 := \mathbb{E}_{X, X'}[h_{p,k}(X, X')^2] = C_k.$$

A similar reasoning to the one of Schrab et al. (2021, Proposition 4) gives us

$$\hat{q}_{1-\alpha}^k \leq C \log\left(\frac{1}{\alpha}\right) \frac{\sqrt{C_k}}{N}$$

provided that the estimated quantile is computed using  $B_1 \geq \frac{3}{\alpha^2} (\ln(\frac{8}{\beta}) + \alpha(1 - \alpha))$  values. The upper bounds involving the  $C_k$  constant can then be combined using the basic inequalities  $2\sqrt{xy} \leq x + y$  and  $\sqrt{x + y} \leq \sqrt{x} + \sqrt{y}$ , together with the fact that

$$\begin{aligned}\text{KSD}_{p,k}^2(q) &= \langle \psi, h_{p,k} \diamond \psi \rangle \\ &= \frac{1}{2} \left( \|\psi\|_2^2 + \|h_{p,k} \diamond \psi\|_2^2 - \|\psi - h_{p,k} \diamond \psi\|_2^2 \right)\end{aligned}$$

to obtain the desired result. A similar 'trick' is also present in the works of Fromont et al. (2013) and Albert et al. (2019).  $\square$

### 3.2. Aggregated test

We can now introduce our aggregated test, which is motivated by the earlier works of [Fromont et al. \(2012; 2013\)](#), [Albert et al. \(2019\)](#), and [Schrab et al. \(2021\)](#).

We compute  $\tilde{K}_k^1, \dots, \tilde{K}_k^{B_2}$  further simulated KSD values from the null hypothesis obtained using either a parametric bootstrap or a wild bootstrap as in Equations (4) or (5), respectively. Consider a finite collection of kernels  $\mathcal{K}$  satisfying the properties presented in Section 2. We construct an aggregated test  $\Delta_\alpha^\mathcal{K}$ , called KSDAGG, which rejects the null hypothesis if one of the single tests  $(\Delta_{u_\alpha w_k}^k)_{k \in \mathcal{K}}$  rejects it, that is

$$\Delta_\alpha^\mathcal{K}(\mathbb{X}_N) := \mathbb{1}(\Delta_{u_\alpha w_k}^k(\mathbb{X}_N) = 1 \text{ for some } k \in \mathcal{K}).$$

The levels of the single tests are adjusted to ensure the aggregated test has the prescribed level  $\alpha$ . This adjustment is performed by introducing positive weights  $(w_k)_{k \in \mathcal{K}}$  satisfying  $\sum_{k \in \mathcal{K}} w_k \leq 1$  and some correction

$$u_\alpha := \sup \left\{ u \in (0, \min_{k \in \mathcal{K}} w_k^{-1}) : \hat{P}_u \leq \alpha \right\} \quad (6)$$

where

$$\hat{P}_u := \frac{1}{B_2} \sum_{b=1}^{B_2} \mathbb{1} \left( \max_{k \in \mathcal{K}} \left( \tilde{K}_k^b - \bar{K}_k^{\bullet \lceil B_1(1-u w_k) \rceil} \right) > 0 \right)$$

is a Monte Carlo approximation of the probability of type I error of our aggregated test with correction  $u$

$$P_u := \mathbb{P}_p \left( \max_{k \in \mathcal{K}} \left( \widehat{\text{KSD}}_{p,k}^2(\mathbb{X}_N) - \hat{q}_{1-u w_k}^k \right) > 0 \right).$$

In practice, we work with a finite collection of kernels, and thus the extrema over  $\mathcal{K}$  are achieved. To compute  $u_\alpha$ , we estimate the supremum in Equation (6) by performing  $B_3$  steps of the bisection method, as proposed by [Schrab et al. \(2021\)](#). Detailed pseudocode for KSDAGG is provided in Algorithm 1.

We verify in the next proposition that performing this correction indeed ensures that our aggregated test  $\Delta_\alpha^\mathcal{K}$  has the prescribed level  $\alpha$ .

**Proposition 3.2.** *For  $\alpha \in (0, 1)$  and  $B_1, B_2, B_3 \in \mathbb{N} \setminus \{0\}$ , the aggregated test  $\Delta_\alpha^\mathcal{K}$  satisfies*

$$\mathbb{P}_p(\Delta_\alpha^\mathcal{K}(\mathbb{X}_N) = 1) \leq \alpha$$

*non-asymptotically when using a parametric bootstrap and asymptotically when using a wild bootstrap.*

*Proof.* This result follows from the proof provided by [Schrab et al. \(2021, Proposition 8\)](#) and the facts that the parametric bootstrap ([Stute et al., 1993](#)) and wild bootstrap ([Chwiałkowski et al., 2016](#)) have non-asymptotic and asymptotic levels  $\alpha$ , respectively.  $\square$

We now provide guarantees for the power of our aggregated test KSDAGG in terms of its uniform separation rate.

---

#### Algorithm 1 KSDAGG

---

**Inputs:** samples  $\mathbb{X}_N = (x_i)_{i=1}^N$ , density  $p$  or score function  $\nabla \log p(\cdot)$ , finite kernel collection  $\mathcal{K}$ , weights  $(w_k)_{k \in \mathcal{K}}$ , level  $\alpha \in (0, e^{-1})$ , estimation parameters  $B_1, B_2, B_3 \in \mathbb{N}$ , parametric or wild bootstrap

**for**  $k \in \mathcal{K}$  **do**

    Compute  $(\tilde{K}_k^b)_{1 \leq b \leq B_1}$  as in Equations (4) or (5)

    Sort in ascending order to obtain  $(\tilde{K}_k^{\bullet b})_{1 \leq b \leq B_1}$

    Compute  $(\tilde{K}_k^b)_{1 \leq b \leq B_2}$  as in Equations (4) or (5)

**end for**

$u_{\min} = 0$

$u_{\max} = \min_{k \in \mathcal{K}} w_k^{-1}$

**for**  $t = 1, \dots, B_3$  **do**

$u = \frac{1}{2}(u_{\min} + u_{\max})$

$\hat{P}_u = \frac{1}{B_2} \sum_{b=1}^{B_2} \mathbb{1} \left( \max_{k \in \mathcal{K}} \left( \tilde{K}_k^b - \bar{K}_k^{\bullet \lceil B_1(1-u w_k) \rceil} \right) > 0 \right)$

**if**  $\hat{P}_u \leq \alpha$  **then**

$u_{\min} = u$

**else**

$u_{\max} = u$

**end if**

**end for**

$u_\alpha = u_{\min}$

**if**  $\max_{k \in \mathcal{K}} \left( \widehat{\text{KSD}}_{p,k}^2(\mathbb{X}_N) - \bar{K}_k^{\bullet \lceil B_1(1-u_\alpha w_k) \rceil} \right) > 0$  **then**

**return** 1 (reject  $\mathcal{H}_0$ )

**else**

**return** 0 (fail to reject  $\mathcal{H}_0$ )

**end if**

---

**Theorem 3.3.** *Let  $\psi := p - q$  denote the difference in densities and assume that  $\max(\|p\|_\infty, \|q\|_\infty) \leq M$ . Consider a collection of kernels  $\mathcal{K}$  with associated positive weights  $(w_k)_{k \in \mathcal{K}}$  satisfying  $\sum_{k \in \mathcal{K}} w_k \leq 1$ . Consider  $C_k$  as defined in Equation (3),  $\alpha \in (0, e^{-1})$ ,  $\beta \in (0, 1)$  and  $B_1, B_2, B_3 \in \mathbb{N}$  satisfying  $B_1 \geq \frac{3}{\alpha^2} (\ln(\frac{8}{\beta}) + \alpha(1 - \alpha))$ ,  $B_2 \geq \frac{8}{\alpha^2} \ln(\frac{2}{\beta})$  and  $B_3 \geq \log_2(\frac{4}{\alpha} \min_{k \in \mathcal{K}} w_k^{-1})$ . There exists a positive constant  $C$  such that if*

$$\|\psi\|_2^2 \geq \min_{k \in \mathcal{K}} \left( \|\psi - h_{p,k} \diamond \psi\|_2^2 + C \log \left( \frac{1}{\alpha w_k} \right) \frac{\sqrt{C_k}}{\beta N} \right)$$

*then the probability of type II error of  $\Delta_\alpha^\mathcal{K}$  is controlled,*

$$\mathbb{P}_q(\Delta_\alpha^\mathcal{K}(\mathbb{X}_N) = 0) \leq \beta.$$

*Proof.* Applying the result of [Schrab et al. \(2021, Theorem 9\)](#), we obtain that the uniform separation rate of our aggregated test  $\Delta_\alpha^\mathcal{K}(\mathbb{X}_N)$  is lower bounded by the minimum over the uniform separation rates of the single tests  $(\Delta_{\alpha w_k}^k)_{k \in \mathcal{K}}$ . Combining this result with Theorem 3.1 concludes the proof.  $\square$

We observe that the aggregation procedure allows to achieve the smallest uniform separation rate of the single tests  $(\Delta_\alpha^k)_{k \in \mathcal{K}}$  up to some logarithmic weighting term  $\log(1/w_k)$ .

### 3.3. Bandwidth selection

A specific application of the setting we have considered is the problem of bandwidth selection for a fixed kernel. Given a kernel  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ , the function

$$k_\lambda(x, y) := k\left(\frac{x}{\lambda}, \frac{y}{\lambda}\right)$$

is also a kernel for any bandwidth  $\lambda > 0$ . A common example is the Gaussian kernel, for which we have  $k(x, y) = \exp(-\|x - y\|_2^2)$  and

$$k_\lambda(x, y) = \exp\left(-\frac{\|x - y\|_2^2}{\lambda^2}\right).$$

As shown by [Gorham & Mackey \(2017\)](#), a more appropriate kernel for goodness-of-fit testing using the KSD is the IMQ (inverse multiquadric) kernel, which is defined with  $k(x, y) = \left(1 + \|x - y\|_2^2\right)^{-\beta_k}$  for a fixed parameter  $\beta_k \in (0, 1)$  as

$$\begin{aligned} k_\lambda(x, y) &= \left(1 + \frac{\|x - y\|_2^2}{\lambda^2}\right)^{-\beta_k} \\ &= \lambda^{2\beta_k} \left(\lambda^2 + \|x - y\|_2^2\right)^{-\beta_k} \\ &\propto \left(\lambda^2 + \|x - y\|_2^2\right)^{-\beta_k} \end{aligned} \quad (7)$$

which is the well-known form of the IMQ kernel with parameters  $\lambda > 0$  and  $\beta_k \in (0, 1)$ . Note that it is justified to consider the kernel up to a multiplicative constant because our single and aggregated tests are invariant under this kernel transformation.

In practice, as suggested by [Gretton et al. \(2012a\)](#) the bandwidth is often set to a heuristic such as the median or the standard deviation of the  $L^2$ -distances between the samples  $(X_i)_{i=1}^N$ , however, these are arbitrary choices with no theoretical guarantees. Another common approach proposed by [Gretton et al. \(2012b\)](#); [Liu et al. \(2020\)](#) is to resort to data splitting in order to select a bandwidth on held-out data, by maximising for a proxy for asymptotic power (see Section 4.1 for details). Both methods were originally proposed for the two-sample problem, but extend straightforwardly to the goodness-of-fit setting.

By considering a kernel collection  $\mathcal{K}_\Lambda = \{k_\lambda : \lambda \in \Lambda\}$  for a collection of bandwidths  $\Lambda$ , we can use our aggregated test KSDAGG to test multiple bandwidths using all the data and without resorting to arbitrary heuristics. We now obtain an expression for the uniform separation rate of  $\Delta_\alpha^{\mathcal{K}_\Lambda}$  in terms of the bandwidth  $\lambda$ .

**Corollary 3.4.** *Consider  $\alpha \in (0, e^{-1})$ ,  $\beta \in (0, 1)$  and  $B_1, B_2, B_3 \in \mathbb{N}$  satisfying the conditions of Theorem 3.3 and assume that  $\max(\|p\|_\infty, \|q\|_\infty) \leq M$ . Given a collection  $\Lambda$  of positive bandwidths with associated positive weights  $(w_\lambda)_{\lambda \in \Lambda}$  satisfying  $\sum_{\lambda \in \Lambda} w_\lambda \leq 1$ , we consider  $\mathcal{K}_\Lambda = \{k_\lambda : \lambda \in \Lambda\}$ . There exists a positive constant  $C$  such that the condition*

$$\|\psi\|_2^2 \geq \min_{\lambda \in \Lambda} \left( \|\psi - h_{p, k_\lambda} \diamond \psi\|_2^2 + C \log\left(\frac{1}{\alpha w_\lambda}\right) \frac{\sqrt{C_{k_\lambda}}}{\beta N} \right)$$

*ensures control over the probability of type II error of the aggregated test*

$$\mathbb{P}_q(\Delta_\alpha^{\mathcal{K}_\Lambda}(\mathbb{X}_N) = 0) \leq \beta.$$

*Proof.* We simply apply Theorem 3.3 to the collection of kernels  $\mathcal{K}_\Lambda$ .  $\square$

We do not impose any restrictions on  $\psi := p - q$  such as assuming it belongs to a specific regularity class. For this reason, our result holds more generally but the dependence on  $\lambda$  in the terms  $\|\psi - h_{p, k_\lambda} \diamond \psi\|_2^2$  and  $-\log(\alpha w_\lambda)(\beta N)^{-1} \sqrt{C_{k_\lambda}}$  is not explicit. For a particular regularity class, one can obtain a uniform separation rate  $N^{-R}$  for some  $R > 0$  by choosing appropriate collections of bandwidths and weights (depending on  $N$ ) such that the two terms have matching orders of  $N$ .

## 4. Implementation and experiments

We consider three different experiments based on a Gamma one-dimensional distribution, a Bernoulli Gaussian Restricted Boltzmann Machine, and a Normalizing Flow for the MNIST dataset. We compare our proposed aggregated test KSDAGG against three alternatives: the KSD test which uses the median bandwidth, a test which splits the data to select an ‘optimal’ bandwidth according to a proxy for asymptotic test power, and a test which uses extra data for bandwidth selection. The ‘extra data’ test is designed simply to provide a best case oracle for the ‘asymptotic test power’ bandwidth selection procedure, and cannot be used in practice (i.e. any extra samples from  $q$  would normally be incorporated into the sample being tested). In order to ensure that our tests always have correct levels for all bandwidth values, dimensions and sample sizes, we use the parametric bootstrap in our experiments.

### 4.1. Alternative bandwidth selection approaches

[Gretton et al. \(2012a\)](#) propose to use the median heuristic as kernel bandwidth, consisting in the median of the  $L^2$ -distances between the samples given by

$$\lambda_{\text{med}} := \text{median}\{\|x_i - x_j\|_2 : 0 \leq i < j \leq N\}.$$

Gretton et al. (2012b) first proposed, for the two-sample problem using a linear-time MMD estimator, to split the data and to use half of it to select an ‘optimal’ bandwidth which maximises a proxy for asymptotic power. This procedure was then extended to quadratic-time estimators and to the goodness-of-fit framework by Jitkrittum et al. (2017), Sutherland et al. (2017) and Liu et al. (2020). These strategies rely on the asymptotic normality of the test statistic under  $\mathcal{H}_a$ . In our setting, the asymptotic power proxy to maximise is the ratio

$$\widehat{\text{KSD}}_{p,k}^2(\mathbb{X}_N) / \widehat{\sigma}_{\mathcal{H}_a}^2 \quad (8)$$

where  $\widehat{\sigma}_{\mathcal{H}_a}^2$  is a regularised positive estimator of the asymptotic variance of  $\widehat{\text{KSD}}_{p,k}^2$  under  $\mathcal{H}_a$  (Liu et al., 2020, Equation 5) which can be computed as

$$\frac{4}{N^3} \sum_{i=1}^N \left( \sum_{j=1}^N H_{i,j} \right)^2 - \frac{4}{N^4} \left( \sum_{i=1}^N \sum_{j=1}^N H_{i,j} \right)^2$$

with  $H_{i,j} := h_{p,k}(X_i, X_j)$ . In our experiments, we also consider a test which has access to  $N$  extra samples drawn from  $q$  to select an ‘optimal’ bandwidth to run the KSD test on the original  $N$  samples  $\mathbb{X}_N$ . This test is interesting to compare to because it uses an ‘optimal’ bandwidth without being detrimental to power.

## 4.2. Experimental details

In our experiments, we use collections of bandwidths of the form

$$\Lambda(\ell_-, \ell_+) := \{2^i \lambda_{\text{med}} : i = \ell_-, \dots, \ell_+\}$$

for the median bandwidth  $\lambda_{\text{med}}$  and integers  $\ell_- < \ell_+$  with uniform weights

$$w_\lambda := \frac{1}{\ell_+ - \ell_- + 1}.$$

For the tests which split the data, we select the bandwidth out of the collection  $\Lambda(\ell_-, \ell_+)$  which maximises the power proxy in Equation (8). All our experiments are run with level  $\alpha = 0.05$  using the IMQ kernel defined in Equation (7) with parameter  $\beta_k = 0.5$ . We use a parametric bootstrap with  $B_1 = B_2 = 500$  simulated KSD values to compute the adjusted test thresholds and  $B_3 = 50$  steps of bisection method to estimate the correction  $u_\alpha$  in Equation (6). To estimate the probability of rejecting the null hypothesis, we average the test outputs across 200 repetitions.

## 4.3. Gamma distribution

For our first experiment, we consider a one-dimensional Gamma distribution with shape parameter 5 and scale

parameter 5 as the model  $p$ . We draw 500 samples from a Gamma distribution with the same scale parameter 5 and with a shifted shape parameter  $5 + s$  for  $s \in \{0, 0.1, 0.2, 0.3, 0.4\}$ . We consider the collection of bandwidths  $\Lambda(0, 10)$ .

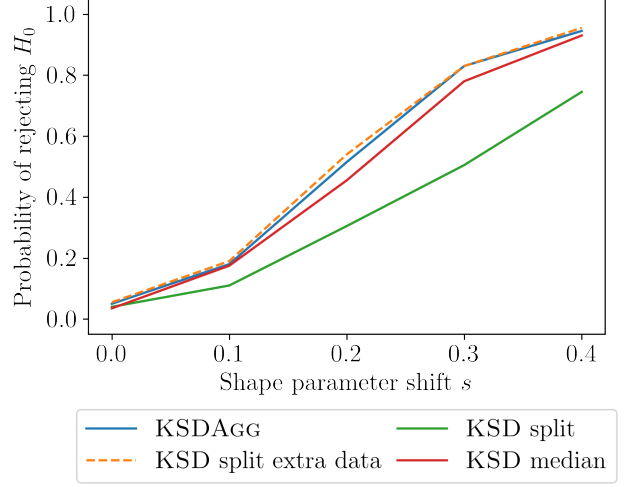


Figure 1. Gamma distribution experiment.

The results we obtained are presented in Figure 1. We observe that all tests have the prescribed level 0.05 under the null hypothesis, corresponding to the case  $s = 0$ . As the shift parameter  $s$  increases, the two densities  $p$  and  $q$  become more different and rejection of the null becomes an easier task, thus the test power increases. Our aggregated test KSDAGG achieves the same power as the ‘best case’ bound on the performance of the asymptotic power heuristic, yielded by the splitting test with extra data. The median test obtains only slightly lower power, the closeness in power can be explained by the fact that this one-dimensional problem is an easy one. We note that the normal splitting test has significantly lower power: this is because, even though it uses an ‘optimal’ bandwidth, it is then run on only half the data, which results in a loss of power.

## 4.4. Bernoulli-Gaussian Restricted Boltzmann Machine

As first considered by Liu et al. (2016) for goodness-of-fit testing using the KSD, we consider a Bernoulli-Gaussian Restricted Boltzmann Machine. It is a graphical model with a binary hidden variable  $h \in \{-1, 1\}^{d_h}$  and a continuous observable variable  $x \in \mathbb{R}^d$ . Those variables have joint density

$$p(x, h) = \frac{1}{Z} \exp \left( \frac{1}{2} x^\top B h + b^\top x + c^\top h - \frac{1}{2} \|x\|^2 \right)$$

where  $Z$  is an unknown normalizing constant. By marginalising over  $h$ , we obtain the density  $p$  of  $x$

$$p(x) = \sum_{h \in \{-1, 1\}^{d_h}} p(x, h).$$

We can sample from it using a Gibbs sampler with 2000 burn-in iterations. We use the dimensions  $d = 50$  and  $d_h = 40$  as considered by Jitkrittum et al. (2017) and Grathwohl et al. (2020). Even though computing  $p$  is intractable for large dimension  $d_h$ , the score function admits a convenient closed form

$$\nabla \log p(x) = b - x + B \frac{\exp(2(B^\top x + c)) - 1}{\exp(2(B^\top x + c)) + 1}.$$

We draw the components of  $b$  and  $c$  from Gaussian standard distributions and sample Rademacher variables taking values in  $\{-1, 1\}$  for the elements of  $B$  for the model  $p$ . We draw 1000 samples from a distribution  $q$  which is constructed in a similar way as  $p$  but with the difference that some Gaussian noise  $\mathcal{N}(0, \sigma)$  is injected into the elements of  $B$ . We consider the standard deviations of the perturbations  $\sigma \in \{0, 0.01, 0.02, 0.03\}$ . We run our experiments with the collection of bandwidths  $\Lambda(-20, 0)$  and provide the results in Figure 2.

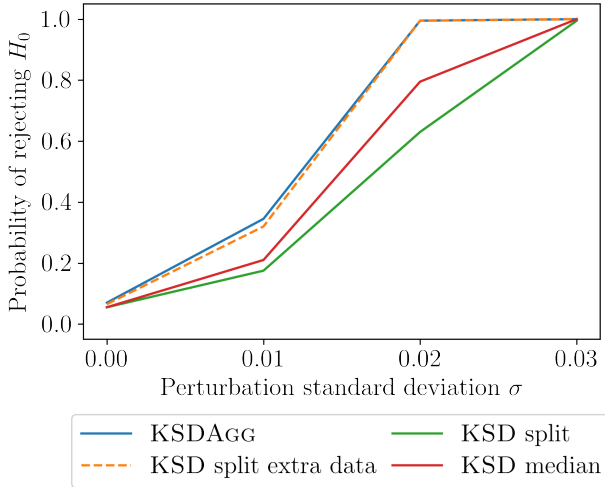


Figure 2. Bernoulli-Gaussian Restricted Boltzmann Machine experiment.

Again, we observe that our aggregated test KSDAGG matches the power obtained by the test which uses extra data to select an ‘optimal’ bandwidth. This means that, in this experiment, KSDAGG obtains the same power as the ‘best’ single test. The difference between KSDAGG and the median heuristic test is now significant, and the splitting test obtains lowest power of the four tests. Again, all tests have well-calibrated levels ( $\sigma = 0$ ) and increasing the noise level  $\sigma$  results in more power for all the tests.



Figure 3. MNIST and Normalizing Flow digits. The top row corresponds to digits from the MNIST dataset. The digits in the middle and bottom rows have been sampled from the Normalizing Flow.

#### 4.5. MNIST Normalizing Flow

Finally, we consider a high-dimensional problem working with images in dimensions  $28^2 = 784$ . We consider a multi-scale Normalizing Flow (Dinh et al., 2017; Kingma & Dhariwal, 2018) which has been trained on the MNIST dataset (LeCun et al., 1998; 2010), it is a generative model which has a probability density  $p$ . As observed in Figure 3, some samples produced by the model can look exactly like MNIST digits, while other do not resemble digits. This Normalizing Flow has been trained to ‘ideally’ produce samples from the MNIST dataset. We are interested in whether or not we can detect the difference in densities. Given some images of digits, are we able to tell if those were sampled from the Normalizing Flow model? We consider the case where the images from  $q$  are sampled from the model (level experiment, confirming performance under  $\mathcal{H}_0$ ), and the case where the samples from  $q$  are drawn from the true MNIST dataset (power experiment). The experiments are run with the collection of bandwidths  $\Lambda(-20, 0)$ . The results are displayed in Table 1 and Figure 4.

Table 1. Normalizing Flow MNIST level experiment.

SAMPLE SIZE	KSDAGG	KSD MEDIAN	KSD SPLIT	KSD SPLIT EXTRA DATA
100	0.075	0.05	0.05	0.05
200	0.05	0.05	0.04	0.05
300	0.055	0.09	0.04	0.09
400	0.05	0.06	0.04	0.06
500	0.04	0.055	0.08	0.055

In Table 1, we observe that the four tests have correct level 0.05 for the five different sample sizes considered (the small fluctuations about the designed test level can be explained by the fact that we are averaging 200 test outputs to estimate these levels). The well-calibrated levels obtained in Table 1 demonstrate the validity of the power results presented in Figure 4.

In Figure 4, we observe that only our aggregated test KSDAGG obtains high power, that is, is able to detect that MNIST samples are not drawn from the Normalizing Flow. The power of the other tests increases only marginally as the sample size increases. We notice that the test which uses extra data to select an ‘optimal’ bandwidth performs poorly when compared to KSDAGG. This could be explained by the fact that this test uses a proxy for the asymptotic power for bandwidth selection, and that in this high-dimensional setting, the asymptotic regime is not attained with sample sizes below 500.

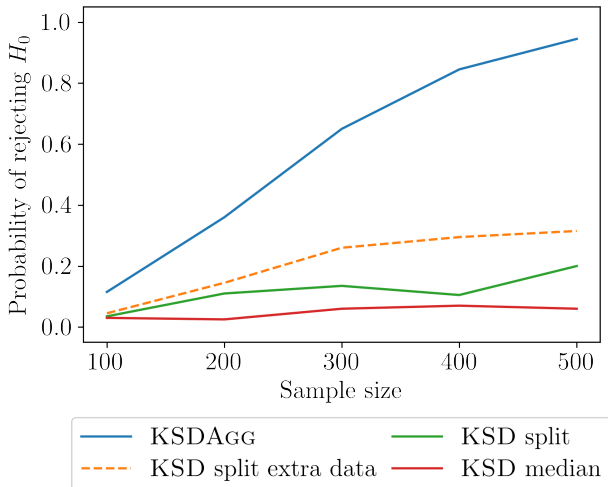


Figure 4. Normalizing Flow MNIST power experiment.

## 5. Discussion

We have introduced KSDAGG, an aggregated goodness-of-fit test based on the Kernel Stein Discrepancy. We have investigated the theoretical properties of this adaptive test. We have shown that it achieves the desired level and have provided conditions to guarantee high power by exhibiting a lower bound on its uniform separation rate. We have observed in our experiments that KSDAGG outperforms alternative state-of-the-art approaches to KSD kernel adaptation for goodness-of-fit testing.

This work covers the problem of KSD adaptivity in the goodness-of-fit framework, this complements the works of Fromont et al. (2013), Kim et al. (2020) and Schrab et al.

(2021) on MMD adaptivity for the two-sample problem and of Albert et al. (2019) and Kim et al. (2020) for HSIC adaptivity for independence testing. A potential future direction of interest could be to tackle the adaptivity problem of the KSD-based linear-time goodness-of-fit test proposed by Jitkrittum et al. (2017). In this setting, the data is split to select feature locations (and the kernel bandwidth), and the KSD test is then run using those adaptive features in data space or in the Fourier domain. A challenging problem would be to obtain those adaptive features using an aggregation procedure which avoids splitting the data.

## References

- Albert, M., Laurent, B., Marrel, A., and Meynaoui, A. Adaptive test of independence based on HSIC measures. *To appear in The Annals of Statistics*, *arXiv preprint arXiv:1902.06441*, 2019.
- Baraud, Y. Non-asymptotic minimax rates of testing in signal detection. *Bernoulli*, 1(8(5):577–606), 2002.
- Carmeli, C., De Vito, E., Toigo, A., and Umanit , V. Vector valued reproducing kernel Hilbert spaces and universality. *Analysis and Applications*, 8(01):19–61, 2010.
- Chwialkowski, K., Sejdinovic, D., and Gretton, A. A wild bootstrap for degenerate kernel tests. In *Advances in neural information processing systems*, pp. 3608–3616, 2014.
- Chwialkowski, K., Strathmann, H., and Gretton, A. A kernel test of goodness of fit. In *International Conference on Machine Learning*, pp. 2606–2615. PMLR, 2016.
- Dinh, L., Sohl-Dickstein, J., and Bengio, S. Density estimation using real NVP. In *International Conference on Learning Representations*, 2017.
- Freidling, T., Poignard, B., Climente-Gonz lez, H., and Yamada, M. Post-selection inference with HSIC-Lasso. In *International Conference on Machine Learning*, pp. 3439–3448. PMLR, 2021.
- Fromont, M., Laurent, B., Lerasle, M., and Reynaud-Bouret, P. Kernels based tests with non-asymptotic bootstrap approaches for two-sample problems. In *Conference on Learning Theory*, PMLR, 2012.
- Fromont, M., Laurent, B., and Reynaud-Bouret, P. The two-sample problem for Poisson processes: Adaptive tests with a nonasymptotic wild bootstrap approach. *The Annals of Statistics*, 41(3):1431–1461, 2013.
- Gorham, J. and Mackey, L. Measuring sample quality with kernels. In *International Conference on Machine Learning*, pp. 1292–1301. PMLR, 2017.

- Grathwohl, W., Wang, K.-C., Jacobsen, J.-H., Duvenaud, D., and Zemel, R. Learning the Stein discrepancy for training and evaluating energy-based models without sampling. In *International Conference on Machine Learning*, pp. 3732–3747. PMLR, 2020.
- Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. Measuring statistical dependence with Hilbert-Schmidt norms. In *International Conference on Algorithmic Learning Theory*. Springer, 2005.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. A kernel two-sample test. *Journal of Machine Learning Research*, 13:723–773, 2012a.
- Gretton, A., Sejdinovic, D., Strathmann, H., Balakrishnan, S., Pontil, M., Fukumizu, K., and Sriperumbudur, B. K. Optimal kernel choice for large-scale two-sample tests. In *Advances in Neural Information Processing Systems*, volume 1, pp. 1205–1213, 2012b.
- Hoeffding, W. A class of statistics with asymptotically normal distribution. In *Breakthroughs in Statistics*, pp. 308–334. Springer, 1992.
- Jitkrittum, W., Xu, W., Szabó, Z., Fukumizu, K., and Gretton, A. A linear-time kernel goodness-of-fit test. In *Advances in Neural Information Processing Systems*, pp. 262–271, 2017.
- Key, O., Fernandez, T., Gretton, A., and Briol, F.-X. Composite goodness-of-fit tests with kernels. *arXiv preprint arXiv:2111.10275*, 2021.
- Kim, I., Balakrishnan, S., and Wasserman, L. Minimax optimality of permutation tests. *To appear in The Annals of Statistics*, *arXiv preprint arXiv:2003.13208*, 2020.
- Kingma, D. P. and Dhariwal, P. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pp. 10236–10245, 2018.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- LeCun, Y., Cortes, C., and Burges, C. MNIST handwritten digit database. AT&T Labs, 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Lee, J. *U-statistics: Theory and Practice*. Citeseer, 1990.
- Leucht, A. and Neumann, M. H. Dependent wild bootstrap for degenerate U- and V-statistics. *Journal of Multivariate Analysis*, 117:257–280, 2013.
- Lim, J. N., Yamada, M., Schölkopf, B., and Jitkrittum, W. Kernel stein tests for multiple model comparison. In *Advances in Neural Information Processing Systems*, pp. 2240–2250, 2019.
- Lim, J. N., Yamada, M., Jitkrittum, W., Terada, Y., Matsui, S., and Shimodaira, H. More powerful selective kernel tests for feature selection. In *International Conference on Artificial Intelligence and Statistics*, pp. 820–830. PMLR, 2020.
- Liu, F., Xu, W., Lu, J., Zhang, G., Gretton, A., and Sutherland, D. J. Learning deep kernels for non-parametric two-sample tests. In *International Conference on Machine Learning*, 2020.
- Liu, Q., Lee, J., and Jordan, M. A kernelized Stein discrepancy for goodness-of-fit tests. In *International Conference on Machine Learning*, pp. 276–284. PMLR, 2016.
- Romano, J. P. and Wolf, M. Exact and approximate step-down methods for multiple hypothesis testing. *Journal of the American Statistical Association*, 100(469):94–108, 2005.
- Schrab, A., Kim, I., Albert, M., Laurent, B., Guedj, B., and Gretton, A. MMD aggregated two-sample test. *arXiv preprint arXiv:2110.15073*, 2021.
- Shao, X. The dependent wild bootstrap. *Journal of the American Statistical Association*, 105(489):218–235, 2010.
- Stute, W., Manteiga, W. G., and Quindimil, M. P. Bootstrap based goodness-of-fit-tests. *Metrika*, 40(1):243–256, 1993.
- Sutherland, D. J., Tung, H.-Y., Strathmann, H., De, S., Ramdas, A., Smola, A., and Gretton, A. Generative models and model criticism via optimized maximum mean discrepancy. In *International Conference on Learning Representations*, 2017.
- Yamada, M., Wu, D., Tsai, Y. H., Ohta, H., Salakhutdinov, R., Takeuchi, I., and Fukumizu, K. Post selection inference with incomplete maximum mean discrepancy estimator. In *International Conference on Learning Representations*, 2019.