

# A kernel Stein test for comparing latent variable models

Heishiro Kanagawa<sup>1</sup> , Wittawat Jitkrittum<sup>2,3</sup>, Lester Mackey<sup>4</sup>, Kenji Fukumizu<sup>5</sup> and Arthur Gretton<sup>1</sup>

<sup>1</sup>Gatsby Computational Neuroscience Unit, University College London, London, UK

<sup>2</sup>Empirical Inference, Max Planck Institute for Intelligent Systems, Tübingen, Germany

<sup>3</sup>Google Research, New York, NY, USA

<sup>4</sup>Microsoft Research New England, Cambridge, MA, USA

<sup>5</sup>The Institute of Statistical Mathematics, Tachikawa, Tokyo, Japan

Address for correspondence: Heishiro Kanagawa, Gatsby Computational Neuroscience Unit, 46 Cleveland Street, London W1T 4JG, UK. Email: [heishiro.kanagawa@gmail.com](mailto:heishiro.kanagawa@gmail.com)

## Abstract

We propose a kernel-based nonparametric test of relative goodness of fit, where the goal is to compare two models, both of which may have unobserved latent variables, such that the marginal distribution of the observed variables is intractable. The proposed test generalizes the recently proposed kernel Stein discrepancy (KSD) tests (Liu et al., *Proceedings of the 33rd international conference on machine learning* (pp. 276–284); Chwialkowski et al., (2016), In *Proceedings of the 33rd international conference on machine learning* (pp. 2606–2615); Yang et al., (2018), In *Proceedings of the 35th international conference on machine learning* (pp. 5561–5570)) to the case of latent variable models, a much more general class than the fully observed models treated previously. The new test, with a properly calibrated threshold, has a well-controlled type-I error. In the case of certain models with low-dimensional latent structures and high-dimensional observations, our test significantly outperforms the relative maximum mean discrepancy test, which is based on samples from the models and does not exploit the latent structure.

**Keywords:** hypothesis testing, kernel methods, mixture models, model selection, Stein's method

## 1 Introduction

A major approach to statistical modeling is the use of variables representing quantities that are unobserved but thought to underlie the observed data: well-known instances include probabilistic PCA (Roweis, 1997; Tipping & Bishop, 1999), factor analysis (see, e.g., Basilevsky, 1994), mixture models (see, e.g., Gilks et al., 1995), topic models for text (Blei et al., 2003), and hidden Markov models (HMMs) (Rabiner, 1989). The hidden structure in these generative models serves multiple purposes: it allows interpretability and understanding of model features [e.g., the topic proportions in a latent Dirichlet allocation (LDA) model of text], and it facilitates modeling by leveraging simple low-dimensional dynamics of phenomena observed in high dimensions (e.g., HMMs with a low-dimensional hidden state). Statistical modelers ultimately use such models to reason about the data; thus, in order to guarantee the validity of the inference, tools for comparing models and evaluating model fit are required.

This article addresses the problem of evaluating and comparing generative probabilistic models, in cases where the models have a latent variable structure, and the marginals over the observed data are intractable. In this scenario, one strategy for evaluating a generative model is to draw samples from it and to compare these samples to the modeled data using a two-sample test: for

Received: June 25, 2021. Revised: April 4, 2023. Accepted: April 6, 2023

© (RSS) Royal Statistical Society 2023.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

instance, [Lloyd and Ghahramani \(2015\)](#) use a test based on the maximum mean discrepancy (MMD) ([Gretton et al., 2012](#)). This approach has two disadvantages, however: it is not computationally efficient due to the sampling step, and it does not take advantage of the information that the model supplies, for instance, the dependence relations among the variables.

Recently, an alternative model evaluation strategy based on Stein's method ([Barbour, 1988](#); [L. H. Y. Chen, 1975](#); [Götze, 1991](#); [Stein, 1972, 1986](#)) has been proposed, which directly employs a closed-form expression for the unnormalized model. Stein's method is a technique from probability theory developed to prove central limit theorems with explicit rates of convergence (see, e.g., [Ross, 2011](#)). The core of Stein's method is that it characterizes a distribution with a *Stein operator*, which, when applied to a function, causes the expectation of the function to be zero under the distribution. For our purposes, we will use the result that a model-specific Stein operator may be defined, to construct a measure of the model's discrepancy. Notably, Stein operators may be obtained without computing the normalizing constant.

Stein operators have been used to design integral probability metrics (IPMs) ([Müller, 1997](#)) to test the goodness of fit of models. IPMs specify a *witness function* which has a large difference in expectation under the sample and model, thereby revealing the difference between the two. When a Stein operator is applied to the IPM function class, the expectation under the model is zero, leaving only the expectation under the sample. A Stein-modified  $W^{2,\infty}$  Sobolev ball was used as the witness function class in [Gorham and Mackey \(2015\)](#) and [Gorham et al. \(2019\)](#). Subsequent work in [Chwialkowski et al. \(2016\)](#), [Liu et al. \(2016\)](#), and [Gorham and Mackey \(2017\)](#) used as the witness function class a Stein-transformed reproducing kernel Hilbert ball, as introduced by [Oates et al. \(2017\)](#): the resulting goodness-of-fit statistic is known as the kernel Stein discrepancy (KSD). Conditions for using the KSD in convergence detection were obtained by [Gorham and Mackey \(2017\)](#). While the foregoing work applies in continuous domains, the approach may also be used for models on a finite domain, where Stein operators ([Bresler & Nagaraj, 2019](#); [Hodgkinson et al., 2020](#); [Ranganath et al., 2016](#); [Reinert & Ross, 2019](#); [Shi et al., 2022](#); [Yang et al., 2018](#)) and associated goodness-of-fit tests ([Yang et al., 2018](#)) have been established. Note that it is also possible to use Stein operators to construct feature dictionaries for comparing models, rather than using an IPM: examples include a test based on Stein features constructed in the sample space so as to maximize test power ([Jitkrittum et al., 2018, 2017](#)) and a test based on Stein-transformed random features ([Huggins & Mackey, 2018](#)). While the aforementioned tests address simple hypotheses, composite tests that use Stein characterizations have been proposed for specific parametric families including gamma ([Betsch & Ebner, 2019b](#); [Henze et al., 2012](#)) and normal distributions ([Betsch & Ebner, 2019c](#); [Henze & Visagie, 2019](#)), and general univariate parametric families ([Betsch & Ebner, 2019a](#)) (note that these tests are not based on IPMs).

While testing goodness of fit alone may be desirable for models of simple phenomena, it will often be the case that in complex domains, no model will fit the data perfectly. In this setting, it is more constructive to ask which model fits better, either within a class of models or in comparing different model classes. A likelihood-ratio test would be an ideal choice for this task, since it is uniformly most powerful ([Lehmann & Romano, 2005](#)). If the models contain latent variables, however, a likelihood-ratio test requires evaluating marginal densities of the models, which are typically intractable. A number of Monte Carlo techniques have been developed to estimate marginal densities or log-density ratios (see, e.g., [Friel & Wyse, 2012](#), for a review). Constructing a test with such techniques is challenging, however; e.g., estimating each marginal density induces a bias in the likelihood ratio, which is difficult to characterize when designing a calibrated threshold (see Section 3.5 for a detailed discussion). Addressing the intractability of the likelihood, [Bounliphone et al. \(2016\)](#) proposed a purely sample-based relative goodness-of-fit test, which compares maximum mean discrepancies between the samples from two rival models with a reference real-world sample. A second relative test was proposed by [Jitkrittum et al. \(2018\)](#), generalizing ([Jitkrittum et al., 2017](#)) and learning the Stein features for which each model outperforms the other; this test requires marginal densities up to normalizing constants and does not apply to latent variable models.

In the present work, we introduce a novel relative goodness-of-fit test for latent variable models (LVMs), which compares models by computing approximate kernel Stein discrepancies. Our contribution is to provide a frequentist test of relative goodness of fit, with an approximate U-statistic

of the kernel Stein discrepancy difference as our test statistic. The statistic is expressed by a posterior expectation of the latent given an observation and is amenable to standard Markov Chain Monte Carlo techniques: in particular, it does not suffer from the challenges in characterizing bias observed in likelihood-ratio estimates for LVMs. Note that our approach differs from Bayesian model selection (Jeffreys, 1961; Kass & Raftery, 1995; Schwarz, 1978; Watanabe, 2013), which reports posterior odds (or Bayes factors) and does not concern controlling frequentist risks such as type-I error rates. To the best of our knowledge, our test represents the first general-purpose, frequentist, relative test for LVMs.

We recall the Stein operator and kernel Stein discrepancy in Section 2, and the notion of relative tests in Section 3. Our main theoretical contributions, also in Section 3, are twofold: first, we derive an appropriate test threshold to account for the randomness in the test statistic caused by sampling the latent variables. Second, we provide guarantees that the resulting test has the correct type-I level (i.e., that the rate of false positives is properly controlled) and that the test is consistent under the alternative: the number of false negatives drops to zero as we observe more data. Finally, in Section 4, we demonstrate our relative test of goodness of fit on a variety of LVMs. Our main point of comparison is the relative MMD test (Bounliphone et al., 2016), where we sample from each model. We demonstrate that the relative Stein test outperforms the relative MMD test in the particular case where the low-dimensional structure of the latent variables can be exploited.

## 2 The kernel Stein discrepancy and LVMs

In this section, we recall the definition of the Stein operator as used in goodness-of-fit testing, as well as the kernel Stein discrepancy, a measure of goodness of fit based on this operator. We will then introduce LVMs, which will bring us to the setting of relative goodness of fit with competing models in Section 3.

Before proceeding, we call attention to our setting: in this article, we treat both continuous- and discrete-valued observations, as formally defined at the outset of Section 2.1. It is our intention to study these two data modalities as they admit the same treatment. The subsequent definitions and analysis of our test are independent of whether a continuous or discrete Stein operator is used, apart from experiments concerning discrete-valued observations. Thus, the detail about discrete models in Section 2.1 may be initially skipped if desired.

### 2.1 Stein operators and kernel Stein discrepancies

Let  $\mathcal{X}$  be the space in which the data take values; for  $D \geq 1$ , the space  $\mathcal{X}$  is either the Euclidean space  $\mathbb{R}^D$  or a finite lattice  $\{0, \dots, L - 1\}^D$  for some  $L > 1$ . Depending on  $\mathcal{X}$ , we shall assume that the densities below are all defined with respect to the Lebesgue measure or the counting measure; i.e., the term *density* includes probability mass functions (pmfs).

**Continuous-valued observations.** Suppose that we are given data  $\{x_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} R$  from an unknown distribution  $R$ , and we wish to test the goodness of fit of a model  $P$ . We first consider the case where the probability distributions  $P, R$  are defined on  $\mathbb{R}^D$  and have respective probability densities  $p, r$ , where all density functions considered in this paper are assumed strictly positive and continuously differentiable. We treat the case of densities defined on bounded domains in the [Online Supplementary Material, Section B.1](#). For differentiable density functions, we define the *score function*,

$$s_p(x) \in \mathbb{R}^D := \frac{\nabla p(x)}{p(x)} = \nabla \log p(x),$$

where the gradient operator is  $\nabla := [\frac{\partial}{\partial x^1}, \dots, \frac{\partial}{\partial x^D}]^\top$ . The score is independent of the normalizing constant for  $p$ , making it computable even when  $p$  is known only up to normalization. Using this score, we define the *Langevin Stein operator* on a space  $\mathcal{F}$  of differentiable functions from  $\mathbb{R}^D$  to  $\mathbb{R}^D$  (Gorham & Mackey, 2015; Oates et al., 2017),

$$[\mathcal{A}_p f](x) = \langle s_p(x), f(x) \rangle + \langle \nabla, f(x) \rangle, \quad f \in \mathcal{F}.$$

A kernel discrepancy may be defined based on the Stein operator (Chwialkowski et al., 2016; Gorham & Mackey, 2017; Liu et al., 2016), which allows us to measure the departure of a distribution  $R$  from a model  $P$ . We define  $\mathcal{F}$  to be a space comprised of  $D$ -dimensional vectors of functions  $f = (f_1, \dots, f_D)$  where the  $d$ th function  $f_d$  is in a reproducing kernel Hilbert space (RKHS) (Aronszajn, 1950; Steinwart & Christmann, 2008, Definition 4.18) with a positive definite kernel  $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  (we use the same kernel for each dimension). The inner product on  $\mathcal{F}$  is  $\langle f, g \rangle_{\mathcal{F}} := \sum_{d=1}^D \langle f_d, g_d \rangle_{\mathcal{F}_k}$ , and  $\mathcal{F}_k$  denotes an RKHS of real-valued functions with kernel  $k$ .

The (Langevin) *kernel Stein discrepancy* (KSD) between  $P$  and  $R$  is defined as

$$\text{KSD}(P \| R) = \sup_{\|f\|_{\mathcal{F}} \leq 1} |\mathbb{E}_{x \sim R} \mathcal{A}_P f(x) - \mathbb{E}_{y \sim P} \mathcal{A}_P f(y)|. \quad (1)$$

Under appropriate conditions on the kernel and measure  $P$ , the expectation  $\mathbb{E}_{y \sim P} \mathcal{A}_P f(y) = 0$  for any  $f \in \mathcal{F}$ . To ensure this property, we will require that  $k \in C^{(1,1)}$ , the set of continuous functions on  $\mathcal{X} \times \mathcal{X}$  with continuous first derivatives and that  $\mathbb{E}_{y \sim P} [\|s_p(y)\|_2] < \infty$  with  $\|\cdot\|_2$  the Euclidean norm. We further assume that the following tail condition holds outside a bounded set:  $p(x) \sqrt{k(x, x)} \leq C \|x\|_2^\delta$  for some constants  $C > 0$  and  $\delta > D - 1$  (see the clarification by South et al., 2021, p. 12, on the tail condition for the Stein's identity). With the vanishing expectation  $\mathbb{E}_{y \sim P} \mathcal{A}_P f(y) = 0$ , the KSD reduces to  $\text{KSD}(P \| R) = \sup_{\|f\|_{\mathcal{F}} \leq 1} |\mathbb{E}_{x \sim R} \mathcal{A}_P f(x)|$ . The use of an RKHS as the function class yields a closed-form expression of the discrepancy by the kernel trick (Chwialkowski et al., 2016; Gorham & Mackey, 2017, Proposition 2),

$$\text{KSD}^2(P \| R) = \mathbb{E}_{x, x' \sim R \otimes R} [h_p(x, x')],$$

if  $\mathbb{E}_{x \sim R} [h_p(x, x)^{1/2}] < \infty$ . Here, the symbol  $R \otimes R$  denotes the product measure of two copies of  $R$  (so  $x$  and  $x'$  are independent random variables identically distributed with the law  $R$ ). The function  $h_p$  (called a *Stein kernel*) is expressed in terms of the RKHS kernel  $k$  and the score function  $s_p$ ,

$$h_p(x, x') = s_p(x)^\top s_p(x') k(x, x') + s_p(x)^\top k_1(x', x) + s_p(x')^\top k_1(x, x') + k_{12}(x, x'),$$

where we have defined

$$k_1(a, b) := \nabla_x k(x, x')|_{x=a, x'=b},$$

$$k_{12}(a, b) := \nabla_x^\top \nabla_{x'} k(x, x')|_{x=a, x'=b}.$$

For a given i.i.d. sample  $\{x_i\}_{i=1}^n \sim R$ , the discrepancy has a simple closed-form finite sample estimate,

$$\text{KSD}^2(P \| R) \approx \frac{1}{n(n-1)} \sum_{i \neq j} h_p(x_i, x_j), \quad (2)$$

which is a U-statistic (Hoeffding, 1948). When the kernel is integrally strictly positive definite (ISPD) (Sriperumbudur et al., 2011, Section 2), and  $R$  admits a density  $r$  that satisfies  $\mathbb{E}_{x \sim R} \|\nabla \log(p(x)/r(x))\|_2 < \infty$ , we have that  $\text{KSD}(P \| R) = 0$  iff  $P = R$  (Barp et al., 2019, Proposition 1). The earlier results of Chwialkowski et al. (2016) and Liu et al. (2016) require more stringent integrability conditions. Gorham and Mackey (2017, Theorem 7) showed that KSD can distinguish any Borel measure  $R$  from  $P$  by assuming conditions such as distant dissipativity (satisfied by finite Gaussian mixtures) (Gorham et al., 2019, Section 3). However, such conditions may be difficult to validate for LVMS. Thus, hereafter, we assume the former condition on the data distribution  $R$ .

**Discrete-valued observations.** We next recall the kernel Stein discrepancy in the discrete setting where distributions are defined on  $\mathcal{X} = \{0, \dots, L-1\}^D$  with  $L > 1$ , as introduced by Yang et al. (2018). In place of derivatives, we specify  $\Delta_k$  as the cyclic forward difference w.r.t.  $k$ th coordinate:  $\Delta_k f(x) = f(x^1, \dots, \tilde{x}^k, \dots, x^D) - f(x^1, \dots, x^k, \dots, x^D)$  where  $\tilde{x}^k = x^k + 1 \bmod L$ , with the

corresponding vector-valued operator  $\Delta = (\Delta_1, \dots, \Delta_D)$ . The inverse operator  $\Delta_k^{-1}$  is given by the backward difference  $\Delta_k^{-1}f(x) = f(x^1, \dots, x^k, \dots, x^D) - f(x^1, \dots, \bar{x}^k, \dots, x^D)$ , where  $\bar{x}^k = x^k - 1 \bmod L$ , and  $\Delta^{-1} = (\Delta_1^{-1}, \dots, \Delta_D^{-1})$ . The score is then  $s_p(x) := p(x)^{-1} \Delta p(x)$ , where it is assumed that the pmf is strictly positive (i.e., it is never zero). The difference Stein operator is then defined as  $\mathcal{A}_p f(x) = \text{tr}[f(x) s_p(x)^\top + \Delta^{-1} f(x)]$ , where it can be shown that  $\mathbb{E}_{x \sim P}[\mathcal{A}_p f(x)] = 0$  (Yang et al., 2018, Theorem 2) (note that we include a trace for consistency with the continuous case—this does not affect the test statistic (Yang et al., 2018, Equation 10)). We have defined the Stein operator and the score function slightly differently from Yang et al. (2018); the change is only in their signs, but this results in the same discrepancy. The difference Stein operator is not the only allowable Stein operator on discrete spaces: other alternatives are given by Yang et al. (2018, Theorem 3), Hodgkinson et al. (2020), and Shi et al. (2022). Although we focus on the Stein operator above, in practice, one might want to consider different Stein operators depending on the application. For instance, the score function  $s_p$  can be numerically unstable, as it contains the reciprocal  $1/p(x)$ ; this can occur when the support of the model is severely mismatched to that of the data. In this particular case, one might choose the Barker–Stein operator proposed by Shi et al. (2022), an instance of the Zanella–Stein operator of Hodgkinson et al. (2020, Example 2). See Online Supplementary Material, Appendix B.2 for details. We compare this operator to the difference operator in an experiment where this mismatch occurs (Section 4.3.3).

As in the continuous case, the KSD can be defined as an IPM, given a suitable choice of reproducing kernel Hilbert space for the discrete domain. An example of kernel is the exponentiated Hamming kernel,  $k(x, x') = \exp(-d_H(x, x'))$ , where  $d_H(x, x') = D^{-1} \sum_{d=1}^D \mathbb{1}(x^d \neq x'^d)$ . The population KSD is again given by the expectation of the Stein kernel,  $\text{KSD}^2(P \| R) = \mathbb{E}_{(x, x') \sim R \otimes R} [h_p(x, x')]$ , where  $h_p$  is defined as

$$h_p(x, x') = s_p(x)^\top s_p(x') k(x, x') + s_p(x)^\top k_1(x', x) + s_p(x')^\top k_1(x, x') + k_{12}(x, x'),$$

and the kernel gradient is replaced by the inverse difference operator, e.g.,  $k_1(x, x') = \Delta_x^{-1} k(x, x')$ , where  $\Delta_x^{-1}$  indicates that the operator  $\Delta^{-1}$  is applied with respect to the argument  $x$ . From Yang et al. (2018, Lemma 8), we have that  $\text{KSD}(P \| R) = 0$  iff  $P = R$ , under the conditions that the probability mass functions for  $P$  and  $R$  are positive and that the Gram matrix defined over all the configurations in  $\mathcal{X}$  is strictly positive definite (i.e., the kernel is integrally strictly positive definite). One can define a kernel satisfying the required condition, for example, by embedding  $\mathcal{X}$  into  $\mathbb{R}^{L \times D}$  with one-hot encoding and using a Taylor-type kernel such as the exponentiated quadratic kernel (Christmann & Steinwart, 2010, Theorem 2.2).

## 2.2 Kernel Stein discrepancies of LVMs

Our objective is to use the KSD to evaluate LVMs, and here we formally specify our target model class. Let  $\mathcal{L}_{\mathcal{X} | \mathcal{Z}} = \{p(\cdot | z) : z \in \mathcal{Z}\}$  be a family of probability density functions on  $\mathcal{X}$  (we call these *likelihood* functions), which are indexed by elements of a set  $\mathcal{Z}$ . A LVM  $P$  is specified by such a family  $\mathcal{L}_{\mathcal{X} | \mathcal{Z}}$  and a (prior) probability measure  $P_{\mathcal{Z}}$  over  $\mathcal{Z}$ . The combination of these defines the marginal density function  $p(x) = \int p(x | z) dP_{\mathcal{Z}}(z)$  and the posterior distribution  $P_{\mathcal{Z}}(dz | x) = \{p(x | z) / p(x)\} P_{\mathcal{Z}}(dz)$ . The distribution  $P$  induced by the former acts as a model of the distribution  $R$  underlying the observation, and the latter enables us to draw an inference over the unobserved variable.

**Remark** In our notation, the variable  $z$  can represent multiple latent variables. The likelihood  $p(x | z)$  often contains parameters, but the dependency on these is suppressed here. If a prior is defined on a parameter, we may treat it as a latent variable; this consideration is relevant to predictive distributions. The likelihood and the prior in a model may be conditioned on some fixed data (e.g., they can be posterior predictive distributions), which we require to be independent of the data used for testing—in such a case, we omit the dependency on the held-out data. For examples, we refer the reader to Section 4.

The definition of the KSD remains the same for LVMs, but an additional difficulty arises in its estimation. Unfortunately, the U-statistic estimator given in (2) requires the score function of the marginal  $p$ , which is challenging to obtain due to the intractability of marginalizing out the latent variable. We will address this challenge by rewriting the score function in terms of the posterior distribution of the latent. In the following, we focus on continuous variable models, but the same conclusion holds for discrete counterparts by replacing gradient operation with cyclic differences.

Under a regularity condition, the score function can be expressed as

$$\mathbf{s}_p(x) = \mathbb{E}_{z|x}[\mathbf{s}_p(x|z)], \quad (3)$$

where  $\mathbf{s}_p(x|z)$  is the score function of the conditional  $p(x|z)$ ; i.e.,  $\mathbf{s}_p(x|z) = p(x|z)^{-1} \nabla_x p(x|z)$  for continuous-valued  $x$ . The reasoning is as follows:

$$\begin{aligned} \frac{\nabla_x p(x)}{p(x)} &= \frac{1}{p(x)} \int \nabla_x p(x|z) dP_Z(z) \\ &= \int \frac{\nabla_x p(x|z)}{p(x|z)} \cdot \frac{p(x|z) dP_Z(z)}{p(x)} = \mathbb{E}_{z|x}[\mathbf{s}_p(x|z)], \end{aligned}$$

where we have assumed the exchangeability of differentiation and integration:  $\nabla_x p(x) = \int \nabla_x p(x|z) dP_Z(z)$ . The identity (3) is an analogue of Fisher's identity (Dempster et al., 1977; Fisher, 1925), which pertinently formed the basis for Stein control variate methodology in Friel et al. (2016), parameter inference for doubly intractable models via score matching (Vértes & Sahani, 2016), and Bayesian model selection with a Hyvärinen score (Dawid & Musio, 2015; Shao et al., 2019). Note that the conditional score  $\mathbf{s}_p(x|z)$  is typically possible to evaluate. For example, consider the following simple form of an exponential family density  $p(x|z) \propto \exp(T(x)\eta(z))$  defined on  $\mathbb{R}^D$  with  $T(x): \mathbb{R}^D \rightarrow \mathbb{R}$  and  $\eta: \mathcal{Z} \rightarrow \mathbb{R}$ ; for this density,  $\mathbf{s}_p(x|z) = \eta(z) \nabla_x T(x)$ . As can be seen in this example, the formula (3) does not require the likelihood  $p(x|z)$  to be normalized. This feature eliminates the need for estimating the normalizing constant of  $p(x|z)$  for each  $z$ , which is required to compute goodness-of-fit measures based on the marginal density  $p(x)$  (Friel & Wyse, 2012); Online Supplementary Material, Section C.2 in the supplementary presents a use case with a truncated model.

With this identity, the KSD is rewritten as follows.

**Lemma 1** Let

$$\begin{aligned} H_p[(x, z), (x', z')] &= \mathbf{s}_p(x|z)^\top \mathbf{s}_p(x'|z') k(x, x') + \mathbf{s}_p(x|z)^\top k_1(x', x) \\ &\quad + k_1(x, x')^\top \mathbf{s}_p(x'|z') + k_{12}(x, x'). \end{aligned} \quad (4)$$

Assume  $\mathbb{E}_{(x,z),(x',z') \sim \tilde{R} \otimes \tilde{R}} |H_p[(x, z), (x', z')]| < \infty$  with the joint distribution  $\tilde{R}(d(x, z)) = P_Z(dz|x)R(dx)$ . If the formula (3) holds, then,

$$\text{KSD}^2(P\|R) = \mathbb{E}_{(x,z),(x',z') \sim \tilde{R} \otimes \tilde{R}} H_p[(x, z), (x', z')].$$

**Proof.** Substituting the formula (3) in the definition of KSD gives the required equation by the Tonelli-Fubini theorem.  $\square$

**Remark** The integrability assumption holds trivially if the input space  $\mathcal{X}$  is finite, while care needs to be taken otherwise. The condition can be checked by examining the absolute integrability of each term in (4). The integrability assumption on the fourth term is mild and is satisfied by common kernels, e.g., the exponentiated quadratic or the inverse multi-quadratic kernels. The condition on the other

terms needs to be checked on a model-by-model basis. It can be shown that the example models in Section 4 satisfy the assumption (please see [Online Supplementary Material, Section B.3](#) in the supplementary material for details).

The new KSD expression is an expectation of a computable symmetric kernel, and constructing an unbiased estimate is straightforward once we obtain a sample. In practice, when the model is complex, sampling from the posterior distribution generally requires simulation, as the posterior is not available in a closed form. Therefore, we propose to approximate the expectation by Markov Chain Monte Carlo (MCMC) methods and construct an approximate U-statistic estimator as follows. Let  $\mathbf{z}_i^{(t)} = (z_{i,1}^{(t)}, \dots, z_{i,m}^{(t)}) \in \mathcal{Z}^m$  be a latent sample of size  $m$  drawn by an MCMC method having  $P_Z(\cdot | x_i)$  as its invariant measure after  $t$  burn-in iterations. Let  $\bar{s}_p(x_i | \mathbf{z}_i^{(t)}) = \frac{1}{m} \sum_{j=1}^m s_p(x_i | z_{i,j}^{(t)})$ . Given a joint sample  $\{(x_i, \mathbf{z}_i^{(t)})\}_{i=1}^n$ , we estimate the KSD by

$$U_n^{(t)}(P) := \frac{1}{n(n-1)} \sum_{i \neq j} \bar{H}_p[(x_i, \mathbf{z}_i^{(t)}), (x_j, \mathbf{z}_j^{(t)})], \tag{5}$$

where

$$\begin{aligned} \bar{H}_p[(x_i, \mathbf{z}_i^{(t)}), (x_j, \mathbf{z}_j^{(t)})] &= \bar{s}_p(x_i | \mathbf{z}_i^{(t)})^\top \bar{s}_p(x_j | \mathbf{z}_j^{(t)}) k(x_i, x_j) + \bar{s}_p(x_i | \mathbf{z}_i^{(t)})^\top k_1(x_j, x_i) \\ &\quad + k_1(x_i, x_j)^\top \bar{s}_p(x_j | \mathbf{z}_j^{(t)}) + k_{12}(x_i, x_j), \end{aligned}$$

and the sum is taken over all distinct sample pairs. If  $P_Z^{(t)}(d\mathbf{z} | x)$  denotes the distribution of an MCMC sample  $\mathbf{z}^{(t)} = (z_1^{(t)}, \dots, z_m^{(t)})$ , then this estimator is indeed a U-statistic, but its expectation is that of kernel  $\bar{H}_p$  with respect to  $P_Z^{(t)}(d\mathbf{z} | x)R(dx)$  instead of  $P_Z(d\mathbf{z} | x)R(dx)$ . Thus, the estimator is biased against the target estimand, the model's KSD, for a finite burn-in period  $t$ , and can therefore be seen an approximation to the *true* U-statistic  $U_n^{(\infty)}$ . Designing a statistical test requires understanding the behavior of the statistic (5), and we will provide its analysis in the next section. Although we focus on MCMC for its approximate unbiasedness in our proposed test, different posterior approximations may be considered in other applications; for example, with a more computationally efficient approach (e.g., variational approximation), the new KSD expression in Lemma 1 might allow us to consider parameter estimation for unnormalized statistical models with latent variables (Barp et al., 2019).

### 3 A relative goodness-of-fit test

We now address the setting of statistical testing for model comparison. We begin this section with our problem settings and notation, and then define a test by showing the asymptotic normality of approximate U-statistics.

#### 3.1 Problem setup

We consider the case where we have two LVMs  $P$  and  $Q$ , and we wish to determine which is a closer approximation of the distribution  $R$  generating our data  $\{x_i\}_{i=1}^n$ . The respective density functions of the models are given by the integrals  $p(x) = \int p(x | z) dP_Z(z)$  and  $q(x) = \int q(x | w) dQ_W(w)$ . As with  $P$ , the latent variable  $w$  is assumed to take values in a set  $\mathcal{W}$  with prior  $Q_W$ . We assume that  $p(x)$  and  $q(x)$  cannot be tractably evaluated, even up to their normalizing constants. Our goal is to determine the *relative* goodness of fit of the models by comparing each model's discrepancy from the data distribution. Our problem is formulated as the following hypothesis test:

$$\begin{aligned} H_0 : \text{KSD}(P \| R) &\leq \text{KSD}(Q \| R) \text{ (null hypothesis),} \\ H_1 : \text{KSD}(P \| R) &> \text{KSD}(Q \| R) \text{ (alternative).} \end{aligned} \tag{6}$$

In other words, the null hypothesis is that the fit of  $P$  to  $R$  (in terms of KSD) is as good as  $Q$ , or better. Note that the KSD in (6) is defined by a particular reproducing kernel, and thus different

kernels yield distinct hypotheses. For kernel selection, we refer the reader to Section 3.4.

We next provide an overview of the formal assumptions made throughout the paper. Let  $(\Omega, \mathcal{S}, \Pi)$  be a probability space, where  $\Omega$  is a sample space,  $\mathcal{S}$  is a  $\sigma$ -algebra, and  $\Pi$  is a probability measure. All random variables (for example, data points  $x_i$  and draws  $\mathbf{z}_i^{(t)}$  from a Markov chain sampler) are understood as measurable functions from the sample space  $\Omega$ . The input space  $\mathcal{X}$  is equipped with the Borel  $\sigma$ -algebra generated by its standard topology. We assume that  $\mathcal{Z}, \mathcal{W}$  are Polish spaces with the Borel  $\sigma$ -algebras defined by their respective topologies, on which the priors  $P_{\mathcal{Z}}, Q_{\mathcal{W}}$  are defined. Finally, we require that the two models are distinct; i.e., their marginal densities disagree on a set of positive  $R$ -measure.

### 3.2 Estimating kernel Stein discrepancies of LVMs

The hypotheses in (6) can be equally stated in terms of the difference of the (squared) KSDs,  $\text{KSD}^2(P\|R) - \text{KSD}^2(Q\|R)$ , which motivates us to design a test statistic by estimating each term. Let  $U_n^{(t)}(P, Q) := U_n^{(t)}(P) - U_n^{(t)}(Q)$  be the difference of KSD estimates, where  $U_n^{(t)}(Q)$  is defined as for  $U_n^{(t)}(P)$  in (5). Note that  $U_n^{(t)}(P, Q)$  is an *approximate* U-statistic (in the sense of the final paragraph in Section 2.2) defined by the difference kernel

$$\bar{H}_{p,q}[(x, \mathbf{z}, \mathbf{w}), (x', \mathbf{z}', \mathbf{w}')] := \bar{H}_p[(x, \mathbf{z}), (x', \mathbf{z}')] - \bar{H}_q[(x, \mathbf{w}), (x', \mathbf{w}')]$$

evaluated on the joint sample  $\{(x_i, \mathbf{z}_i^{(t)}, \mathbf{w}_i^{(t)})\}_{i=1}^n$ . The statistic takes as input random variables with evolving laws, and defining a test require us to understand the behavior of such statistics. This section delivers an analysis in a general setting.

We first characterize the asymptotic distribution of an approximate U-statistic. The following theorem shows that such a statistic is asymptotically normal around the expectation of the true U-statistic provided its bias vanishes fast.

**Theorem 1** (Asymptotic normality). Let  $\{\gamma_t\}_{t=1}^\infty$  be a sequence of Borel probability measures on a Polish space  $\mathcal{Y}$  and  $\gamma$  be another Borel probability measure. Let  $\{Y_i^{(t)}\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \gamma_t$ , and for a symmetric function  $h: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ , define a U-statistic and its mean by

$$U_n^{(t)} = \frac{1}{n(n-1)} \sum_{i \neq j} h(Y_i^{(t)}, Y_j^{(t)}), \quad \theta_t = \mathbb{E}_{(Y, Y') \sim \gamma_t \otimes \gamma_t} [h(Y, Y')].$$

Let  $\theta = \mathbb{E}_{(Y, Y') \sim \gamma \otimes \gamma} [h(Y, Y')]$ . Let  $v_t := \mathbb{E}_{(Y, Y') \sim \gamma_t \otimes \gamma_t} [|\tilde{h}_t(Y, Y')|^3]^{1/3}$  with  $\tilde{h}_t = h - \theta_t$ , and assume  $\limsup_{t \rightarrow \infty} v_t < \infty$ . Assume that  $\sigma_t^2 = 4\text{Var}_{Y \sim \gamma_t} [\mathbb{E}_{Y' \sim \gamma_t} [h(Y, Y')]]$  converges to a constant  $\sigma^2$ . Assume that we have  $\theta_t \rightarrow \theta$  as  $t \rightarrow \infty$ . Then, in the limit of large  $n$  and of  $t$  growing as a function of  $n$  such that  $\sqrt{n}(\theta_t - \theta) \rightarrow 0$ , the following two statements hold: if  $\sigma > 0$ , we have

$$\sqrt{n}(U_n^{(t)} - \theta) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$$

where  $\xrightarrow{d}$  denotes convergence in distribution; in the case  $\sigma = 0$ ,  $\sqrt{n}(U_n^{(t)} - \theta) \rightarrow 0$  in probability.

The proof is in [Online Supplementary Material, Section A](#) in the supplement. Note that in the preceding and following results, the limit of  $n$  and  $t$  is taken simultaneously rather than sequentially, such that the condition  $\sqrt{n}(\theta_t - \theta) \rightarrow 0$  holds: see discussion below and in Section 3.3. By letting  $Y_i^{(t)} = (x_i, \mathbf{z}_i^{(t)}, \mathbf{w}_i^{(t)})$  and  $h = \bar{H}_{p,q}$  in the foregoing theorem, we obtain the same conclusion for the difference estimate  $U_n^{(t)}(P, Q)$ .













As PPCA models have tractable marginals, we also compare our test with the KSD test using exact score functions (i.e., no MCMC simulation), which serves as the performance upper-bound. The MCMC sampler we use is HMC; more precisely, we use the NumPyro (Phan et al., 2019) implementation of No-U-Turn Sampler (NUTS) (Hoffman & Gelman, 2014); we take  $t = 200$  burn-in samples and  $m = 500$  consecutive draws for computing a score estimate  $\bar{s}_p$ .

We use two kernel functions: (a) the exponentiated quadratic (EQ) kernel  $k(x, x') = \exp\{-\|x - x'\|_2^2/(2\lambda^2)\}$ , and (b) the IMQ kernel (9) with  $\beta = 0.5$ ,  $c = 1$ , and  $\Lambda = \lambda^2 I$ . All three tests use the same kernel function, which allows us to investigate the effect of using the Stein-modified kernel. The length scale parameter  $\lambda$  is set to the median of the pairwise (Euclidean) distances of holdout samples from  $R$  so that the parameter (and thus the hypothesis) is fixed across trials. We include the EQ kernel in our comparison, as the population MMD is possible to compute, allowing us to verify the hypothesis in advance.

We simulate null and alternative cases by perturbing the weight parameter  $A$ ; we add a positive value  $\delta > 0$  to the  $(1, 1)$ -entry of  $A$ . Let us denote a perturbed weight by  $A_\delta$ . Note that the data PPCA model has a Gaussian marginal  $\mathcal{N}(0, AA^\top + \psi^2 I_x)$ . Therefore, this perturbation gives a model  $\mathcal{N}(0, A_\delta A_\delta^\top + \psi^2 I_x)$ , where the first row and column of  $A_\delta A_\delta^\top$  deviate from those of  $AA^\top$ . The perturbation is additive and increasing in  $\delta$ , as each element of  $A$  is positive. We create a problem by specifying perturbation parameters  $(\delta_P, \delta_Q)$  for  $(P, Q)$ . For the EQ-kernel MMD, we numerically confirmed that the perturbation gives a worse model for a larger perturbation. While the population KSD is not analytically tractable, this perturbation affects the score function through the covariance matrix, and the same behavior is expected for KSD; see [Online Supplementary Material, Section B.6](#) in the supplement for details.

**Problem 1 (null).** We create a null scenario by choosing  $(\delta_P, \delta_Q) = (1, 1 + 10^{-5})$  ( $P$  has a smaller covariance perturbation and is closer to  $R$  than  $Q$ ). For different null settings, we refer the reader to [Online Supplementary Material, Section C.4](#) in the supplement. We run the tests with significance levels  $\alpha = 0.01, 0.05$ . [Table 1](#) reports the finite-sample size of the three tests for significance level  $\alpha = 0.05$ . The result for  $\alpha = 0.01$  is omitted as none of the tests rejected the hypotheses. The size of the proposed LKSD test is indeed controlled. The extremely small type-I errors of the KSD tests are caused by the sensitivity of KSD to this perturbation; the population KSD value is negative and far from zero, and the test statistics easily fall in the acceptance region. The other two tests also have their error rates lower than the significance level. Note that their test thresholds are determined by treating the population discrepancy differences as zero, resulting in conservative tests.

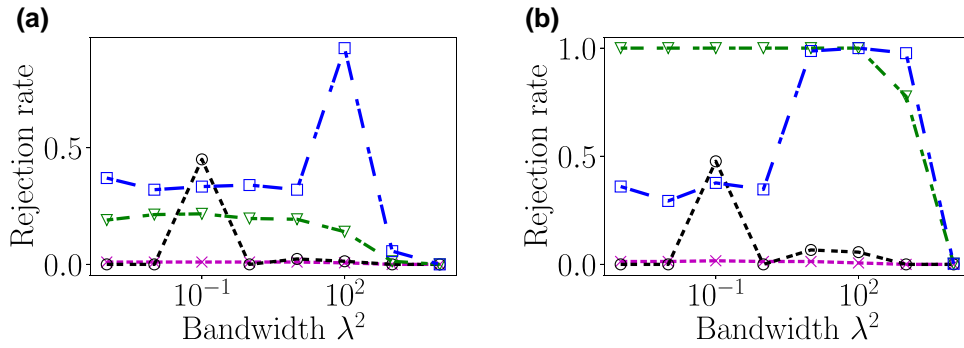
**Problem 2 (alternative).** We investigate the power of the proposed test. We set up an alternative scenario by fixing  $\delta_P = 2$  for  $P$  and  $\delta_Q = 1$  for  $Q$ . For comparison with different parameter settings, please see [Online Supplementary Material, Appendix C.1](#). The significance level  $\alpha$  is fixed at 0.05. All the other parameters are chosen as in Problem 1. [Figure 1](#) shows the plot of the test power against the sample size in each problem. The KSD reaches a near 100 percent rejection rate relatively quickly, indicating that information from the score function is helpful for these problems.

**Table 1.** Type-I errors the MMD test of Bounliphone et al. (2016), the proposed LKSD test, and the KSD test in PPCA problem 1

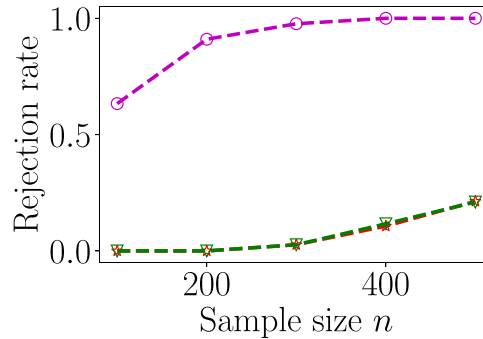
| Sample size $n$ | Rejection rates |       |       |         |       |       |
|-----------------|-----------------|-------|-------|---------|-------|-------|
|                 | EQ-med          |       |       | IMQ-med |       |       |
|                 | MMD             | LKSD  | KSD   | MMD     | LKSD  | KSD   |
| 100             | 0.000           | 0.013 | 0.000 | 0.000   | 0.010 | 0.000 |
| 200             | 0.000           | 0.000 | 0.000 | 0.000   | 0.000 | 0.000 |
| 300             | 0.003           | 0.007 | 0.000 | 0.003   | 0.003 | 0.000 |
| 400             | 0.003           | 0.007 | 0.000 | 0.003   | 0.000 | 0.000 |
| 500             | 0.007           | 0.013 | 0.000 | 0.007   | 0.007 | 0.000 |

*Note.* Rejection rates are computed on 300 trials with significance level  $\alpha = 0.05$ . The columns EQ-med and IMQ-med denote EQ and IMQ kernels with the median bandwidth, respectively.





**Figure 2.** Power curves of the proposed LKSD test and the MMD test in PCCA Problem 2. The perturbation parameters are set as  $(\delta_p, \delta_Q = 2, 1)$ , each result is computed on 300 trials. The significance level  $\alpha = 0.05$ . Markers:  $\nabla$  (LKSD test with IMQ kernel);  $\square$  (LKSD test with EQ kernel);  $\circ$  (MMD test with IMQ kernel);  $\times$  (MMD test with EQ kernel). (a)  $n = 100$  and (b)  $n = 300$ .



**Figure 3.** Power curves of the MMD test, the proposed LKSD test, and the KSD test in PCCA Problem 2. All the test use the covariance-preconditioned IMQ kernel. The perturbation parameters are set as  $(\delta_p, \delta_Q = 2, 1)$ . Each result is computed on 300 trials. The significance level  $\alpha = 0.05$ . Markers:  $\nabla$  (the LKSD test);  $\star$  (the KSD test);  $\circ$  (the relative MMD test).

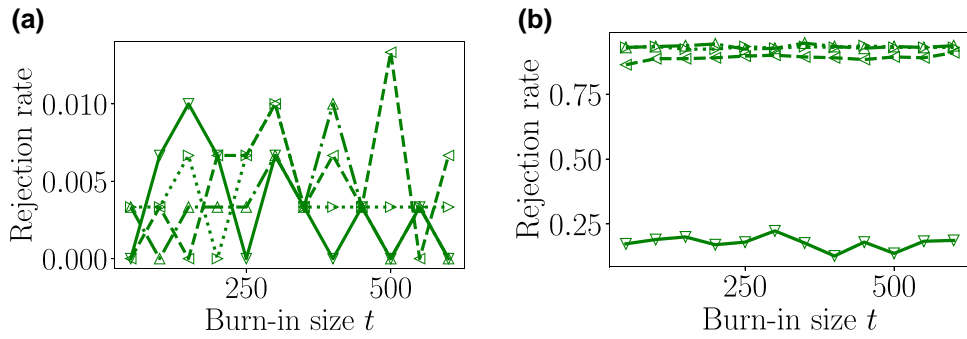
We next consider a slow-converging sampler for which the burn-in length  $t$  becomes crucial. We consider the null case (Problem 1) and replace the sampler for the first model  $P$  with MALA. We set the step size for the MALA sampler to make its convergence slow; we use the step size of  $10^{-4}D_z^{-1/3}$ . We initialize the two samplers differently to make sure that the resulting distributions differ when the samplers have not converged: the MALA sampler for  $P$  is initialized with samples from a Gaussian  $\mathcal{N}\{(1, \dots, 1), I_z\}$  and the NUTS sampler for  $Q$  a uniform distribution  $U[-2, 2]^{D_z}$ . Figure 5 demonstrates the relation between type-I error rates and choices of  $t$  and  $m$ . In contrast to the previous experiment, the burn-in has a clear effect on the type-I error: insufficient burn-in leads to uncontrolled error rates. The right panel ( $n = 300$ ) shows that the test has substantially higher type-I error rates than in the left ( $n = 100$ ). Comparison between these cases illustrates that a larger sample size  $n$  requires more intensive burn-in, as the test becomes more confident to reject. A large value of  $m$  improves the test as in the previous experiment. It can be understood that the contribution of burn-in samples is negligible in the score approximation. Although our analysis in Corollary 1 requires long burn-in, taking large  $m$  appears to be more important in practice, especially under a computational budget constraint. This experiment thus confirms the importance of the quality of the sampler.

## 4.2 Dirichlet process mixtures

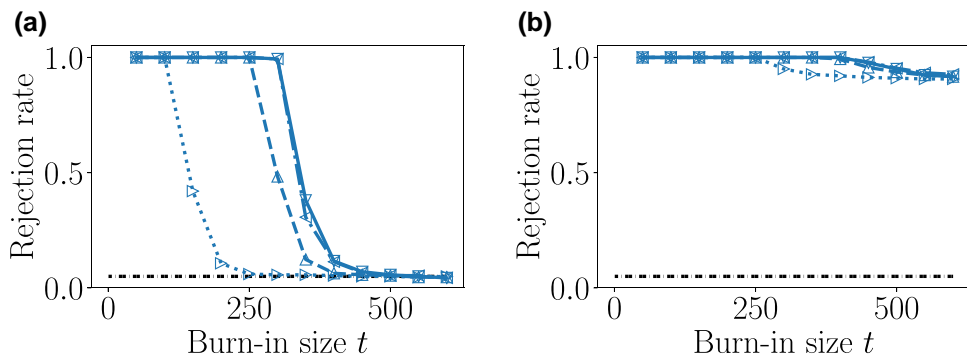
Our next experiment applies our test to a Dirichlet process mixtures (DPM) model. Let  $\psi(x | z)$  be a probability density function on  $\mathbb{R}^D$ . We consider a mixture density

$$\int \psi(x | z) d\rho(z), \quad (10)$$





**Figure 4.** The effect of MCMC quality on the test’s performance. Rejection rates against burn-in size  $t$  with varying Markov chain sample size  $m$ . PPCA Problems 1 and 2 with  $\alpha = 0.05$ . Both samplers use NUTS. Markers:  $\nabla$  ( $m = 1$ );  $\triangleleft$  ( $m = 10$ );  $\triangle$  ( $m = 100$ );  $\triangleright$  ( $m = 1,000$ ). (a) Problem 1 (null  $H_0$  is true) and (b) Problem 2 (alternative  $H_1$  is true).



**Figure 5.** The effect of a poor MCMC sampler on the test. Type-I error rates against the burn-in size  $t$  with varying Markov chain sample size  $m$ . PCCA Problem 1 (the null  $H_0$  is true). The dark dashed line indicates the significance level  $\alpha = 0.05$ . The samplers for  $P$  and  $Q$  are respectively MALA and NUTS. Markers:  $\nabla$  ( $m = 1$ );  $\triangleleft$  ( $m = 10$ );  $\triangle$  ( $m = 100$ );  $\triangleright$  ( $m = 1,000$ ). (a)  $n = 100$  and (b)  $n = 300$ .

where  $\rho$  is a Borel probability measure on a Polish space  $\mathcal{Z}$ . A DPM model (Ferguson, 1983) places a Dirichlet process prior  $DP(a)$  on the mixing distribution  $\rho$ . Thus, a DPM model  $DPM(a)$  assumes the following generative process:

$$x_i | z_i, \phi, F \stackrel{\text{i.i.d.}}{\sim} \psi(x | z_i), \quad z_i | F \stackrel{\text{i.i.d.}}{\sim} F, \quad F \sim DP(a).$$

Here,  $a$  is a finite Borel measure on  $\mathcal{Z}$ . Note that the marginal density (on a single observation) is given by

$$\mathbb{E}_F \left[ \int \psi(x | z) dF(z) \right].$$

Although the prior has an infinite-dimensional component, the required conditional score function is simply  $s_\psi(x | z, \phi)$ ; thus we only need to sample from a finite-dimensional posterior  $P_Z(dz | x)$ . If a model is conditioned on held-out data  $\mathcal{D}$ , then the predictive density  $p(x | \mathcal{D})$  is  $\mathbb{E}_{F | \mathcal{D}}[\int \psi(x | z) dF(z)]$ , which may be used to estimate the density (10). The score function is given by the expectation of  $s_\psi(x | z)$  with respect to the posterior

$$\frac{\psi(x | z, \phi)}{p(x | \mathcal{D})} \bar{F}_D(dz)$$

with  $\bar{F}_D$  being the mean measure of  $P_F(dF | \mathcal{D})$ . Sampling from the posterior can be performed with a combination of the Metropolis–Hastings algorithm and Gibbs sampling (see, e.g., Ghosal & van

der Vaart, 2017, Chapter 5). For the score formula and the MCMC procedure, we refer the reader to [Online Supplementary Material, Section B.8](#) in the supplement. By setting  $\psi$  to an isotropic normal density, for example, we can guarantee the integrability assumption in Lemma 1 (see [Online Supplementary Material, Section B.3](#)).

Our problem below considers comparing the predictive densities defined by two models with different Dirichlet process priors. Note that since candidate models  $P, Q$  here are point estimates derived from their respective posterior means of  $F$ , we discard some aspects of uncertainty in estimating the target (10); our setting only concerns evaluating the quality of those point estimates in approximating the data generating distribution  $R$ .

**Experiment details.** For the data distribution  $R$ , we use the mixture (10) defined by  $\psi(x|z) = \mathcal{N}(x; z, 2I)$  and  $\rho = \mathcal{N}(0, I)$ ; this choice yields  $R = \mathcal{N}(0, 3I)$ . We consider the following simple Gaussian DPM model GDPM( $\mu$ ):

$$x_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(z_i, 2I), \quad z_i \stackrel{\text{i.i.d.}}{\sim} F, \quad F \sim \text{DP}(a), \quad a = \mathcal{N}(\mu, I),$$

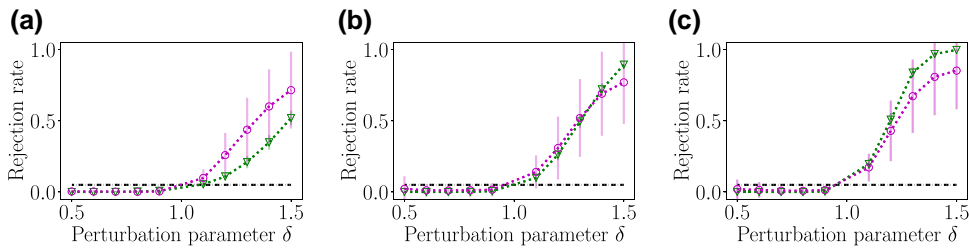
where  $\mu \in \mathbb{R}^D$ . Note that without conditioning on observations, the model's marginal density is simply a Gaussian distribution  $\mathcal{N}(\mu, 3I)$ , which does not require approximation.

We therefore compare predictive distributions, i.e., we compare two GDPM models conditioned on *training data*  $\mathcal{D}_{\text{tr}} = \{\tilde{x}_i\}_{i=1}^{n_{\text{tr}}} \stackrel{\text{i.i.d.}}{\sim} R$ . We consider two GDPM models with *wrong* priors, where their prior means are shifted. Specifically, we take two models chosen as  $Q = \text{GDPM}(\bar{\mathbf{1}})$  and  $P = \text{GDPM}(\delta\bar{\mathbf{1}})$  with  $\bar{\mathbf{1}} = \mathbf{1}/\sqrt{D}$ . Unlike the preceding experiments, we condition the two models on the training data and obtain the predictive distributions, denoted by  $P_{\mathcal{D}_{\text{tr}}}$  and  $Q_{\mathcal{D}_{\text{tr}}}$ , respectively; our problem is thus the comparison between  $P_{\mathcal{D}_{\text{tr}}}$  and  $Q_{\mathcal{D}_{\text{tr}}}$ . The distributions now require simulating their posterior, and we use a random-scan Gibbs sampler and the Metropolis algorithm with a burn-in period  $t = 1,000$  and the size of the latents  $m = 500$ . For sampling observables from the models, we use a random-scan Gibbs sampler with a burn-in period 2,000. We expect that if the training sample size  $n_{\text{tr}}$  is small, a larger perturbation would give a worse model as the effect of the prior is still present; we thus set  $n_{\text{tr}} = 5$ . Due to the small sample size, the expected model relation might not hold, depending on the draw of  $\mathcal{D}_{\text{tr}}$ . Therefore, we examine the rejection rates of the LKSD and MMD tests, averaged over 50 draws; for each draw of  $\mathcal{D}_{\text{tr}}$ , we estimate the rejection rates based on 100 trials. Our problem is formed by varying the perturbation scale  $\delta$  for  $P_{\mathcal{D}_{\text{tr}}}$ , which is chosen from a regular grid  $\{0.5, 0.6, \dots, 0.9, 1.1, \dots, 1.5\}$ . This construction gives a null case when  $\delta < 1$ , the alternative otherwise. We set the dimension  $D$  to 10 and the significance level  $\alpha$  to 0.05. As in Section 4.1, we use the IMQ kernel with median scaling.

Figure 6 reports the rejection rates of the two tests for each of  $n \in \{50, 100, 200\}$ . Note that the curves in the graph do not represent type-I errors nor power, as they are rejection rates *averaged* over draws  $\mathcal{D}_{\text{tr}}$ , each of which forms a different problem. It can be seen that on average, both tests have correct sizes ( $\delta < 1$ ). In the alternative regime ( $\delta > 1$ ), the LKSD test underperforms the MMD with a small sample size ( $n = 50$ ); however, its improvement in power is faster and exceeds the MMD at  $n = 200$ . These results imply that the LKSD estimate has a large variance for a small sample size, whereas its estimand (the population difference) is also larger, and thus the mean of the test statistic diverges faster. Thus, it may be understood that the KSD is more sensitive to model differences in this setting.

### 4.3 LDA

Our final experiment studies the behavior of the LKSD test on discrete data using LDA models. LDA is a mixed-membership model (Airoldi et al., 2014) for grouped discrete data such as text corpora. We follow Blei et al. (2003) and use the terminology of text data for ease of exposition. Accordingly, the following terms are defined using our notation. A word is an element in a discrete set (a vocabulary)  $\{0, \dots, L-1\}$  of size  $L$ . A document  $x$  is a sequence of  $D$  words, i.e.,  $x \in \{0, \dots, L-1\}^D$  is a  $D$ -dimensional discrete vector. A prominent feature of LDA is that it groups similar words assuming that they come from a shared latent *topic*, which serves as a mixture component. An LDA model assumes the following generative process on a corpus of documents  $\{x_i\}_{i=1}^n$ :



**Figure 6.** Comparison in Gaussian Dirichlet mixture models. Rejection rates plotted against the perturbation parameter  $\delta$ . The sample size  $n$  is chosen from  $\{50, 100, 200\}$ . The rejection rates are averaged over draws of  $\mathcal{D}_{tr}$ . The supposed null and alternative regimes are  $\delta < 1$  and  $\delta > 1$ , respectively. Markers:  $\nabla$  (the LKSD test);  $\circ$  (the relative MMD test). The dark dashed line indicates the significance level  $\alpha = 0.05$ . The errorbars indicate the standard deviations of the estimated rejection rates. (a)  $n = 50$ . (b)  $n = 100$  and (c)  $n = 200$ .

1. For each document  $i \in \{1, \dots, n\}$ , generate a distribution over  $K$  topics  $\theta_i \stackrel{i.i.d.}{\sim} \text{Dir}(a)$  (the Dirichlet distribution), where  $\theta_i$  is a probability vector of size  $K \geq 1$ .
2. For the  $j$ th word  $x_i^j, j \in \{1, \dots, D\}$  in a document  $i$ ,
  - (a) Choose a topic  $z_i^j \stackrel{i.i.d.}{\sim} \text{Cat}(\theta_i)$ .
  - (b) Draw a word from  $x_i^j \stackrel{i.i.d.}{\sim} \text{Cat}(b_k)$ , where  $b_k$  is the distribution over words for topic  $k$ , and the topic assignment  $z_i^j = k$ .

Here,  $a = (a_1, \dots, a_K)$  is a vector of positive real numbers, and  $b = (b_1, \dots, b_K)^T \in [0, 1]^{K \times L}$  represents a collection of  $K$  distributions over  $L$  words. In summary, an LDA model  $P = \text{LDA}(a, b)$  assumes the factorization

$$\prod_{i=1}^n p(x_i | z_i, \theta_i; a, b) p(z_i, \theta_i; a, b) = \prod_{i=1}^n \left\{ \prod_{j=1}^D p(x_i^j | z_i^j, b) p_z(z_i^j | \theta_i) \right\} p_\theta(\theta_i | a),$$

where  $z_i$  and  $\theta_i$  act as latent variables.

Because of the independence structure over words, the conditional score function is simply given as

$$s_p(x | z, \theta, a, b) = s_p(x | z, b) = \left( \frac{p(\tilde{x}^j | z^j, b)}{p(x^j | z^j, b)} - 1 \right)_{j=1, \dots, D}, \quad \text{where } \tilde{x}^j = x^j + 1 \text{ mod } L.$$

Score approximation requires the posterior distribution  $p(z | x; a, b)$  with respect to  $z$ . Marginalization of  $\theta$  renders latent topics dependent on each other, and thus the posterior is intractable. A latent topic is conjugate to the corresponding topic distribution given all other topics. Therefore, an MCMC method such as collapsed Gibbs sampling allows us to sample from  $p(z | x; a, b)$ . As the observable and the latent are supported on finite sets, the use of Lemma 1 is justified; the finite moment assumptions in Corollary 1 are guaranteed; and the consistency of the population mean and variance of the test statistic follows from the convergence of  $\mathbb{E}_{z|x}^{(t)}[\bar{s}_p(x | z)]$  and  $\mathbb{E}_{w|x}^{(t)}[\bar{s}_q(x | w)]$  for each  $x \in \mathcal{X}$ .

#### 4.3.1 Synthetic data—prior sparsity perturbation

In the below two problems, we observe a sample  $\{x_i\}_{i=1}^n$  from an LDA model  $R = \text{LDA}(a, b)$ . The number of topics is  $K = 3$ . The hyper-parameter  $a$  is chosen as  $a = (a_0, a_0, a_0)$ ; for model  $R$ , we set  $a_0 = 0.1$ . Each of three rows in  $b = (b_1, b_2, b_3)^T \in [0, 1]^{3 \times L}$  is fixed at a value drawn from the symmetric Dirichlet distribution with all the concentration parameters one, and the vocabulary size is  $L = 10,000$ . Each  $x_i \in \{0, \dots, L - 1\}^D$  is a document consisting of  $D = 50$  words.



**Table 2.** Rejection rates of the MMD test and the LKSD test in LDA experiments

(a) Type-I errors of the KSD and MMD tests in LDA Problem 1;  $(\delta_P, \delta_Q) = (0.5, 0.6)$ . The significance level  $\alpha = 0.05$ .

| Sample size $n$ | Rejection rates |       |
|-----------------|-----------------|-------|
|                 | MMD             | LKSD  |
| 100             | 0.003           | 0.013 |
| 200             | 0.010           | 0.007 |
| 300             | 0.007           | 0.003 |
| 400             | 0.003           | 0.007 |
| 500             | 0.007           | 0.010 |

(b) Power of the KSD and MMD tests in LDA Problem 2;  $(\delta_P, \delta_Q) = (1.0, 0.5)$ . The significance level  $\alpha$  is chosen from  $\{0.01, 0.05\}$ .

| Sample size $n$ | Rejection rates       |       |                       |       |
|-----------------|-----------------------|-------|-----------------------|-------|
|                 | Level $\alpha = 0.01$ |       | Level $\alpha = 0.05$ |       |
|                 | MMD                   | LKSD  | MMD                   | LKSD  |
| 100             | 0.000                 | 0.010 | 0.007                 | 0.070 |
| 200             | 0.003                 | 0.030 | 0.010                 | 0.183 |
| 300             | 0.000                 | 0.097 | 0.003                 | 0.283 |
| 400             | 0.000                 | 0.197 | 0.010                 | 0.463 |
| 500             | 0.000                 | 0.280 | 0.007                 | 0.570 |

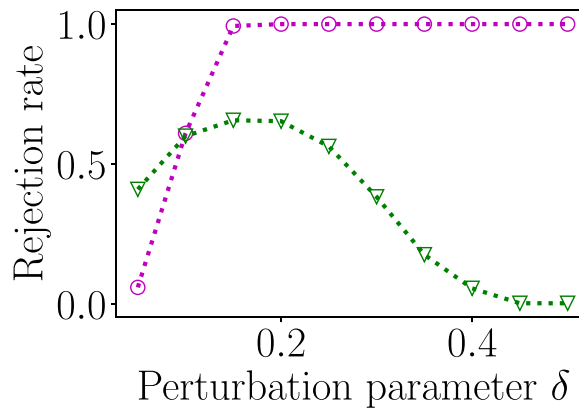
*Note.* Each result is based on 300 trials.

service arXiv. We treat the abstract of an article as a document and use paper categories to set up a problem. Specifically, we construct a problem by choosing three paper categories for model  $P$ ,  $Q$  and the data distribution  $R$ . Unlike the preceding experiments, for a model category, we fit an LDA model to the dataset of abstracts in the category. As the KSD requires the number of words to be fixed, then for a given data category, we extract abstracts of length no less than  $D = 100$  and subsample excess words. This process yields a dataset of articles of equal length  $D$ ; for each trial, we obtain the data  $\{x_i\}_{i=1}^n$  by subsampling from the larger set of articles. Thus, our problem is to compare two LDA models trained on different article sets and assess their fit to the dataset.

In the following experiments, we examine the power of LKSD and MMD tests. We vary the sample size  $n$  from 100 to 500. We fix the dataset category to stat.TH (statistics theory) and inspect two combinations of model categories. To train an LDA model  $LDA(a, b)$ , we use the Gensim implementation (Rehurek & Sojka, 2011) of the variational algorithm of M. Hoffman et al. (2010). For sparsity parameters  $a$ , we use the parameter returned by this algorithm; we point-estimate topics  $b$  using the mean of the topics under the variational distribution. The number of topics is set to 100. The vocabulary set is comprised of words that appear in the abstracts of three chosen categories. As in the previous experiments, we use the IMQ-BoW kernel for both tests. We fix the significance level  $\alpha$  at 0.05.

As we have seen the numerical instability issue in the previous section, we also consider an alternative KSD that is stable but computationally more expensive, as mentioned in Section 2.1 and the supplement (Online Supplementary Material, Section B.2). For this, we take a burn-in size  $t = 500$  and a Markov chain size  $m = 1,000$ . We denote this method by LKSD-stable.

**Probability theory vs statistical methodology.** We choose math.PR (mathematics probability theory) for  $P$  and stat.ME (statistics methodology) for  $Q$ . In addition to the taxonomic proximity



**Figure 7.** Power estimates plotted against perturbation parameters  $\delta$ . The significance level  $\alpha = 0.05$ ; the sample size  $n = 300$ . Markers:  $\nabla$  (the LKSD test);  $\circ$  (the MMD test).

to stat.TH, the category stat.ME has a larger proportion of articles shared with the target category: 3, 121 of 18, 973 (stat.ME) vs. 2, 884 of 46, 769 (math.PR). Thus, we expect  $Q$  to outperform  $P$ . This combination results in a vocabulary set of size  $L = 126, 190$ . For score estimation, we set the burn-in length  $t$  to 500 and the Markov chain sample size  $m$  to 5,000. Additionally, we run the LKSD test with  $m = 15,000$  (labelled LKSD-extra) and the MMD test with the model sample size  $n_{\text{model}} = 10,000$  (labelled MMD-extra). The sample size  $n_{\text{model}}$  is thresholded at 10,000 as the computational cost exceeds that of the LKSD test (in fact, sampling in this case makes the MMD by an order of magnitude slower due to the large vocabulary size).

**Table 3** summarizes the result. The MMD test underperforms all the KSD-based tests; extra sampling did not lead to a significant improvement. We can see that increasing the Markov chain size  $m$  boosts the LKSD test, as it reduces the variance of the score estimator. The low power of the MMD test indicates that the model difference is too subtle to discern from the word compositions of generated documents; the LKSD tests offer a different viewpoint based on the model information.

**Machine learning vs statistical methodology.** Our second experiment uses cs.LG (computer science machine learning) for  $P$ , while  $Q$  uses the same category as the previous experiment. With this combination, the vocabulary size  $L$  is 208, 671. By the same reasoning as above, the second model  $Q$  is expected to be better than  $P$ . We run the same tests as above and compare their performances.

**Table 4** summarizes the result. This experiment serves as a negative case study for the LKSD test: the MMD tests achieved power 1 for all sample-size choices (MMD-extra is omitted here), whereas the power of the LKSD test does not exceed even the significance level  $\alpha$  for most sample size settings (LKSD-extra is omitted as increasing the Markov chain size did not improve the power). We attribute this failure to the unmatched support of the model  $P$  in the test distribution. This reasoning is supported by the high power of the MMD, as the BoW feature easily detects deviation of document patterns in this case. Thus, as we noted in the synthetic experiment in Section 4.3.2, the LKSD test fails when there is a severe mismatch in data and model support. The stable LKSD test approaches the same level as the MMD at  $n = 500$ , but still underperforms. While stable, the KSD used for this test can also suffer from the mismatch of the support, since it depends on the same density ratio as in the unstable counterpart.

## 5 Conclusion

We have developed a test of relative goodness of fit for latent variable models based on the kernel Stein discrepancy. The proposed test applies to a wide range of models, since the requirements of the test are mild: (a) models have MCMC samplers for inferring their latent variables, and (b) likelihoods have evaluable score functions. The proposed test complements existing model









