

Discrepancy-based Inference for Intractable Generative Models using Quasi-Monte Carlo

Ziang Niu^{1,*}, Johanna Meier^{2,*}, François-Xavier Briol^{3,†}

¹University of Pennsylvania, ²Leibniz Universität Hannover, ³University College London,

*contributed equally, †corresponding author.

July 5, 2022

Abstract

Intractable generative models are models for which the likelihood is unavailable but sampling is possible. Most approaches to parameter inference in this setting require the computation of some discrepancy between the data and the generative model. This is for example the case for minimum distance estimation and approximate Bayesian computation. These approaches require sampling a high number of realisations from the model for different parameter values, which can be a significant challenge when simulating is an expensive operation. In this paper, we propose to enhance this approach by enforcing “sample diversity” in simulations of our models. This will be implemented through the use of quasi-Monte Carlo (QMC) point sets. Our key results are sample complexity bounds which demonstrate that, under smoothness conditions on the generator, QMC can significantly reduce the number of samples required to obtain a given level of accuracy when using three of the most common discrepancies: the maximum mean discrepancy, the Wasserstein distance, and the Sinkhorn divergence. This is complemented by a simulation study which highlights that an improved accuracy is sometimes also possible in some settings which are not covered by the theory.

1 Introduction

A particular challenge for statistics is the growing complexity of phenomena modelled by scientists, and as a result the growing complexity of the models themselves. This can often lead to cases where a closed form of the likelihood is not available anymore. As a result, classical parameter estimation tools such as maximum likelihood estimation or Bayesian inference cannot be used. Within these so-called intractable likelihood models, intractable generative models are parametric families of probability distributions which are specified through a generative process, so that it is possible to obtain realisations for any value of the parameter [26]. These models are widely used throughout the sciences including genetics [9], astronomy [20] and ecology [8]. In machine learning, one of the main applications is for the simulation of realistic looking images [56]; see the recent line of work on generative adversarial networks [40].

Denote by \mathbb{P}_θ any element of a parametric family of interest with parameter θ , and let \mathcal{X} be the space of realisations from this model. The generative process of \mathbb{P}_θ is usually summarised through a pair (\mathbb{U}, G_θ) which includes a relatively simple probability distribution \mathbb{U} (such as a Gaussian or uniform) on some space \mathcal{U} and a parametric map $G_\theta : \mathcal{U} \rightarrow \mathcal{X}$ called a generator or simulator. To obtain n independent and identically distributed (IID) realisations $\{x_i\}_{i=1}^n$ from the model for some fixed parameter θ , one can simply sample IID realisations $u_i \sim \mathbb{U}$, then map

these samples through the generator $x_i = G_\theta(u_i)$. The main advantage of generative models is that one can model ever more complex phenomena by increasing the flexibility of the generator, as long as the map G_θ can be evaluated pointwise.

Since simulating data is the only option available in the case of generative models, many inference methods for this class are based on simulating synthetic data for various parameter values, then comparing the simulated data to the observations to select a “good” parameter value. The latter usually requires defining some notion of distance, or discrepancy, between the two datasets. Once a discrepancy is defined, one possible approach is the framework of *minimum distance estimation (MDE)* [71], where an estimator is constructed as the minimiser (over the set of model parameters) of the discrepancy between datasets. In the Bayesian literature, an alternative approach called *approximate Bayesian computation (ABC)* [9] consists of constructing a pseudo-posterior distribution over parameters by selecting parameter values simulated from a prior distribution for which the discrepancy between simulated and actual data is small. In all of the cases above, thinking of the actual data as an approximation to the data-generating process of interest, the main computational challenge can be summarised as having to efficiently estimate some discrepancy given access to realisations of two distributions.

There is a vast literature on possible discrepancies, each with competing advantages for parameter estimation including efficiency, robustness to model misspecification, computational cost and sample complexity. In this paper, we will not aim to be exhaustive, but will focus on a small subset of discrepancies which are popular in the literature because they lend themselves to efficient implementations. The first discrepancy is the maximum mean discrepancy (MMD) [41], which compares embeddings of probability distributions into reproducing kernel Hilbert spaces, and can be straightforwardly computed through evaluations of a kernel. This was studied by [16, 21, 22, 4, 27] in the context of MDE, and by [53, 31, 51, 82, 14] for the case where G_θ is a neural network in particular. It was also used by [61, 70, 55, 46, 13] in the context of ABC. The two other discrepancies we will consider are the Wasserstein distance, as well as its relaxation called the Sinkhorn divergence. These can be efficiently implemented thanks to algorithmic advances in computational optimal transport [73]. They were considered for MDE by [6, 10, 38, 28, 85, 60, 63, 78] and for ABC by [11, 39, 58].

Clearly, any algorithmic development improving our ability to estimate these discrepancies will significantly reduce the overall computational cost of implementing all of the algorithms described above. We propose to tackle this problem through the use of quasi-Monte Carlo (QMC) point sets [30]. In particular, we focus on the case where \mathbb{U} is a uniform distribution¹ and replace independent and identical distributed (IID) realisations by some QMC point set. This is a rather simple algorithmic trick, which we will call QMC sampling and which has been explored for a wide range of models; see for example [19] for copula models, or [43, 44] for neural networks. Once again, a full review of QMC sampling is out of scope for this paper. Intuitively, this approach consists of generating a more “diverse” set of samples from the model. This can be observed visually through the example in Figure 1 which compares realisations from a Gaussian distribution obtained through Monte Carlo (MC) and QMC. Clearly, the realisations obtained through QMC provide an improved approximation of \mathbb{P}_θ in the intuitive sense that they provide a more uniform coverage of areas of high-probability under \mathbb{P}_θ .

The main contribution of this paper is a set of theoretical results demonstrating the advantages of QMC sampling for performing inference with discrepancies. In particular, Theorem 1, Theorem 2 and Theorem 3 provide sample complexity results with respect to the MMD, Wasserstein and Sinkhorn divergence respectively. In each case, the theorem provides sufficient conditions for estimating the discrepancy at a rate which is linear (up to log factors) in the num-

¹The assumption that \mathbb{U} is uniform is relatively minor due to Sklar’s theorem, which states that any multivariate distribution can be obtained through a transformation of a uniform distribution.

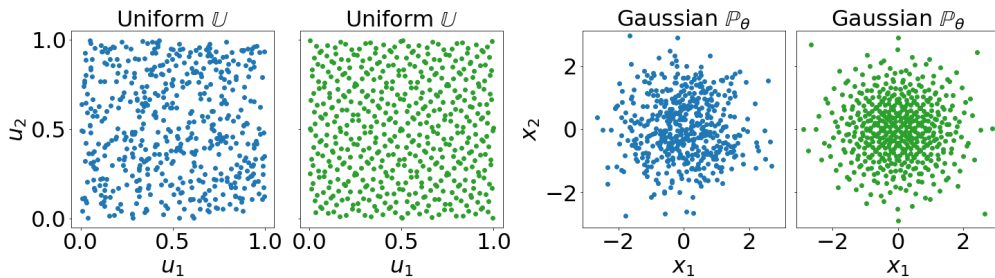


Figure 1: Realisations from $\mathbb{P}_\theta = \mathcal{N}(0, I)$, a zero mean Gaussian with identity covariance matrix. We compare realisations from a $\text{Unif}([0, 1]^2)$ (left plot) against a QMC point set (center left plot), together with their projections through the generator $x = G_\theta(u) = (\Phi_{\theta_1}^{-1}(u_1), \Phi_{\theta_2}^{-1}(u_2))$ (center right and right plot), where Φ_θ denotes the cumulative distribution function of \mathbb{P}_θ .

ber of realisations n . This is a significant improvement upon the usual MC rate which decreases at a root- n speed. Of course, such speed-ups do come at the cost of the generality of the method as they require certain regularity conditions on G_θ and \mathcal{X} . Despite this drawback, we show through an extensive simulation study that faster rates than MC (although not necessarily linear) can still be obtained for QMC in some settings not covered by our theory. We therefore see this paper as an initial step in the study of the use of QMC sampling for discrepancy estimation.

The remainder of this paper is structured as follows. In Section 2, we introduce intractable generative models and the most common distances used for inference, including the MMD, Wasserstein distance and Sinkhorn divergence. In Section 3, we derive our novel sample complexity results. Finally, the performance of the performance of these novel estimators is studied numerically in Section 4. We conclude by discussing potential future directions in Section 5.

2 Background

This section will recall background material on inference for intractable generative models (in Section 2.1), then introduce the main discrepancies considered in the literature (in Section 2.2).

2.1 Inference for Intractable Generative Models through Discrepancies

Throughout this paper, we will consider settings where the base space is $\mathcal{U} = [0, 1]^s$, the data space satisfies $\mathcal{X} \subseteq \mathbb{R}^d$ and the parameter space satisfies $\Theta \subseteq \mathbb{R}^p$ for $s, p, d \in \mathbb{N}_+ = \{1, 2, 3, \dots\}$. We will denote by $\mathcal{P}(\mathcal{X})$ the set of all Borel probability distributions on \mathcal{X} .

The inference task of interest can be summarised as follows. Given IID realisations $\{y_j\}_{j=1}^m$ from some unknown $\mathbb{Q} \in \mathcal{P}(\mathcal{X})$, we would like to find the parameter $\theta^* \in \Theta$ such that \mathbb{P}_{θ^*} is closest to \mathbb{Q} in some sense. In particular, if $\mathbb{Q} \in \{\mathbb{P}_\theta \in \mathcal{P}(\mathcal{X}) : \theta \in \Theta\}$ (i.e. the model is well-specified), our task is to recover the parameter value θ^* which was used to simulate the observations $\{y_j\}_{j=1}^m$. One approach is to use a discrepancy, which we will define to be any function $D : \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X}) \rightarrow [0, \infty)$. Specific examples will be provided in Section 2.2, but for now we will only assume such a discrepancy has been selected, and describe how it can be used for inference. Firstly, we may construct an estimator through the framework of MDE [71]:

$$\hat{\theta}_m^D \in \arg \min_{\theta \in \Theta} D(\mathbb{P}_\theta, \mathbb{Q}^m),$$

where $\mathbb{Q}^m(dx) = \frac{1}{m} \sum_{j=1}^m \delta_{y_j}(dx)$ is an empirical measure, and δ_{y_j} a Dirac measure at y_j . Of course, this is usually an intractable optimisation problem since it requires evaluating D pointwise at \mathbb{P}_θ , which is itself unknown. As a result, a common approach is to solve the optimisation problem through evaluations of $D(\mathbb{P}_\theta^n, \mathbb{Q}^m)$, or of its gradient, where $\mathbb{P}_\theta^n(dx) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}(dx)$ and is obtained from realisations $\{x_i\}_{i=1}^n$ from \mathbb{P}_θ . For all discrepancies considered in this paper, $D(\mathbb{P}_\theta^n, \mathbb{Q}^m)$ is a biased estimate of $D(\mathbb{P}_\theta, \mathbb{Q}^m)$. This leads to the use of stochastic optimisation methods with biased gradient estimates, which leads to a bias in the estimated parameter [83, 48]. However, any approach leading to more efficient estimation of $D(\mathbb{P}_\theta, \mathbb{Q}^m)$ may be able to significantly reduce this bias.

Secondly, we may use ABC, which aims to construct a pseudo-posterior which closely approximates the exact Bayesian posterior [9]. This can be achieved by sampling parameter values $\{\theta_k\}_{k=1}^K$ (for some $K \in \mathbb{N}$) from a prior distribution π , then for each of these values simulating a dataset $\{x_i^k\}_{i=1}^n$. Each of these parameter values is then accepted as a realisation from the pseudo-posterior if $D(\mathbb{P}_{\theta_k}^n, \mathbb{Q}^m) \leq \varepsilon$ holds for some threshold parameter $\varepsilon > 0$. This straightforward procedure allows us to sample from the following pseudo-posterior:

$$\Pi_\varepsilon^D(d\theta|y_1, \dots, y_m) \propto \Pi(d\theta) \mathbb{E} \left[\mathbf{1}_{\{D(\mathbb{P}_\theta^n, \mathbb{Q}^m) \leq \varepsilon\}}(d\theta) \right],$$

where $\mathbf{1}_A$ is an indicator function for the event A , and the expectation is with respect to the randomness in the simulated data. Note that this sampling procedure is only necessary due to the intractability of $D(\mathbb{P}_{\theta_k}, \mathbb{Q}^m)$ for intractable generative models; if this quantity was tractable, we would instead want to verify whether $D(\mathbb{P}_{\theta_k}, \mathbb{Q}^m) \leq \varepsilon$ instead of $D(\mathbb{P}_{\theta_k}^n, \mathbb{Q}^m) \leq \varepsilon$.

QMC has previously been used for ABC [18], but this was used to improve sampling of parameters instead of simulating the data. Finally, we also point out that recent generalised Bayesian procedures for generative models are also discrepancy-based; see for example [77, 69].

Clearly MDE and ABC critically rely on $D(\mathbb{P}_\theta^n, \mathbb{Q}^m)$ approaching $D(\mathbb{P}_\theta, \mathbb{Q}^m)$ at a fast rate in n . Whether this is possible will depend on the discrepancy D .

2.2 Examples of Discrepancies for Inference

Recall that any discrepancy D such that $\forall \mathbb{P}_1, \mathbb{P}_2, \mathbb{P}_3 \in \mathcal{P}(\mathcal{X})$: (i) $D(\mathbb{P}_1, \mathbb{P}_2) = 0$ if and only if $\mathbb{P}_1 = \mathbb{P}_2$, (ii) $D(\mathbb{P}_1, \mathbb{P}_2) = D(\mathbb{P}_2, \mathbb{P}_1)$, and (iii) $D(\mathbb{P}_1, \mathbb{P}_2) \leq D(\mathbb{P}_1, \mathbb{P}_3) + D(\mathbb{P}_3, \mathbb{P}_2)$, is called a probability metric on $\mathcal{P}(\mathcal{X})$. If only (i) holds, D is called a (statistical) divergence. The discrepancies in this paper closely relate to integral probability metrics (IPMs) [57]. Given a set of functions \mathcal{F} , an IPM is a probability metric which takes the form:

$$D_{\mathcal{F}}(\mathbb{P}, \mathbb{Q}) := \sup_{f \in \mathcal{F}} \left| \int_{\mathcal{X}} f(x) \mathbb{P}(dx) - \int_{\mathcal{X}} f(x) \mathbb{Q}(dx) \right|.$$

In practice, \mathcal{F} needs to be large enough to be able to differentiate \mathbb{P} from \mathbb{Q} , but also small enough so that $D_{\mathcal{F}}(\mathbb{P}, \mathbb{Q})$ can be computed, or at least approximated up to high accuracy. It should also not be too large since we might otherwise have $D_{\mathcal{F}}(\mathbb{P}, \mathbb{Q}) = \infty$ for all $\mathbb{P} \neq \mathbb{Q}$. This can significantly restrict the choices of \mathcal{F} available for inference. In the case where \mathcal{F} is finite-dimensional, the discrepancy above can be thought of as comparing a finite number of summary statistics of \mathbb{P} and \mathbb{Q} , as commonly done for the method of simulated moments or in ABC. For this case, the use of QMC was previously studied in [35]. In contrast, our work will focus on the most common discrepancies based on infinite-dimensional \mathcal{F} , which we introduce below.

Maximum Mean Discrepancy Let $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathbb{R} : \|f\|_{\mathcal{H}_k} \leq 1\}$, the unit-ball of a reproducing kernel Hilbert space (RKHS) \mathcal{H}_k with kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. In this case, the

IPM is called the (kernel) *maximum mean discrepancy* [41]. We will assume that the kernel is characteristic, which guarantees that the discrepancy is a metric on the set

$$\mathcal{P}_k(\mathcal{X}) := \{\mathbb{P} \in \mathcal{P}(\mathcal{X}) : \int_{\mathcal{X}} \sqrt{k(x,x)} \mathbb{P}(dx) < \infty\} \subseteq \mathcal{P}(\mathcal{X}),$$

see [79] for more details. The name MMD originates from the fact that the IPM can be expressed as $\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\int_{\mathcal{X}} k(\cdot, y) \mathbb{P}(dy) - \int_{\mathcal{X}} k(\cdot, y) \mathbb{Q}(dy)\|_{\mathcal{H}_k}$, which is the size of the difference between \mathbb{P} and \mathbb{Q} when embedded in \mathcal{H}_k . The squared-MMD can alternatively be expressed as

$$\begin{aligned} \text{MMD}^2(\mathbb{P}, \mathbb{Q}) &:= \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) \mathbb{P}(dx) \mathbb{P}(dy) - 2 \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) \mathbb{P}(dx) \mathbb{Q}(dy) \\ &\quad + \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) \mathbb{Q}(dx) \mathbb{Q}(dy). \end{aligned} \quad (1)$$

Note that this expression does not require the computation of a supremum anymore. Given two empirical measures $\mathbb{P}^n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ and $\mathbb{Q}^m = \frac{1}{m} \sum_{j=1}^m \delta_{y_j}$ approximating \mathbb{P} and \mathbb{Q} respectively, this expression lends itself naturally to the following approximation:

$$\text{MMD}^2(\mathbb{P}^n, \mathbb{Q}^m) = \frac{\sum_{i \neq j}^n k(x_i, x_j)}{n^2} - \frac{2 \sum_{i=1}^n \sum_{j=1}^m k(x_i, y_j)}{nm} + \frac{\sum_{i \neq j}^m k(y_i, y_j)}{m^2}. \quad (2)$$

The use of a U-statistic may also be preferred in some case; see for example [16]. One of the main advantages of the MMD is the fact that it can be easily approximated, but also that the kernel choice allows for significant flexibility. The most common example is the Gaussian (or squared-exponential) kernel $k(x, x') = \lambda^2 \exp(-\|x - x'\|_2^2 / \sigma^2)$ where $\lambda, \sigma > 0$. QMC point sets were already used with the MMD in [43, 44] in the context of neural network generators, but those papers do not study the sample complexity of the approach from a theoretical viewpoint.

Wasserstein Distance Let $c : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$ be a metric (called cost function), $p \geq 1$ and $\Gamma(\mathbb{P}, \mathbb{Q}) \subset \mathcal{P}(\mathcal{X} \times \mathcal{X})$ be the set of distributions with marginals $\mathbb{P} \in \mathcal{P}(\mathcal{X})$ and $\mathbb{Q} \in \mathcal{P}(\mathcal{X})$ in the first and second coordinate respectively. The Wasserstein distance can be expressed as:

$$W_{c,p}(\mathbb{P}, \mathbb{Q}) := \left(\min_{\gamma \in \Gamma(\mathbb{P}, \mathbb{Q})} \int_{\mathcal{X} \times \mathcal{X}} c^p(x, y) \gamma(dx, dy) \right)^{\frac{1}{p}},$$

A common choice for c is the Euclidean distance, but other metrics can be used. The Wasserstein distance is a probability metric on the set

$$\mathcal{P}_{c,p}(\mathcal{X}) = \{\mathbb{P} \in \mathcal{P}(\mathcal{X}) : \int_{\mathcal{X}} c^p(x, y) \mathbb{P}(dx) < \infty \forall y \in \mathcal{X}\} \subseteq \mathcal{P}(\mathcal{X}).$$

Although computing the Wasserstein distance for general \mathbb{P} and \mathbb{Q} is usually not possible, it is straightforward to do so for empirical measures \mathbb{P}^n and \mathbb{Q}^m (see for example Chapter 3 in [73]):

$$W_{c,p}(\mathbb{P}^n, \mathbb{Q}^m) = \left(\min_P \sum_{i=1}^n \sum_{j=1}^m c^p(x_i, y_j) P_{ij} \right)^{\frac{1}{p}},$$

where the minimisation is performed over all $n \times m$ matrices such that $P_{ij} \geq 0 \forall i, j$, $\sum_{i=1}^n P_{ij} = \frac{1}{m}$ and $\sum_{j=1}^m P_{ij} = \frac{1}{n}$. To approximate $W_{c,p}(\mathbb{P}, \mathbb{Q})$, a natural approach is to use $W_{c,p}(\mathbb{P}^n, \mathbb{Q}^m)$, but this is known to have a slow convergence rate as n increases whenever $d > 1$ [36]. In the special case where $p = 1$, the Wasserstein distance is an IPM which corresponds to taking \mathcal{F} to be the set of functions with Lipschitz constant 1: $\{f : \mathcal{X} \rightarrow \mathbb{R} \text{ s.t. } \forall x, y \in \mathcal{X}, |f(x) - f(y)| \leq c(x, y)\}$. This is therefore another setting of infinite-dimensional \mathcal{F} where the supremum does not need to be computed numerically.

Sinkhorn Divergence A common relaxation of the Wasserstein distance is the following:

$$\begin{aligned}\bar{W}_{c,p,\lambda}(\mathbb{P}, \mathbb{Q}) &:= \min_{\gamma \in \Gamma(\mathbb{P}, \mathbb{Q})} \int_{\mathcal{X} \times \mathcal{X}} c^p(x, y) \gamma(dx, dy) + \lambda H(\gamma \| \mathbb{P} \otimes \mathbb{Q}), \\ H(\gamma \| \mathbb{P} \otimes \mathbb{Q}) &:= \int_{\mathcal{X} \times \mathcal{X}} \log \left(\frac{\gamma(dx, dy)}{\mathbb{P}(dx) \mathbb{Q}(dy)} \right) \gamma(dx, dy)\end{aligned}$$

where $H(\pi \| \mathbb{P} \otimes \mathbb{Q})$ is called the relative entropy, and $\mathbb{P} \otimes \mathbb{Q}$ is the product measure. Since this discrepancy is not normalised, it is common to work instead with the *Sinkhorn divergence* [37]:

$$S_{c,p,\lambda}(\mathbb{P}, \mathbb{Q}) = \bar{W}_{c,p,\lambda}(\mathbb{P}, \mathbb{Q}) - \frac{1}{2} (\bar{W}_{c,p,\lambda}(\mathbb{P}, \mathbb{P}) + \bar{W}_{c,p,\lambda}(\mathbb{Q}, \mathbb{Q})),$$

which guarantees the resulting value is greater or equal to zero. The Sinkhorn divergence is also symmetric, but does not satisfy the triangle inequality and so is not a metric. However, it does interpolate between the two IPMs we have seen so far: as $\lambda \rightarrow 0$, $S_{c,p,\lambda}(\mathbb{P}, \mathbb{Q}) \rightarrow W_{c,p}(\mathbb{P}, \mathbb{Q})$, whereas when $\lambda \rightarrow \infty$, $S_{c,p,\lambda}(\mathbb{P}, \mathbb{Q}) \rightarrow \text{MMD}(\mathbb{P}, \mathbb{Q})$ with kernel $k = -c$ [33]. Once again, it is straightforward to compute $S_{c,p,\lambda}(\mathbb{P}^n, \mathbb{Q}^m)$ in the case of empirical measures, and this can be used to estimate the exact Sinkhorn divergence: $S_{c,p,\lambda}(\mathbb{P}, \mathbb{Q}^m)$. From a computational viewpoint, one particular advantage of the Sinkhorn divergence over the Wasserstein distance is that it has better sample complexity when using Monte Carlo points in multiple dimensions [37]. We will return to this point in the next section on QMC sample complexity.

Sliced Discrepancies A final example of discrepancies commonly used for inference are the so-called *sliced discrepancies* [50]. The main motivation for these is to construct discrepancies which will be useful for high-dimensional problems. This is done by projecting probability distributions on \mathcal{X} to probability distributions on some lower dimensional space \mathcal{Y} (usually one dimension) using a map $\mathcal{S}_\xi : \mathcal{P}(\mathcal{X}) \rightarrow \mathcal{P}(\mathcal{Y})$, then comparing these projections using any discrepancy $D : \mathcal{P}(\mathcal{Y}) \times \mathcal{P}(\mathcal{Y}) \rightarrow [0, \infty)$, such as those discussed above. The corresponding sliced discrepancy consists of an average over possible projections:

$$SD(\mathbb{P}, \mathbb{Q}) = \int_{\Xi} D(\mathcal{S}_\xi \mathbb{P}, \mathcal{S}_\xi \mathbb{Q}) d\xi,$$

where $\mathcal{S}_\xi \mathbb{P}, \mathcal{S}_\xi \mathbb{Q}$ are the projections of \mathbb{P}, \mathbb{Q} along the direction $\xi \in \Xi$, and Ξ is the space of directions considered. In order to compute the discrepancy, an MC estimator is used: $\widehat{SD}(\mathbb{P}, \mathbb{Q}) = \frac{1}{L} \sum_{l=1}^L D(\mathcal{S}_{\xi_l} \mathbb{P}, \mathcal{S}_{\xi_l} \mathbb{Q})$ where $\{\xi_l\}_{l=1}^L$ are MC realisations from a uniform distribution over Ξ . The most common sliced-discrepancy is the sliced-Wasserstein distance $SW_{c,p}$ [28, 85, 60, 63], in which case D is $W_{c,p}$ and the projections are constructed using the Radon transform.

3 Sample Complexity with Quasi-Monte Carlo

Now that we have introduced the main discrepancies which will be considered in this paper, we are ready to introduce our novel sample complexity results based on QMC and RQMC. We first introduce the methodology in Section 3.1, then provide theoretical results demonstrating improved sample complexity for MMD in Section 3.3 and for the Wasserstein distance and its Sinkhorn approximation in Section 3.4 and 3.5 respectively. These results all build upon the work of [7], which considered the use of QMC for integrating compositions of functions.

Notation For two sequences $\{f_n\}_{n \in \mathbb{N}}$ and $\{g_n\}_{n \in \mathbb{N}}$, $f_n = O(g_n) \Leftrightarrow \limsup_{n \rightarrow \infty} |f_n/g_n| < \infty$. For some $f : \mathcal{X} \rightarrow \mathbb{R}$ and multi-index $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}^d$, we will denote by $\partial^\alpha f$ the partial derivative $\partial^{|\alpha|} f / \partial^{\alpha_1} x_1 \dots \partial^{\alpha_d} x_d$. The space $\mathcal{C}^m(\mathcal{X})$ of m -continuously differentiable

functions ($m \in \mathbb{N}$) corresponds to functions such that $\partial^\alpha f$ is continuous $\forall \alpha \in \mathbb{N}^d$ such that $|\alpha| = \alpha_1 + \dots + \alpha_d \leq m$. Similarly, $\mathcal{C}^{m,m}(\mathcal{X} \times \mathcal{X})$ will denote functions $f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that $\partial^{\alpha,\alpha} f$ exists and is continuous $\forall \alpha \in \mathbb{N}^d$ with $|\alpha| \leq m$. Relatedly, if we have a set $\beta \subseteq 1 : d$, we write ∂_β to denote the (first-order) mixed partial derivatives of f with respect to the coordinates in the set β . Finally, we will write $L^p(\mathcal{X})$ to denote the p -integrable functions; i.e. $f : \mathcal{X} \rightarrow \mathbb{R}$ satisfying $\|f\|_{L^p(\mathcal{X})} := (\int_{\mathcal{X}} f^p(x) dx)^{\frac{1}{p}} < \infty$ (where we will use the common abuse of terminology to avoid technicalities with equivalence classes).

3.1 Enhancing Sample Diversity through quasi-Monte Carlo

Recall that to obtain realisations $\{x_i\}_{i=1}^n$ from \mathbb{P}_θ , the generative approach consist of obtaining realisations $\{u_i\}_{i=1}^n \sim \text{Unif}([0, 1]^s)$, then mapping these through the generator: $x_i = G_\theta(u_i)$. Under sufficient regularity conditions on G_θ , we would expect two realisations x_1, x_2 to be far from one another whenever u_1, u_2 are also far from one another. The main idea in this paper is that we may improve sample diversity by selecting $\{u_i\}_{i=1}^n$ according to a QMC point set. This notion of diversity is usually measured through the *star-discrepancy* of a point set:

$$D^*(\{u_i\}_{i=1}^n) := \sup_{v \in [0, 1]^s} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{[0, v)}(u_i) - \prod_{i=1}^s v_i \right|.$$

We will call a point set $\{u_i\}_{i=1}^n$ such that $D^*(\{u_i\}_{i=1}^n) = O(n^{-1}(\log n)^{\alpha_s})$ for some $\alpha_s > 0$ as $n \rightarrow \infty$ a *QMC point set*, and α_s will usually depend on the dimensionality s of the domain \mathcal{U} . This is also sometimes referred to as a low-discrepancy point set, but we will avoid this terminology to avoid any confusion between discrepancies on probability distributions and the star discrepancy. Popular constructions [30] include Hammersley point sets, which are based on infinite van der Corput sequence, and can achieve $\alpha_s = s - 1$. Alternatively, lattice point sets achieve $\alpha_2 = 2$ and $\alpha_s = s$ for $s \geq 3$, (t, m, s) -nets in base b achieve $\alpha_s = s - 1$, and the Halton sequence achieves $\alpha_s = s$.

Bounds on $D^*(\{u_i\}_{i=1}^n)$ are particularly useful since they provide bounds on the integration error for an estimate $\frac{1}{n} \sum_{i=1}^n f(u_i)$ of some real-valued function $f : [0, 1]^s \rightarrow \mathbb{R}$ whenever it has bounded Hardy-Krause variation, which will be denoted by $V_{\text{HK}}(f)$. Since the notation for the Hardy-Krause variation is rather involved, we refer the reader to Appendix A for details.

Related constructions are the *randomized QMC (RQMC) point sets*, which are sets of points $\{u_i\}_{i=1}^n$ with distribution $\text{Unif}([0, 1]^s)$ such that $\exists N, B > 0$ such that for $\forall n \geq N$, $D^*(\{u_i\}_{i=1}^n) \leq B(\log n)^{\alpha_s} n^{-1}$ with probability 1 for some $\alpha_s > 0$. The most common approach to construct these consists of “scrambling” a QMC point set, which consists of applying random transformations which preserve the low discrepancy structure. This allows those point sets to be used to obtain unbiased estimates of integrals of some functions against $[0, 1]^s$. Details on the construction of the scrambled points can be found in Chapter 17 in [67].

In the remainder, we will provide technical conditions on \mathcal{X} and G_θ so that for any D amongst the discrepancies previously mentioned and assuming we use n QMC points, we have

$$|D(\mathbb{P}_\theta, \mathbb{Q}^m) - D(\mathbb{P}_\theta^n, \mathbb{Q}^m)| = O(n^{-1}(\log n)^{\alpha_s}).$$

This is an improvement on the MC rate for which the rate would be $O(n^{-\frac{1}{2}})$. Since the cost of generating MC or QMC realisations is linear in the number of samples, a natural approach to balance the error in n and m of estimating $D(\mathbb{P}_\theta, \mathbb{Q})$ is to take n growing with \sqrt{m} . Note however that this optimal scaling is asymptotic and relies on a number of unknown constants dependent on the QMC point set used and the cost of evaluating the generator. This scaling will be studied further in the experiments.

3.2 Technical Assumptions

Before stating our sample complexity results, we introduce and discuss the assumptions that will be required. Our first assumption concerns the domain of the generator and the point sets:

Assumption 1. *Given a model \mathbb{P}_θ with generative process $(\text{Unif}([0, 1]^s), G_\theta)$, we assume we have access to $x_i = G_\theta(u_i)$ for $i = 1, \dots, n$ where $\{u_i\}_{i=1}^n \subset [0, 1]^s$ form a QMC or RQMC point set for some $\alpha_s > 0$. Furthermore, we write $\mathbb{P}_\theta^n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$.*

This assumption is very mild since it only assumes we can write the generative model in terms of a generator mapping from $[0, 1]^s$ (which is always possible due to Sklar’s theorem) and that we have access to a QMC or RQMC point set such as those mentioned above. Such point sets are widely available, for example in Python through the packages SciPy [84] and QMCPy [23].

For the MMD and Sinkhorn divergence results, we will also require a second assumption on the generator. For this, we will use the notation $a_v : b_{-v}$ to represent a point $u \in [a, b]^s$ with $u_j = a_j$ for $j \in v$, and $u_j = b_j$ for $j \notin v$; see Appendix A for more details.

Assumption 2. *The generator is a map $G_\theta : [0, 1]^s \rightarrow \mathcal{X}$ where:*

1. $\partial^{(1, \dots, 1)}(G_\theta)_j \in \mathcal{C}([0, 1]^s)$ for all $j = 1, \dots, d$.
2. $\partial^v(G_\theta)_j(\cdot : 1_{-v}) \in L^{p_j}([0, 1]^{|v|})$ for all $j = 1, \dots, d$ and $v \in \{0, 1\}^s \setminus (0, \dots, 0)$, where $p_j \in [1, \infty]$ and $\sum_{j=1}^d p_j^{-1} \leq 1$.

Assumption 2.1 is fairly straightforward and simply requires that the mixed partial derivative of the generator with respect to each coordinate is a continuous function, which is usually a condition which should be easy to verify (this needs to be done on a case-by-case basis). For example, in the case of neural network-based generators, the chain rule guarantees that this assumption will be satisfied whenever the activation functions are smooth enough. This is for example the case for the logistic, hyperbolic tangent, Gaussian, softplus and softmax activation functions which are all infinitely differentiable. However, neural generators with less regular activation functions such as the rectified linear unit will not satisfy the condition.

Assumption 2.2 requires certain integrability conditions for derivatives of the generator. When \mathcal{X} is compact, it follows directly from the first condition. However, this is not true when \mathcal{X} is not bounded and the requirement that $\sum_{j=1}^d p_j^{-d} \leq 1$ is slightly harder to satisfy in that case, especially for high dimensional problems. One straightforward, but rather restrictive, way of guaranteeing the condition is to enforce that derivatives of the form $\partial^v(G_\theta)_j(\cdot : 1_{-v})$ are all bounded. Alternatively, we could require that $\partial^v(G_\theta)_j \in L^{p_j}([0, 1]^{|v|})$ for all $j = 1, \dots, d$ and $v \in \{0, 1\}^s \setminus (0, \dots, 0)$, where $p_j \in [1, \infty]$ and $\sum_{j=1}^d p_j^{-1} \leq 1/2$; see Corollary 7 of [7] for a more detailed discussion. This holds for example when the generator has bounded derivatives.

3.3 Sample Complexity for Maximum Mean Discrepancy

We are now ready to present our sample complexity results. Our first set of results will provide sufficient conditions on k and G_θ to guarantee improved sample complexity by the use of (R)QMC point sets. We say that a kernel k is bounded if $\exists C > 0$ such that $\sup_{x, x' \in \mathcal{X}} |k(x, x')| \leq C$. Before presenting this result, we briefly recall a result using IID samples which will be used as a reference.

Proposition 1 (Lemma 1 in [16]). *Assume that k is bounded and let $\mathbb{P} \in \mathcal{P}_k(\mathcal{X})$. Let $\mathbb{P}^n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ where $\{x_i\}_{i=1}^n$ are IID realisations from \mathbb{P} . Then, with probability $1 - \delta$:*

$$\text{MMD}(\mathbb{P}, \mathbb{P}^n) = O(n^{-\frac{1}{2}}) \sqrt{\log(\delta^{-1})}.$$

We also only provide a simplified version of the statement which does not make the constants explicit for simplicity. It is also possible to obtain similar results for convergence of the MMD in the case of dependent realisations; see [22]. Although the rate in n is independent of dimensions, we will require a large number of samples in order to converge to zero due to the small exponent. The original statement is valid for finite n , but we present it in this asymptotic form for ease of comparison with the QMC/RQMC result below.

We now present a new sample complexity for MMD using QMC sequences. To do so, we need to show that the space of functions of the form $f \circ G_\theta$ for $f \in \mathcal{H}_k$ is continuously embedded into a space for which QMC can provide fast convergence rates. This is a challenging task, as was highlighted by [52], and we provide an auxiliary theorem for this (Theorem 4) in Appendix B.1. For this theorem to hold, we show that sufficient conditions can be obtained by ensuring that the generator G_θ and domain \mathcal{X} are regular enough.

Theorem 1. *Let $\mathbb{P}_\theta \in \mathcal{P}_k(\mathcal{X})$ and suppose Assumption 1 and 2 hold. Further assume that $k \in \mathcal{C}^{s,s}(\mathcal{X})$ and $\forall t \in \mathbb{N}_0^d, |t| \leq s, \sup_{x \in \mathcal{X}} \partial^{t,t} k(x, x) < C_k$ where C_k is some universal constant depending only on k . Then,*

$$\text{MMD}(\mathbb{P}_\theta, \mathbb{P}_\theta^n) = O(n^{-1}(\log n)^{\alpha_s}).$$

A direct implication is the following corollary, which follows from the triangle inequality.

Corollary 1. *Suppose the conditions in Theorem 1 hold. Then,*

$$|\text{MMD}(\mathbb{P}_\theta, \mathbb{Q}^m) - \text{MMD}(\mathbb{P}_\theta^n, \mathbb{Q}^m)| = O(n^{-1}(\log n)^{\alpha_s}).$$

The proof is in Appendix B.2. When using QMC, our result is only valid asymptotically in n , whereas for MC the result is also valid for finite n , although it only holds with probability $1 - \delta$. When using a RQMC point set the result above holds with probability 1 for finite but large enough n . As compared to Proposition 1, this theorem requires additional regularity from the generator (as per Assumption 2), but also smoothness for k . It does however provide a significantly faster convergence rate. The smoothness condition for the kernel is always satisfied in the case of the Gaussian kernel since it is infinitely differentiable; see Section 4.4. of [80]. Note that Theorem 1 has direct implications for the work of [43, 44], which considered the use of QMC sampling in the context of MMD generative adversarial networks.

3.4 Sample Complexity for the Wasserstein Distance

The main competitor to MMD for inference in generative models is the Wasserstein distance. An interesting question is therefore whether QMC can also lead to improved sample complexity results in this setting. We first recall a result for the case of MC realisations. Extensions of this result to dependent realisations can also be found in [36].

Proposition 2 (Theorem 1 in [36], simplified). *Let $\mathcal{X} = \mathbb{R}^d$, $p > 0$, c be a metric and $\mathbb{P} \in \mathcal{P}_{c,q}(\mathcal{X})$ for $q > p$. Let \mathbb{P}^n be the empirical measure obtained from n IID realisations of \mathbb{P} . Then,*

$$\mathbb{E} [W_{c,p}^p(\mathbb{P}, \mathbb{P}^n)] = \begin{cases} O\left(n^{-\frac{1}{2}} + n^{-\frac{(q-p)}{q}}\right) & \text{if } p > \frac{d}{2} \text{ and } q \neq 2p. \\ O\left(n^{-\frac{1}{2}} \log(1+n) + n^{-\frac{(q-p)}{q}}\right) & \text{if } p = \frac{d}{2} \text{ and } q \neq 2p. \\ O\left(n^{-\frac{p}{d}} + n^{-\frac{(q-p)}{q}}\right) & \text{if } p \in [1, \frac{d}{2}) \text{ and } q \neq \frac{d}{(d-p)}. \end{cases}$$

The result above is in expectation, but leads directly to a result in probability using Markov's inequality. This result shows a significant disadvantage of using the Wasserstein distance for

inference in generative models from a computational viewpoint: it suffers from the curse of dimensionality when p is small relative to d (the scenario most common in practice). Indeed, in the third case considered above the n required to estimate the distance accurately increases exponentially quickly with d .

The case most commonly considered in practice for inference in generative models is $p = 1$ (see for example [10, 11]), in which case the rate is $O(n^{-1/2})$ if $d = 1$, $O(n^{-1/2} \log(1+n))$ if $d = 2$, and $O(n^{-1/d})$ for $d \geq 3$. In the next result, we derive a novel result to show the impact of the use of QMC point sets to estimate the Wasserstein distance when $q = 1$, in which case the Wasserstein is an IPM. The proof is in Appendix B.3.

Theorem 2. *Let $\mathbb{P}_\theta \in \mathcal{P}_{c,1}(\mathcal{X})$ where $c(x, y) = \|x - y\|$ for some norm $\|\cdot\|$ on $\mathcal{X} \subseteq \mathbb{R}^d$. Suppose that Assumption 1 holds with $s = d = 1$, and assume that $V_{\text{HK}}(G_\theta) < \infty$. Then,*

$$W_{c,1}(\mathbb{P}_\theta, \mathbb{P}_\theta^n) = O(n^{-1}(\log n)^{\alpha_s}).$$

Since our goal is to approximate $W_{c,1}(\mathbb{P}_\theta, \mathbb{Q}^m)$ with $W_{c,1}(\mathbb{P}_\theta^n, \mathbb{Q}^m)$, we also consider:

Corollary 2. *Suppose the conditions in Theorem 2 hold. Then,*

$$|W_{c,1}(\mathbb{P}_\theta, \mathbb{Q}^m) - W_{c,1}(\mathbb{P}_\theta^n, \mathbb{Q}^m)| = O(n^{-1}(\log n)^{\alpha_s}).$$

We note that the assumption that $V_{\text{HK}}(G_\theta) < \infty$ is weaker than that imposed in Assumption 2, so that the discussion about sufficient conditions also holds here. This result shows that the convergence rate can be improved by a $O(n^{-1/2})$ term (up to logarithms) when using a RQMC/QMC point set instead of MC samples in $d = 1$ (once again, QMC results are only valid asymptotically). This is significant since in $d = 1$, the computational cost for the Wasserstein distance is $O(n \log n)$, which is significantly faster than the $O(n^2)$ cost for the MMD distance. For $d > 1$, the optimal rate for approximating an arbitrary distribution with a deterministic point set is $W_{c,1}(\mathbb{P}, \mathbb{P}^n) = O(n^{-1/d})$; see Theorem 2 in [64]. We therefore cannot hope to obtain an improved sample complexity result in this case.

Fortunately, this is not the end of the story. First, the $d = 1$ rate also transfers to sliced-Wasserstein distances in $d > 1$ using Theorem 2 in [59]. As we will see in the next section, the use of QMC and RQMC for the sliced-Wasserstein distance leads to very favourable computational costs, and warrants further study. Second, the next section will show that the Sinkhorn divergence can also be approximated at a fast rate even for $d > 1$.

3.5 Sample Complexity for the Sinkhorn Divergence

As for the other discrepancies, we will first review an existing result about the sample complexity of the Sinkhorn divergence with MC samples. Note that the result, which was proved in [37], is in terms of distance between estimated Sinkhorn divergence and the exact Sinkhorn divergence. Results of this form can be obtained from our theorems for the MMD and Wasserstein distance since they are both metrics and hence satisfy the triangle inequality, but here we are working with a divergence instead of a metric and so directly present the result in this form.

Proposition 3 (Corollary 1 in [37]). *Let $\mathbb{P}, \mathbb{Q} \in \mathcal{P}_{c,p}(\mathcal{X})$ on some bounded $\mathcal{X} \subset \mathbb{R}^d$, and suppose $c \in \mathcal{C}^{\infty,\infty}(\mathcal{X} \times \mathcal{X})$ is a Lipschitz continuous cost function. Let \mathbb{P}_n and \mathbb{Q}_n consist of n IID realisations from \mathbb{P} and \mathbb{Q} respectively. Then, with probability $1 - \delta$:*

$$|S_{c,p,\lambda}(\mathbb{P}, \mathbb{Q}) - S_{c,p,\lambda}(\mathbb{P}_n, \mathbb{Q}_n)| = O(n^{-\frac{1}{2}}) \sqrt{\log(\delta^{-1})}.$$

The constant in this rate depends on λ and d , and more detailed can be found in Theorem 3 of [37]. Most strikingly, the dependence on λ is exponential as $\lambda \rightarrow \infty$. See also [54] for a more refined result when using the squared Euclidean metric as cost function. Given a fixed value of λ and d , the rate in n is the MC rate. As we will see in the next results, this can be improved upon using QMC/RQMC point sets. Note here we need to restrict the domain to be compact, which is more restrictive than for our results for the MMD or Wasserstein distance, but is similar to the requirement in 3.

Assumption 3. *Assume that the domain $\mathcal{X} \subset \mathbb{R}^d$ is a compact space.*

This restriction for the domain is necessary since our proof builds on [37], which requires this assumption to hold. Although there are some results that allow the support of the distribution to be unbounded, for example in [54] where the compactness assumption was relaxed to distributions with sub-Gaussian tails on unbounded domains, this proof technique require us to enforce stronger regularity conditions for the generator which would limit the applicability of the proof.

Theorem 3. *Let $c \in \mathcal{C}^{\infty, \infty}(\mathcal{X} \times \mathcal{X})$ and suppose $\mathbb{P}_\theta, \mathbb{Q} \in \mathcal{P}_{c,p}(\mathcal{X})$. Furthermore, suppose Assumptions 1, 2 and 3 hold. Then*

$$|S_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{Q}^m) - S_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{Q}^m)| = O(n^{-1}(\log n)^{\alpha_s}).$$

The proof is available in Appendix B.4. Note that the rate is now the same as that possible when using QMC/RQMC for the MMD, and it significantly improves on what is possible when working with the Wasserstein distance.

4 Numerical Experiments

In this section, we will return to the uniform and Gaussian models first studied in Figure 1, then consider inference for intractable generative models including the multivariate g-and-k distributions, a flexible class of bivariate Beta distributions, and the deep neural network generator of a variational autoencoder. The aims of this section are two-fold. First, we will verify that the theoretical results in the previous section hold in practice. Second, we will look at QMC sampling in settings where Assumption 2 and 3 are violated. The requirements on the smoothness of G_θ and the assumption that \mathcal{X} is compact are rather restrictive but necessary to transfer existing theoretical results from the QMC theory to the setting of generative models. Thankfully, we will see that there are many settings where these assumptions are not satisfied but the approach nevertheless provides significant speed-ups. As such, our paper provides further evidence complementing the extensive discussion of this issue in Chapters 15, 16 and 17 of Art Owen’s book [67], and opens the way for further extensions of our theoretical results in Section 3.

Our simulation study uses the SciPy [84], JAX [15], QMCPy [23], POT [34] TensorFlow [1] libraries. The code can be found at

<https://github.com/johannnamr/Discrepancy-based-inference-using-QMC>.

Unless stated otherwise, all the RQMC results are based on generalised Halton or Sobol sequences which have been randomised using the scrambling factors of [32]. The approximation of sliced-distances are based on randomly sampled slices as described in Section 2. Additional results are provided in Appendix C.

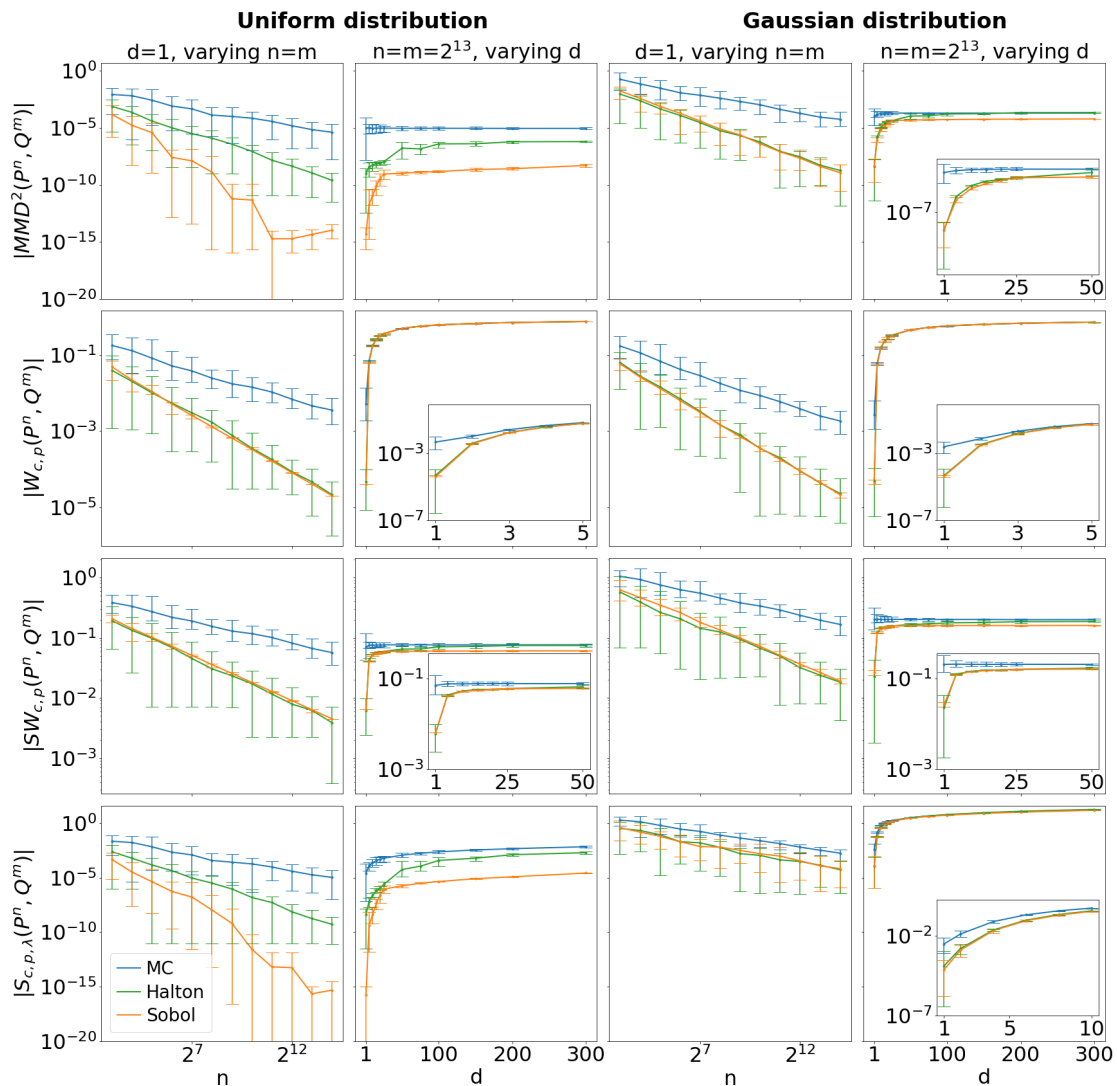


Figure 2: Sample Complexity for a Uniform and Gaussian model in MMD, Wasserstein distance, sliced-Wasserstein distance and Sinkhorn divergence.

4.1 Sample Complexity for Uniform and Gaussian Models

We first revisit the examples in Figure 1 which considered uniform and Gaussian distributions. These examples are of course very simple and do not require inference tools for generative models, but their simplicity allows us to study the sample complexity of QMC/RQMC in a wide range of scenarios. For the uniform distribution $\mathbb{P}_\theta = \text{Unif}([0, 1]^d)$, we will use $\mathbb{U} = \text{Unif}([0, 1]^s)$ with $G_\theta(u) = u$. For the Gaussian distribution $\mathbb{P}_\theta = \mathcal{N}(0, I)$, we use $\mathbb{U} = \text{Unif}([0, 1]^s)$ together with the inverse CDF of the univariate standard Gaussian Φ element-wise: $G_\theta(u) = (\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_d))$. The simulator G_θ does not depend on θ here since we only study the sample complexity results for a fixed distribution.

For these examples, we have $s = d$, $\mathbb{P}^n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$, $\mathbb{Q}^m = \frac{1}{m} \sum_{j=1}^m \delta_{y_j}$ and $n = m$, where

$\{x_i\}_{i=1}^n$ and $\{y_j\}_{j=1}^m$ are obtained through the generator G_θ . Our main results are presented in Figure 2, and include simulations with MC (in blue), RQMC with Halton sequences (in green) and RQMC with Sobol sequences (in orange). All the experiments have been repeated 25 times. The lines provide the average, and the error bars also represent intervals for the range of values observed. The smaller windows provide a zoomed-in plot for the cases where the gains in performance quickly reduce with d .

The first row computes $|\text{MMD}^2(\mathbb{P}^n, \mathbb{Q}^m)|$ when using a squared-exponential kernel with lengthscale $l = 1.5d^{1/2}$. This quantity should decrease as $O(n^{-1/2})$ (see Proposition 1) when using MC, and as $O(n^{-1}(\log n)^{\alpha_s})$ (see Theorem 1) when using RQMC. These rates clearly hold for both models when $d = 1$, and we see that RQMC quickly provides orders of magnitude improvements as n grows. For the uniform example, the Sobol sequence significantly outperforms the Halton sequence, but this is not the case for the Gaussian example. This is in line with theoretical results showing that the root-mean squared error for Sobol sequences can decrease as $O(n^{-3/2}(\log n)^{\alpha_s})$ [66], and could motivate further theoretical work extending the results in this paper. Significant improvements are also observed for larger values of d , although the gains (if any) are limited for $d > 100$ in the Gaussian case. This is not surprising since the Gaussian model does not satisfy the necessary conditions of Theorem 1 since \mathcal{X} is unbounded and the generator is not sufficiently regular (G_θ is unbounded, and as a result has infinite Hardy-Krause variation; see [67] Section 15.11). The lengthscale is adapted so as to increase with dimension; this is necessary as the distance between points grows exponentially with d due to the curse of dimensionality.

Additional experiments with the Matérn kernel with smoothness parameter $\nu = 3/2, 5/2$ and $7/2$ are also provided in Figure 12. We observe that the performance is significantly improved when using QMC points sets regardless of the choice of kernel, although this advantage decreases when d increases, and is larger for smoother kernels. This is interesting to see since the Matérn kernel does not satisfy the conditions of Theorem 1 when d is large. Finally, we notice from Figure 14 that the results are not very sensitive to the choice of QMC point set.

The second row of Figure 2 illustrates $|W_{c,p}(\mathbb{P}^n, \mathbb{Q}^m)|$ with $c(x, y) = \|x - y\|_2$ and $p = 1$. The QMC point sets lead to significant gains when $d = 1$, but not for larger d (a small advantage is seen until $d = 5$, but this is very limited). Further experiments for alternative choices of c and p can also be found in Figure 15, where similar results are observed. All of these results are consistent with what we would expect from Theorem 2, even though the regularity conditions of the theorem are not satisfied in the Gaussian case. The third row of Figure 2 illustrates $|SW_{c,p}(\mathbb{P}^n, \mathbb{Q}^m)|$ with $L = 100$ random slices when $c = \|x - y\|_2$ and $p = 1$. Clearly, we are able to obtain a gain in accuracy when using RQMC, and this is the case even for large d , which is a significant improvement on what is possible with the exact Wasserstein distance. Although the rate is $O(n^{-1}(\log n)^{\alpha_s})$ regardless of the value of d , the gains from using RQMC do become smaller in higher dimensions because the constant in this rate does still depend on d .

Finally, the fourth row of Figure 2 looks at the value of $|S_{c,p,\lambda}(\mathbb{P}^n, \mathbb{Q}^m)|$ with $\lambda = 2d$, $p = 2$ and $c(x, y) = \|x - y\|_2$. Once again, we observe that RQMC provides significant gains in performance in $d = 1$, but also for $d > 1$ in the case of the uniform. For the Gaussian, although the performance is improved to some extent for $d > 1$, these gains are really small. Interestingly, Figure 16 shows that the gains crucially depend on p and c , but also on the choice of λ . In particular, although Figure 2 could lead us to believe that there are close to no gains for the Gaussian case, this is clearly not the case when using an increased regularisation level λ .

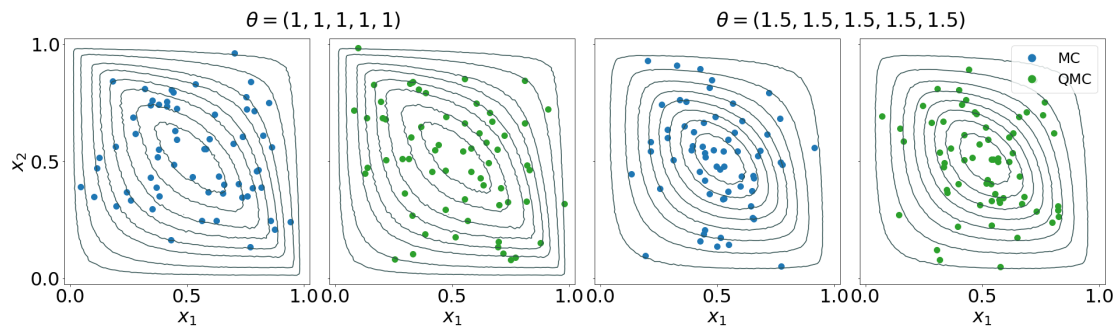


Figure 3: Realisation from the bivariate Beta model using MC and QMC point sets for $\theta = (1, 1, 1, 1, 1)$ and $\theta = (1.5, 1.5, 1.5, 1.5, 1.5)$.

4.2 Inference for Bivariate Beta Distributions

We now move on to studying discrepancy-based inference for intractable generative models with QMC and RQMC. Ever since the work of [65], there has been an interest in designing flexible classes of multivariate distributions which generalise the Beta distribution (as an indicator, [65] has over 180 citations to date). One popular approach is that of [5], which has been used by [25] to model household purchasing habits, and by [76] to model indicators of well-being. Although flexible, this does lead to an intractable density which makes inference challenging. We will focus on the $d = 2$ and $p = 5$ version of the model previously considered by [45, 62], and whose marginals are $\text{Beta}(\theta_1 + \theta_3, \theta_4 + \theta_5)$ and $\text{Beta}(\theta_2 + \theta_4, \theta_3 + \theta_5)$ distributed respectively in the first and second coordinate. In particular, denoting by $\lfloor x \rfloor$ the integer part of some $x \in \mathbb{R}$:

$$G_\theta^1(u) := \frac{\tilde{u}_1 + \tilde{u}_3}{\tilde{u}_1 + \tilde{u}_3 + \tilde{u}_4 + \tilde{u}_5}, \quad G_\theta^2(u) := \frac{\tilde{u}_2 + \tilde{u}_4}{\tilde{u}_2 + \tilde{u}_3 + \tilde{u}_4 + \tilde{u}_5}, \quad \tilde{u}_i = -\sum_{k=1}^{\lfloor \theta_i \rfloor} \ln(u_{ik}) + u_{i0},$$

where $u_{i0} \sim \text{Gamma}(\theta_i - \lfloor \theta_i \rfloor, 1)$, $u = (u_{11}, \dots, u_{1\theta_1}, u_{21}, \dots, u_{5\theta_5}) \sim \text{Unif}([0, 1]^s)$ and $s = \sum_{i=1}^5 \lfloor \theta_i \rfloor$. Note that the dimension s of the base space now depends on the value of θ .

In the special case where θ_i is an integer, $u_{i0} = 0$ is fixed (as opposed to sampled from a Gamma). In this case, both \mathcal{X} and G_θ satisfy the conditions in Assumptions 3 and 2. When this is not the case, u_{i0} can be generated through rejection sampling (see Appendix C.2). In that case, the generator does not satisfy Assumption 2 anymore, and also has a much higher-dimensional domain; i.e. $s = \sum_{i=1}^5 \lfloor \theta_i \rfloor + 15$. Here, the first term comes from the simulation of Gamma random variables with integer parameters $\lfloor \theta_1 \rfloor, \dots, \lfloor \theta_5 \rfloor$, and the second term is the dimensionality required to simulate five Gamma random variables with scalar parameters in $(0, 1)$ (that is, the simulation of a Gamma through rejection sampling requires a three-dimensional point). Despite these challenges, we will see below that certain gains in performance are still possible.

Figure 3 provides realisations from this model through MC and QMC sampling (in blue and green respectively). As observed, the QMC point set provides a slightly better coverage of the distribution, although the difference is not very large visually. In those cases, $s = 5$ for the left-hand side plot, whereas $s = 20$ for the right-hand side plot. We note that this case significantly differs from the examples in the previous section since we have $s \gg d$, which may partly explain why the difference is not as striking visually. However, looking at Figure 4 (which is the equivalent of Figure 2 for this model), we can see that QMC leads to a significant improvement in terms of sample complexity, especially in the case of integer parameter and to a lesser extent with scalar parameter values. Once again, this difference between left-hand side and right-hand side plot is

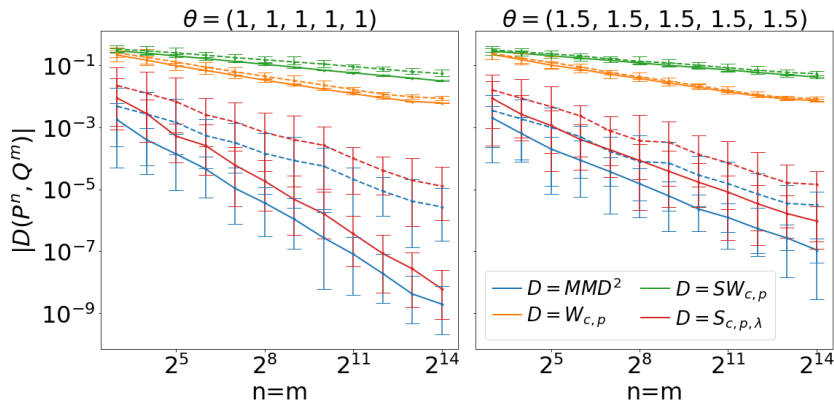


Figure 4: Sample complexity results for the bivariate Beta distribution with integer and scalar parameters. In both cases $d=2$, but the left-hand side plot uses $s = 5$ whereas the right-hand side plot uses $s = 20$. These differences in dimensionality seem to impact the convergence rate obtained through QMC point sets. Each solid line corresponds to RQMC, while the dashed lines correspond to MC point sets.

most likely due the difference in value of s , and the fact that Assumption 2 is not satisfied in the latter case. We also note that the advantage provided by QMC is particularly significant for the MMD and the sliced Wasserstein distance.

For the last part of this experiment, we perform inference for the parameter θ using an MDE approach with the MMD, the Wasserstein distance, the sliced Wasserstein distance and the Sinkhorn divergence. The generator G_θ is not differentiable in θ and we therefore propose to use a gradient-free global optimisation algorithm. We utilise the differential evolution algorithm due to [81], which is implemented as a sub-routine of the `optimize` function in the python library `SciPy` [84]. The dataset consists of $m = 2^{16}$ points, from which a minibatch of 2^{10} points is sampled at random at every iteration. Depending on the considered experiment, either n samples are generated using MC or $n, n^{\frac{3}{4}}, n^{\frac{2}{3}}$ or $n^{\frac{1}{2}}$ are simulated using RQMC at every iteration. The optimisation algorithm is run for 3,000 iterations for every setting. For the MMD, a squared-exponential kernel with lengthscale $l = 1.5d^{1/2}$ is used. The Wasserstein distance is computed with $c(x, y) = \|x - y\|_2$ and $p = 1$ as is the sliced Wasserstein distance based on 100 projections. The Sinkhorn divergence is considered with $\lambda = 5d$, $c(x, y) = \|x - y\|_2$ and $p = 2$. For the experiments, we focused on the case where $\theta^* = (\theta_1^*, \theta_2^*, \theta_3^*, \theta_4^*, \theta_5^*) = (1, 1, 1, 1, 1)$ as this was studied by [45, 62]. Therefore, the bounds, within each parameter is optimised by the differential evolution algorithm, are set to $[0, 2]$. We note that although the true parameter is integer valued, the optimisation algorithm will have to simulate data for parameter values which are scalar-valued. As a result, the dimensionality of the domain of the generator will generally be $s \in [15, 25]$ (assuming that the optimisation routine does not explore regions of the parameter space with large parameter values relative to θ^*).

The results of our experiments are presented in Figure 5, where we studied the computational cost and the accuracy of the estimates in l_2 norm for each choice of discrepancy. In each of these settings, we compared an MC method based on n points with an RQMC method with $n, n^{\frac{3}{4}}, n^{\frac{2}{3}}$ and $n^{\frac{1}{2}}$ points. As could be reasonably expected, the RQMC-based estimator with n points is significantly more expensive than an MC with n points, but it is also much more accurate in l_2 error. Similarly, the RQMC-based estimators with $n^{\frac{2}{3}}$ or $n^{\frac{1}{2}}$ are less accurate in l_2 error but usually cheaper than MC with n points. More interestingly, we see that the RQMC estimator

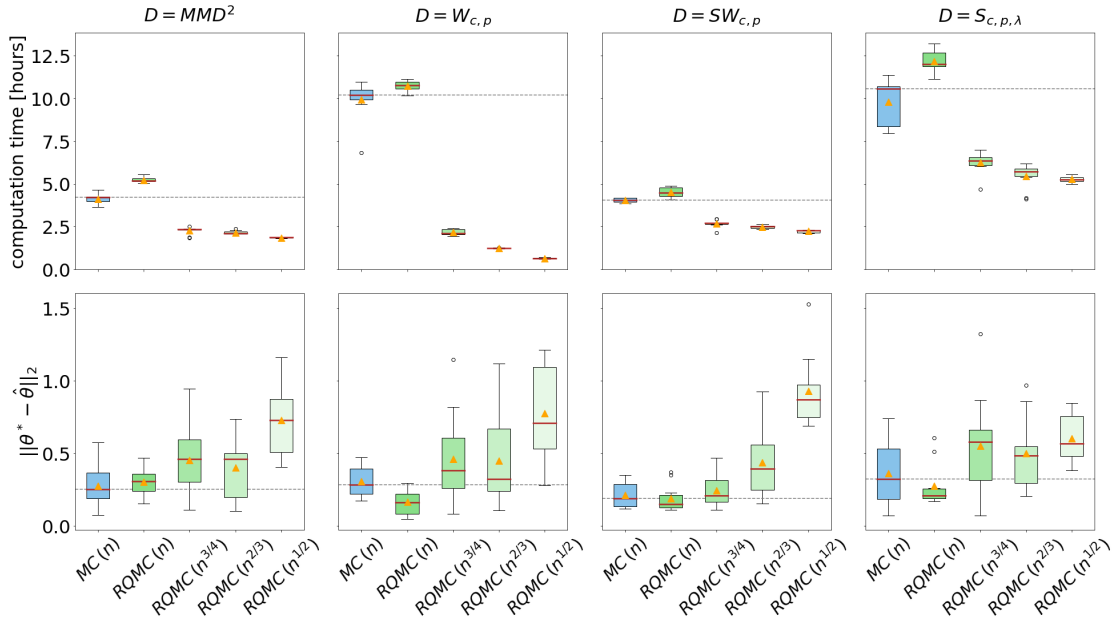


Figure 5: Minimum distance estimation for the parameters of the bivariate Beta distribution with the MMD, Wasserstein, sliced Wasserstein and Sinkhorn divergence. Each box-plot correspond to the result of 10 repetitions of an identical experiment, and the red bars provide the sample median. The dotted horizontal line correspond to the median value for the MC-based estimators with n points.

with $n^{\frac{3}{4}}$ points is both cheaper and more accurate than the MC estimator with n points for the Wasserstein distance, whilst for the Sinkhorn divergence it is cheaper and provides roughly the same level of accuracy. This clearly highlights that RQMC point sets can provide advantages even in cases not necessarily covered by our theoretical results. Surprisingly, this is not the case for the MMD, for which the performance of the RQMC-based estimator with $n^{\frac{3}{4}}$ is slightly worse than for the MC-based estimator in this experiment. We speculate that this may be due to a poor choice of kernel or an issue with the optimisation method since we obtained encouraging sample complexity results in Figure 4.

4.3 Inference for Multivariate g-and-k Models

Next, we will consider is the multivariate extension of the g-and-k distribution considered in [45, 62]. This parametric class is very flexible as it contains four parameters controlling the mean, variance, skewness and kurtosis of the marginals, as well as a fifth parameter controlling correlations across neighbouring coordinates. Unfortunately, inference is made challenging by the fact that the density is not available in closed-form. It is however straightforward to sample from this distribution, and it has recently become one of the most common target problems to assess the performance of inference schemes for generative models; see e.g. [74, 10, 45, 11, 16, 62, 27] for a small subset of recent papers using this model. The g-and-k has been applied to a range of applied problems, including (amongst others) insurance modelling [72], ranking and selection [42], and modelling of the prices of short-term rentals [75].

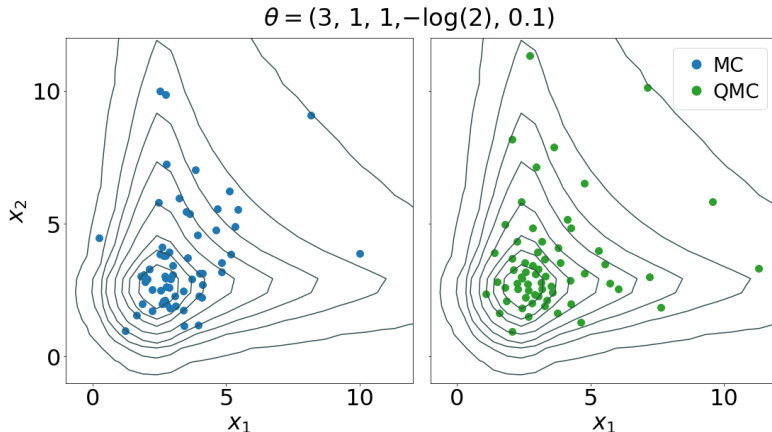


Figure 6: Realisations from the multivariate g-and-k distribution for $d = 2$. The scatter plots for MC and QMC are both based on $n = 2^6$ realisations, of which two fall outside the plotted interval for each point set. These points were omitted to allow for a zoomed-in view of the mode of the distribution.

The generator for this model is:

$$G_\theta(u) := \theta_1 + \theta_2 \left(1 + 0.8 \frac{(1 - \exp(-\theta_3 z))}{(1 + \exp(-\theta_3 z))} \right) (1 + z^2)^{\theta_4} z,$$

where $U = \text{Unif}([0, 1]^d)$, $z = \Sigma^{\frac{1}{2}} \Phi^{-1}(u)^\top$ where $\Phi^{-1}(u)$ is the inverse CDF of the univariate standard Gaussian distribution applied element-wise and $\Sigma \in \mathbb{R}^{d \times d}$ is a symmetric tri-diagonal Toeplitz matrix with diagonal entries all equal to 1 and off-diagonal entries equal to θ_5 . Its square-root can be obtained in closed form and is provided in Appendix C.3. Note that $s = d$ and we can straightforwardly replace MC realisations with a QMC or RQMC point set. Another important remark is that this generator does not satisfy the conditions of Assumption 2 since we are using the inverse CDF of a standard Gaussian. As parameter of interest, we consider $\theta = (\theta_1, \theta_2, \theta_3, \exp(\theta_4), \theta_5)$, where the rescaling of is used to avoid numerical instabilities during optimisation.

Figure 6 presents a scatter plot of two point sets of size $n = 2^6$ obtained through MC and RQMC in the case where $\theta = (3, 1, 1, -\log(2), 0.1)$. We can observe that the RQMC-based point set provides a better coverage of areas of high probability than the MC-based point set. These visual results also bare out in the estimates of the discrepancies in Figure 7, where we plot the sample complexity as a function of n for different values of d for the MMD (with squared-exponential kernel and lengthscale $l = 1.5d^{1/2}$), the sliced Wasserstein distance (with $c(x, y) = \|x - y\|_2$, $p = 1$ and 100 projections) and the Sinkhorn divergence (with $\lambda = 5d$, $c(x, y) = \|x - y\|_2$ and $p = 2$). Here, the Wasserstein distance is omitted due to the prohibitive computational cost when d is large.

In each case, the RQMC algorithms significantly outperform their MC counterpart, although this improved performance is limited for higher values of d . For example, in the case of MMD, the RQMC rates were of the form $n^{-\alpha}$ with α equal to 0.73, 0.65, 0.58 and 0.54 in dimensions 5, 10, 25 and 50 respectively, whereas α was approximately 0.5 for MC in all cases. This is in line with what we would expect following the results of Section 4.1 where the use of the inverse CDF of a Gaussian was studied in detail.

In the last part of our experiments for the multivariate g-and-k distribution, we adapt a

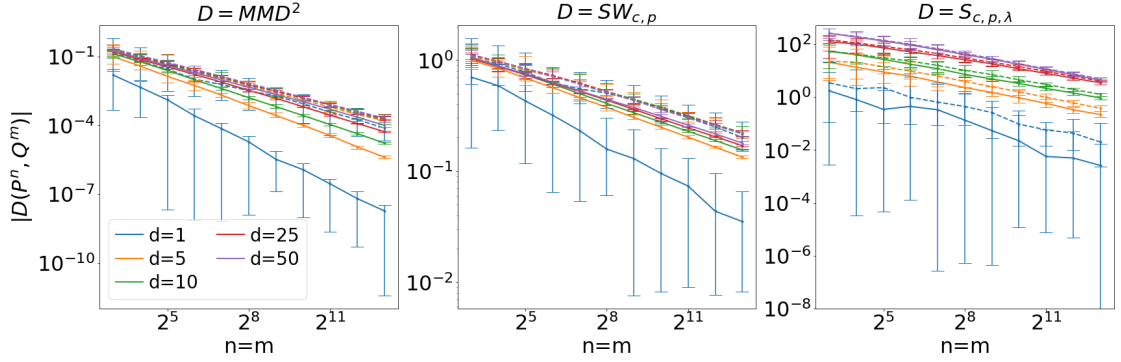


Figure 7: Sample complexity results for the multivariate g-and-k distribution in various dimensions d . The Wasserstein distance is omitted due to the fact that the discrepancy is prohibitively expensive when $d > 1$ and n is large. Each solid line corresponds to RQMC, whereas the dashed lines correspond to MC point sets.

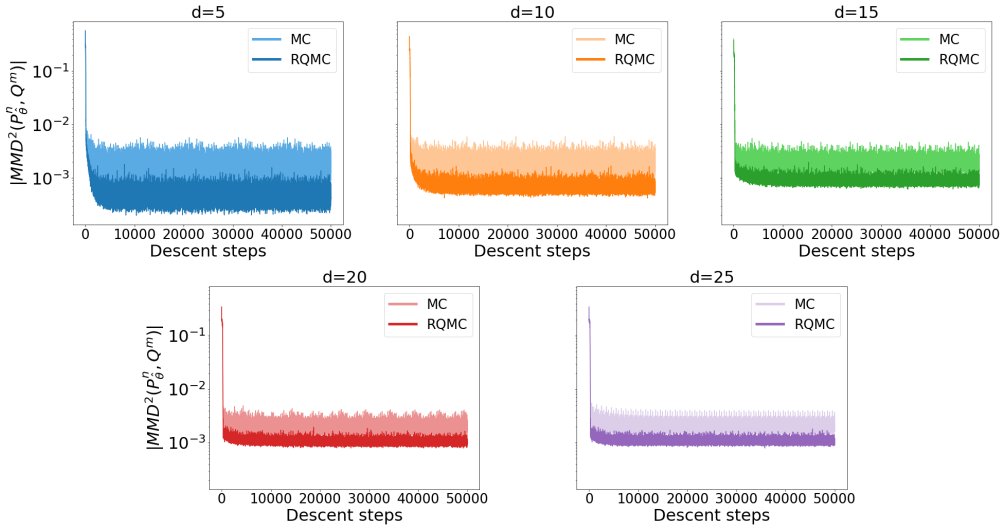


Figure 8: Minimum distance estimation with the MMD for the multivariate g-and-k distribution. The figure plots the estimated MMD between the model with the estimated parameter and the data as a function of the number of stochastic gradient descent steps.

gradient-based optimisation method to perform inference for the parameter θ using an MDE approach that builds on the MMD. The considered stochastic gradient descent (SGD) algorithm is similar to the one of [16], but uses an approximation of the squared MMD using empirical measures as in (2) instead of a U-statistic approximation. From $m = 2^{16}$ data points, a minibatch of 2^{11} points is sampled for every descent step. Using either the MC or QMC approach, $n = 2^9$ data points are simulated for each descent step. The step size of the SGD algorithm is fixed at 0.2 (for both MC and QMC) and the optimisation is run for 50,000 descent steps. The squared MMD and its gradient are computed based on the squared-exponential kernel with lengthscale $l = 1.5d^{1/2}$. To obtain the gradient of the multivariate g-and-k distribution, we make

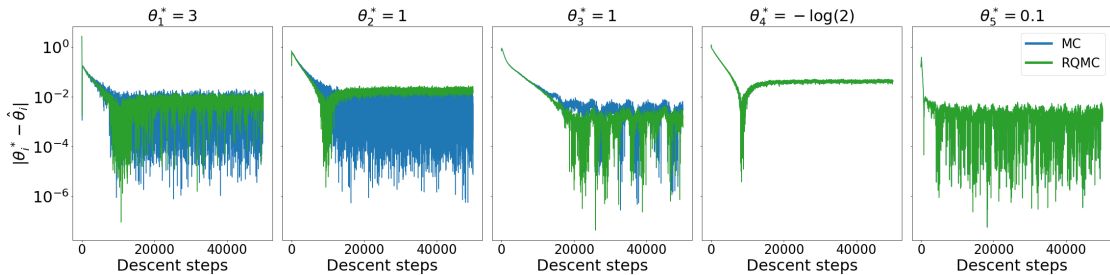


Figure 9: Minimum distance estimation with the MMD for the multivariate g-and-k distribution. The figure plots the l_1 error between the estimated parameter and the true value as a function of the number of stochastic gradient descent steps.

use of automatic differentiation provided by the python library JAX [15]. The experiments aim at retrieving the true parameter $\theta^* = (\theta_1^*, \theta_2^*, \theta_3^*, \exp(\theta_4^*), \theta_5^*) = (3, 1, 1, -\log(2), 0.1)$ and start the SGD algorithm at $\theta^0 = (0.3, 0.3, 0.3, 0.3, 0.3)$.

Figure 8 illustrates the results for $|\text{MMD}^2(\mathbb{P}_\theta^n, \mathbb{Q}^m)|$ as a function of the number of descent steps of the stochastic gradient descent method, where \mathbb{Q}^m here corresponds to the entire original data (i.e. $m = 2^{16}$). The experiment is repeated for a range of values of d between 5 and 25. As we can observe, the estimated MMD is much more accurate when sampling is done with RQMC. In Figure 9, we then look at the case of $d = 5$ in more details. In particular, the figure shows the l_1 distance between the true and estimated parameters of the g-and-k as a function of the number of descent steps. It is overall unclear which of RQMC and MC outperforms the other, and this depends on which parameters are of most interest. RQMC seems to outperform MC for θ_3 , performs equally well as MC for θ_4 and θ_5 (the curves overlap), and tends to do worse for θ_2 . For all parameters, the jumps in l_1 error between descent steps is much larger for MC than RQMC, highlighting that RQMC estimates have a much smaller variance. The contrast between Figure 8 and Figure 9 highlights that minimisation of a discrepancy does not necessarily mean that the estimates for all parameter values will be accurate. In fact, Figure 21 in Appendix C.3 actually shows that RQMC actually has a worse performance than MC as d grows when looking at the results in terms of l_2 errors instead of MMD (as in Figure 8). In this case, we expect that such a counter-intuitive result is due to the value of m being too small relative to that of n for RQMC, which could lead to over-fitting.

4.4 Inference for Generative Neural Networks

Our final model is a generative neural network which was trained using the Sinkhorn divergence by [37]. More precisely, this model is the decoder network of a variational autoencoder (VAE) given by $G_\theta : \mathcal{U} \rightarrow \mathcal{X}$ with $\mathcal{U} = [0, 1]^2$ and $\mathcal{X} = [0, 1]^{784}$ (i.e. $s = 2$ and $d = 784$) where:

$$G_\theta(u) = \phi_2(\phi_1(\phi_1(u^\top W^1 + b^1)^\top W^2 + b^2)^\top W^3 + b^3)$$

and θ is a vector containing all entries of the weight matrices $W^1 \in \mathbb{R}^{2 \times 500}$, $W^2 \in \mathbb{R}^{500 \times 500}$, $W^3 \in \mathbb{R}^{500 \times 784}$ and biases $b^1 \in \mathbb{R}^{500}$, $b^2 \in \mathbb{R}^{500}$, $b^3 \in \mathbb{R}^{784}$ so that $p = 644784$. Additionally, $\phi_1(x) = \log(\exp(x) + 1)$ (a softplus activation function) and $\phi_2(x) = (1 + \exp(-x))^{-1}$ (a logistic activation function), and the output of the generator is a 784-dimensional vector which can be rescaled to form a 28×28 pixel image. Since G_θ is the composition of smooth functions, it is itself smooth. Furthermore, since \mathcal{X} is bounded, the derivatives of G_θ (which are continuous) must also be bounded, and G_θ therefore satisfies Assumption 2, and hence our theorems hold.

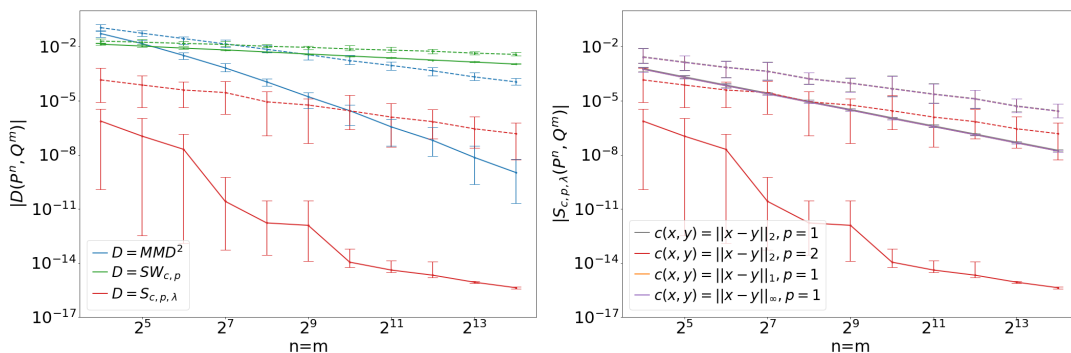


Figure 10: Sample complexity results for the generative neural network. Each solid line corresponds to RQMC, whereas the dashed lines correspond to MC.

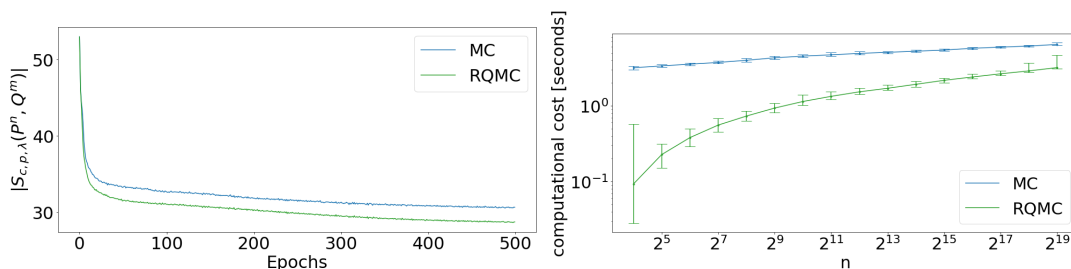


Figure 11: Left: Minimum distance estimation with the Sinkhorn divergence for the VAE. The figure plots the Sinkhorn divergence between the model with the estimated parameters and the data as a function of the number of training epochs. Right: Comparison of the computational cost for simulating from the generative neural network using MC and RQMC. Each line represents the average of 500 repetitions and the error bars give the minimum and maximum values observed.

In the right-hand side plot of Figure 10, the sample complexity is plotted as a function of n for the MMD (with squared-exponential kernel and lengthscale $l = 0.01$), the sliced Wasserstein distance (with $c(x, y) = \|x - y\|_2$, $p = 1$ and 100 projections), and the Sinkhorn divergence (with $c(x, y) = \|x - y\|_2$, $p = 2$ and $\lambda = 1$). Here, the Wasserstein distance is omitted due to the prohibitive computational cost in high dimensions. We observe that QMC leads to significant improvements in sample complexity, especially for the MMD and Sinkhorn divergence.

Comparing the sample complexity for the Sinkhorn divergence with different choices of cost c and order p in the right-hand side plot of Figure 10, we find that the choice of squared Euclidean cost, i.e. $c(x, y) = \|x - y\|_2$ and $p = 2$, significantly outperforms the other considered choices. It is therefore used in the following experiments.

In the final experiment, the generative neural network is trained as the decoder network of a VAE on the MNIST dataset for 500 epochs with mini-batches of size 300 using the Adam optimizer [49]. This MDE approach is based on the Sinkhorn divergence with parameters $\lambda = 1$, $c(x, y) = \|x - y\|_2$ and $p = 2$ and a dataset of size $m = n = 55,000$. This setup corresponds to the one used by [38]. The implementation of this experiment uses the python library TensorFlow [1] and SciPy[84] to generate Sobol points. Using RQMC sampling, we observe in the left-hand side plot of Figure 11 that the training loss decreases significantly faster in the number of training epochs than when using MC sampling.

The right-hand side plot of Figure 11 compares the computational cost of simulating from the generative neural network, which implies sampling with $d = 2$. We observe that RQMC sampling is much cheaper than MC for all considered n .

5 Conclusion

This paper focused on the use of QMC and RQMC point sets for discrepancy-based inference in intractable generative models. We showed (in Theorems 1, 2 and 3) that the sample complexity becomes $O(n^{-1}(\log n)^{\alpha_s})$ instead of $O(n^{-1/2})$ for the MMD and Sinkhorn divergence in arbitrary dimension d . These faster rates can provide significant improvements on the current state-of-the-art with no significant increases in computational cost (since QMC point sets can be pre-computed). Unfortunately, the rate for the Wasserstein-1 distance can only be improved when $d = 1$, and is otherwise gated at $O(n^{-1/d})$ due to a well-known curse of dimensionality. However, we showed that the recently introduced sliced-Wasserstein distance can obtain the optimal $O(n^{-1}(\log n)^{\alpha_s})$ rate regardless of d .

One significant drawback of our results is that they not only require the generator to satisfy certain regularity conditions (see Assumption 2), but also that \mathcal{X} is compact (see Assumption 3). These are common assumptions for the QMC literature (see the discussion in [67]), but these nonetheless exclude many cases of practical interest. Despite these limitations, we showed in Section 4 that QMC/RQMC can still provide significant gains when the assumptions do not hold; for example when using the inverse transform approach to sampling from Gaussian distributions (which has an unbounded generator) and sampling Gamma random variables through rejection sampling. This is in line with work in the QMC literature (see for example [66]) and future work could explore these cases from a theoretical viewpoint in more detail.

Another potential line of future research would be to explore the use of other point sets, including weighted point sets, for inference in generative models. This was recently studied in the context of the Sinkhorn divergence by [12], who use quantization to improve sample qualities. However, alternative approaches could also be used. For example, Bayesian quadrature [17] is known to provide optimally weighted point set for the MMD, and could lead to faster sample complexity results. We expect that such approaches could provide significant improvements in performance, particularly in cases of computationally expensive generators. Higher-order digital nets could also be used to provide dimension-independent convergence rates, albeit with further assumptions on the generator. In this respect, one could think of adapting the architecture of deep generative models as well as the choice of discrepancy so as to ensure that such fast rates can be obtained.

Acknowledgments

The authors are grateful to Chris Oates and two anonymous reviewers for helpful comments and suggestions on this paper. FXB was supported by the Lloyd’s Register Foundation programme on data-centric engineering at The Alan Turing Institute under the EPSRC grant [EP/N510129/1].

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu,

- and X. Zheng. “TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems”. In: *arXiv:1603.04467* (2015). Software available from tensorflow.org.
- [2] R. A. Adams and J. J. F. Fournier. *Sobolev spaces*. 2nd edition. Vol. 140. Academic Press, 2006.
- [3] J. H. Ahrens and U. Dieter. “Computer methods for sampling from Gamma, Beta, Poisson and Binomial distributions”. In: *Computing* 12.3 (1974), pp. 223–246.
- [4] P. Alquier and M. Gerber. “Universal robust regression via maximum mean discrepancy”. In: *arXiv:2006.00840* 1 (2020).
- [5] B. C. Arnold and H. K. T. Ng. “Flexible bivariate beta distributions”. In: *Journal of Multivariate Analysis* 102.8 (2011), pp. 1194–1202.
- [6] F. Bassetti, A. Bodini, and E. Regazzini. “On minimum Kantorovich distance estimators”. In: *Statistics & Probability Letters* 76.12 (2006), pp. 1298–1302.
- [7] K. Basu and A. B. Owen. “Transformations and Hardy–Krause variation”. In: *SIAM Journal on Numerical Analysis* 54.3 (2016), pp. 1946–1966.
- [8] M. A. Beaumont. “Approximate Bayesian computation in evolution and ecology”. In: *Annual Review of Ecology, Evolution, and Systematics* 41.1 (2010), pp. 379–406.
- [9] M. A. Beaumont, W. Zhang, and D. J. Balding. “Approximate Bayesian computation in population genetics”. In: *Genetics* 162.4 (2002), pp. 2025–2035.
- [10] E. Bernton, P. E. Jacob, M. Gerber, and C. P. Robert. “Inference in generative models using the Wasserstein distance”. In: *Information and Inference* 8.4 (2017), pp. 657–676.
- [11] E. Bernton, P. E. Jacob, M. Gerber, and C. P. Robert. “Approximate Bayesian computation with the Wasserstein distance”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 81.2 (2019), pp. 235–269.
- [12] G. Beugnot, A. Genevay, K. Greenewald, and J. Solomon. “Improving approximate optimal transport distances using quantization”. In: *arXiv:2102.12731* (2021).
- [13] A. Bharti, F.-X. Briol, and T. Pedersen. “A general method for calibrating stochastic radio channel models with kernels”. In: *arXiv:2012.09612*. To appear in *IEEE Transactions in Antennas and Propagation* (2020).
- [14] M. Bińkowski, D. J. Sutherland, M. Arbel, and A. Gretton. “Demystifying MMD GANs”. In: *International Conference on Learning Representations*. 2018.
- [15] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. *JAX: composable transformations of Python+NumPy programs*. Version 0.2.5. 2018.
- [16] F.-X. Briol, A. Barp, A. B. Duncan, and M. Girolami. “Statistical inference for generative models with maximum mean discrepancy”. In: *arXiv:1906.05944* (2019).
- [17] F.-X. Briol, C. J. Oates, M. Girolami, M. A. Osborne, and D. Sejdinovic. “Probabilistic integration: a role in statistical computation?” In: *Statistical Science* 34.1 (2019), pp. 1–22.
- [18] A. Buchholz and N. Chopin. “Improving approximate Bayesian computation via quasi-Monte Carlo”. In: *Journal of Computational and Graphical Statistics* 28.1 (2019), pp. 205–219.
- [19] M. Cambou, M. Hofert, and C. Lemieux. “Quasi-random numbers for copula models”. In: *Statistics and Computing* 27.5 (2017), pp. 1307–1329.

- [20] E. Cameron and A. N. Pettitt. “Approximate Bayesian computation for astronomical model analysis: a case study in galaxy demographics and morphological transformation at high redshift”. In: *Monthly Notices of the Royal Astronomical Society* 425.1 (2012), pp. 44–65.
- [21] B.-E. Cherief-Abdellatif and P. Alquier. “MMD-Bayes: robust Bayesian estimation via maximum mean discrepancy”. In: *Proceedings of The 2nd Symposium on Advances in Approximate Bayesian Inference*. Vol. 118. PMLR, 2020, pp. 1–21.
- [22] B.-E. Chérif-Abdellatif and P. Alquier. “Finite sample properties of parametric MMD estimation: robustness to misspecification and dependence”. In: *Bernoulli (to appear)* (2021).
- [23] S. C. Choi, F. J. Hickernell, M. McCourt, J. Rathinavel, and A. Sorokin. *QMCPy: a quasi-Monte Carlo Python Library*. 2020.
- [24] G. M. Constantine and T. H. Savits. “A multivariate Faa di Bruno formula with applications”. In: *Transactions of the American Mathematical Society* 348.2 (1996), pp. 503–520.
- [25] R. Crackel and J. Flegal. “Bayesian inference for a flexible class of bivariate beta distributions”. In: *Journal of Statistical Computation and Simulation* 87.2 (2017), pp. 295–312.
- [26] K. Cranmer, J. Brehmer, and G. Louppe. “The frontier of simulation-based inference”. In: *Proceedings of the National Academy of Sciences* 117.48 (2020).
- [27] C. Dellaporta, J. Knoblauch, T. Damoulas, and F.-X. Briol. “Robust Bayesian inference for simulator-based models via the MMD posterior bootstrap”. In: *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*. 2022, pp. 943–970.
- [28] I. Deshpande, Z. Zhang, and A. G. Schwing. “Generative modeling using the sliced Wasserstein distance”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [29] L. Devroye. *Non-uniform random variate generation*. New York et al.: Springer, 1986.
- [30] J. Dick and F. Pillichshammer. *Digital nets and sequences: discrepancy theory and quasi-Monte Carlo integration*. Cambridge: Cambridge University Press, 2010.
- [31] G. K. Dziugaite, D. M. Roy, and Z. Ghahramani. “Training generative neural networks via maximum mean discrepancy optimization”. In: *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*. 2015, pp. 258–267.
- [32] H. Faure and C. Lemieux. “Generalized Halton sequences in 2008: A comparative study”. In: *ACM Transactions on Modeling and Computer Simulation* 19.4 (2009), pp. 1–31.
- [33] J. Feydy, T. Séjourné, F.-X. Vialard, S.-I. Amari, A. Trounev, and G. Peyré. “Interpolating between optimal transport and MMD using Sinkhorn divergences”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR 89. 2019, pp. 2681–2690.
- [34] R. Flamary, N. Courty, A. Gramfort, M. Z. Alaya, A. Boisbunon, S. Chambon, L. Chapel, A. Corenflos, K. Fatras, N. Fournier, L. Gautheron, N. T. Gayraud, H. Janati, A. Rakotomamonjy, I. Redko, A. Rolet, A. Schutz, V. Seguy, D. J. Sutherland, R. Tavenard, A. Tong, and T. Vayer. “POT: Python Optimal Transport”. In: *Journal of Machine Learning Research* 22.78 (2021), pp. 1–8.
- [35] J.-J. Forneron. “A scrambled method of moments”. In: *arXiv:1911.09128* (2019).
- [36] N. Fournier and A. Guillin. “On the rate of convergence in Wasserstein distance of the empirical measure”. In: *Probability Theory and Related Fields* 162.3-4 (2015), pp. 707–738.

- [37] A. Genevay, L. Chizat, F. Bach, M. Cuturi, and P. Gabriel. “Sample complexity of Sinkhorn divergences”. In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. 2019, pp. 1574–1583.
- [38] A. Genevay, G. Peyre, and M. Cuturi. “Learning generative models with Sinkhorn divergences”. In: *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*. 2018, pp. 1608–1617.
- [39] P.-O. Goffard and P. J. Laub. “Approximate Bayesian computations to fit and compare insurance loss models”. In: *arXiv:2007.03833* (2020).
- [40] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. “Generative adversarial nets”. In: *Advances in Neural Information Processing Systems*. Vol. 27. 2014, pp. 2672–2680.
- [41] A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola. “A kernel method for the two-sample-problem”. In: *Advances in Neural Information Processing Systems*. Vol. 19. 2006, pp. 513–520.
- [42] M. Haynes, H. MacGillivray, and K. Mengersen. “Robustness of ranking and selection rules using generalised g-and-k distributions”. In: *Journal of Statistical Planning and Inference* 65 (1997), pp. 45–66.
- [43] M. Hofert, A. Prasad, and M. Zhu. “Quasi-random sampling for multivariate distributions via generative neural networks”. In: *Journal of Computational and Graphical Statistics* (2021).
- [44] M. Hofert, A. Prasad, and M. Zhu. “Applications of multivariate quasi-random sampling with neural networks”. In: *arXiv:2012.08036* (2020).
- [45] B. Jiang, T. Y. Wu, and W. H. Wong. “Approximate Bayesian computation with Kullback-Leibler divergence as data discrepancy”. In: *International Conference on Artificial Intelligence and Statistics*. 2018, pp. 1711–1721.
- [46] T. Kajihara, K. Yamazaki, M. Kanagawa, and K. Fukumizu. “Kernel recursive ABC: Point estimation with intractable likelihood”. In: *International Conference on Machine Learning*. 2018, pp. 2400–2409.
- [47] M. Kanagawa, P. Hennig, D. Sejdinovic, and B. K. Sriperumbudur. “Gaussian processes and kernel methods: A review on connections and equivalences”. In: *arXiv:1807.02582* (2018).
- [48] B. Karimi, B. Miasojedow, E. Moulines, and H.-T. Wai. “Non-asymptotic analysis of biased stochastic approximation scheme”. In: *Conference on Learning Theory*. 2019.
- [49] D. P. Kingma and J. Ba. “Adam: A method for stochastic optimization”. In: *arXiv:1412.6980* (2014).
- [50] S. Kolouri, K. Nadjahi, U. Simsekli, and S. Shahrampour. “Generalized sliced distances for probability distributions”. In: *arxiv:2002.12537* (2020).
- [51] C.-L. Li, W.-C. Chang, Y. Cheng, Y. Yang, and B. Poczos. “MMD GAN: towards deeper understanding of moment matching network”. In: *Advances in Neural Information Processing Systems*. Vol. 30. 2017, pp. 2203–2213.
- [52] Y. Li, L. Kang, and F. J. Hickernell. “Is a transformed low discrepancy design also low discrepancy?” In: *Contemporary experimental design, multivariate analysis and data mining*. Springer, 2020, pp. 69–92.

- [53] Y. Li, K. Swersky, and R. Zemel. “Generative moment matching networks”. In: *Proceedings of the 32nd International Conference on Machine Learning*. Vol. 37. 2015, pp. 1718–1727.
- [54] G. Mena and J. Niles-Weed. “Statistical bounds for entropic optimal transport: sample complexity and the central limit theorem”. In: *Advances in Neural Information Processing Systems*. Vol. 32. 2019, pp. 4541–4551.
- [55] J. Mitrovic, D. Sejdinovic, and Y. W. Teh. “DR-ABC: Approximate Bayesian computation with kernel-based distribution regression”. In: *International Conference on Machine Learning* 3 (2016), pp. 2209–2218.
- [56] S. Mohamed and B. Lakshminarayanan. “Learning in implicit generative models”. In: *arXiv:1610.03483* (2016).
- [57] A. Muller. “Integral probability metrics and their generating classes of functions”. In: *Advances in Applied Probability* 29.2 (1997), pp. 429–443.
- [58] K. Nadjahi, V. De Bortoli, A. Durmus, R. Badeau, and U. Şimşekli. “Approximate bayesian computation with the sliced-wasserstein distance”. In: *EEE International Conference on Acoustics, Speech and Signal Processing* (2020), pp. 5470–5474.
- [59] K. Nadjahi, A. Durmus, L. Chizat, S. Kolouri, S. Shahrampour, and U. Şimşekli. “Statistical and topological properties of sliced probability divergences”. In: *Neural Information Processing Systems*. 2020.
- [60] K. Nadjahi, A. Durmus, U. Simsekli, and R. Badeau. “Asymptotic guarantees for learning generative models with the sliced-Wasserstein distance”. In: *Advances in Neural Information Processing Systems*. Vol. 32. 2019, pp. 250–260.
- [61] S. Nakagome, K. Fukumizu, and S. Mano. “Kernel approximate Bayesian computation in population genetic inferences”. In: *Statistical Applications in Genetics and Molecular Biology* 12.6 (2013), pp. 667–678.
- [62] H. D. Nguyen, J. Arbel, H. Lu, and F. Forbes. “Approximate Bayesian computation via the energy statistic”. In: *IEEE Access* 8 (2020), pp. 131683–131698.
- [63] K. Nguyen, N. Ho, T. Pham, and H. Bui. “Distributional sliced-Wasserstein and applications to generative modeling”. In: *International Conference on Learning Representations*. 2021.
- [64] E. Novak. “Some results on the complexity of numerical integration”. In: *Monte Carlo and quasi-Monte Carlo methods*. Vol. 163. Springer Proceedings in Mathematics & Statistics. Springer, 2016, pp. 161–183.
- [65] I. Olkin and R. Liu. “A bivariate beta distribution”. In: *Statistics and Probability Letters* 62.4 (2003), pp. 407–412.
- [66] A. B. Owen. “Halton sequences avoid the origin”. In: *SIAM Review* 48.3 (2006), pp. 487–503.
- [67] A. B. Owen. *Monte Carlo Theory, Methods and Examples*. 2013.
- [68] A. B. Owen. “Multidimensional variation for quasi-Monte Carlo”. In: *International Conference on Statistics in honour of Professor Kai-Tai Fang’s 65th birthday*. Ed. by J. Fan and G. Li. 2005, pp. 49–74.
- [69] L. Pacchiardi and R. Dutta. “Generalized Bayesian likelihood-free inference using scoring rules estimators”. In: *arXiv:2104.03889* (2021).

- [70] M. Park, W. Jitkrittum, and D. Sejdinovic. “K2-ABC: approximate Bayesian computation with kernel embeddings”. In: *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*. Vol. 51. PMLR, 2016, pp. 398–407.
- [71] W. C. Parr and W. R. Schucany. “Minimum distance and robust estimation”. In: *Journal of the American Statistical Association* 75.371 (1980), pp. 616–624.
- [72] G. Peters, W. Chen, and R. Gerlach. “Estimating quantile families of loss distributions for non-life insurance modelling via L-moments”. In: *Risks* 4.2 (2016), p. 14.
- [73] G. Peyré and M. Cuturi. *Computational optimal transport: with applications to data science*. Foundations and Trends in Machine Learning. 2019.
- [74] D. Prangle. “gk: An R Package for the g-and-k and generalised g-and-h Distributions”. In: *arXiv:1706.06889* (2017).
- [75] G. S. Rodrigues, D. J. Nott, and S. A. Sisson. “Likelihood-free approximate Gibbs sampling”. In: *Statistics and Computing* 30.4 (2020), pp. 1057–1073.
- [76] J. M. Sarabia, F. Prieto, and V. Jordá. “Bivariate beta-generated distributions with applications to well-being data”. In: *Journal of Statistical Distributions and Applications* 1.15 (2014).
- [77] S. M. Schmon, P. W. Cannon, and J. Knoblauch. “Generalized posteriors in Approximate Bayesian Computation”. In: *3rd symposium on Advances in Approximate Bayesian*. 2020, pp. 1–11.
- [78] Z. Shen, Z. Wang, A. Ribeiro, and H. Hassani. “Sinkhorn natural gradient for generative models”. In: *Advances In Neural Information Processing Systems*. 2020, pp. 1646–1656.
- [79] B. K. Sriperumbudur, A. Gretton, K. Fukumizu, B. Schölkopf, and G. R. G. Lanckriet. “Hilbert space embeddings and metrics on probability measures”. In: *Journal of Machine Learning Research* 11 (2010).
- [80] I. Steinwart, A. Christmann, M. Jordan, J. Kleinberg, and B. Schölkopf. *Support vector machines*. Information Science and Statistics. Dordrecht: Springer, 2008.
- [81] R. Storn and K. Price. “Differential evolution – a simple and efficient heuristic for global optimization over continuous spaces”. In: *Journal of Global Optimization* 11.4 (1997), pp. 341–359.
- [82] D. J. Sutherland, H.-Y. Tung, H. Strathmann, S. De, A. Ramdas, A. Smola, and A. Gretton. “Generative models and model criticism via optimized maximum mean discrepancy”. In: *International Conference on Learning Representations*. 2017.
- [83] V. B. Tadic and A. Doucet. “Asymptotic bias of stochastic gradient search”. In: *Annals of Applied Probability* 27.6 (2017), pp. 3255–3304.
- [84] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors. “SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python”. In: *Nature Methods* 17 (2020), pp. 261–272.
- [85] J. Wu, Z. Huang, D. Acharya, W. Li, J. Thoma, D. P. Paudel, and L. van Gool. “Sliced Wasserstein generative models”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019.

Appendix

First, in Appendix A, we recall relevant background material on QMC. Then, in Appendix B, we provide all the proofs for the results in the main text. In Appendix C, we provide additional numerical experiments to complement the results in the main text.

A Additional Background

For completeness, we first recall several definitions and results which are relevant for QMC. Our presentation closely follows [68], and we refer the reader to this paper for further details. For some vector $u \in \mathbb{R}^s$, we will denote its j 'th component as u_j , so that $u = (u_1, \dots, u_s)$. We first introduce the s -fold alternating sum of f over $[a, b] \subset \mathbb{R}$:

$$\Delta(f; a, b) = \sum_{v \subseteq \{1, \dots, s\}} (-1)^{|v|} f(a_v : b_{-v})$$

where a, b are two s dimensional vectors. We write $|v|$ for the cardinality of the multi-index v , and $-v$ for the sequence $\{1, \dots, s\} \setminus v$ which contains all elements of $\{1, \dots, s\}$ not in v . Furthermore, a_v denotes a $|v|$ -tuple of real values representing the components a_j for $j \in v$. The symbol $a_v : b_{-v}$ represents the point $u \in [a, b]^s$ with $u_j = a_j$ for $j \in v$, and $u_j = b_j$ for $j \notin v$.

Let $\mathcal{Y} = \{u \in [a, b]^s \mid 0 < u_1 < u_2 < \dots < u_s = 1\}$ be a ladder on $[a, b]$. For $j = 1, 2, \dots, s$, denote by \mathcal{Y}^j a ladder on $[a_j, b_j]$. A (multi-dimensional) ladder on $[a, b]$ has the form $\mathcal{Y} = \prod_{j=1}^s \mathcal{Y}^j$. For $y \in \mathcal{Y}$, the successor point y_+ is defined by taking $(y_+)_j$ to be the successor of y_j in \mathcal{Y}^j . The variation of f over \mathcal{Y} is then given by:

$$V_{\mathcal{Y}}(f) = \sum_{y \in \mathcal{Y}} |\Delta(f; y, y_+)|.$$

Let \mathbb{Y}^j denote the set of all ladders on $[a_j, b_j]$ and put $\mathbb{Y} = \prod_{j=1}^s \mathbb{Y}^j$. Then, the *variation of f in the sense of Hardy and Krause* is given by:

$$V_{HK}(f) = \sum_{v \subseteq \{1, \dots, s\}} V_{[a_{-v}, b_{-v}]} f(u_{-v} : b_v).$$

We can now finally present the Koksma-Hwlaka inequality, which decouples the quadrature error into a term depending on the function, the Hardy-Krause variation, and a term depending on the point set, the star discrepancy.

Lemma 1 (Theorem 15.5 in [67]). *Let $\mathcal{U} = [0, 1]^s$, $f : \mathcal{U} \rightarrow \mathbb{R}$ and $\{u_i\}_{i=1}^n \subset \mathcal{U}$. Then, if $V_{HK}(f) < \infty$, we have:*

$$\left| \int_{[0,1]^s} f(u) du - \frac{1}{n} \sum_{i=1}^n f(u_i) \right| \leq V_{HK}(f) D^*(\{u_i\}_{i=1}^n).$$

Combining this result with the definition of QMC or RQMC point set allows us to provide results on the convergence of QMC/RQMC estimators for functions with bounded Hardy-Krause variation.

B Proof of Theoretical Results

In this appendix, we provide proofs of all the theoretical results in the main text. Firstly, in Section B.1, we provide some useful preliminary results. Then, Section B.2 contains the proof of our results on MMD, Section B.3 the proof of our results on the Wasserstein distance, and Section B.4 the proof of our results on the Sinkhorn divergence.

B.1 Preliminary Results

Before stating our main result for this section, Theorem 4, we recall a preliminary results which will be used in its proof.

Lemma 2 (Generalised Hölder's Inequality; Corollary 2.6 of [2]). *Suppose that $p, p_1, \dots, p_r \in (0, \infty]$ and $\sum_{i=1}^r p_i^{-1} = p^{-1}$. Then, if $\|f_i\|_{L^{p_i}(\mathcal{X})} < \infty$ for all $i \in \{1, \dots, r\}$, we have:*

$$\|\prod_{i=1}^r f_i\|_{L^p(\mathcal{X})} \leq \prod_{i=1}^r \|f_i\|_{L^{p_i}(\mathcal{X})}.$$

We now provide an intermediate result which upper bounds the norm of the composition of two functions. For this, we will need to introduce an ordering on multi-indices. Let a, b be two multi-indices, then $a \prec b$ means that $|a| < |b|$ or, $|a| = |b|$ and $a_i < b_i$ for the smallest i such that $a_i \neq b_i$. The proof closely follows [7], but allows for additional smoothness of $g \circ h$.

Theorem 4. *Let $\mathcal{X} \subset \mathbb{R}^d$ be an open set and let \mathcal{H}_k be an RKHS with kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying $k \in C^{s,s}(\mathbb{R}^d \times \mathbb{R}^d)$ with $\sup_{x \in \mathcal{X}} \partial^{t,t} k(x, x) < C_k, \forall t \in \mathbb{N}_0^d$ such that $|t| \leq s$ where C_k is some universal constant only depending on kernel. Suppose $g : \mathcal{X} \rightarrow \mathbb{R}$ satisfies $g \in \mathcal{H}_k$ and $h : [0, 1]^s \rightarrow \mathcal{X}$. Then, assuming h is sufficiently regular for all norms to exist:*

$$V_{\text{HK}}(g \circ h) \leq C \|g\|_{\mathcal{H}_k(\mathcal{X})} \times \sum_{\alpha \neq \emptyset, \alpha \subseteq 1:s} \sum_{1 \leq |t| \leq |\alpha|} \sum_{l=1}^{|\alpha|} \sum_{(\ell_r, k_r) \in S(l, \alpha, t)} \prod_{r=1}^l \|\partial_{\ell_r} h_{k_r}(\cdot : 1-\alpha)\|_{L^{p_r}([0,1]^{|\alpha|})}$$

for any $\sum_{r=1}^s p_r^{-1} \leq 1$ and where

$$S(l, \alpha, t) = \left\{ (\ell_r, k_r), r = 1, \dots, l \mid \ell_r \in 1 : s, k_r \in 1 : d, \cup_{r=1}^l \ell_r = \alpha, \right. \\ \left. \ell_r \cap \ell_{r'} = \emptyset \text{ for } r \neq r', \text{ and } |\{j \in 1 : l \mid k_j = i\}| = t_i \right\}.$$

Proof. Starting with Equation 3 in [7] and recalling that $h : [0, 1]^s \rightarrow \mathcal{X}$ and $g : \mathcal{X} \rightarrow \mathbb{R}$:

$$V_{\text{HK}}(g \circ h) \leq \sum_{\alpha \neq \emptyset, \alpha \subseteq 1:s} \|\partial_\alpha (g \circ h)(\cdot : 1-\alpha)\|_{L^1([0,1]^{|\alpha|})} \quad (3)$$

where we recall that $V_{\text{HK}}(g \circ h)$ denotes the variation of $g \circ h$ in the sense of Hardy and Krause. In order to express the norm of $g \circ h$, we first need an expression for its partial derivatives. We will use Theorem 1 in [24] which a Faa di Bruno formula for mixed partial derivative. In particular, for $\alpha \subseteq 1 : s$:

$$\partial_\alpha (g \circ h)(u : 1-\alpha) = \sum_{1 \leq |t| \leq |\alpha|} \partial^t g(h(u : 1-\alpha)) \sum_{l=1}^{|\alpha|} \sum_{(\ell_r, k_r) \in S(l, \alpha, t)} \prod_{r=1}^l \partial_{\ell_r} h_{k_r}(u : 1-\alpha)$$

To clarify, here the first sum is over all multi-indices $t \in \mathbb{N}_0^d$, and $\partial_\alpha (g \circ h)$ denotes mixed partial derivatives where we differentiate at most once per coordinate. Taking the L^p norm of these derivatives, we get that for $\alpha \subseteq 1 : s$:

$$\begin{aligned} & \|\partial_\alpha (g \circ h)(\cdot : 1-\alpha)\|_{L^p([0,1]^{|\alpha|})} \\ &= \left\| \sum_{1 \leq |t| \leq |\alpha|} \partial^t g(h(\cdot : 1-\alpha)) \sum_{l=1}^{|\alpha|} \sum_{(\ell_r, k_r) \in S(l, \alpha, t)} \prod_{r=1}^l \partial_{\ell_r} h_{k_r}(\cdot : 1-\alpha) \right\|_{L^p([0,1]^{|\alpha|})} \\ &\leq \sum_{1 \leq |t| \leq |\alpha|} \left\| \partial^t g(h(\cdot : 1-\alpha)) \sum_{l=1}^{|\alpha|} \sum_{(\ell_r, k_r) \in S(l, \alpha, t)} \prod_{r=1}^l \partial_{\ell_r} h_{k_r}(\cdot : 1-\alpha) \right\|_{L^p([0,1]^{|\alpha|})} \\ &\leq \sum_{1 \leq |t| \leq |\alpha|} \|\partial^t g(h(\cdot : 1-\alpha))\|_{L^\infty([0,1]^{|\alpha|})} \left\| \sum_{l=1}^{|\alpha|} \sum_{(\ell_r, k_r) \in S(l, \alpha, t)} \prod_{r=1}^l \partial_{\ell_r} h_{k_r}(\cdot : 1-\alpha) \right\|_{L^p([0,1]^{|\alpha|})} \\ &\leq \sum_{1 \leq |t| \leq |\alpha|} \|\partial^t g(h(\cdot : 1-\alpha))\|_{L^\infty([0,1]^{|\alpha|})} \sum_{l=1}^{|\alpha|} \sum_{(\ell_r, k_r) \in S(l, \alpha, t)} \left\| \prod_{r=1}^l \partial_{\ell_r} h_{k_r}(\cdot : 1-\alpha) \right\|_{L^p([0,1]^{|\alpha|})} \end{aligned} \quad (4)$$

Here, the first inequality follows by the triangle inequality. The second inequality follows from Hölder's inequality (Lemma 2 with $p_1 = \infty$ and $p_2 = p$). Finally, the third inequality once again follows from the triangle inequality. The rest of the proof will consist of bounding each of the remaining norms separately.

For the first norm, we will use the fact that $g \in \mathcal{H}_k$ and we can therefore bound the norm of its derivatives. Since $k \in \mathcal{C}^{m \times m}(\mathbb{R}^d \times \mathbb{R}^d)$, we have $k \in \mathcal{C}^{m \times m}(\mathcal{X} \times \mathcal{X})$. Following Corollary 4.36 of [80], we have that $g \in \mathcal{H}_k$ implies $g \in \mathcal{C}^m(\mathcal{X})$, and $\forall t \in \mathbb{N}_0^d$ with $|t| \leq m$ and $\forall x \in \mathcal{X}$,

$$\partial^t g(x) \leq \|g\|_{\mathcal{H}_k(\mathcal{X})} (\partial^{t,t} k(x, x))^{\frac{1}{2}}$$

Given the assumption that $\sup_{x \in \mathcal{X}} \partial^{t,t} k(x, x) \leq C_k$, we obtain the following inequality combining above results with $m = s$

$$\|\partial^t g(h(\cdot : 1_{-\alpha}))\|_{L^\infty([0,1]^{|\alpha|})} \leq \|\partial^t g\|_{L^\infty(\mathcal{X})} = \sup_{x \in \mathcal{X}} \partial^t g(x) \leq C_k^{1/2} \|g\|_{\mathcal{H}_k(\mathcal{X})}. \quad (5)$$

For the second norm, we can once again make use of Hölder's inequality (Lemma 2) which guarantees that if $\partial_{\ell_r} h_{k_r}(\cdot : 1_{-\alpha}) \in L^{p_r}([0,1]^{|\alpha|})$ for $\alpha \subseteq 1 : s$ and $\sum_{r=1}^l p_r^{-1} \leq p^{-1}$, then

$$\left\| \prod_{r=1}^l \partial_{\ell_r} h_{k_r}(\cdot : 1_{-\alpha}) \right\|_{L^p([0,1]^{|\alpha|})} \leq \prod_{r=1}^l \|\partial_{\ell_r} h_{k_r}(\cdot : 1_{-\alpha})\|_{L^{p_r}([0,1]^{|\alpha|})}. \quad (6)$$

Plugging the inequalities in 5 and 6 into 4, we get that for $\alpha \subseteq 1 : s$:

$$\begin{aligned} & \|\partial_\alpha (g \circ h)(\cdot : 1_{-\alpha})\|_{L^p([0,1]^{|\alpha|})} \\ & \leq C \|g\|_{\mathcal{H}_k(\mathcal{X})} \sum_{1 \leq |t| \leq |\alpha|} \sum_{l=1}^{|\alpha|} \sum_{(\ell_r, k_r) \in S(l, \alpha, t)} \prod_{r=1}^l \|\partial_{\ell_r} h_{k_r}(\cdot : 1_{-\alpha})\|_{L^{p_r}([0,1]^{|\alpha|})}. \end{aligned}$$

Plugging this bound with $p = 1$ in Equation 3 concludes the proof. \square

B.2 Proof of Theorem 1 and Corollary 1

B.2.1 Proof of Theorem 1

Proof. First, we notice that under our assumptions, we may directly apply Theorem 4 in order to get that $\exists C_\theta > 0$ such that for any $f \in \mathcal{H}_k$:

$$V_{\text{HK}}(f \circ G_\theta) \leq C_\theta \|f\|_{\mathcal{H}_k}.$$

More precisely, G_θ takes the place of h in Theorem 4, and all norms depending on G_θ are bounded thanks to Assumption 2. We can then directly combine this result with the Koksma-Hlawka inequality (Lemma 1) to get a bound on the MMD:

$$\begin{aligned} \text{MMD}(\mathbb{P}_\theta, \mathbb{P}_\theta^n) &= \sup_{\|f\|_{\mathcal{H}_k(\mathcal{X})} \leq 1} \left| \int_{\mathcal{X}} f(x) \mathbb{P}_\theta(dx) - \int_{\mathcal{X}} f(x) \mathbb{P}_\theta^n(dx) \right| \\ &= \sup_{\|f\|_{\mathcal{H}_k(\mathcal{X})} \leq 1} \left| \int_{[0,1]^s} f(G_\theta(u)) du - \frac{1}{n} \sum_{i=1}^n f(G_\theta(u_i)) \right| \\ &\leq \sup_{\|f\|_{\mathcal{H}_k(\mathcal{X})} \leq 1} V_{\text{HK}}(f \circ G_\theta) D^*(\{x_i\}_{i=1}^n) \\ &\leq \sup_{\|f\|_{\mathcal{H}_k(\mathcal{X})} \leq 1} C_\theta \|f\|_{\mathcal{H}_k(\mathcal{X})} D^*(\{x_i\}_{i=1}^n) = C_\theta D^*(\{x_i\}_{i=1}^n) \end{aligned}$$

By definition, we know that whenever $\{u_i\}_{i=1}^n$ is a QMC point set, we have $D^*(\{u_i\}_{i=1}^n) = O(n^{-1}(\log n)^{\alpha_s})$. This concludes the proof. \square

B.2.2 Proof of Corollary 1

Proof. The proof is trivial by using the fact that MMD is a distance and thus the triangle inequality holds $|\text{MMD}(\mathbb{P}_\theta, \mathbb{Q}^m) - \text{MMD}(\mathbb{P}_\theta^n, \mathbb{Q}^m)| \leq \text{MMD}(\mathbb{P}_\theta, \mathbb{P}_\theta^n)$. The rate therefore follows from Theorem 1. \square

B.3 Proof of Theorem 2 and Corollary 2

B.3.1 Proof of Theorem 2

Proof. Using the Kantorovich-Rubinstein duality theorem, we may express the Wasserstein distance as an integral probability metric associated to the class of Lipschitz continuous functions when $p = 1$:

$$\begin{aligned} W_{c,1}(\mathbb{P}_\theta, \mathbb{P}_\theta^n) &:= \sup_{\|f\|_{\text{L}} \leq 1} \left| \int_{\mathcal{X}} f(x) \mathbb{P}_\theta(dx) - \int_{\mathcal{X}} f(x) \mathbb{P}_\theta^n(dx) \right| \\ &= \sup_{\|f\|_{\text{L}} \leq 1} \left| \int_{[0,1]^s} f(G_\theta(u)) du - \frac{1}{n} \sum_{i=1}^n f(G_\theta(u_i)) \right| \end{aligned} \quad (7)$$

where $\|g\|_{\text{L}} := |g(x) - g(y)|/c(x, y)$. Then, using the Koksma-Hwlaka inequality in Lemma 1, we get:

$$\left| \int_{[0,1]^s} f(G_\theta(u)) du - \frac{1}{n} \sum_{i=1}^n f(G_\theta(u_i)) \right| \leq V_{\text{HK}}(f \circ G_\theta) D^*(\{u_i\}_{i=1}^n). \quad (8)$$

Let \mathcal{Y}^N be the ladder $\{v \in [0, 1]^N | 0 < v_1 < v_2 < \dots < v_N = 1\}$. Assuming f is Lipschitz and G_θ has bounded variation in the sense of Hardy and Krause, we have

$$\begin{aligned} V_{\text{HK}}(f \circ G_\theta) &= \sup_{N \geq 1} \sup_{v \in \mathcal{Y}^N} \sum_{i=1}^N |f(G_\theta(v_i)) - f(G_\theta(v_{i-1}))| \\ &= \|f\|_{\text{L}} \sup_{N \geq 1} \sup_{v \in \mathcal{Y}^N} \sum_{i=1}^N c(G_\theta(v_i), G_\theta(v_{i-1})) \\ &\leq M \|f\|_{\text{L}} \sup_{N \geq 1} \sup_{v \in \mathcal{Y}^N} \sum_{i=1}^N |G_\theta(v_i) - G_\theta(v_{i-1})| \\ &= M \|f\|_{\text{L}} V_{\text{HK}}(G_\theta). \end{aligned} \quad (9)$$

where the first equality follows by definition of the Hardy-Krause variation, the second equality from the definition of the Lipschitz norm, and the first inequality from the fact that all norms are equivalent on \mathbb{R} so that $\exists M > 0$ such that $c(x, y) \leq M|x - y|$ for all $x, y \in \mathbb{R}$. Combining the results in Equations 7, 8 and 9, we get:

$$\begin{aligned} W_{c,1}(\mathbb{P}_\theta, \mathbb{P}_\theta^n) &= \sup_{\|f\|_{\text{L}} \leq 1} \left| \int_{[0,1]^s} f(G_\theta(u)) du - \frac{1}{n} \sum_{i=1}^n f(G_\theta(u_i)) \right| \\ &\leq \sup_{\|f\|_{\text{L}} \leq 1} M \|f\|_{\text{L}} V_{\text{HK}}(G_\theta) D^*(\{u_i\}_{i=1}^n) = M V_{\text{HK}}(G_\theta) D^*(\{u_i\}_{i=1}^n). \end{aligned}$$

The proof of the theorem is concluded by noting the rate of convergence for the star discrepancy in the case of QMC or RQMC point sets. \square

B.3.2 Proof of Corollary 2

Proof. The proof is simple by noticing when $p \geq 1$, the Wasserstein distance $W_{c,p}$ with distance function c is indeed a distance satisfying the triangle inequality (Proposition 2.3 in [73]). \square

B.4 Proof of Theorem 3

We will now prove Theorem 3. The bounds follow the main approach in [37], but need to be significantly modified to accommodate QMC or RQMC point sets instead of IID realisations.

Proof. Since the Sinkhorn divergence is a normalised version of the regularised optimal transport problem, we can use this definition together with the triangle inequality to get:

$$\begin{aligned} & |S_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{Q}) - S_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{Q})| \\ &= |\bar{W}_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{Q}) - \bar{W}_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{Q}) - \frac{1}{2} (\bar{W}_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{P}_\theta) - \bar{W}_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{P}_\theta^n))| \\ &\leq |\bar{W}_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{Q}) - \bar{W}_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{Q})| + \frac{1}{2} |\bar{W}_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{P}_\theta) - \bar{W}_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{P}_\theta^n)|. \end{aligned} \quad (10)$$

We will now focus on bounding these terms. To do so, we first recall that the regularized optimal transport problem can be expressed as follows:

$$\bar{W}_{c,p,\lambda}(\mathbb{P}, \mathbb{Q}) = \max_{g \in \mathcal{C}(\mathcal{X}), h \in \mathcal{C}(\mathcal{Y})} \mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(g, h) \right] + \lambda,$$

where

$$F_\lambda^{x,y}(g, h) = g(x) + h(y) - \lambda \exp\left(\frac{g(x) + h(y) - c^p(x,y)}{\lambda}\right).$$

We will denote by (g^*, h^*) the optimal potentials for $\bar{W}_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{Q})$ (i.e. the functions g and h attaining the maximum), by (\bar{g}, \bar{h}) the optimal potentials for $\bar{W}_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{Q})$, by (\tilde{g}, \tilde{h}) the optimal potentials for $\bar{W}_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{P}_\theta)$, and by (\hat{g}, \hat{h}) the optimal potentials for $\bar{W}_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{P}_\theta^n)$.

Now we can upper bound the first term in (10) using the triangle inequality as follows:

$$\begin{aligned} & |\bar{W}_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{Q}) - \bar{W}_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{Q})| \\ &= \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(g^*, h^*) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(\bar{g}, \bar{h}) \right] \right| \\ &\leq \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(g^*, h^*) \right] - \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(\bar{g}, \bar{h}) \right] \right| \\ &\quad + \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(\bar{g}, \bar{h}) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(\bar{g}, \bar{h}) \right] \right| \end{aligned} \quad (11)$$

We will now bound the remaining terms. The first term of (11) can be upper bounded by:

$$\begin{aligned} & \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(g^*, h^*) \right] - \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(\bar{g}, \bar{h}) \right] \right| \\ &= \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(g^*, h^*) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(g^*, h^*) \right] + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(g^*, h^*) \right] \right. \\ &\quad \left. - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(\bar{g}, \bar{h}) \right] + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(\bar{g}, \bar{h}) \right] - \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(\bar{g}, \bar{h}) \right] \right| \\ &\leq \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(g^*, h^*) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(g^*, h^*) \right] \right| \\ &\quad + \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(\bar{g}, \bar{h}) \right] - \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(\bar{g}, \bar{h}) \right] \right| \end{aligned} \quad (12)$$

where the inequality follows from triangle inequality and the fact that

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(g^*, h^*) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(\bar{g}, \bar{h}) \right] \leq 0$$

because of the optimality of (\tilde{g}, \tilde{h}) . We will now turn to the second term in (10), which can be similarly upper-bounded as follows:

$$\begin{aligned}
& \left| \bar{W}_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{P}_\theta) - \bar{W}_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{P}_\theta^n) \right| \\
&= \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{P}_\theta} \left[F_\lambda^{X,Y}(\tilde{g}, \tilde{h}) \right] - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n F_\lambda^{x_i, y_j}(\hat{g}, \hat{h}) \right| \\
&= \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{P}_\theta} \left[F_\lambda^{X,Y}(\tilde{g}, \tilde{h}) \right] - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n F_\lambda^{x_i, y_j}(\hat{g}, \hat{h}) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{P}_\theta} \left[F_\lambda^{x_i, Y}(\tilde{g}, \tilde{h}) \right] \right. \\
&\quad \left. + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{P}_\theta} \left[F_\lambda^{x_i, Y}(\tilde{g}, \tilde{h}) \right] + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n F_\lambda^{x_i, y_j}(\tilde{g}, \tilde{h}) - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n F_\lambda^{x_i, y_j}(\hat{g}, \hat{h}) \right| \\
&\leq \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{P}_\theta} \left[F_\lambda^{X,Y}(\tilde{g}, \tilde{h}) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{P}_\theta} \left[F_\lambda^{x_i, Y}(\tilde{g}, \tilde{h}) \right] \right| \\
&\quad + \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{P}_\theta} \left[F_\lambda^{x_i, Y}(\tilde{g}, \tilde{h}) \right] - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n F_\lambda^{x_i, y_j}(\tilde{g}, \tilde{h}) \right| \tag{13}
\end{aligned}$$

where the last inequality holds since

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n F_\lambda^{x_i, y_j}(\tilde{g}, \tilde{h}) - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n F_\lambda^{x_i, y_j}(\hat{g}, \hat{h}) \leq 0$$

due to the definition of $\tilde{g}, \tilde{h}, \hat{g}$ and \hat{h} .

Combining (11), (12) and (13), we end up with several terms which take the form of absolute integration errors for integrating against \mathbb{P}_θ for various choices of potentials. From Theorem 2 in [37], we know that if $c \in \mathcal{C}^{\infty, \infty}(\mathcal{X} \times \mathcal{X})$, then all of these potentials are in $W^{m,2}(\mathcal{X})$ for $m = d/2 + 1$. We will now obtain an upper bound on the integration error for any arbitrary potentials $g, h \in W^{m,2}(\mathcal{X})$. Firstly, using the definition of $F_\lambda^{x,y}(g, h)$ and the triangle inequality:

$$\begin{aligned}
& \left| \mathbb{E}_{X \sim \mathbb{P}_\theta, Y \sim \mathbb{Q}} \left[F_\lambda^{X,Y}(g, h) \right] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim \mathbb{Q}} \left[F_\lambda^{x_i, Y}(g, h) \right] \right| \\
&\leq \left| \int_{\mathcal{X}} g(x) \mathbb{P}_\theta(dx) - \frac{1}{n} \sum_{i=1}^n g(x_i) \right| \\
&\quad + \lambda \left| \int_{\mathcal{X}} \int_{\mathcal{X}} \exp\left(\frac{g(x)+h(y)-c^p(x,y)}{\lambda}\right) \mathbb{Q}(dy) \mathbb{P}_\theta(dx) - \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{X}} \exp\left(\frac{g(x_i)+h(y)-c^p(x_i,y)}{\lambda}\right) \mathbb{Q}(dy) \right| \\
&= \left| \int_{[0,1]^s} g(G_\theta(u)) du - \frac{1}{n} \sum_{i=1}^n g(G_\theta(u_i)) \right| \tag{14}
\end{aligned}$$

where the equality holds due to the duality equation (see e.g. Equation 6 in [37]):

$$\exp\left(-\frac{g(x)}{\lambda}\right) = \int_{\mathcal{X}} \exp\left(\frac{h(y)-c^p(x,y)}{\lambda}\right) \mathbb{Q}(dy).$$

To bound the expression above, we may then use the Koksma-Hwlaka inequality in Lemma 1:

$$\left| \int_{[0,1]^s} g(G_\theta(u)) du - \frac{1}{n} \sum_{i=1}^n g(G_\theta(u_i)) \right| \leq V_{\text{HK}}(g \circ G_\theta) D^* (\{u_i\}_{i=1}^n).$$

To conclude the proof, our approach will be to upper bound $V_{\text{HK}}(g \circ G_\theta)$ using Theorem 4 for some sufficiently smooth kernel k which we will take to be Matérn kernel k of smoothness $m - d/2 = 1$ (see Appendix C.1 for a definition). This will require that G_θ is sufficiently regular to satisfy the assumptions in Theorem 4, but this is true thanks to Assumption 2. As a result $\exists C_\theta > 0$ such that $V_{\text{HK}}(g \circ G_\theta) \leq C_\theta \|g\|_{\mathcal{H}_k}$, which leads to a bound of the form:

$$\left| \int_{[0,1]^s} g(G_\theta(u)) du - \frac{1}{n} \sum_{i=1}^n g(G_\theta(u_i)) \right| \leq C_\theta \|g\|_{\mathcal{H}_k} D^* (\{u_i\}_{i=1}^n). \tag{15}$$

It has been proven in Theorem 2 of [38] that the potentials $g, h \in W^{m,2}(\mathcal{X})$, where $\mathcal{X} \subseteq \mathbb{R}^d$ is a compact space and $m \in \mathbb{N}$. Conveniently, when $m > d/2$, we know that $W^{m,2}(\mathcal{X})$ is norm-equivalent to the RKHS \mathcal{H}_k with Matérn kernel k of smoothness $m - d/2 = 1$ (see Example 2.6

in [47]), so that $\exists C_1, C_2 > 0$ such that:

$$C_1 \|g\|_{\mathcal{H}_k(x)} \leq \|g\|_{W^{m,2}(x)} \leq C_2 \|g\|_{\mathcal{H}_k(x)}$$

We can then combine this result with Equation 15 to get a bound of the form

$$\left| \int_{[0,1]^s} g(G_\theta(u)) du - \frac{1}{n} \sum_{i=1}^n g(G_\theta(u_i)) \right| \leq \frac{C_\theta}{C_1} \|g\|_{W^{m,2}(x)} D^* (\{u_i\}_{i=1}^n). \quad (16)$$

Putting all of the pieces together we end up with

$$|S_{c,p,\lambda}(\mathbb{P}_\theta, \mathbb{Q}) - S_{c,p,\lambda}(\mathbb{P}_\theta^n, \mathbb{Q})| \leq \tilde{C}_\theta D^* (\{u_i\}_{i=1}^n).$$

where the bound follows from combining Equations 10 and 11 to obtain an upper bound in terms of integration error, then Equation 16 to upper bound such error, and finally combining all of the constants. This concludes our proof. \square

C Additional Numerical Experiments

In this section, we provide additional details on the numerical experiments presented in the main text, and also complement these with additional results to provide a more complete picture of the impact of QMC and RQMC point sets. First, in Section C.1, we provide additional experiments on the sample complexity for the uniform and Gaussian models. Sections C.2 and C.3 then provide additional details on the experiments with the bivariate Beta and multivariate g-and-k distributions respectively.

C.1 Uniform and Gaussian Models

The first set of additional experiments focuses on the sample complexity of MMD. These experiments were once again performed with generalised Halton sequences randomised using the scrambling factors of [32], and with a lengthscale of $l = 1.5d^{1/2}$. In Figure 12, we compare the sample complexity of MMD when different kernels are used. In particular, we compare a squared-exponential kernel with Matérn kernels of smoothness $3/2, 5/2$ and $7/2$. The Matérn kernels take the form:

$$k_\nu(x, x') = \frac{\lambda^{2\nu} 2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|x-x'\|_2}{\sigma} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu} \|x-x'\|_2}{\sigma} \right),$$

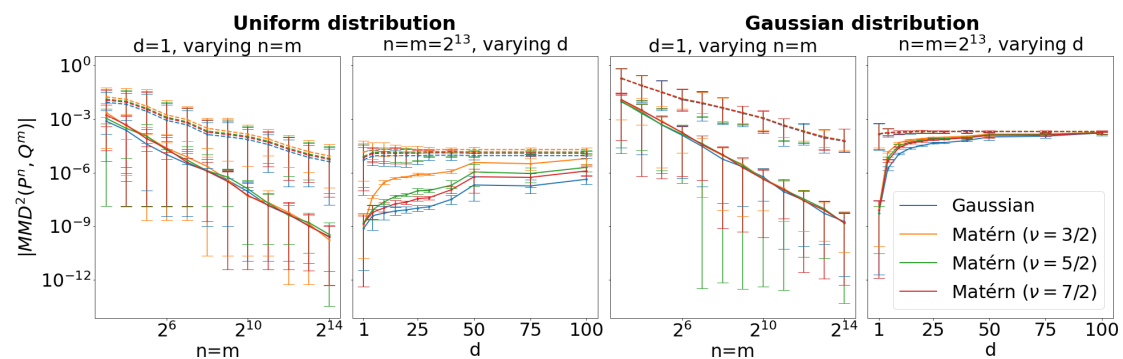


Figure 12: Sample complexity results for the MMD squared with different choices of kernels. The smoother the kernel, the more QMC point sets improve performance for $d > 1$.

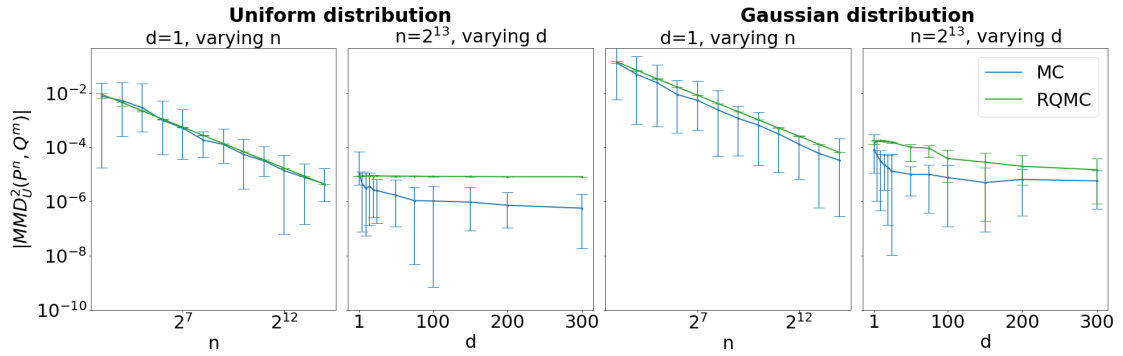


Figure 13: Sample complexity results for the U-statistic approximation of the maximum mean discrepancy squared for the uniform and Gaussian distributions. We compare MC realisations with realisations obtained through a RQMC point set. The setup is identical to that of Figure 2 (top row), except we use a U-statistic approximation instead of the squared MMD with empirical measures.

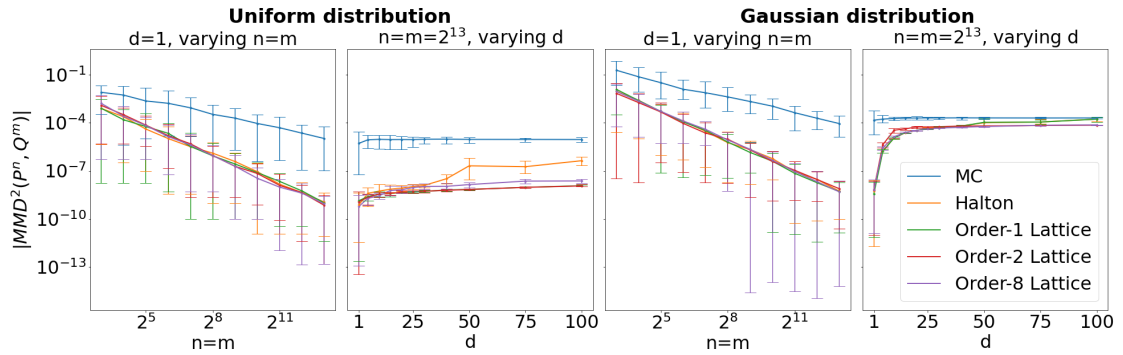


Figure 14: Sample complexity results for the MMD squared for different QMC point sets. The performance does not seem to be significantly impacted by the choice of QMC point set.

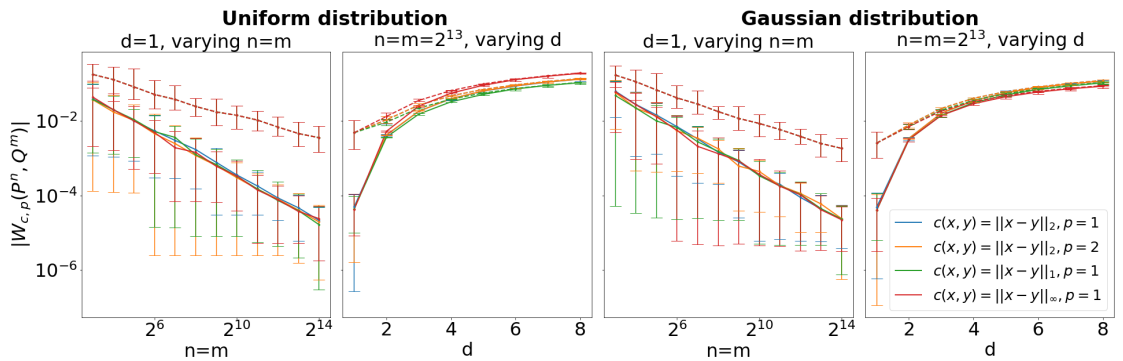


Figure 15: Sample complexity results for the Wasserstein distance with various choices of the cost c and order p .

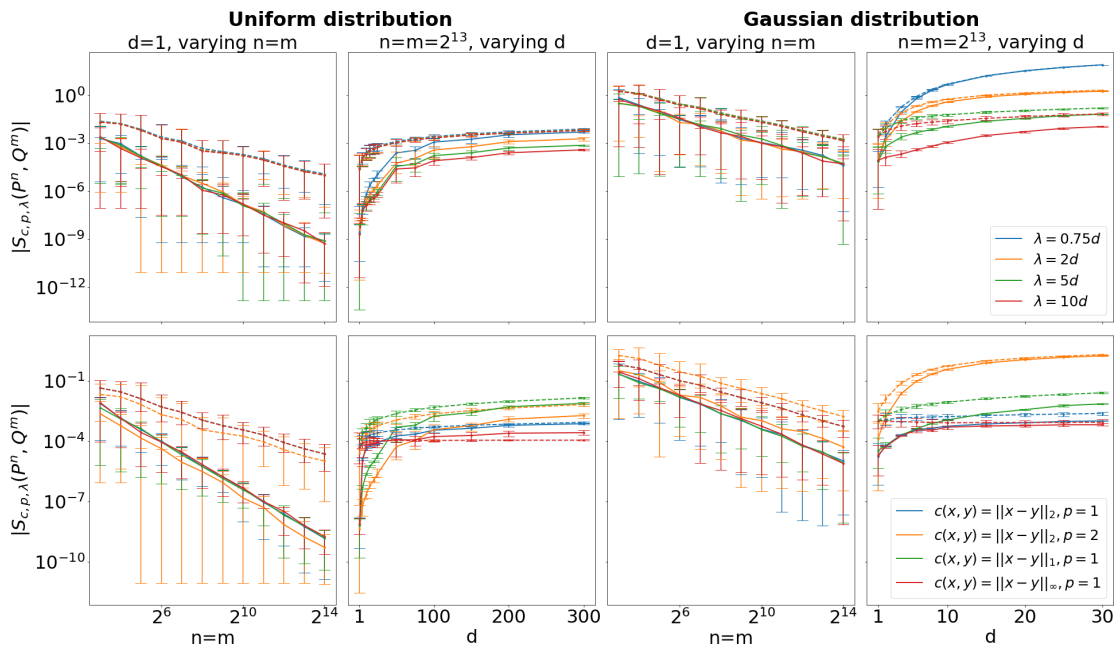


Figure 16: Sample complexity results for the Sinkhorn divergence with various choices of cost c , order p and regularisation λ . The top row corresponds to $c = \|x - y\|_2$ and $p = 2$, whereas the second row corresponds to $\lambda = 2d$.

where $\nu > 0$ is the smoothness parameter, Γ is the Gamma function, and K_ν is the modified Bessel function of the second kind of order ν . In dimension $d = 1$, all kernels lead to similar sample complexity results for either MC and QMC point sets. However, for $d > 1$, we see a clear improvement when using a smoother kernel and QMC point sets, with the squared-exponential kernel providing the best overall performance. This clearly supports our choice of squared-exponential kernel for the experiments in the main text, and also shows the importance of the smoothness requirements on the kernel in Theorem 1.

Another choice we made in the main text was to focus on the MMD with empirical measures. However, many papers in the literature use a U-statistic approximation instead:

$$\text{MMD}_U^2(\mathbb{P}^n, \mathbb{Q}^m) = \frac{\sum_{i \neq j}^n k(x_i, x_j)}{n(n-1)} - \frac{2 \sum_{i=1}^n \sum_{j=1}^m k(x_i, y_j)}{nm} + \frac{\sum_{i \neq j}^m k(y_i, y_j)}{m(m-1)},$$

see for example [16, 70]. The main advantage of the U-statistic is that it is unbiased, but it does have a larger variance. This turns out to have a significant impact when using QMC point sets in which case we cannot obtain an improved convergence rate. This is illustrated in Figure 13 where we reproduced the sample complexity plots in the top row of Figure 2 using the U-statistic. As can be observed, we are not able to obtain a faster convergence rate, and this is the case even in $d = 1$. In fact, the results are significantly worse than MC when $d > 1$.

These experiments were complemented by a study of the impact of the QMC point sets in Figure 14 where we compare an order-1, order-2, and order-8 lattice which were shifted to obtain randomised point sets. As observed, there is only negligible differences in the performance of the different QMC point sets when d is small, but further gains can be obtained when d is large in the case of the uniform distribution.

Next, we studied the impact of the choice of c and p on sample complexity results for the Wasserstein distance. Note that the result in Theorem 2 is only valid for $p = 1$. As we can see in Figure 15, the performance is similar across various choices of c and p . In each case, a faster rate is obtained for $d = 1$ indicating that the result of our theorem could potentially be extended to $p \neq 1$. However, in all cases this gain in performance quickly vanishes as d increases. A similar study was performed for the Sinkhorn divergence in Figure 16 (bottom row). In this case, we may rely on Theorem 3 which is also valid for $d > 1$. As we can see, there seems to be a larger impact due to the choice of cost function or of p , and this should warrant further study.

C.2 Bivariate Beta Model

As mentioned in Section 4.2 in the main text, it is possible to sample from the bivariate beta model using uniform random variables whenever all parameter take integer values. However, in the more general setting where the parameters may take scalar values, we will also require realisations from a Gamma random variable.

In order to make this model amenable to realisations from QMC point sets, we therefore need an approach to sampling from Gamma random variables using uniform random variables. A number of approaches are highlighted in Chapter IX.3. of [29], but we will focus specifically on the rejection sampling algorithm by Ahrens and Dieter [3] which we recall in Algorithm 1.

Algorithm 1 Rejection sampling for n realisations of a $\text{Gamma}(\alpha, 1)$ where $\alpha \in (0, 1)$

Require: A 3-dimensional point set of size n : $\{u_i = (u_{i1}, u_{i2}, u_{i3})\}_{i=1}^n \subset [0, 1]^3$

```

Set  $b = (\alpha + e)/e$ 
for  $i$  in  $1, \dots, n$  do
  Set  $p_i = b \times u_{i1}$ .
  if  $p_i \leq 1$  then
    Let  $x_i = p_i^{1/\alpha}$ 
    If  $u_{i2} \leq \exp(-x_i)$ , accept  $x_i$ , else reject  $x_i$ .
  else
    Let  $x_i = -\log((b - p_i)/\alpha)$ 
    If  $u_{i3} \leq x_i^{\alpha-1}$ , accept  $x_i$ , else reject  $x_i$ .
  end if
end for

```

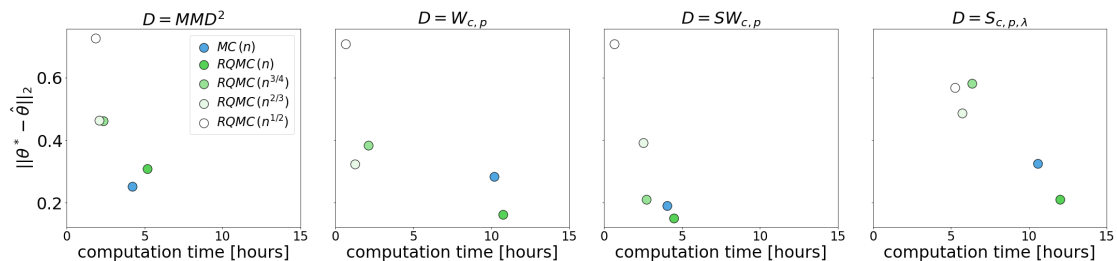


Figure 17: Minimum distance estimation for the parameters of the bivariate Beta distribution with the MMD, Wasserstein, sliced Wasserstein and Sinkhorn divergence. Each point corresponds to the median of 10 repetitions of an identical experiment.

Alternative Representation Figure 17 gives an alternative representation of the results in Figure 5 in the main text. It highlights the relationship between the computation time and the l_2 error of the estimates in the different setups.

C.3 Univariate and Multivariate g-and-k Models

In this final subsection, we provide additional details for the g-and-k models.

Generator In order to simulate from the multivariate g-and-k distribution studied in this paper, we will simply need to simulate some uniform random variables and transform these. In order to do so, one quantity of interest will be the matrix-square root of $\Sigma \in \mathbb{R}^{d \times d}$. We recall that Σ is a symmetric tri-diagonal Toeplitz matrix with diagonal entries all equal to 1 and off-diagonal entries equal to θ_5 , i.e.:

$$\Sigma = \begin{bmatrix} 1 & \theta_5 & 0 & \dots & 0 \\ \theta_5 & 1 & \theta_5 & \ddots & \vdots \\ 0 & \theta_5 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & \theta_5 \\ 0 & \dots & 0 & \theta_5 & 1 \end{bmatrix}.$$

For such matrices, the square-root is known in closed form and its computation does not require the use of an algorithm. It has entries given by

$$\left(\Sigma^{\frac{1}{2}}\right)_{ij} = \frac{2}{d+1} \sum_{k=1}^d \sqrt{1 + 2\theta_5 \cos\left(\frac{k\pi}{d+1}\right)} \sin\left(\frac{ik\pi}{d+1}\right) \sin\left(\frac{jk\pi}{d+1}\right).$$

This can be used directly in the expression for the generator of this model.

Computational Cost Figure 19 describes the computational cost of simulating n realisations of the g-and-k distribution using our implementation. In particular, it compares MC and RQMC for a range of values of d . When n is less than 2^{12} , the cost is usually slightly smaller with RQMC, but as n goes beyond this point the cost of using MC was significantly smaller.

Additional Numerical Results To complement the results in the main text, we first provide results for parameter estimation in the case $d = 1$, which is the most common in the literature. In this case, $p = 4$ since the parameter θ_5 does not enter the generator.

The results were obtained without sub-sampling the dataset and are provided in Figure 18. As observed in the top left plot, the stochastic optimisation algorithm is able to attain low values of the MMD squared in a much smaller number of steps when using RQMC as opposed to MC. This then leads to an improved parameter estimate as measure in terms of l_2 -norm between the estimated parameter and the true parameter θ^* ; see the top right plot. The bottom row of the figure gives the error for each of the four parameters as the number of step increases. In each case, the RQMC estimates provide significant improvements over the MC estimates, although the gains are limited for the second parameter (which controls the variance).

To complement these results, we provide a histogram obtained by sampling $n = 2^{11}$ from the model at θ^* from MC and QMC, and compare these to a histogram of \mathbb{P}_{θ^*} (obtained in practice by sampling a number of samples order of magnitude larger). The results are provided in Figure 20. The RQMC-based realisations provide a much better approximation of the distribution near

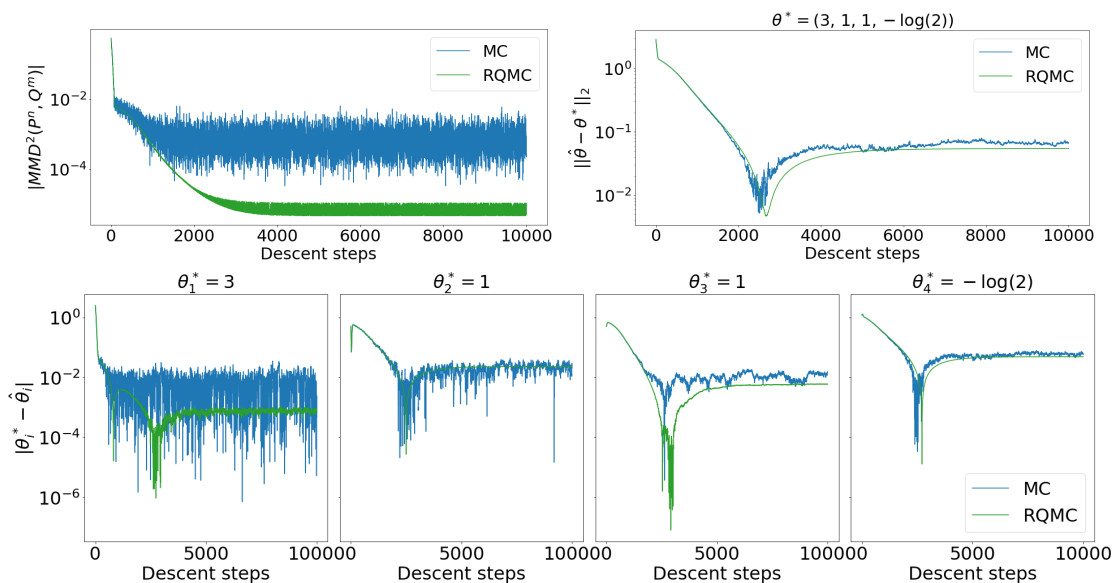


Figure 18: Minimum MMD estimation of the parameters of the univariate g-and-k distribution using stochastic optimisation.

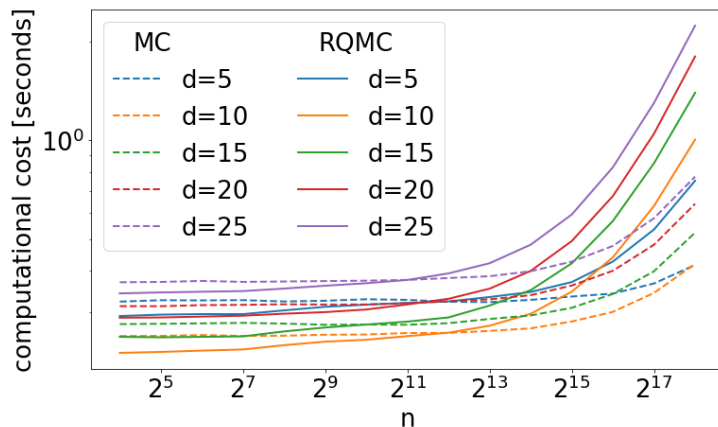


Figure 19: Comparison of the computational cost for simulating from the multivariate g-and-k distribution using MC and RQMC. Each line represents the average of 1000 repetitions.

the mode. This is confirmed by the table which provides the distance between the MC-based histogram or the RQMC-based histogram and the truth in terms of various choices of distance including the Kullback-Leibler divergence, the l_2 norm, or the Hellinger distance. We also notice that both MC and RQMC provide relatively poor approximation at the tail of the distribution. This is most likely due to the small number of realisations used to create the histogram.

Finally, Figure 21 provides the l_2 error between true and estimated parameters for the experiment presented in Figure 8. Clearly, a smaller value of the estimated MMD does not necessarily guarantee a better parameter estimate.

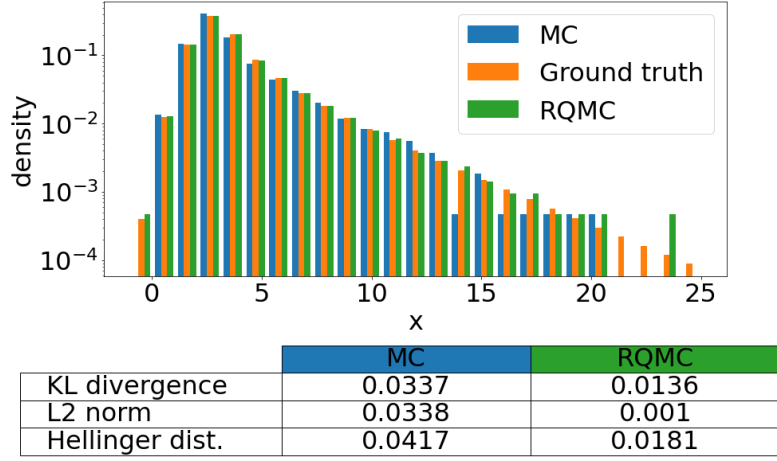


Figure 20: Histograms of the univariate g-and-k distribution at θ^* , together with MC and RQMC-based approximations with $n = 2^{11}$.

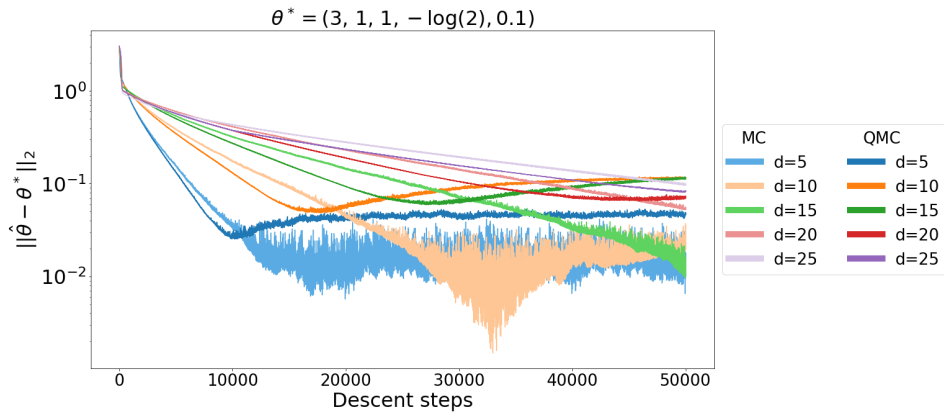


Figure 21: Minimum distance estimation with the MMD for the multivariate g-and-k distribution. The figure plots the l_2 error between the estimated parameter and the true value as the number of stochastic gradient descent steps increases.