



**HAL**  
open science

## On two ways to use determinantal point processes for Monte Carlo integration – Long version

Guillaume Gautier, Rémi Bardenet, Michal Valko

### ► To cite this version:

Guillaume Gautier, Rémi Bardenet, Michal Valko. On two ways to use determinantal point processes for Monte Carlo integration – Long version. NeurIPS 2019 - Thirty-third Conference on Neural Information Processing Systems, Jun 2019, Vancouver, Canada. hal-02277739

**HAL Id: hal-02277739**

**<https://hal.science/hal-02277739v1>**

Submitted on 12 Dec 2019

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

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

---

# On two ways to use determinantal point processes for Monte Carlo integration

---

Guillaume Gautier<sup>†\*</sup>

g.gautier@inria.fr

Rémi Bardenet<sup>†</sup>

remi.bardenet@gmail.com

Michal Valko<sup>‡\*†</sup>

valkom@deepmind.com

<sup>†</sup>Univ. Lille, CNRS, Centrale Lille, UMR 9189 – CRIStAL, 59651 Villeneuve d’Ascq, France

<sup>\*</sup>Inria Lille-Nord Europe, 40 avenue Halley 59650 Villeneuve d’Ascq, France

<sup>‡</sup>DeepMind Paris, 14 Rue de Londres, 75009 Paris, France

## Abstract

When approximating an integral by a weighted sum of function evaluations, determinantal point processes (DPPs) provide a way to enforce repulsion between the evaluation points. This negative dependence is encoded by a kernel. Fifteen years before the discovery of DPPs, Ermakov & Zolotukhin (EZ, 1960) had the intuition of sampling a DPP and solving a linear system to compute an unbiased Monte Carlo estimator of the integral. In the absence of DPP machinery to derive an efficient sampler and analyze their estimator, the idea of Monte Carlo integration with DPPs was stored in the cellar of numerical integration. Recently, Bardenet & Hardy (BH, 2019) came up with a more natural estimator with a fast central limit theorem (CLT). In this paper, we first take the EZ estimator out of the cellar, and analyze it using modern arguments. Second, we provide an efficient implementation<sup>1</sup> to sample exactly a particular multidimensional DPP called *multivariate Jacobi ensemble*. The latter satisfies the assumptions of the aforementioned CLT. Third, our new implementation lets us investigate the behavior of the two unbiased Monte Carlo estimators in yet unexplored regimes. We demonstrate experimentally good properties when the kernel is adapted to basis of functions in which the integrand is sparse or has fast-decaying coefficients. If such a basis and the level of sparsity are known (e.g., we integrate a linear combination of kernel eigenfunctions), the EZ estimator can be the right choice, but otherwise it can display an erratic behavior.

## 1 Introduction

Numerical integration is a core task of many machine learning applications, including most Bayesian methods (Robert, 2007). Both deterministic (Davis & Rabinowitz, 1984; Dick & Pillichshammer, 2010) and random (Robert & Casella, 2004) algorithms have been proposed; see also (Evans & Swartz, 2000) for a survey. All methods require evaluating the integrand at carefully chosen points, called *quadrature nodes*, and combining these evaluations to minimize the approximation error.

Recently, a stream of work has made use of prior knowledge on the smoothness of the integrand using kernels. Oates et al. (2017) and Liu & Lee (2017) used kernel-based control variates, splitting the computational budget into regressing the integrand and integrating the residual. Bach (2017) looked for the best way to sample i.i.d. nodes and combine the resulting evaluations. Finally, Bayesian quadrature (O’Hagan, 1991; Huszár & Duvenaud, 2012; Briol et al., 2015), herding (Chen et al., 2010; Bach et al., 2012), or the biased importance sampling estimate of Delyon & Portier (2016) all favor *dissimilar* nodes, where dissimilarity is measured by a kernel. Our work falls in this last cluster.

We build on the particular approach of Bardenet & Hardy (2019) for Monte Carlo integration based on projection *determinantal point processes* (DPPs, Hough et al., 2006; Kulesza & Taskar, 2012). DPPs are a repulsive distribution over configurations of points, where repulsion is again parametrized by a kernel. In a sense, DPPs are the kernel machines of point processes.

---

<sup>1</sup> [github.com/guilgautier/DPPy](https://github.com/guilgautier/DPPy)

Fifteen years before [Macchi \(1975\)](#) even formalized DPPs, [Ermakov & Zolotukhin \(EZ, 1960\)](#) had the intuition to use a determinantal structure to sample quadrature nodes, followed by solving a linear system to aggregate the evaluations of the integrand into an unbiased estimator. This linear system yields a simple and interpretable characterization of the variance of their estimator. [Ermakov & Zolotukhin](#)'s result did not diffuse much<sup>2</sup> in the Monte Carlo community, partly because the mathematical and computational machinery to analyze and implement it was not available. Seemingly unaware of this previous work, [Bardenet & Hardy \(2019\)](#) came up with a more natural estimator of the integral of interest, and they could build upon the thorough study of DPPs in random matrix theory ([Johansson, 2006](#)) to obtain a fast central limit theorem (CLT). Since then, DPPs with stationary kernels have also been used by [Mazoyer et al. \(2019\)](#) for Monte Carlo integration. In any case, these DPP-based Monte Carlo estimators crucially depend on efficient sampling procedures for the corresponding (potentially multidimensional) DPP.

**Our contributions.** First, we reveal the close link between DPPs and the approach of [Ermakov & Zolotukhin \(1960\)](#). Second, we provide a simple proof of their result and survey the properties of the EZ estimator with modern arguments. In particular, when the integrand is a linear combination of the eigenfunctions of the kernel of the underlying DPP, the corresponding Fourier-like coefficients can be estimated with zero variance. In other words, one sample of the corresponding DPP yields perfect interpolation of the underlying integrand, by solving a linear system. Third, we propose an efficient Python implementation for exact sampling of a particular DPP, called *multivariate Jacobi ensemble*. The code<sup>1</sup> is available in the DPPy toolbox of [Gautier et al. \(2019\)](#). This implementation allows to numerically investigate the behavior of the two Monte Carlo estimators derived by [Bardenet & Hardy \(2019\)](#) and [Ermakov & Zolotukhin \(1960\)](#), in regimes yet unexplored for any of the two. Fourth, important theoretical properties of both estimators, like the CLT of ([Bardenet & Hardy, 2019](#)), are technically involved. A CLT for EZ promises to be even more difficult to establish. The current empirical investigation provides a motivation and guidelines for more theoretical work. Our point is not to compare DPP-based Monte Carlo estimators to the wide choice of numerical integration algorithms, but to get a fine understanding of their properties so as to fine-tune their design and guide theoretical developments.

## 2 Quadrature, DPPs, and the multivariate Jacobi ensemble

In this section, we quickly survey classical quadrature rules. Then, we define DPPs and give the key properties that make them useful for Monte Carlo integration. Finally, among so-called *projection* DPPs, we introduce the multivariate Jacobi ensemble used by [Bardenet & Hardy \(2019\)](#) to generate quadrature nodes, and on which we base our experimental work.

### 2.1 Standard quadrature

Following [Briol et al. \(2015, Section 2.1\)](#), let  $\mu(dx) = \omega(x) dx$  be a positive Borel measure on  $\mathbb{X} \subset \mathbb{R}^d$  with finite mass and density  $\omega$  w.r.t. the Lebesgue measure. This paper aims to compute integrals of the form  $\int f(x)\mu(dx)$  for some test function  $f : \mathbb{X} \rightarrow \mathbb{R}$ . A quadrature rule approximates such integrals as a weighted sum of evaluations of  $f$  at some *nodes*  $\{x_1, \dots, x_N\} \subset \mathbb{X}$ ,

$$\int f(x)\mu(dx) \approx \sum_{n=1}^N \omega_n f(x_n), \quad (1)$$

where the weights  $\omega_n \triangleq \omega_n(x_1, \dots, x_N)$  do not need to be non-negative nor sum to one.

Among the many quadrature designs mentioned in the introduction ([Evans & Swartz, 2000, Section 5](#)), we pay special attention to the textbook example of the (deterministic) Gauss-Jacobi rule. This scheme applies to dimension  $d = 1$ , for  $\mathbb{X} \triangleq [-1, 1]$  and  $\omega(x) \triangleq (1-x)^a(1+x)^b$  with  $a, b > -1$ . In this case, the nodes  $\{x_1, \dots, x_N\}$  are taken to be the zeros of  $p_N$ , the orthonormal Jacobi polynomial of degree  $N$ , and the weights  $\omega_n \triangleq 1/K(x_n, x_n)$  with  $K(x, x) \triangleq \sum_{k=0}^{N-1} p_k(x)^2$ . In particular, this specific quadrature rule allows to perfectly integrate polynomials up to degree  $2N - 1$  ([Davis & Rabinowitz, 1984, Section 2.7](#)). In a sense, the DPPs of [Bardenet & Hardy \(2019\)](#) are a random, multivariate generalization of Gauss-Jacobi quadrature, as we shall see in Section 3.1.

<sup>2</sup> Many thanks to Mathieu Gerber of Univ. Bristol, UK, for digging up this result from his human memory.

Monte Carlo integration can be defined as random choices of nodes in (1). Importance sampling, for instance, corresponds to i.i.d. nodes, while Markov chain Monte Carlo corresponds to nodes drawn from a carefully chosen Markov chain; see, e.g., [Robert & Casella \(2004\)](#) for more details. Finally, quasi-Monte Carlo (QMC, [Dick & Pillichshammer, 2010](#)) applies to  $\mu$  uniform over a compact subset of  $\mathbb{R}^d$ , and constructs deterministic nodes that spread uniformly, as measured by their *discrepancy*.

## 2.2 Projection DPPs

DPPs can be understood as a parametric class of point processes, specified by a base measure  $\mu$  and a kernel  $K : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{C}$ . The latter is commonly assumed to be Hermitian and trace-class. For the resulting process to be well defined, it is necessary and sufficient that the kernel  $K$  is positive semi-definite with eigenvalues in  $[0, 1]$ , see, e.g., [Soshnikov \(2000, Theorem 3\)](#). When the eigenvalues further belong to  $\{0, 1\}$ , we speak of a *projection* kernel and a *projection* DPP. One practical feature of projection DPPs is that they almost surely produce samples with fixed cardinality, equal to the rank  $N$  of the kernel. More generally, they are the building blocks of DPPs. Indeed, under general assumptions, all DPPs are mixtures of projection DPPs ([Hough et al., 2006, Theorem 7](#)). Hereafter, unless specifically stated,  $K$  is assumed to be a real-valued, symmetric, projection kernel.

One way to define a projection DPP with  $N$  points is to take  $N$  functions  $\phi_0, \dots, \phi_{N-1}$  orthonormal w.r.t.  $\mu$ , i.e.,  $\langle \phi_k, \phi_\ell \rangle \triangleq \int \phi_k(x)\phi_\ell(x)\mu(dx) = \delta_{k\ell}$ , and consider the kernel  $K_N$  associated to the orthogonal projector onto  $\mathcal{H}_N \triangleq \text{span}\{\phi_k, 0 \leq k \leq N-1\}$ , i.e.,

$$K_N(x, y) \triangleq \sum_{k=0}^{N-1} \phi_k(x)\phi_k(y). \quad (2)$$

We say that the set  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathbb{X}$  is drawn from the projection DPP with base measure  $\mu$  and kernel  $K_N$ , denoted by  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \sim \text{DPP}(\mu, K_N)$ , when  $(\mathbf{x}_1, \dots, \mathbf{x}_N)$  has joint distribution

$$\frac{1}{N!} \det(K_N(x_p, x_n))_{p,n=1}^N \mu^{\otimes N}(dx). \quad (3)$$

$\text{DPP}(\mu, K_N)$  indeed defines a probability measure over sets since (3) is invariant by permutation and the orthonormality of the  $\phi_k$ s yields the normalization. See also [Appendix A.1](#) for more details on the construction of projection DPPs from sets of linearly independent functions.

The repulsion of projection DPPs may be understood geometrically by considering the Gram formulation of the kernel (2), namely

$$K_N(x, y) = \Phi(x)^\top \Phi(y), \quad \text{where} \quad \Phi(x) \triangleq (\phi_0(x), \dots, \phi_{N-1}(x))^\top. \quad (4)$$

This allows to rewrite the joint distribution (3) as

$$\frac{1}{N!} \underbrace{\det \Phi(x_{1:N}) \Phi(x_{1:N})^\top}_{=(\det \Phi(x_{1:N}))^2} \mu^{\otimes N}(dx), \quad \text{where} \quad \Phi(x_{1:N}) \triangleq \begin{pmatrix} \phi_0(x_1) & \dots & \phi_{N-1}(x_1) \\ \vdots & & \vdots \\ \phi_0(x_N) & \dots & \phi_{N-1}(x_N) \end{pmatrix}. \quad (5)$$

Thus, the larger the determinant of the *feature matrix*  $\Phi(x_{1:N})$ , i.e., the larger the volume of the parallelotope spanned by the *feature vectors*  $\Phi(x_1), \dots, \Phi(x_N)$ , the more likely  $x_1, \dots, x_N$  co-occur.

## 2.3 The multivariate Jacobi ensemble

In this part, we specify a projection kernel. We follow [Bardenet & Hardy \(2019\)](#) and take its eigenfunctions to be multivariate orthonormal polynomials. In dimension  $d = 1$ , letting  $(\phi_k)_{k \geq 0}$  in (2) be the orthonormal polynomials w.r.t.  $\mu$  results in a projection DPP called an *orthogonal polynomial ensemble* (OPE, [König, 2004](#)). When  $d > 1$ , orthonormal polynomials can still be uniquely defined by applying the Gram-Schmidt procedure to a set of monomials, provided the base measure is not pathological. However, there is no natural order on multivariate monomials: an ordering  $\mathfrak{b} : \mathbb{N}^d \rightarrow \mathbb{N}$  must be picked before we apply Gram-Schmidt to the monomials in  $L^2(\mu)$ . We follow [Bardenet & Hardy \(2019, Section 2.1.3\)](#) and consider multi-indices  $k \triangleq (k^1, \dots, k^d) \in \mathbb{N}^d$  ordered by their maximum degree  $\max_i k^i$ , and for constant maximum degree, by the usual lexicographic order. We still denote the corresponding multivariate orthonormal polynomials by  $(\phi_k)_{k \in \mathbb{N}^d}$ .

By multivariate OPE we mean the projection DPP with base measure  $\mu(dx) \triangleq \omega(x) dx$  and orthogonal projection kernel  $K_N(x, y) \triangleq \sum_{\mathfrak{b}(k)=0}^{N-1} \phi_k(x)\phi_k(y)$ . When the base measure is separable, i.e.,  $\omega(x) = \omega^1(x^1) \times \dots \times \omega^d(x^d)$ , multivariate orthonormal polynomials are products of univariate ones, and the kernel (2) reads

$$K_N(x, y) = \sum_{\mathfrak{b}(k)=0}^{N-1} \prod_{i=1}^d \phi_{k^i}^i(x^i)\phi_{k^i}^i(y^i), \quad (6)$$

where  $(\phi_\ell^i)_{\ell \geq 0}$  are the orthonormal polynomials w.r.t.  $\omega^i(z) dz$ . For  $\mathbb{X} = [-1, 1]^d$  and  $\omega^i(z) = (1-z)^{a^i}(1+z)^{b^i}$ , with  $a^i, b^i > -1$ , the resulting DPP is called a *multivariate Jacobi ensemble*.

### 3 Monte Carlo integration with projection DPPs

Our goal is to design random quadrature rules (1) on  $\mathbb{X} \triangleq [-1, 1]^d$  with desirable properties. We focus on computing  $\int f(x)\mu(dx)$  with the two unbiased DPP-based Monte Carlo estimators of Bardenet & Hardy (BH, 2019) and Ermakov & Zolotukhin (EZ, 1960). We start by presenting the natural BH estimator which, when associated to the multivariate Jacobi ensemble, comes with a CLT with a faster rate than classical Monte Carlo. Then, we survey the properties of the less obvious EZ estimator. Using a generalization of the Cauchy-Binet formula we provide a slight improvement of the key result of EZ. Despite the lack of result illustrating a fast convergence rate, the EZ estimator has a practical and interpretable variance. In particular, this estimator turns a single DPP sample into a perfect integrator as well as a perfect interpolator of functions that are linear combinations of eigenfunctions of the associated kernel. Finally, we detail our exact sampling procedure for multivariate Jacobi ensemble, which allows to exploit the best of both the BH and EZ estimators.

#### 3.1 A natural estimator

For  $f \in L^1(\mu)$ , Bardenet & Hardy (2019) consider

$$\widehat{I}_N^{\text{BH}}(f) \triangleq \sum_{n=1}^N \frac{f(\mathbf{x}_n)}{K_N(\mathbf{x}_n, \mathbf{x}_n)}, \quad (7)$$

as an unbiased estimator of  $\int f(x)\mu(dx)$ , with variance (see, e.g., Lavancier et al., 2012, Section 2.1)

$$\mathbb{V}\text{ar}[\widehat{I}_N^{\text{BH}}(f)] = \frac{1}{2} \int \left( \frac{f(x)}{K_N(x, x)} - \frac{f(y)}{K_N(y, y)} \right)^2 K_N(x, y)^2 \mu(dx)\mu(dy), \quad (8)$$

which clearly captures a notion of smoothness of  $f$  w.r.t.  $K_N$  but its interpretation is not obvious.

For  $\mathbb{X} = [-1, 1]^d$ , the interest in multivariate Jacobi ensemble among DPPs comes from the fact that (7) can be understood as a (randomized) multivariate counterpart of the Gauss-Jacobi quadrature introduced in Section 2.1. Moreover, for  $f$  essentially  $\mathcal{C}^1$ , Bardenet & Hardy (2019, Theorem 2.7) proved a CLT with faster-than-classical-Monte-Carlo decay,

$$\sqrt{N^{1+1/d}} \left( \widehat{I}_N^{\text{BH}}(f) - \int f(x)\mu(dx) \right) \xrightarrow[N \rightarrow \infty]{\text{law}} \mathcal{N}(0, \Omega_{f, \omega}^2), \quad (9)$$

with  $\Omega_{f, \omega}^2 \triangleq \frac{1}{2} \sum_{k \in \mathbb{N}^d} (k^1 + \dots + k^d) \mathcal{F}_{\frac{f\omega}{\omega_{\text{eq}}}}(k)^2$ , where  $\mathcal{F}_g$  denotes the Fourier transform of  $g$ , and  $\omega_{\text{eq}}(x) \triangleq 1 / \prod_{i=1}^d \pi \sqrt{1 - (x^i)^2}$ . In the fast CLT (9), the asymptotic variance is governed by the smoothness of  $f$  since  $\Omega_{f, \omega}$  is a measure of the decay of the Fourier coefficients of the integrand.

#### 3.2 The Ermakov-Zolotukhin estimator

We start by stating the main finding of Ermakov & Zolotukhin (1960), see also Evans & Swartz (2000, Section 6.4.3) and references therein. To the best of our knowledge, we are the first to make the connection with projection DPPs, as defined in Section 2.2. This allows us to give a slight improvement and provide a simpler proof of the original result, based on a generalization of the Cauchy-Binet formula (Johansson, 2006). Finally, we apply Ermakov & Zolotukhin's (1960) technique to build an unbiased estimator of  $\int f(x)\mu(dx)$ , which comes with a practical and interpretable variance.

**Theorem 1.** Consider  $f \in L^2(\mu)$  and  $N$  functions  $\phi_0, \dots, \phi_{N-1} \in L^2(\mu)$  orthonormal w.r.t.  $\mu$ . Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \sim \text{DPP}(\mu, K_N)$ , with  $K_N(x, y) = \sum_{k=0}^{N-1} \phi_k(x)\phi_k(y)$ . Consider the linear system

$$\begin{pmatrix} \phi_0(\mathbf{x}_1) & \dots & \phi_{N-1}(\mathbf{x}_1) \\ \vdots & & \vdots \\ \phi_0(\mathbf{x}_N) & \dots & \phi_{N-1}(\mathbf{x}_N) \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_N) \end{pmatrix}. \quad (10)$$

Then, the solution of (10) is unique,  $\mu$ -almost surely, with coordinates  $y_k = \frac{\det \Phi_{\phi_{k-1}, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})}$ , where  $\Phi_{\phi_{k-1}, f}(\mathbf{x}_{1:N})$  is the matrix obtained by replacing the  $k$ -th column of  $\Phi(\mathbf{x}_{1:N})$  by  $f(\mathbf{x}_{1:N})$ . Moreover, for all  $1 \leq k \leq N$ , the coordinate  $y_k$  of the solution vector satisfies

$$\mathbb{E}[y_k] = \langle f, \phi_{k-1} \rangle, \quad \text{and} \quad \text{Var}[y_k] = \|f\|^2 - \sum_{\ell=0}^{N-1} \langle f, \phi_\ell \rangle^2. \quad (11)$$

We improved the original result by showing that  $\text{Cov}[y_j, y_k] = 0$ , for all  $1 \leq j \neq k \leq N$ .

Before we provide the proof, also detailed in Appendix A.2, several remarks are in order. We start by considering a function  $f \triangleq \sum_{k=0}^{M-1} \langle f, \phi_k \rangle \phi_k$ ,  $1 \leq M \leq \infty$ , where  $(\phi_k)_{k \geq 0}$  forms an orthonormal basis of  $L^2(\mu)$ , e.g., the Fourier basis or wavelet bases (Mallat & Peyré, 2009). Next, we build the orthogonal projection kernel  $K_N$  onto  $\mathcal{H}_N \triangleq \text{span}\{\phi_0, \dots, \phi_{N-1}\}$  as in (2). Then, Theorem 1 shows that solving (10), with points  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \sim \text{DPP}(\mu, K_N)$ , provides unbiased estimates of the  $N$  Fourier-like coefficients  $(\langle f, \phi_k \rangle)_{k=0}^{N-1}$ . Remarkably, these estimates are uncorrelated and have the same variance (11) equal to the residual  $\sum_{k=N}^{\infty} \langle f, \phi_k \rangle^2$ . The faster the decay of the coefficients, the smaller the variance. In particular, for  $M \leq N$ , i.e.,  $f \in \mathcal{H}_N$ , the estimators have zero variance. Put differently,  $f$  can be reconstructed perfectly from only one sample of  $\text{DPP}(\mu, K_N)$ .

*Proof.* First, the joint distribution (5) of  $(\mathbf{x}_1, \dots, \mathbf{x}_N)$  is proportional to  $(\det \Phi(x_{1:N}))^2 \mu^{\otimes N}(x)$ . Thus, the matrix  $\Phi(\mathbf{x}_{1:N})$  defining the linear system (10) is invertible,  $\mu$ -almost surely, and the expression of the coordinates follows from Cramer's rule. Then, we treat the case  $k = 1$ , the others follow the same lines. The proof relies on the orthonormality of the  $\phi_k$ s and a generalization of the Cauchy-Binet formula (A.1), cf. Lemma A. The expectation in (11) reads

$$\begin{aligned} \mathbb{E} \left[ \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})} \right] &\stackrel{(5)}{=} \frac{1}{N!} \int \det \Phi_{\phi_0, f}(x_{1:N}) \det \Phi(x_{1:N}) \mu^{\otimes N}(dx) \\ &\stackrel{(A.1)}{=} \det \begin{pmatrix} \langle f, \phi_0 \rangle & (\langle f, \phi_\ell \rangle)_{\ell=1}^{N-1} \\ 0_{N-1,1} & I_{N-1} \end{pmatrix} = \langle f, \phi_0 \rangle. \end{aligned} \quad (12)$$

Similarly, the second moment reads

$$\begin{aligned} \mathbb{E} \left[ \left( \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})} \right)^2 \right] &\stackrel{(5)}{=} \frac{1}{N!} \int \det \Phi_{\phi_0, f}(x_{1:N}) \det \Phi_{\phi_0, f}(x_{1:N}) \mu^{\otimes N}(dx) \\ &\stackrel{(A.1)}{=} \det \begin{pmatrix} \|f\|^2 & (\langle f, \phi_\ell \rangle)_{\ell=1}^{N-1} \\ (\langle f, \phi_k \rangle)_{k=1}^{N-1} & I_{N-1} \end{pmatrix} = \|f\|^2 - \sum_{k=1}^{N-1} \langle f, \phi_k \rangle^2. \end{aligned} \quad (13)$$

Finally, the variance in (11) = (13) - (12)<sup>2</sup>. The covariance is treated in Appendix A.2.  $\square$

In the setting of the multivariate Jacobi ensemble described in Section 2.3, the first orthonormal polynomial  $\phi_0$  is constant, equal to  $\mu([-1, 1]^d)^{-1/2}$ . Hence, a direct application of Theorem 1 yields

$$\widehat{I}_N^{\text{EZ}}(f) \triangleq \frac{y_1}{\phi_0} = \mu([-1, 1]^d)^{1/2} \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})}, \quad (14)$$

as an unbiased estimator of  $\int_{[-1, 1]^d} f(x) \mu(dx)$ , see Appendix A.3. We also show that (14) can be viewed as a quadrature rule (1) with weights summing to  $\mu([-1, 1]^d)$ . Unlike the variance of  $\widehat{I}_N^{\text{BH}}(f)$  in (8), the variance of  $\widehat{I}_N^{\text{EZ}}(f)$  clearly reflects the accuracy of the approximation of  $f$  by its projection onto  $\mathcal{H}_N$ . In particular, it allows to integrate and interpolate polynomials up to “degree”  $b^{-1}(N-1)$ , perfectly. Nonetheless, its limiting theoretical properties, like a CLT, look hard to establish. In particular, the dependence of each quadrature weight on all quadrature nodes makes the estimator a peculiar object that doesn't fit the assumptions of traditional CLTs for DPPs (Soshnikov, 2000).

### 3.3 How to sample from projection DPPs and the multivariate Jacobi ensemble

To perform Monte Carlo integration with DPPs, it is crucial to sample the points and evaluate the weights efficiently. However, sampling from continuous DPPs remains a challenge. In this part, we review briefly the main technique for projection DPP sampling before we develop our method to generate samples from the multivariate Jacobi ensemble. The code<sup>1</sup> is available in the DPPy toolbox of [Gautier et al. \(2019\)](#), the associated documentation<sup>3</sup> contains a lot more details on DPP sampling.

In both finite and continuous cases, except for some specific instances, exact DPP sampling usually requires the spectral decomposition of the underlying kernel ([Lavancier et al., 2012](#), Section 2.4). However, for projection DPPs, prior knowledge of the eigenfunctions is not necessary, only an oracle to evaluate the kernel is required. Next, we describe the generic exact sampler of [Hough et al. \(2006](#), Algorithm 18). It is based on the chain rule and valid exclusively for projection DPPs.

For simplicity, consider a projection DPP  $(\mu, K_N)$  with  $\mu(dx) = \omega(x) dx$  and  $K_N$  as in (2). This DPP has exactly  $N$  points,  $\mu$ -almost surely ([Hough et al., 2006](#), Lemma 17). To get a valid sample  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , it is enough to apply the chain rule to sample  $(\mathbf{x}_1, \dots, \mathbf{x}_N)$  and forget the order the points were selected. The chain rule scheme can be derived from two different perspectives. Either using Schur complements to express the determinant in the joint distribution (3),

$$\frac{K_N(x_1, x_1)}{N} \omega(x_1) dx_1 \prod_{n=2}^N \frac{K_N(x_n, x_n) - \mathbf{K}_{n-1}(x_n)^\top \mathbf{K}_{n-1}^{-1} \mathbf{K}_{n-1}(x_n)}{N - (n-1)} \omega(x_n) dx_n, \quad (15)$$

where  $\mathbf{K}_{n-1}(\cdot) = (K_N(x_1, \cdot), \dots, K_N(x_{n-1}, \cdot))^\top$ , and  $\mathbf{K}_{n-1} = (K_N(x_p, x_q))_{p,q=1}^{n-1}$ . Or geometrically using the base $\times$ height formula to express the squared volume in the joint distribution (5),

$$\frac{\|\Phi(x_1)\|^2}{N} \omega(x_1) dx_1 \prod_{n=2}^N \frac{\text{distance}^2(\Phi(x_n), \text{span}\{\Phi(x_p)\}_{p=1}^{n-1})}{N - (n-1)} \omega(x_n) dx_n. \quad (16)$$

Note that the numerators in (15) correspond to the incremental posterior variances of a noise-free Gaussian process model with kernel  $K_N$  ([Rasmussen & Williams, 2006](#)), giving yet another intuition for repulsion. When using the chain rule, the practical challenge is twofold: find efficient ways to (i) evaluate the conditional densities, (ii) sample exactly from the conditionals.

In this work, we take  $\mathbb{X} = [-1, 1]^d$  and focus on sampling the multivariate Jacobi ensemble with parameters  $|a^i|, |b^i| \leq 1/2$ , cf. Section 2.3. We remodeled the original implementation<sup>4</sup> of the multivariate Jacobi ensemble sampler accompanying the work of [Bardenet & Hardy \(BH, 2019\)](#) in a more Pythonic way. In particular, we address the previous challenges in the following way:

(i) contrary to BH, we leverage the Gram structure to vectorize the computations and consider (16) instead of (15), and evaluate  $K_N(x, y)$  via (4) instead of (6). The overall procedure is akin to a sequential Gram-Schmidt orthogonalization of the feature vectors  $\Phi(x_1), \dots, \Phi(x_N)$ . Moreover we pay special attention to avoiding unnecessary evaluations of orthogonal polynomials (OP) when computing a feature vector  $\Phi(x)$ . This is done by coupling the slicing feature of the Python language with the dedicated method `scipy.special.eval_jacobi`, used to evaluate the three-term recurrence relations satisfied by each of the univariate Jacobi polynomials. Given the chosen ordering  $\mathbf{b}$ , the computation of  $\Phi(x)$  requires the evaluation of  $d$  recurrence relations up to depth  $\sqrt[4]{N}$ .

(ii) like BH, we sample each conditional in turn using a rejection sampling mechanism with the same proposal distribution. But BH take as proposal  $\omega_{\text{eq}}(x) dx$ , which corresponds to the limiting marginal of the multivariate Jacobi ensemble as  $N$  goes to infinity; see ([Simon, 2011](#), Section 3.11). On our side, we use a two-layer rejection sampling scheme. We rather sample from the  $n$ -th conditional using the marginal distribution  $N^{-1} K_N(x, x) \omega(x) dx$  as proposal and rejection constant  $N/(N - (n-1))$ . This allows us to reduce the number of (costly) evaluations of the acceptance ratio. The marginal distribution itself is sampled using the same proposal  $\omega_{\text{eq}}(x) dx$  and rejection constant as BH. The rejection constant, of order  $2^d$ , is derived from [Chow et al. \(1994\)](#) and [Gautschi \(2009\)](#). We further reduced the number of OP evaluations by considering  $N^{-1} K_N(x, x) \omega(x) dx$  as a mixture, where each component in (6) involves only one OP. In the end, the expected total number of rejections is of order  $2^d N \log N$ . Finally, we implemented a specific rejection free method for the univariate Jacobi ensemble; a special continuous projection DPP which can be sampled exactly in  $\mathcal{O}(N^2)$ , by computing the eigenvalues of a random tridiagonal matrix ([Killip & Nenciu, 2004](#), Theorem 2).

<sup>3</sup> [dppy.readthedocs.io](http://dppy.readthedocs.io) <sup>4</sup> [github.com/rbardenet/dppmc](https://github.com/rbardenet/dppmc)

All these improvements resulted in dramatic speedups. For example, on a modern laptop, sampling a 2D Jacobi ensemble with  $N = 1000$  points, see Figure 1(a), takes less than a minute, compared to hours previously. For more details on the sampling procedure, we refer to Appendix A.4.

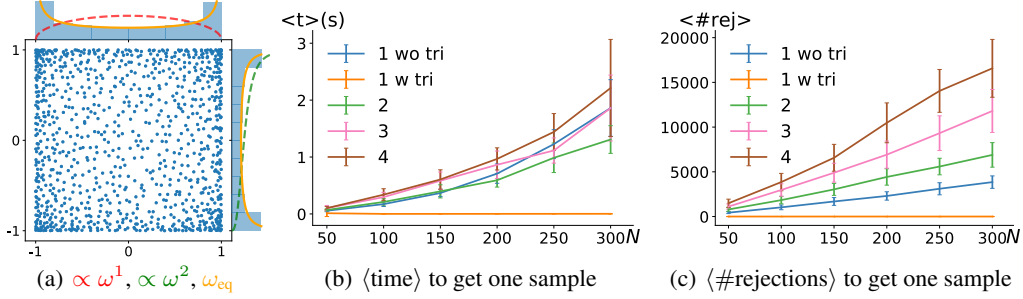


Figure 1: (a) A sample from a 2D Jacobi ensemble with  $N = 1000$  points. (b)-(c)  $a^i, b^i = -1/2$ , the colors and numbers correspond to the dimension. For  $d = 1$ , the tridiagonal model (tri) of Killip & Nenciu offers tremendous time savings. (c) The total number of rejections grows as  $2^d N \log(N)$ .

## 4 Empirical investigation

We perform three main sets of experiments to investigate the properties of the BH (7) and EZ (14) estimators of the integral  $\int f(x)\mu(dx)$ . We add the baseline vanilla Monte Carlo, where points are drawn i.i.d. proportionally to  $\mu$ . The two estimators are built from the multivariate Jacobi ensemble, cf. Section 2.3. First, we extend, for larger  $N$ , the experiments of Bardenet & Hardy (2019) illustrating their fast CLT (9) on a smooth function. Then, we illustrate Theorem 1 by considering polynomial functions that can be either fully or partially decomposed in the eigenbasis of the DPP kernel. Finally, we compare the potential of both estimators on various non smooth functions.

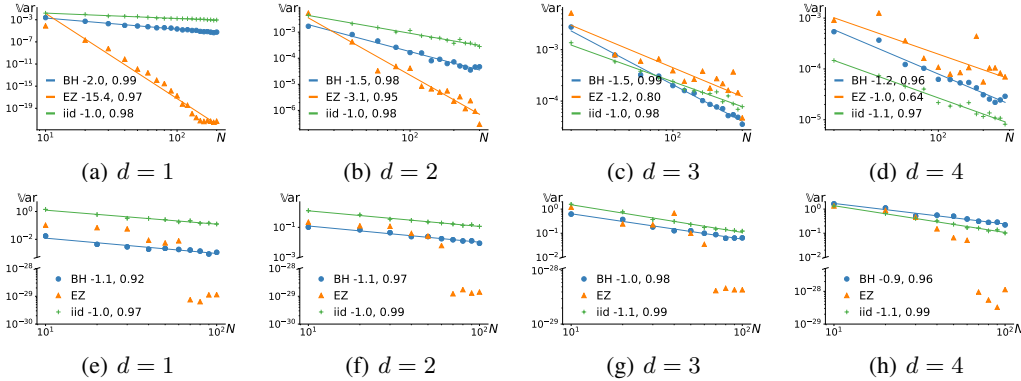


Figure 2: (a)-(d) cf. Section 4.1, the numbers in the legend are the slope and  $R^2$  (e)-(h) cf. Section 4.2.

### 4.1 The bump experiment

Bardenet & Hardy (2019, Section 3) illustrate the behavior of  $\hat{I}_N^{\text{BH}}$  and its CLT (9) on a unimodal, smooth *bump* function; see Appendix B.1. The expected variance decay is of order  $1/N^{1+1/d}$ . We reproduce their experiment in Figure 2 for larger  $N$ , and compare with the behavior of  $\hat{I}_N^{\text{EZ}}$ . In short,  $\hat{I}_N^{\text{EZ}}$  dramatically outperforms  $\hat{I}_N^{\text{BH}}$  in  $d \leq 2$ , with surprisingly fast empirical convergence rates. When  $d \geq 3$ , performance decreases, and  $\hat{I}_N^{\text{BH}}$  shows both faster and more regular variance decay.

To know whether we can hope for a CLT for  $\hat{I}_N^{\text{EZ}}$ , we performed Kolmogorov-Smirnov tests for  $N = 300$ , which yielded small  $p$ -values across dimensions, from 0.03 to 0.24. This is compared to the same  $p$ -values for  $\hat{I}_N^{\text{BH}}$ , which range from 0.60 to 0.99. The results are presented in Appendix B.1. The lack of normality of  $\hat{I}_N^{\text{EZ}}$  is partly due to a few outliers. Where these outliers come from is left for future work; ill-conditioning of the linear system (10) is an obvious candidate. Besides, contrary to  $\hat{I}_N^{\text{BH}}$ , the estimator  $\hat{I}_N^{\text{EZ}}$  has no guarantee to preserve the sign of integrands having constant sign.



## 4.2 Integrating sums of eigenfunctions

In the next series of experiments, we are mainly interested in testing the variance decay of  $\widehat{I}_N^{\text{EZ}}(f)$  prescribed by Theorem 1. To that end, we consider functions of the form

$$f(x) = \sum_{\mathbf{b}(k)=0}^{M-1} \frac{1}{\mathbf{b}(k) + 1} \phi_k(x), \quad (17)$$

whose integral w.r.t.  $\mu$  is to be estimated based on realizations of the multivariate Jacobi ensemble with kernel  $K_N(x, y) = \sum_{\mathbf{b}(k)=0}^{N-1} \phi_k(x)\phi_k(y)$ , where  $N \neq M$  a priori. This means that the function  $f$  can be either fully ( $M \leq N$ ) or partially ( $M > N$ ) decomposed in the eigenbasis of the kernel. In both cases, we let the number of points  $N$  used to build the two estimators vary from 10 to 100 in dimensions  $d = 1$  to 4. In the first setting, we set  $M = 70$ . Thus,  $N$  eventually reaches the number of functions used to build  $f$  in (17), after what  $\widehat{I}_N^{\text{EZ}}$  is an exact estimator, see Figure 2(e)-(h). The second setting has  $M = N + 1$ , so that the number of points  $N$  is never enough for the variance in (11) to be zero. The results of both settings are to be found in Appendix B.2.

In the first case, for each dimension  $d$ , we indeed observe a drop in the variance of  $\widehat{I}_N^{\text{EZ}}$  once the number of points of the DPP hits the threshold  $N = M$ . This is in perfect agreement with Theorem 1: once  $f \in \mathcal{H}_M \subseteq \mathcal{H}_N$ , the variance in (11) is zero. In the second setting, as  $N$  increases the contribution of the extra mode  $\phi_{\mathbf{b}^{-1}(N)}$  in (17) decreases as  $\frac{1}{N}$ . Hence, from Theorem 1 we expect a variance decay of order  $\frac{1}{N^2}$ , which we observe in practice.

## 4.3 Further experiments

In Appendices B.3-B.6 we test the robustness of both BH and EZ estimators, when applied to functions presenting discontinuities or which do not belong to the span of the eigenfunctions of the kernel. Although the conditions of the CLT (9) associated to  $\widehat{I}^{\text{BH}}$  are violated, the corresponding variance decay is smooth but not as fast. For  $\widehat{I}^{\text{EZ}}$ , the performance deteriorates with the dimension. Indeed, the cross terms arising from the Taylor expansion of the different functions introduce monomials, associated to large coefficients, that do not belong to  $\mathcal{H}_N$ . Sampling more points would reduce the variance (11). But more importantly, for EZ to excel, this suggests to adapt the kernel to the basis where the integrand is known to be sparse or to have fast-decaying coefficients. In regimes where BH and EZ do not shine, vanilla Monte Carlo becomes competitive for small values of  $N$ .

## 5 Conclusion

Ermakov & Zolotukhin (EZ, 1960) proposed a non-obvious unbiased Monte Carlo estimator using projection DPPs. It requires solving a linear system, which in turn involves evaluating both the  $N$  eigenfunctions of the corresponding kernel and the integrand at the  $N$  points of the DPP sample. This is yet another connection between DPPs and linear algebra. In fact, solving this linear system provides unbiased estimates of the Fourier-like coefficients of the integrand  $f$  with each of the  $N$  eigenfunctions of the DPP kernel. Remarkably, these estimators have identical variance, and this variance measures the accuracy of the approximation of  $f$  by its projection onto these eigenfunctions. With modern arguments, we have provided a much shorter proof of these properties than in the original work of (Ermakov & Zolotukhin, 1960). Beyond this, little is known on the EZ estimator. While coming with a less interpretable variance, the more direct estimator proposed by Bardenet & Hardy (BH, 2019) has an intrinsic connection with the classical Gauss quadrature and further enjoys stronger theoretical properties when using multivariate Jacobi ensemble.

Our experiments highlight the key features of both estimators when the underlying DPP is a multivariate Jacobi ensemble, and further demonstrate the known properties of the BH estimator in yet unexplored regimes. Although EZ shows a *surprisingly fast* empirical convergence rate for  $d \leq 2$ , its behavior is more erratic for  $d \geq 3$ . Ill-conditioning of the linear system is a potential source of outliers in the distribution of the estimator. Regularization may help but would introduce a stability/bias trade-off. More generally, EZ seems worth investigating for integration or even interpolation tasks where the function is known to be decomposable in the eigenbasis of the kernel, i.e., in a setting similar to the one of Bach (2017). Finally, the new implementation of an exact sampler for multivariate Jacobi ensemble unlocks more large-scale empirical investigations and asks for more theoretical work. The associated code<sup>1</sup> is available in the DPPy toolbox of Gautier et al. (2019).

## References

- Bach, F. [On the Equivalence between Kernel Quadrature Rules and Random Feature Expansions](#). *Journal of Machine Learning Research*, 2017. arXiv:1502.06800.
- Bach, F., Lacoste-Julien, S., and Obozinski, G. [On the Equivalence between Herding and Conditional Gradient Algorithms](#). In *International Conference on Machine Learning (ICML)*, 2012. arXiv:1203.4523.
- Bardenet, R. and Hardy, A. [Monte Carlo with Determinantal Point Processes](#). *Annals of Applied Probability*, in press, 2019. arXiv:1605.00361.
- Briol, F.-X., Oates, C. J., Girolami, M., and Osborne, M. A. [Frank-Wolfe Bayesian Quadrature: Probabilistic Integration with Theoretical Guarantees](#). In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015. arXiv:1506.02681.
- Chen, Y., Welling, M., and Smola, A. [Super-Samples from Kernel Herding](#). In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2010. arXiv:1203.3472.
- Chow, Y., Gatteschi, L., and Wong, R. [A Bernstein-type inequality for the Jacobi polynomial](#). *Proceedings of the American Mathematical Society*, 1994.
- Davis, P. J. and Rabinowitz, P. [Methods of numerical integration](#). Academic Press. 1984.
- Delyon, B. and Portier, F. [Integral approximation by kernel smoothing](#). *Bernoulli*, 2016. arXiv:1409.0733.
- Dick, J. and Pillichshammer, F. [Digital nets and sequences : discrepancy and quasi-Monte Carlo integration](#). Cambridge University Press. 2010.
- Ermakov, S. M. and Zolotukhin, V. G. [Polynomial Approximations and the Monte-Carlo Method](#). *Theory of Probability & Its Applications*, 1960.
- Evans, M. and Swartz, T. [Approximating integrals via Monte Carlo and deterministic methods](#). Oxford University Press. 2000.
- Gautier, G., Polito, G., Bardenet, R., and Valko, M. [DPPy: DPP Sampling with Python](#). *Journal of Machine Learning Research - Machine Learning Open Source Software (JMLR-MLOSS)*, in press, 2019. arXiv:1809.07258.
- Gautschi, W. [How sharp is Bernstein's Inequality for Jacobi polynomials?](#) *Electronic Transactions on Numerical Analysis*, 2009.
- Hough, J. B., Krishnapur, M., Peres, Y., and Virág, B. [Determinantal Processes and Independence](#). In *Probability Surveys*. 2006. arXiv:math/0503110.
- Huszár, F. and Duvenaud, D. [Optimally-Weighted Herding is Bayesian Quadrature](#). In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2012. arXiv:1204.1664.
- Johansson, K. [Random matrices and determinantal processes](#). *Les Houches Summer School Proceedings*, 2006.
- Killip, R. and Nenciu, I. [Matrix models for circular ensembles](#). *International Mathematics Research Notices*, 2004. arXiv:math/0410034.
- König, W. [Orthogonal polynomial ensembles in probability theory](#). *Probability Surveys*, 2004. arXiv:math/0403090.
- Kulesza, A. and Taskar, B. [Determinantal Point Processes for Machine Learning](#). *Foundations and Trends in Machine Learning*, 2012. arXiv:1207.6083.
- Lavancier, F., Møller, J., and Rubak, E. [Determinantal point process models and statistical inference : Extended version](#). *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, 2012. arXiv:1205.4818.

- Liu, Q. and Lee, J. D. **Black-Box Importance Sampling**. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017. arXiv:1610.05247.
- Macchi, O. **The coincidence approach to stochastic point processes**. *Advances in Applied Probability*, 1975.
- Mallat, S. and Peyré, G. **A wavelet tour of signal processing : the sparse way**. Elsevier/Academic Press. 2009.
- Mazoyer, A., Coeurjolly, J.-F., and Amblard, P.-O. **Projections of determinantal point processes**. *ArXiv e-prints*, 2019. arXiv:1901.02099v3.
- Oates, C. J., Girolami, M., and Chopin, N. **Control functionals for Monte Carlo integration**. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2017. arXiv:1410.2392.
- O’Hagan, A. **Bayes–Hermite quadrature**. *Journal of Statistical Planning and Inference*, 1991.
- Rasmussen, C. E. and Williams, C. K. I. **Gaussian processes for machine learning**. MIT Press. 2006.
- Robert, C. P. **The Bayesian choice : from decision-theoretic foundations to computational implementation**. Springer. 2007.
- Robert, C. P. and Casella, G. **Monte Carlo statistical methods**. Springer-Verlag New York. 2004.
- Simon, B. **Szegő’s theorem and its descendants: Spectral theory for  $l^2$  perturbations of orthogonal polynomials**. M. B. Porter Lecture Series, Princeton Univ. Press, Princeton, NJ. 2011.
- Soshnikov, A. **Determinantal random point fields**. *Russian Mathematical Surveys*, 2000. arXiv:math/0002099.

## A Methodology

### A.1 The generalized Cauchy-Binet formula: a modern argument

Johansson (2006, Section 2.2) developed a natural way to build DPPs associated to projection (potentially non-Hermitian) kernels. In this part, we focus on the generalization of the Cauchy-Binet formula (Johansson, 2006, Proposition 2.10). Its usefulness is twofold for our purpose. First, it serves to justify the fact that the normalization constant of the joint distribution (3) is one, i.e., it is indeed a probability distribution. Second, we use it as a modern and simple argument to prove a slight improvement of the result of Ermakov & Zolotukhin (1960), cf. Theorem 1. An extended version of the proof is given in Appendix A.2.

**Lemma A.** (Johansson, 2006, Proposition 2.10) *Let  $(\mathbb{X}, \mathcal{B}, \mu)$  be a measurable space and consider measurable functions  $\phi_0, \dots, \phi_{N-1}$  and  $\psi_0, \dots, \psi_{N-1}$ , such that  $\phi_k \psi_\ell \in L^1(\mu)$ . Then,*

$$\det(\langle \phi_k, \psi_\ell \rangle)_{k,\ell=0}^{N-1} = \frac{1}{N!} \int \det \Phi(x_{1:N}) \det \Psi(x_{1:N}) \mu^{\otimes N}(dx), \quad (\text{A.1})$$

where

$$\Phi(x_{1:N}) = \begin{pmatrix} \phi_0(x_1) & \dots & \phi_{N-1}(x_1) \\ \vdots & & \vdots \\ \phi_0(x_N) & \dots & \phi_{N-1}(x_N) \end{pmatrix} \quad \text{and} \quad \Psi(x_{1:N}) = \begin{pmatrix} \psi_0(x_1) & \dots & \psi_{N-1}(x_1) \\ \vdots & & \vdots \\ \psi_0(x_N) & \dots & \psi_{N-1}(x_N) \end{pmatrix}$$

### A.2 Proof of Theorem 1

First, we recall the result of Ermakov & Zolotukhin (1960), cf. Theorem 1. Then, we provide a modern proof based on the generalization of the Cauchy-Binet formula, cf. Lemma A, where we exploit the orthonormality of the eigenfunctions of the kernel.

**Theorem B.** *Consider  $f \in L^2(\mu)$  and  $N$  functions  $\phi_0, \dots, \phi_{N-1} \in L^2(\mu)$  orthonormal w.r.t.  $\mu$ , i.e.,*

$$\langle \phi_k, \phi_\ell \rangle \triangleq \int \phi_k(x) \phi_\ell(x) \mu(dx) = \delta_{k\ell}, \quad \forall 0 \leq k, \ell \leq N-1. \quad (\text{A.2})$$

Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \sim \text{DPP}(\mu, K_N)$ , with projection kernel  $K_N(x, y) = \sum_{k=0}^{N-1} \phi_k(x) \phi_k(y)$ . That is to say  $(\mathbf{x}_1, \dots, \mathbf{x}_N)$  has joint distribution

$$\frac{1}{N!} \det(K_N(x_p, x_q))_{p,q=1}^N \mu^{\otimes N}(dx) = \frac{1}{N!} (\det \Phi(x_{1:N}))^2 \mu^{\otimes N}(dx). \quad (\text{A.3})$$

Consider the linear system  $\Phi(\mathbf{x}_{1:N})y = f(\mathbf{x}_{1:N})$ , that is,

$$\begin{pmatrix} \phi_0(\mathbf{x}_1) & \dots & \phi_{N-1}(\mathbf{x}_1) \\ \vdots & & \vdots \\ \phi_0(\mathbf{x}_N) & \dots & \phi_{N-1}(\mathbf{x}_N) \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_N) \end{pmatrix}. \quad (\text{A.4})$$

Then, the solution of (A.4) is unique,  $\mu$ -almost surely, with coordinates

$$y_k = \frac{\det \Phi_{\phi_{k-1}, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})}, \quad (\text{A.5})$$

where  $\Phi_{\phi_{k-1}, f}(\mathbf{x}_{1:N})$  is the matrix obtained by replacing the  $k$ -th column of  $\Phi(\mathbf{x}_{1:N})$  by  $f(\mathbf{x}_{1:N})$ . Moreover, for all  $1 \leq k \leq N$ , the coordinate  $y_k$  of the solution vector satisfies

$$\mathbb{E}[y_k] = \langle f, \phi_{k-1} \rangle, \quad \text{and} \quad \text{Var}[y_k] = \|f\|^2 - \sum_{\ell=0}^{N-1} \langle f, \phi_\ell \rangle^2. \quad (\text{A.6})$$

We improved the original result by showing that  $\text{Cov}[y_j, y_k] = 0$ , for all  $1 \leq j \neq k \leq N$ .

*Proof of Theorem B.* First, the joint distribution (A.3) of  $(\mathbf{x}_1, \dots, \mathbf{x}_N)$  is proportional to  $(\det \Phi(x_{1:N}))^2 \mu^{\otimes N}(x)$ . Thus,  $\det \Phi(\mathbf{x}_{1:N}) \neq 0$ ,  $\mu$ -almost surely. Hence, the matrix  $\Phi(\mathbf{x}_{1:N})$  defining the linear system (A.4) is invertible,  $\mu$ -almost surely.

The expression of the coordinates (A.5) follows from Cramer's rule.

Then, we treat the case  $k = 1$ , the others follow the same lines. The proof relies on Lemma A where we exploit the orthonormality of the  $\phi_k$ s (A.2). The expectation (A.6) reads

$$\begin{aligned}
\mathbb{E} \left[ \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})} \right] &\stackrel{(A.3)}{=} \frac{1}{N!} \int \det \Phi_{\phi_0, f}(x_{1:N}) \det \Phi(x_{1:N}) \mu^{\otimes N}(dx) \\
&\stackrel{(A.1)}{=} \det \begin{pmatrix} \langle f, \phi_0 \rangle & (\langle f, \phi_\ell \rangle)_{\ell=1}^{N-1} \\ (\langle \phi_k, \phi_0 \rangle)_{k=1}^{N-1} & (\langle \phi_k, \phi_\ell \rangle)_{k, \ell=1}^{N-1} \end{pmatrix} \\
&\stackrel{(A.2)}{=} \det \begin{pmatrix} \langle f, \phi_0 \rangle & (\langle f, \phi_\ell \rangle)_{\ell=1}^{N-1} \\ 0_{N-1,1} & I_{N-1} \end{pmatrix} \\
&= \langle f, \phi_0 \rangle. \tag{A.7}
\end{aligned}$$

Similarly, the second moment reads

$$\begin{aligned}
\mathbb{E} \left[ \left( \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})} \right)^2 \right] &\stackrel{(A.3)}{=} \frac{1}{N!} \int \det \Phi_{\phi_0, f}(x_{1:N}) \det \Phi_{\phi_0, f}(x_{1:N}) \mu^{\otimes N}(dx) \\
&\stackrel{(A.1)}{=} \det \begin{pmatrix} \langle f, f \rangle & (\langle f, \phi_\ell \rangle)_{\ell=1}^{N-1} \\ (\langle \phi_k, f \rangle)_{k=1}^{N-1} & (\langle \phi_k, \phi_\ell \rangle)_{k, \ell=1}^{N-1} \end{pmatrix} \\
&\stackrel{(A.2)}{=} \det \begin{pmatrix} \|f\|^2 & (\langle f, \phi_\ell \rangle)_{\ell=1}^{N-1} \\ (\langle f, \phi_k \rangle)_{k=1}^{N-1} & I_{N-1} \end{pmatrix} \\
&= \|f\|^2 - \sum_{k=1}^{N-1} \langle f, \phi_k \rangle^2. \tag{A.8}
\end{aligned}$$

Finally, the variance in (A.6) = (A.8) - (A.7)<sup>2</sup>.

With the same arguments, for  $j \neq k$ , we can compute the covariance  $\text{Cov}[y_j, y_k]$ . For simplicity, we treat only the case  $j = 1, k = 2$ , the general case follows the same lines.

$$\begin{aligned}
\text{Cov}[y_1, y_2] &= \mathbb{E} \left[ \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N}) \det \Phi_{\phi_1, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})} \right] - \mathbb{E} \left[ \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})} \right] \mathbb{E} \left[ \frac{\det \Phi_{\phi_1, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})} \right] \\
&\stackrel{(A.3), (A.7)}{=} \frac{1}{N!} \int \det \Phi_{\phi_0, f}(x_{1:N}) \det \Phi_{\phi_1, f}(x_{1:N}) \mu^{\otimes N}(dx) - \langle f, \phi_0 \rangle \langle f, \phi_1 \rangle \\
&\stackrel{(A.1)}{=} \det \begin{pmatrix} \langle f, \phi_0 \rangle & \langle f, f \rangle & (\langle f, \phi_\ell \rangle)_{\ell=2}^{N-1} \\ (\langle \phi_k, \phi_0 \rangle)_{k=1}^{N-1} & (\langle \phi_k, f \rangle)_{k=1}^{N-1} & (\langle \phi_k, \phi_\ell \rangle)_{k=1, \ell=2}^{N-1} \end{pmatrix} - \langle f, \phi_0 \rangle \langle f, \phi_1 \rangle \\
&\stackrel{(A.2)}{=} \det \begin{pmatrix} \langle f, \phi_0 \rangle & \|f\|^2 & (\langle f, \phi_\ell \rangle)_{\ell=2}^{N-1} \\ 0 & \langle \phi_1, f \rangle & 0 \\ 0_{N-2,1} & (\langle \phi_k, f \rangle)_{k=2}^{N-1} & I_{N-2} \end{pmatrix} - \langle f, \phi_0 \rangle \langle f, \phi_1 \rangle \\
&= \langle f, \phi_0 \rangle \langle f, \phi_1 \rangle - \langle f, \phi_0 \rangle \langle f, \phi_1 \rangle = 0.
\end{aligned}$$

□

### A.3 The EZ estimator as a quadrature rule

In this part, we consider Theorem B in the setting where one of the eigenfunctions of the kernel, say  $\phi_0$  is constant. In this case, we show that  $\widehat{I}_N^{\text{EZ}}(f)$  defined by (14) provides an unbiased estimate of  $\int_{\mathbb{X}} f(x) \mu(dx)$  with known variance. In addition, it can be seen as a quadrature rule in the sense of (1), with weights a priori non negative weights  $\omega_n$  that sum to  $\mu(\mathbb{X})$ . This is a non obvious fact, judging from the expression (14) of the estimator.

**Proposition 1.** Consider  $\phi_0$  constant in Theorem B. Then, solving the corresponding linear system (A.4) allows to construct

$$\widehat{I}_N^{\text{EZ}}(f) \triangleq \frac{y_1}{\phi_0} = \mu(\mathbb{X})^{1/2} \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N})}{\det \Phi(\mathbf{x}_{1:N})}, \quad (\text{A.9})$$

as an unbiased estimator of  $\int_{\mathbb{X}} f(x) \mu(dx)$ , with variance equal to  $\mu(\mathbb{X}) \times (\text{A.6})$ . In addition, it can be seen as a random quadrature rule (1) with weights summing to  $\mu(\mathbb{X})$ .

*Proof.* Since  $\phi_0$  is constant with unit norm we have  $\phi_0 = \mu(\mathbb{X})^{-1/2}$ , so that

$$\mathbb{E}[\widehat{I}_N^{\text{EZ}}(f)] = \frac{1}{\phi_0} \mathbb{E}[y_1] = \frac{1}{\phi_0} \langle f, \phi_0 \rangle = \int_{\mathbb{X}} f(x) dx,$$

and

$$\text{Var}[\widehat{I}_N^{\text{EZ}}(f)] = \frac{1}{\phi_0^2} \text{Var}[y_1] = \mu(\mathbb{X}) \times (\text{A.6}).$$

In addition, (A.9) can be written

$$\widehat{I}_N^{\text{EZ}}(f) = \mu(\mathbb{X})^{1/2} \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N})}{\phi_0 \det \Phi_{\phi_0, 1}(\mathbf{x}_{1:N})} = \mu(\mathbb{X}) \frac{\det \Phi_{\phi_0, f}(\mathbf{x}_{1:N})}{\det \Phi_{\phi_0, 1}(\mathbf{x}_{1:N})},$$

and the expansion of the numerator w.r.t. the first column yields

$$\widehat{I}_N^{\text{EZ}}(f) = \sum_{n=1}^N f(\mathbf{x}_n) \underbrace{\frac{\mu(\mathbb{X})}{\det \Phi_{\phi_0, 1}(\mathbf{x}_{1:N})} (-1)^{1+n} \det(\phi_k(x_p))_{k=1, p=1 \neq n}^{N-1, N}}_{\triangleq \omega_n(\mathbf{x}_{1:N})}. \quad (\text{A.10})$$

Note that there is a priori no reason for the weights to be nonnegative. Finally,

$$\sum_{n=1}^N \omega_n(\mathbf{x}_{1:N}) = \frac{\mu(\mathbb{X})}{\det \Phi_{\phi_0, 1}(\mathbf{x}_{1:N})} \underbrace{\sum_{n=1}^N (-1)^{1+n} \det(\phi_k(x_p))_{k=1, p=1 \neq n}^{N-1, N}}_{= \det \Phi_{\phi_0, 1}(\mathbf{x}_{1:N})} = \mu(\mathbb{X}).$$

This concludes the proof.  $\square$

#### A.4 Sampling from the multivariate Jacobi ensemble

We mention that the code<sup>1</sup> and the documentation<sup>3</sup> associated to this work are available in the DPPy toolbox of Gautier et al. (2019).

In dimension  $d = 1$ , we implemented the random tridiagonal matrix model of Killip & Nenciu (2004, Theorem 2) to sample from the univariate Jacobi ensemble, with base measure  $\mu(dx) = (1-x)^a (1+x)^b dx$ , where  $a, b > -1$ . That is to say, this one dimensional continuous projection DPP with  $N$  points can be sampled in  $\mathcal{O}(N^2)$ , by computing the eigenvalues of random tridiagonal matrix with i.i.d. coefficients of size  $N \times N$ .

Next, for  $d \geq 2$ , we detail the procedure described in Section 3.3 for sampling exactly from the multivariate Jacobi ensemble with parameters  $|a^i|, |b^i| \leq \frac{1}{2}$ , for all  $1 \leq i \leq d$ .

More specifically, we consider sampling exactly from the projection DPP( $\mu, K_N$ ) where

- $\mu(dx) = \omega(x) dx$ , with

$$\omega(x) = \prod_{i=1}^d \omega^i(x^i), \quad \text{where} \quad \omega^i(z) = \prod_{i=1}^d (1-z)^{a^i} (1+z)^{b^i}, \quad \text{and} \quad |a^i|, |b^i| \leq \frac{1}{2}. \quad (\text{A.11})$$

- $K_N(x, y) = \sum_{b(b)=0}^{N-1} \phi_k(x) \phi_k(y)$ , with

$$\phi_k(x) = \prod_{i=1}^d \phi_{k^i}^i(x^i), \quad \text{where} \quad \int_{-1}^1 \phi_u^i(z) \phi_v^i(z) \omega^i(z) dz = \delta_{uv}. \quad (\text{A.12})$$

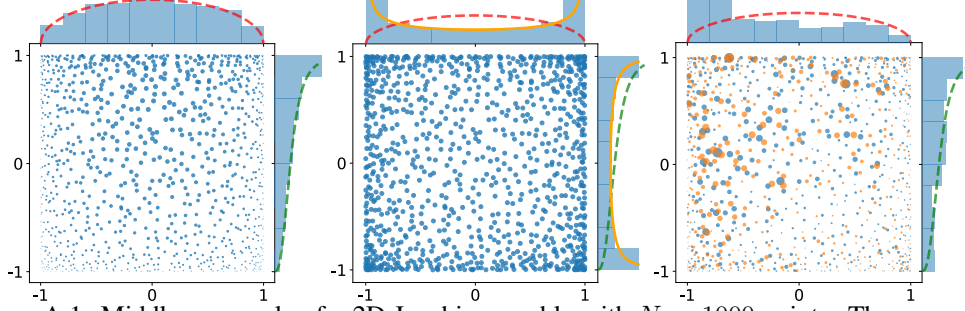


Figure A.1: Middle: a sample of a 2D Jacobi ensemble with  $N = 1000$  points. The normalized reference densities, proportional to  $(1-x)^{a^1}(1+x)^{b^1}$  and  $(1-y)^{a^2}(1+y)^{b^2}$ , are displayed in dashed lines. The empirical marginal densities which converges to the arcsine density  $\omega_{\text{eq}}(x) = \frac{1}{\pi\sqrt{1-x^2}}$  is plotted in solid line. Left: we plot the same sample where the disk centered at  $\mathbf{x}_n$  now has now an area proportional to the weight  $1/K_N(\mathbf{x}_n, \mathbf{x}_n)$  in  $\hat{I}_N^{\text{BH}}(f)$  in (7). Observe that these weights serve as a proxy for the reference measure, like Gaussian quadrature. Right: the weight in  $\hat{I}_N^{\text{EZ}}(f)$  given by (A.10); observe that they can be either positive or negative. The histogram of the absolute value of the weights is plotted on the marginal axes

As an illustration, Figure A.1 displays a sample of a  $d = 2$  Jacobi ensemble with  $N = 1000$  points. Our sampling scheme is an instance of the generic chain-rule-based procedure of Hough et al. (2006, Algorithm 18) where the knowledge of the eigenfunctions can be leveraged, see also Lavancier et al. (2012, Algorithm 1). In our case, sampling  $N$  points in dimension  $d$ , requires an expected total number of rejections of order  $2^d N \log(N)$ . As mentioned in Section 3.3, to sample from  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \sim \text{DPP}(\mu, K_N)$  it is enough to sample  $(\mathbf{x}_1, \dots, \mathbf{x}_N)$  and forget the order the points were selected. Starting from the two formulations (3) and (5) of the joint distribution, the chain rule scheme can be derived from two different perspectives. Either by expressing the determinant  $\det(K_N(x_p, x_n))_{p,n=1}^N$  using Schur complements

$$\begin{aligned}
 (3) &= \frac{1}{N!} \det(K_N(x_p, x_n))_{p,n=1}^N \prod_{n=1}^N \omega(x_n) dx_n & (\text{A.13}) \\
 &= \frac{K_N(x_1, x_1)}{N} \omega(x_1) dx_1 \prod_{n=2}^N \omega(x_n) dx_n \frac{K_N(x_n, x_n) - \mathbf{K}_{n-1}(x_n)^\top \mathbf{K}_{n-1}^{-1} \mathbf{K}_{n-1}(x_n)}{N - (n-1)} \omega(x_n) dx_n,
 \end{aligned}$$

where  $\mathbf{K}_{n-1}(\cdot) = (K_N(x_1, \cdot), \dots, K_N(x_{n-1}, \cdot))^\top$ , and  $\mathbf{K}_{n-1} = (K_N(x_p, x_q))_{p,q=1}^{n-1}$ . Or geometrically using the base $\times$ height formula to express  $(\det \Phi(x_{1:N}))^2$  as the squared volume of the parallelepiped spanned by  $\Phi(x_1), \dots, \Phi(x_N)$

$$\begin{aligned}
 (5) &= \frac{1}{N!} \text{volume}^2(\Phi(x_1), \dots, \Phi(x_N)) \prod_{n=1}^N \omega(x_n) dx_n \\
 &= \frac{\|\Phi(x_1)\|^2}{N} \omega(x_1) dx_1 \prod_{n=2}^N \frac{\text{distance}^2(\Phi(x_n), \text{span}\{\Phi(x_p)\}_{p=1}^{n-1})}{N - (n-1)} \omega(x_n) dx_n. & (\text{A.14})
 \end{aligned}$$

Note that, contrary to (A.14), the formulation (A.13) does not require a priori knowledge of the eigenfunctions of the projection kernel  $K_N$ .

Like Bardenet & Hardy (2019), we sample each conditional in turn using rejection sampling with the same proposal distribution and rejection bound. But where Bardenet & Hardy (2019) use the formulation (A.13) of the chain rule we consider the geometrical perspective (A.14). This allows for a implementation that is simpler (no need to update  $\mathbf{K}_{n-1}^{-1}$ ), fully vectorized, and more interpretable: akin to a sequential Gram-Schmidt orthogonalization of the feature vectors  $\Phi(x_1), \dots, \Phi(x_N)$ .

Moreover, contrary to Bardenet & Hardy (2019) who take  $\omega_{\text{eq}}(x) dx$  as proposal to sample from the each of the conditionals, we use a two-layer rejection sampling scheme. We rather sample from the  $n$ -th conditional using the marginal distribution  $N^{-1} K_N(x, x) \omega(x) dx$ . This choice of proposal allows us to reduce the number of (costly) evaluations of the acceptance ratio.

The rejection constant associated to the  $n$ -th conditional in (A.13) reads

$$\begin{aligned} & \frac{(N - (n - 1))^{-1} (K_N(x, x) - \mathbf{K}_{n-1}(x)^\top \mathbf{K}_{n-1}^{-1} \mathbf{K}_{n-1}(x)) \omega(x)}{N^{-1} K_N(x, x) \omega(x)} \\ &= \frac{N}{N - (n - 1)} \frac{K_N(x, x) - \mathbf{K}_{n-1}(x)^\top \mathbf{K}_{n-1}^{-1} \mathbf{K}_{n-1}(x)}{K_N(x, x)} \leq \frac{N}{N - (n - 1)}. \end{aligned} \quad (\text{A.15})$$

The marginal distribution itself is sampled using the same proposal  $\omega_{\text{eq}}(x) dx$  and rejection constant as [Bardenet & Hardy \(2019\)](#). However, we further reduce the number of computations by considering  $N^{-1} K_N(x, x) \omega(x) dx$  as a mixture, see Section A.4.1

#### A.4.1 Generate samples from the marginal distribution

First, observe that the marginal density can be written as a mixture of  $N$  probability densities where each component is assigned the same weight  $1/N$

$$\frac{1}{N} K_N(x, x) \omega(x) = \frac{1}{N} \sum_{\mathfrak{b}(k)=0}^{N-1} \phi_k(x)^2 \omega(x). \quad (\text{A.16})$$

Thus, sampling from (A.16) can be done in two steps:

- (i) select a multi-index  $k = \mathfrak{b}^{-1}(n)$  with  $n$  drawn uniformly at random in  $\{0, \dots, N - 1\}$
- (ii) sample from  $\phi_k(x)^2 \omega(x) dx$

We perform Step (ii) using rejection sampling with proposal distribution

$$\omega_{\text{eq}}(x) dx = \prod_{i=1}^d \frac{1}{\pi \sqrt{1 - (x^i)^2}} dx^i, \quad (\text{A.17})$$

which corresponds to the limiting marginal distribution of the multivariate Jacobi ensemble as  $N$  goes to infinity; see [\(Simon, 2011, Section 3.11\)](#) and [Figure A.1](#). The acceptance ratio writes

$$\begin{aligned} & \frac{\phi_k(x)^2 \omega(x)}{\omega_{\text{eq}}(x)} \stackrel{(\text{A.12})(\text{A.11})}{=} \prod_{i=1}^d \frac{\phi_{k^i}^i(x^i)^2 \times (1 - x^i)^{a^i} (1 + x^i)^{b^i}}{\pi^{-1} (1 - x^i)^{-\frac{1}{2}} (1 + x^i)^{-\frac{1}{2}}} \\ &= \prod_{i=1}^d \pi (1 - x^i)^{a^i + \frac{1}{2}} (1 + x^i)^{b^i + \frac{1}{2}} \phi_{k^i}^i(x^i)^2. \end{aligned} \quad (\text{A.18})$$

Each of the terms that appear in (A.18) can be bounded using the following recipe:

- (a) For  $k^i = 0$ ,  $\phi_0^i$  is constant and the orthonormality w.r.t.  $(1 - x)^{a^i} (1 + x)^{b^i} dx$  yields

$$(\phi_0^i)^2 \int_{-1}^1 (1 - x)^{a^i} (1 + x)^{b^i} dx = 1 \iff (\phi_0^i)^2 = \frac{1}{2^{a^i + b^i + 1} B(a^i + 1, b^i + 1)}, \quad (\text{A.19})$$

so that the corresponding term in (A.18) becomes

$$\frac{\pi (1 - x)^{a^i + \frac{1}{2}} (1 + x)^{b^i + \frac{1}{2}}}{2^{a^i + b^i + 1} B(a^i + 1, b^i + 1)} \leq \frac{\pi (1 - m)^{a^i + \frac{1}{2}} (1 + m)^{b^i + \frac{1}{2}}}{2^{a^i + b^i + 1} B(a^i + 1, b^i + 1)} \triangleq C_{k^i=0} \leq 2, \quad (\text{A.20})$$

$$\text{where } m = \operatorname{argmax}_{-1 \leq x \leq 1} (1 - x)^{a^i + \frac{1}{2}} (1 + x)^{b^i + \frac{1}{2}} = \begin{cases} 0, & \text{if } a^i = b^i = -\frac{1}{2}, \\ \frac{b^i - a^i}{a^i + b^i + 1}, & \text{otherwise.} \end{cases}$$

- (b) For  $k^i \geq 1$ , we use the bound  $C_{k^i \geq 1}$  (A.22) provided originally by [Chow et al. \(1994\)](#). As mentioned by [Gautschi \(2009\)](#), this bound is probably maximal for  $k^i = 1$  and parameters  $a^i \approx -0.0691$ ,  $b^i = 1/2$ , with value  $\approx 0.64297807\pi \approx 2.02$ .



Finally, the expected number of rejections to perform Step (ii) is equal to  $\prod_{i=1}^d C_{k^i}$  which is of order  $2^d$ , and the expected total number of rejections of the chain rule (A.13) is of order

$$\sum_{n=1}^N 2^d \frac{N}{N - (n - 1)} = 2^d N \sum_{n=1}^N \frac{1}{n} \approx 2^d N \log(N). \quad (\text{A.21})$$

**Proposition 2.** (Gautschi, 2009, Equation 1.3) Let  $(\phi_k)_{k \geq 0}$  be the (univariate) orthonormal polynomials w.r.t.  $(1 - x)^a (1 + x)^b dx$  with  $|a| \leq \frac{1}{2}$ ,  $|b| \leq \frac{1}{2}$ . Then, for any  $x \in [-1, 1]$  and  $k \geq 1$ ,

$$\pi(1 - x)^{a + \frac{1}{2}} (1 + x)^{b + \frac{1}{2}} \phi_k(x)^2 \leq \frac{2 \Gamma(k + a + b + 1) \Gamma(k + \max(a, b) + 1)}{k! (k + \frac{a+b+1}{2})^{2 \max(a, b)} \Gamma(k + \min(a, b) + 1)}. \quad (\text{A.22})$$

#### A.4.2 Empirical timing and number of rejections

In Figure A.2 we illustrate the following observations. Computing the acceptance ratio (A.15) requires to propagate the recurrence relations up to order  $\sqrt[d]{N}$ . Thus, for a given number of points  $N$ , the larger the dimension, the smaller the depth of the recurrence. This could hint that, evaluating the kernel (6) becomes cheaper as  $d$  increases. However, the rejection rate also increases, so that in practice, it is not cheaper to sample in larger dimensions because the number of rejections dominates. In the particular case of dimension  $d = 1$ , samples are generated using the fast and rejection-free tridiagonal matrix model of Killip & Nenciu (2004, Theorem 2). This grants huge time savings compared to the acceptance-rejection method.

Finally, some remarks are in order. Sampling from the  $n$ -th conditional distribution using rejection sampling is common practice (Lavancier et al., 2012, Section 2.4.2). However, tailored proposals with tight rejection constants are required (Lavancier et al., 2012, Appendices E-F). Taking the marginal distribution  $N^{-1} K_N(x, x) \omega(x) dx$  as proposal yields a  $N/(N - (n - 1))$  rejection constant and applies in the general case. Nevertheless, it remains to sample from this marginal distribution. Rejection sampling might be a first option to sample from  $N^{-1} K_N(x, x) \omega(x) dx$ , but when the eigenfunctions are available it could be another option to see it as a mixture (cf. Section A.4.1), where good proposals for each  $\phi_k(x)^2 w(x) dx$  are required.

In the case of (multivariate) orthogonal polynomial ensembles (cf. Section 2.3), evaluations of  $K_N(x, y)$  (6) can be performed using the Gram representation (4),  $K_N(x, y) = \Phi(x)^T \Phi(y)$  and one can leverage the three-term recurrence relations satisfied by each of the univariate Jacobi polynomials  $(\phi_\ell^i)_\ell$ . This is what we do in our special case, we use the dedicated function `scipy.special.eval_jacobi` to evaluate, up to depth  $\sqrt[d]{N}$ , the three-term recurrence relations satisfied by each of the univariate Jacobi. Instead of calling the recursive routine internally to evaluate  $\Phi(x)$ , the corresponding  $d \sqrt[d]{N}$  univariate polynomials or  $N$  multivariate polynomials could be stored in some way and evaluated pointwise on the fly. The preprocessing time and the memory required would increase but it might accelerate the evaluation of  $\Phi(x)$ .

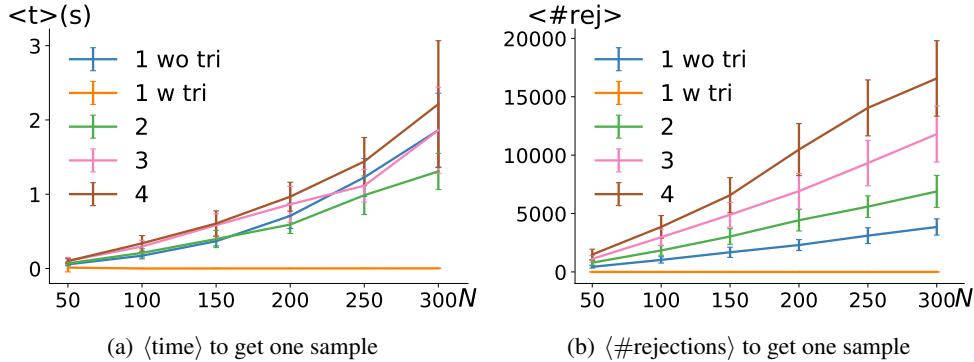


Figure A.2:  $a^i, b^i = -1/2$ , the colors and numbers correspond to the dimension. For  $d = 1$ , the tridiagonal model (tri) of Killip & Nenciu offers tremendous time savings. (b) The total number of rejections grows as  $2^d N \log(N)$  (A.21).

## B Experiments

### B.1 Reproducing the bump example

In Section 4.1, we reproduce the experiment of [Bardenet & Hardy \(2019, Section 3\)](#) where they illustrate the behavior of  $\widehat{I}_N^{\text{BH}}$  on a unimodal, smooth bump function:

$$f(x) = \prod_{i=1}^d \exp\left(-\frac{1}{1-\varepsilon-(x^i)^2}\right) \mathbb{1}_{(-\sqrt{1-\varepsilon}, \sqrt{1-\varepsilon})}(x^i). \quad (\text{B.1})$$

We take  $\varepsilon = 0.05$ . For each value of  $N$ , we sample 100 times from the same multivariate Jacobi ensemble with i.i.d. uniform parameters on  $[-1/2, 1/2]$ , compute the resulting 100 values of each estimator, and plot the two resulting sample variances. In addition, in Figure B.2 we test the potential hope for a CLT for  $\widehat{I}_N^{\text{EZ}}$  and compare with  $\widehat{I}_N^{\text{BH}}$  for which the CLT (9) holds, in the regime  $N = 300$ .

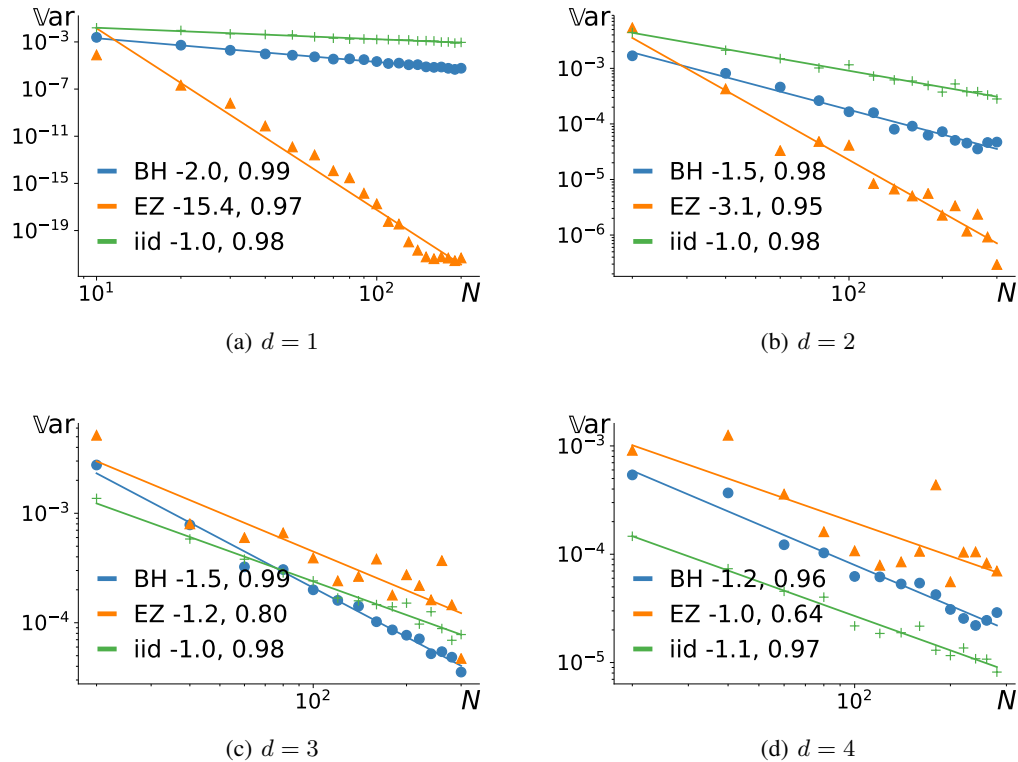


Figure B.1: Reproducing the bump function ( $\varepsilon = 0.05$ ) experiment of [Bardenet & Hardy \(2019\)](#), cf. Section 4.1. Observe the expected variance decay of order  $1/N^{1+1/d}$  for BH. Although vanilla Monte Carlo becomes competitive for small  $N$  as  $d$  increases, its variance decay is of order  $1/N \geq 1/N^{1+1/d}$ . Thus, there will always be meeting point, for some  $N^*$ , after which the variance of BH will be smaller. For  $d = 1$ , EZ has almost no variance for  $N \geq 100$ : the bump function is extremely well approximated by a polynomials of degree  $N \geq 100$ .

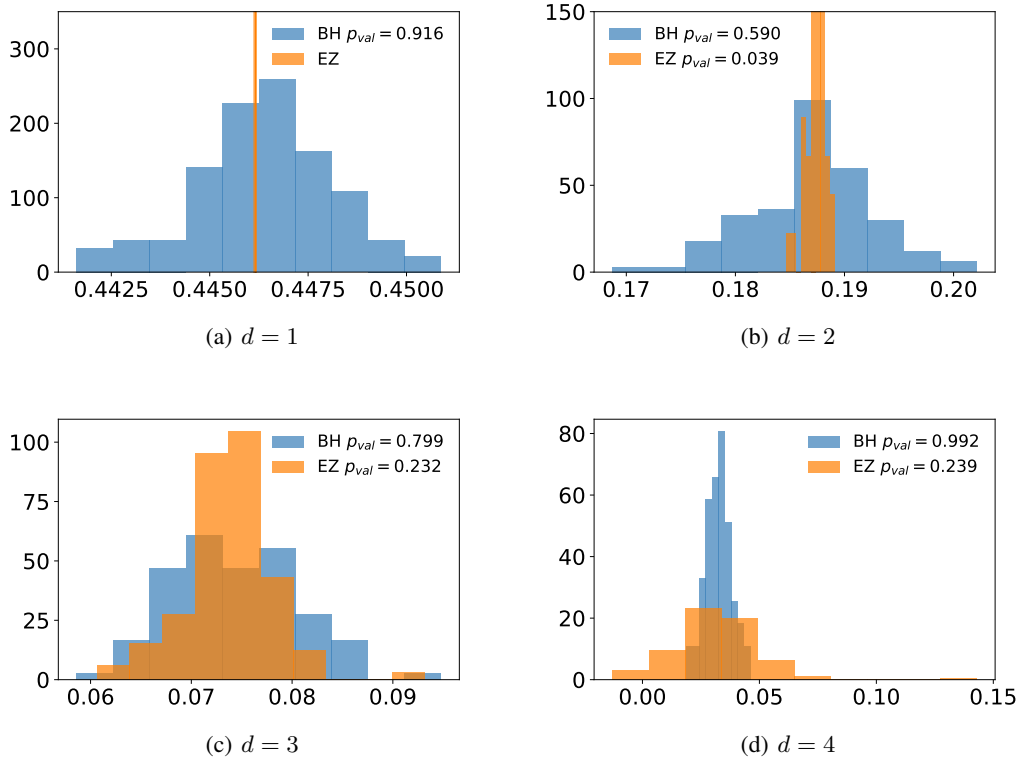


Figure B.2: Histogram of 100 independent estimates  $\hat{I}_N^{\text{BH}}$  and  $\hat{I}_N^{\text{EZ}}$  of the integral of the bump function ( $\varepsilon = 0.05$ ) with  $N = 300$  and associated p-value of Kolmogorov-Smirnov test, cf. Section 4.1. The fluctuations of BH confirm to be Gaussian (cf. CLT (9)). (a) the bump function is extremely well approximated by a polynomial of degree 300 hence  $\hat{I}_N^{\text{EZ}}$  has almost no variance. (b)-(c)-(d) A few outliers seem to break the potential Gaussianity of  $\hat{I}_N^{\text{EZ}}$ . (d)  $\hat{I}_N^{\text{EZ}}(f)$  does not preserve the sign of the integrand.

## B.2 Integrating sums of eigenfunctions

Figure B.3 gives the results of the first setting set in Section 4.2, where we integrate a sum of  $M = 70$  kernel eigenfunctions. In this case,  $EZ$  has zero variance once  $N \geq M$ , a performance that can be reached neither by BH nor vanilla Monte Carlo.

Figure B.4 illustrates the second setting, where the sum always has one more eigenfunction than there are points in the DPP samples. In this case, the conditions for the CLT of BH, cf (9), are not met; there is no  $1/N^{1+1/d}$  guarantee on the variance decay for BH estimator. The performance of BH and vanilla Monte Carlo are comparable. By construction, the variance of EZ decays as  $1/N^2 \leq 1/N$ . Thus, there will always be meeting point, for some  $N^*$ , after which the variance of EZ will be smaller than vanilla Monte Carlo.

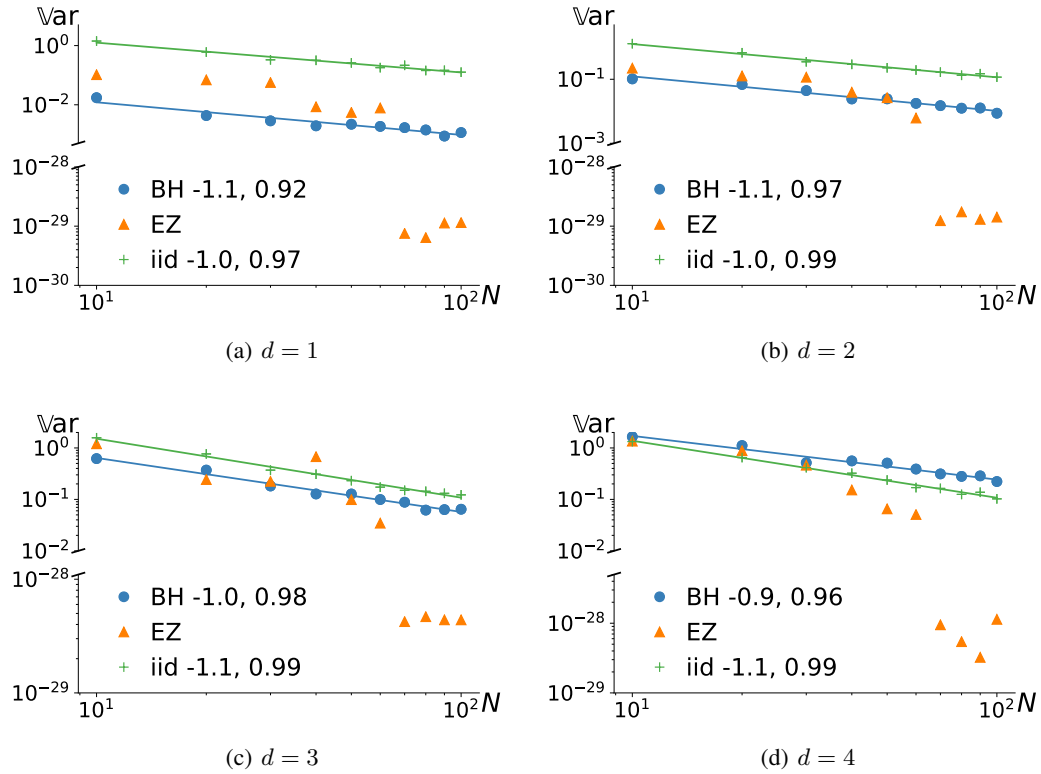
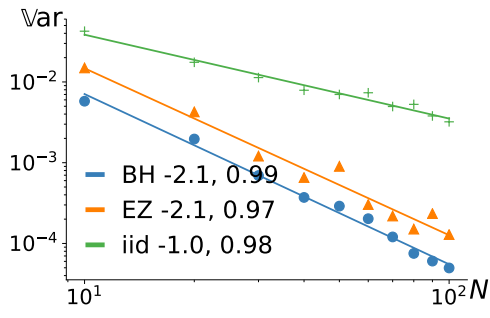
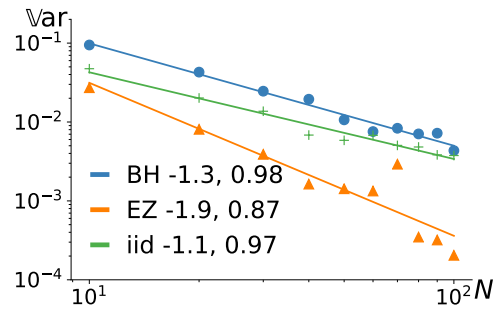


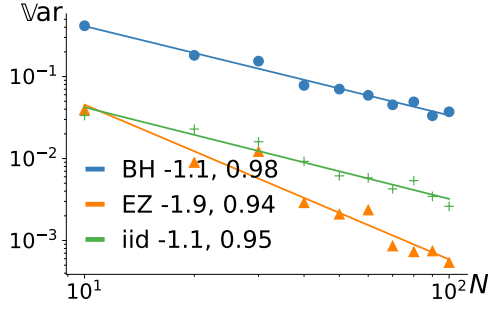
Figure B.3: Comparison of  $\widehat{I}_N^{\text{BH}}$  and  $\widehat{I}_N^{\text{EZ}}$  integrating a finite sum of 70 eigenfunctions of the DPP kernel as in (17), cf. Section 4.2.



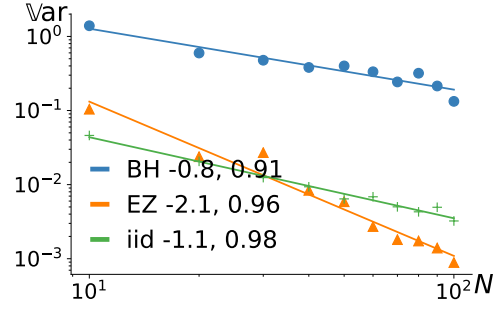
(a)  $d = 1$



(b)  $d = 2$



(c)  $d = 3$



(d)  $d = 4$

Figure B.4: Comparison of  $\widehat{I}_N^{\text{BH}}$  and  $\widehat{I}_N^{\text{EZ}}$  for a linear combination of  $N + 1$  eigenfunctions of the DPP kernel as in (17), cf. Section 4.2.

We now consider cases where the guarantees of BH not EZ are unknown.

### B.3 Integrating absolute value

We consider estimating the integral

$$\int_{[-1,1]^d} \prod_{i=1}^d |x^i| (1-x^i)^{a^i} (1+x^i)^{b^i} dx^i \quad (\text{B.2})$$

where  $a^1, b^1 = -\frac{1}{2}$  and  $a^i, b^i$  i.i.d. uniformly in  $[-\frac{1}{2}, \frac{1}{2}]$ , using BH (7) and EZ (14) estimators.

Results are given in Figure B.5. In dimension  $d = 1$ , the absolute value is well approximated by its truncated Taylor series of low order and EZ performs very well, but as the dimension increases, its performance is more erratic. For  $d \leq 2$ , the performance of BH is smooth and better than vanilla Monte Carlo. In particular, for  $d \leq 2$ , the rate  $1/N^{1+1/d}$  seems to hold for BH while the conditions for the CLT (9) are not satisfied. But it seems no longer true in larger dimension.

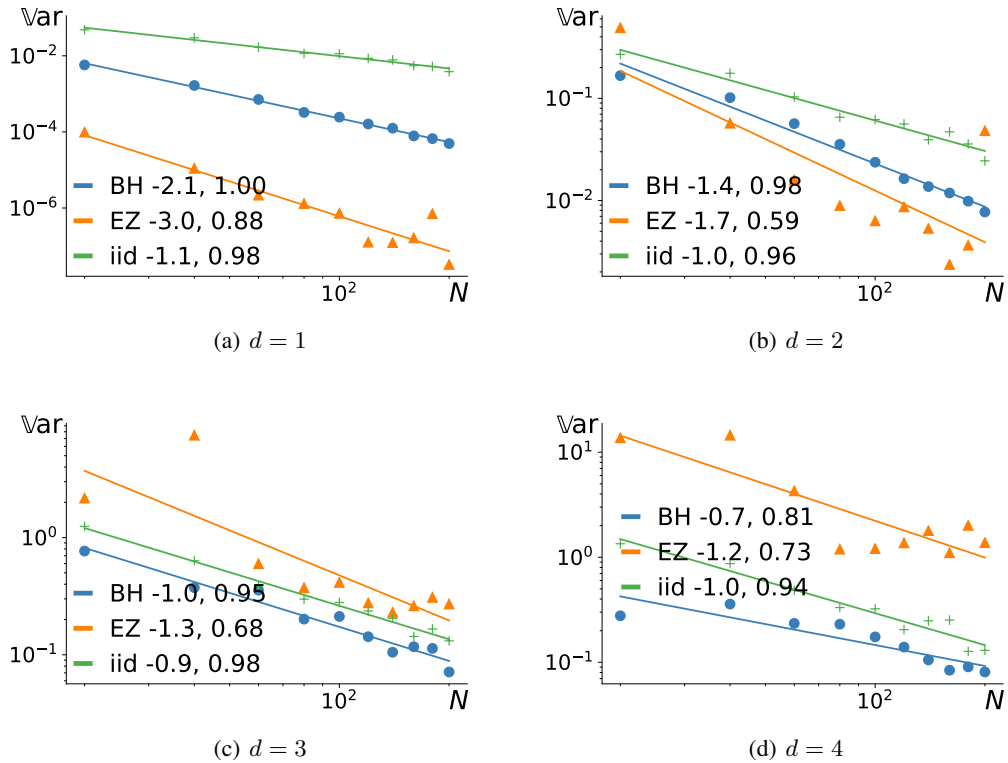


Figure B.5: Comparison of  $\widehat{I}_N^{\text{BH}}$  and  $\widehat{I}_N^{\text{EZ}}$  for absolute value, cf. Section 4.3.

#### B.4 Integrating Heaviside

Let  $H(x) = \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{otherwise} \end{cases}$ . We consider estimating the integral

$$\int_{[-1,1]^d} \prod_{i=1}^d 2 \left( H(x^i) - \frac{1}{2} \right) (1-x^i)^{a^i} (1+x^i)^{b^i} dx^i \quad (\text{B.3})$$

where  $a^1, b^1 = -\frac{1}{2}$  and  $a^i, b^i$  i.i.d. uniformly in  $[-\frac{1}{2}, \frac{1}{2}]$ , using BH (7) and EZ (14) estimators.

Results are given in Figure B.6. The EZ estimator behaves in a very erratic way; it does not seem robust to the discontinuity we have introduced. This can be explained by considering  $H(x) = \frac{1}{2} \lim_{\epsilon \rightarrow 0} 1 + \tanh \frac{x}{\epsilon}$  and taking the product of the Taylor series expansions of  $\tanh$ ; the square of the coefficients in front of the monomials in such expansion become very large as  $\epsilon \rightarrow 0$ . One could expect better behavior for very large  $N$ . The performance of BH is smooth and the rate  $1/N^{1+1/d}$  seems to hold despite the conditions for the CLT (9) are not satisfied.

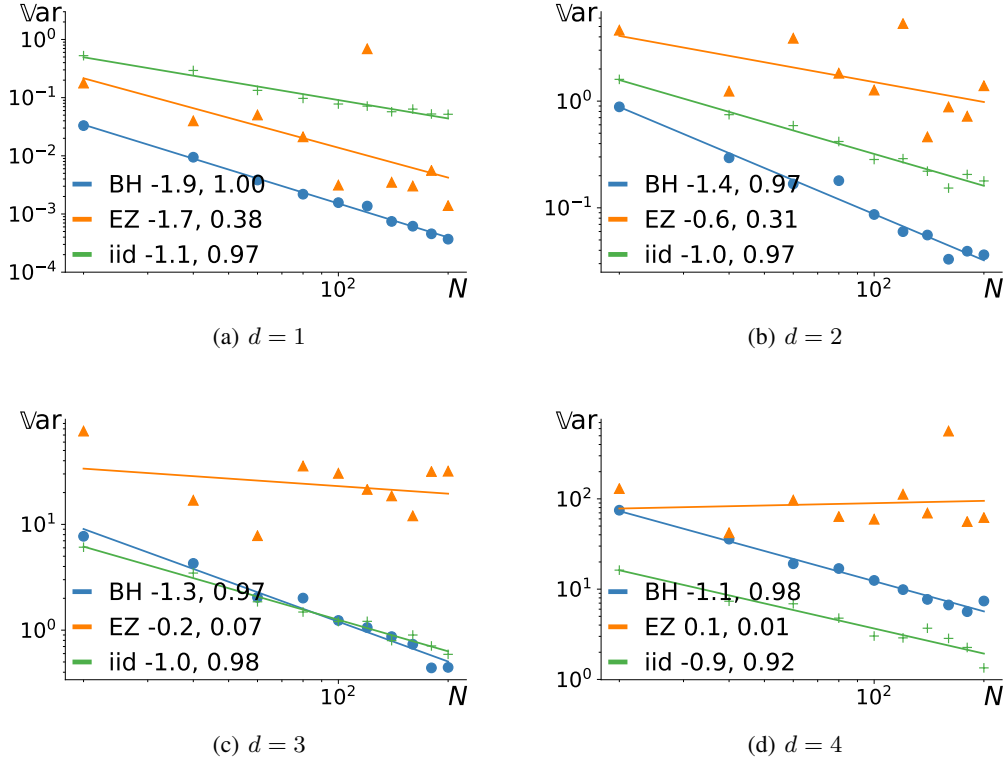


Figure B.6: Comparison of  $\widehat{I}_N^{\text{BH}}$  and  $\widehat{I}_N^{\text{EZ}}$  for Heaviside function, cf. Section 4.3.

## B.5 Integrating cosine

We consider estimating the integral

$$\int_{[-1,1]^d} \prod_{i=1}^d \cos(\pi x^i) (1-x^i)^{a^i} (1+x^i)^{b^i} dx^i \quad (\text{B.4})$$

where  $a^1, b^1 = -\frac{1}{2}$  and  $a^i, b^i$  i.i.d. uniformly in  $[-\frac{1}{2}, \frac{1}{2}]$ , using BH (7) and EZ (14) estimators.

Results are given in Figure B.7. The EZ estimator behaves well for  $d \leq 2$  but its performance deteriorates for  $d \geq 3$ . Indeed, the cross terms arising from the Taylor expansion of the different  $\cos(\pi x^i)$  introduce monomials, associated to large coefficients, that do not belong to  $\mathcal{H}_N$ . One could expect better behavior for very large  $N$ . For  $d \leq 2$ , the rate  $1/N^{1+1/d}$  for BH seems to hold despite the conditions for the CLT (9) are not satisfied. For  $d \geq 3$ , BH and vanilla Monte Carlo behave similarly.

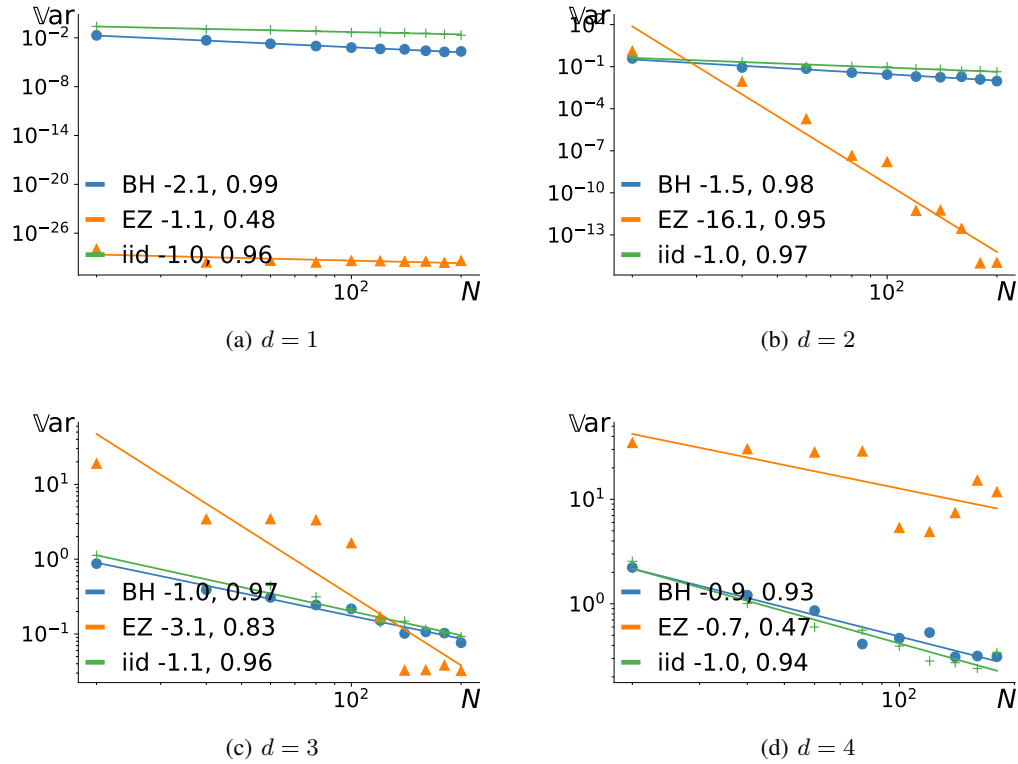


Figure B.7: Comparison of  $\widehat{I}_N^{\text{BH}}$  and  $\widehat{I}_N^{\text{EZ}}$  for cosine, cf. Section 4.3.



### B.6 Integrating a mixture of smooth and non smooth functions

Let  $f(x) = H(x)(\cos(\pi x) + \cos(2\pi x) + \sin(5\pi x))$ . We consider estimating the integral

$$\int_{[-1,1]^d} \prod_{i=1}^d f(x^i)(1-x^i)^{a^i}(1+x^i)^{b^i} dx^i \quad (\text{B.5})$$

where  $a^1, b^1 = -\frac{1}{2}$  and  $a^i, b^i$  i.i.d. uniformly in  $[-\frac{1}{2}, \frac{1}{2}]$ , using BH (7) and EZ (14) estimators.

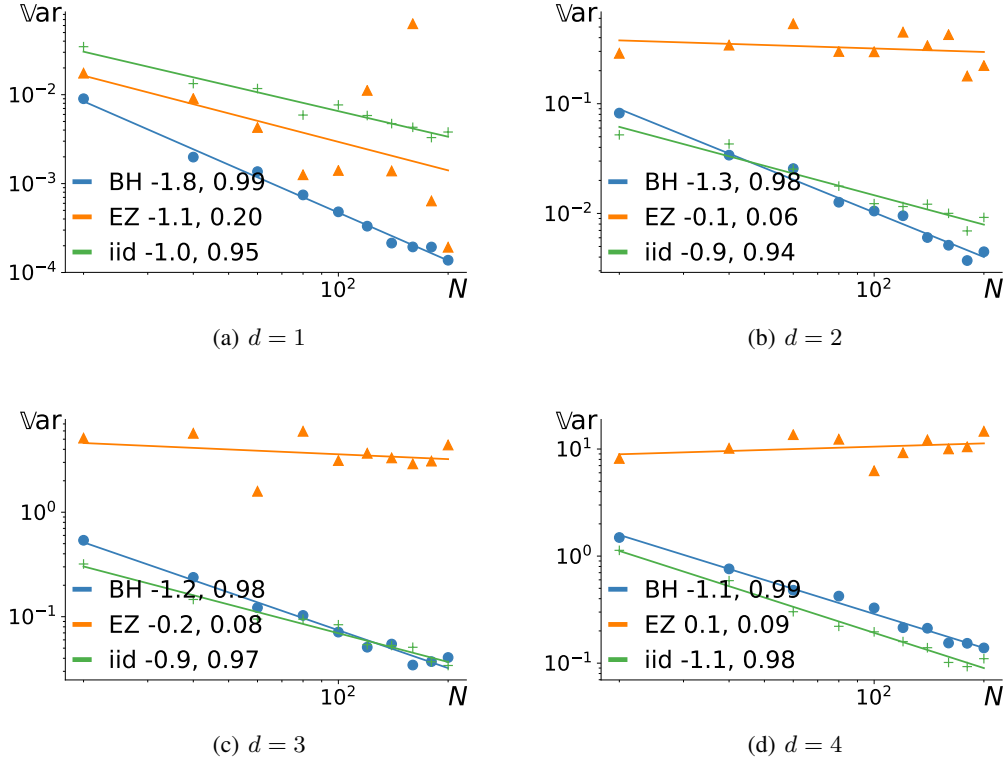


Figure B.8: Comparison of  $\hat{I}_N^{\text{BH}}$  and  $\hat{I}_N^{\text{EZ}}$ , cf. Section 4.3.