

Density Estimation in Uncertainty Propagation Problems Using a Surrogate Model*

Adi Ditkowski[†], Gadi Fibich[†], and Amir Sagiv[†]

Abstract. The effect of uncertainties and noise on a quantity of interest (model output) is often better described by its probability density function (PDF) than by its moments. Although density estimation is a common task, the adequacy of approximation methods (surrogate models) for density estimation has not been analyzed before in the uncertainty-quantification literature. In this paper, we first show that standard surrogate models (such as generalized polynomial chaos), which are highly accurate for moment estimation, might completely fail to approximate the PDF, even for one-dimensional noise. This is because density estimation requires that the surrogate model accurately approximate the gradient of the quantity of interest and not just the quantity of interest itself. Hence, we develop a novel spline-based algorithm for density estimation whose convergence rate in L^q is polynomial in the sampling resolution. This convergence rate is better than that of standard statistical density estimation methods (such as histograms and kernel density estimators) at dimensions $1 \leq d \leq \frac{5}{2}m$, where m is the spline order. Furthermore, we obtain the convergence rate for density estimation with any surrogate model that approximates the quantity of interest and its gradient in L^∞ . Finally, we demonstrate our algorithm for problems in nonlinear optics and fluid dynamics.

Key words. uncertainty quantification, density estimation, probability density function, nonlinear dynamics, spline, surrogate model

AMS subject classifications. 65D07, 65Z05, 62G07, 78A60

DOI. 10.1137/18M1205959

1. Introduction. Uncertainties and noise are prevalent in mathematical models in all branches of science. In such cases, the solution of the (otherwise deterministic) model becomes random, and so one is interested in computing its statistics. This problem, sometimes known as *forward uncertainty propagation*, arises in various areas such as biochemistry [32, 34], fluid dynamics [6, 21, 30, 34], structural engineering [47], hydrology [7], and nonlinear optics [41].

In many applications, one is interested in computing the *probability density function* (PDF) of a certain “quantity of interest” (output) of the model [1, 6, 7, 21, 32, 41, 53]. Often, density estimation is performed using standard uncertainty propagation methods and surrogate models [22, 47], such as stochastic finite element and generalized polynomial chaos (gPC) [23, 35, 46, 59], hp-gPC [56], and Wiener–Haar expansion [31], since these methods can

*Received by the editors August 15, 2018; accepted for publication (in revised form) November 5, 2019; published electronically February 4, 2020.

<https://doi.org/10.1137/18M1205959>

Funding: The research of the first author was supported by the United States–Israel Binational Science Foundation under grant 2016197. The research of the second and third authors was partially supported by the Israel Science Foundation (ISF) under grant 177/13.

[†]School of Mathematical Sciences, Tel Aviv University, Tel Aviv 6997801, Israel (adid@tauex.tau.ac.il, fibich@tauex.tau.ac.il, asagiv88@gmail.com).

approximate moments with spectral accuracy [60, 61]. In this paper we show, however, that methods which are robust and highly accurate for moment approximation are not necessarily so for density estimation. To the best of our knowledge, this observation has not been made before in the UQ literature.

Why is it then that robust moment approximation does not imply robust density estimation? This is because the quantity of interest $f(\alpha)$ and its PDF $p_f(\alpha)$ are explicitly related by (see Lemma 4.1)

$$p_f(y) = \sum_{\alpha \in f^{-1}(y)} \frac{c(\alpha)}{|f'(\alpha)|},$$

where α is the one-dimensional random parameter and $c(\alpha)d\alpha$ is its distribution. This formula and its multidimensional counterpart (Lemma 5.2) show that even if f is well approximated by a function g in L^q , the corresponding density p_g might not be a good approximation of p_f . Indeed, for p_g to approximate p_f , then g' needs to be close to f' , and $g'(\alpha)$ should also vanish if and only if $f'(\alpha)$ does. These conditions might not be satisfied by some of the above-mentioned standard UQ methods. In contrast, spline interpolation approximates both the function and its gradient [3, 25, 40, 44] and does not tend to produce spurious extremal points. Therefore, we construct a novel algorithm for density estimation in uncertainty propagation problems using splines as our surrogate model. With cubic splines, our density-estimation algorithm has a *guaranteed* convergence rate of at least h^3 , where h is the maximal sampling spacing (resolution). More generally, with splines of order m , the convergence rate is at least h^m . These rates are superior to those of the standard kernel density estimators (KDEs) [51, 58] for noise dimension $1 \leq d \leq \frac{5}{2}m$. Our choice of splines is motivated by the availability and efficiency of one- and multidimensional spline toolboxes. Nevertheless, other surrogate models can be used in this algorithm, and indeed this paper lays the theoretical framework for deriving the convergence rate of such methods (Corollaries 4.8 and 5.5). We show, essentially, that density estimation convergence can be performed with any surrogate model for which the L^∞ error of both the function and its gradient converge to zero as the spacing resolution h vanishes. Because we only rely on solving the underlying deterministic model (i.e., our method is nonintrusive), and because interpolation by spline is a standard numerical procedure, our proposed method is very easy to implement.

While the focus of this paper is on density estimation, we also consider the problem of moment estimation using a small sample size. Traditionally, the error bounds of moment estimation for spectral methods (e.g., gPC) are obtained asymptotically as N , the number of samples, goes to infinity. In some applications, however, each solution of the deterministic model is computationally expensive and so the number of samples is limited to, e.g., $N < 100$. Hence, spectrally convergent methods might fail to attain the desired accuracy due to insufficient sampling resolution, *even for one-dimensional noise*. In contrast, the spline-based method approximates moments accurately even when the sample size is small. In addition, high derivatives and discontinuities have little effect on our method's accuracy, due to the fact that spline interpolation is predominantly local (see section 4). Another advantage over gPC is that splines are not limited to a specific choice of sampling points.

The paper is organized as follows. Section 2 introduces the general settings and notation and presents several density-estimation applications from the forward uncertainty propagation

literature. Section 3 reviews standard statistical density-estimation methods (histogram, KDEs) and the gPC method for moment and density estimation. In section 4 we present our spline-based algorithm for moment and density estimation in the one-dimensional case. We then prove that the density-estimation error scales as N^{-m} , where N is the number of samples and m is the order of the splines (Theorem 4.7). Section 5 generalizes our algorithm to d -dimensional noises using tensor-product splines of order m . This section also contains our key theoretical result (Theorem 5.3), that the density-estimation error in the d -dimensional case scales as $N^{-\frac{m}{d}}$.

In section 7 we compare numerically the moment-estimation and density-estimation accuracy of our spline-based method with that of gPC and KDE in one dimension. In addition, in section 6.4 we show that both gPC and our spline-based method can approximate moments and the PDF of certain nonsmooth quantities of interest. We conclude this section with two- and three-dimensional numerical examples (section 6.5). In all cases, the density-estimation errors are consistent with our error estimates (Theorems 4.7 and 5.3). We use our method to compute the PDF of the rotation angle of the polarization ellipse in nonlinear optics (section 7) and the PDF of the shock location in the Burgers equation (section 8). In all these cases, we confirm that the spline-based density estimation converges at least at a cubic rate and observe that the spline-based moments are more accurate than the gPC ones for small sample sizes. Section 9 concludes with open questions and future research directions.

2. Settings and computational goals. We consider initial value problems of the form

$$(2.1) \quad u_t(t, \mathbf{x}; \boldsymbol{\alpha}) = Q(u, \mathbf{x}; \boldsymbol{\alpha})u, \quad u(t = 0, \mathbf{x}; \boldsymbol{\alpha}) = u_0(\mathbf{x}; \boldsymbol{\alpha}),$$

where $\mathbf{x} \in \mathbb{R}^d$, Q is a possibly nonlinear differential operator, and $\boldsymbol{\alpha} \in \Omega \subset \mathbb{R}^m$ is a random variable which is distributed according to a continuous weight function $c(\boldsymbol{\alpha})$, the PDF of the input parameters, such that $\int_{\Omega} c(\boldsymbol{\alpha}) d\boldsymbol{\alpha} = 1$. The randomness of $u(t, \mathbf{x}; \boldsymbol{\alpha})$ is due to the dependence of Q and/or u_0 on $\boldsymbol{\alpha}$.

For a given *quantity of interest* $f(\boldsymbol{\alpha}) := f(u(t, \mathbf{x}); \boldsymbol{\alpha})$, we may wish to perform the following:

1. Moment estimation. Compute the mean, variance, or standard deviation of $f(\boldsymbol{\alpha})$:

$$(2.2) \quad \mathbb{E}_{\boldsymbol{\alpha}}[f] := \int_{\Omega} f(\boldsymbol{\alpha}) c(\boldsymbol{\alpha}) d\boldsymbol{\alpha}, \quad \text{Var}[f] := |\mathbb{E}_{\boldsymbol{\alpha}}[f]|^2 - \mathbb{E}_{\boldsymbol{\alpha}}[|f|^2], \quad \sigma(f) := \sqrt{\text{Var}[f]}.$$

2. Density estimation. Compute the PDF of $f(\boldsymbol{\alpha})$:

$$(2.3a) \quad p(y) := \frac{dP(y)}{dy}, \quad y \in \mathbb{R},$$

where P is the cumulative distribution function (CDF),

$$(2.3b) \quad P(y) := \text{Prob}\{f(\boldsymbol{\alpha}) < y\}.$$

2.1. Applications. Two examples of density estimation in UQ which will be discussed in this paper are the effect of shot-to-shot variation in nonlinear optics (section 7) and hydrodynamical shock formation (section 8). We briefly present two other examples of density estimation in the UQ literature, for which our method can also be applied:

1. Out-of-equilibrium chemical reactions. Belosouov–Zhabotinsky type systems model out-of-equilibrium chemical reactions. One concrete system is the Oregonator [18],

$$\begin{aligned}\frac{dX}{dt} &= k_1 Y - k_2 XY + k_3 X - k_4 X^2, \\ \frac{dY}{dt} &= -k_1 Y - k_2 XY + k_5 Z, \\ \frac{dZ}{dt} &= k_3 X - k_5 Z,\end{aligned}$$

where X , Y , and Z are the concentrations of three different chemical species, and $\{k_i\}_{i=1}^5$ are the rate parameters, often estimated empirically [34]. For large values of t , this system exhibits sustained, temporal oscillations with a frequency $F = F(k_1, \dots, k_5)$. To deal with an uncertainty in the parameters k_4 and k_5 , the authors of [32] computed the moments of X, Y, Z , and the PDF of the oscillations frequency F . This is an example of (2.1)–(2.3) with $\boldsymbol{\alpha} = (k_4, k_5)$ and $f = X, Y, Z$, and F .

2. Heat convection. Consider the flow of a fluid in a two-dimensional box $\mathbf{x} = (x, y) \in [x_1, x_2] \times [y_1, y_2]$, which is modeled by the Navier–Stokes like equations

$$\begin{aligned}\nabla \cdot \mathbf{u} &= 0, & \frac{\partial \theta}{\partial t} + \mathbf{u} \cdot \nabla \theta &= \nabla^2 \theta, \\ \frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u} &= -\nabla p + \text{Pr} \nabla^2 \mathbf{u} + F(\mathbf{u}, \theta),\end{aligned}$$

where $\mathbf{u}(t, \mathbf{x}; \boldsymbol{\alpha})$ is the fluid velocity, $p(t, \mathbf{x}; \boldsymbol{\alpha})$ is the pressure, $\theta(t, \mathbf{x}; \boldsymbol{\alpha})$ is the temperature, Pr is the Prandtl number, and F is the buoyant force [21]. The temperature is a known constant θ_0 on one side of the box but is random on the other side, i.e.,

$$\theta(t, x_1, y) \equiv \theta_0, \quad \theta(t, x_2, y) = \theta_1(y; \boldsymbol{\alpha}).$$

The PDF of the pressure and of the velocity were computed in [53] when $\theta_1(y; \boldsymbol{\alpha}) = \theta_1(\boldsymbol{\alpha})$ and $\boldsymbol{\alpha}$ is uniformly distributed in $[\alpha_{\min}, \alpha_{\max}]$ and in [21] when $\theta_1(y; \boldsymbol{\alpha})$ is a Gaussian random process.

3. Review of existing methods. We briefly present the standard methods in the literature for (2.1)–(2.3).

3.1. Monte Carlo method, the histogram method, and kernel density estimators. Given N independently and identically distributed (i.i.d.) samples $\{\boldsymbol{\alpha}_j\}_{j=1}^N$, the simplest moment estimator is the Monte Carlo approximation $E_{\boldsymbol{\alpha}}[f] \approx \frac{1}{N} \sum_{n=1}^N f(\boldsymbol{\alpha}_n)$. The Monte Carlo method is intuitive and easy to implement. The main drawback of this method is its slow convergence rate of $O(N^{-1/2})$. In cases where each computation of $f(\boldsymbol{\alpha}_j)$ is expensive (e.g., when it requires to solve numerically (2.1) with $\boldsymbol{\alpha} = \boldsymbol{\alpha}_j$), this slow convergence rate can make the Monte Carlo method impractical.

Density estimation using N i.i.d. samples of $f(\boldsymbol{\alpha})$, denoted by $\{f_j\}_{j=1}^N$, is a fundamental problem in nonparametric statistics. A widely used method for density estimation is the histogram method, in which one partitions the range of $f(\boldsymbol{\alpha})$ into L disjoint bins $\{B_\ell\}_{\ell=1}^L$ and approximates the PDF p with the *histogram estimator*

$$(3.1) \quad p_{\text{hist}}(y) := \frac{1}{N} \sum_{\ell=1}^L (\# \text{ of samples for which } f_j \in B_\ell) \cdot \mathbb{1}_{B_\ell}(y),$$

where $\mathbb{1}_{B_\ell}$ is the characteristic function of bin B_ℓ [58]. An alternative family of estimators, unlike the histogram method that can provide a smooth PDF, is KDEs

$$(3.2) \quad p_{\text{kde}}(y) := \frac{1}{Nh} \sum_{j=1}^N K\left(\frac{y - f_j}{h}\right),$$

where $h > 0$ is the “window size” and K is the kernel function (e.g., $K(t) = (2\pi)^{-1/2}e^{-t^2/2}$); see [51, 58]. The L^1 error of the KDE method asymptotically scales as $N^{-2/5}$ [13].¹ As with the Monte Carlo method, this rate is too slow when each evaluation of f_j is computationally expensive.

3.2. Generalized polynomial chaos. The Monte Carlo method, the histogram method, and KDE are all statistical methods, in the sense that they only rely on the sampled values $\{f_j\}_{j=1}^N$. Much more information can be extracted from $\{f_j\}_{j=1}^N$ if the two following conditions hold:

1. The “original” $\{\boldsymbol{\alpha}_j\}_{j=1}^N$ for which $f(\boldsymbol{\alpha}_j) = f_j$ are known.
2. $f(\boldsymbol{\alpha})$ is smooth.²

These two conditions often hold in the general settings of section 2. In such cases, a powerful numerical approach, gPC, can be applied [22, 23, 35, 59]. For clarity, we review the gPC method for a one-dimensional random variable α , i.e., $\Omega \subseteq \mathbb{R}$.

We define the set of orthogonal polynomials $\{p_n(x)\}_{n=0}^\infty$ with respect to $c(\alpha)$ by the conditions [48]

$$(3.3) \quad \text{Deg}(p_n) = n, \quad \int_{\Omega} p_n^*(\alpha)p_m(\alpha)c(\alpha)d\alpha = \delta_{n,m},$$

where p_n^* denotes the complex conjugation of p_m . This family of orthogonal polynomials constitutes an orthonormal basis of the space of square integrable functions, i.e., for all $f \in L^2(\Omega, c)$,

$$(3.4) \quad f(\alpha) = \sum_{n=0}^\infty \hat{f}(n)p_n(\alpha), \quad \hat{f}(n) := \int_{\Omega} f(\alpha)p_n(\alpha)c(\alpha)d\alpha, \quad n = 0, 1, \dots$$

This expansion *converges spectrally* for the classical families of orthogonal polynomials, e.g., the normalized Hermite and Legendre polynomials.³ Specifically, if f is analytic, the truncated expansion (3.4) has the exponential accuracy

¹The mean L^2 error (the squared root of the “MISE”) also asymptotically scales as $N^{-\frac{2}{5}}$ [51, 58].

²In section 6.4 we show how our method can be extended to nonsmooth functions.

³That is, if f is in C^r , then $\{\hat{f}(n)\} \leq cn^{-r}$, and if f is analytic, then $|\hat{f}(n)| \leq ce^{-\gamma n}$, for some $c, \gamma > 0$.

$$(3.5) \quad \left\| f(\alpha) - \sum_{n=0}^{N-1} \hat{f}(n) p_n(\alpha) \right\|_2 \sim C e^{-\gamma N}, \quad N \gg 1,$$

for some constants $C, \gamma > 0$ [50, 57, 59].

The expansion coefficients $\{\hat{f}(n)\}$ (see (3.4)) can be approximated using the Gauss quadrature formula $\mathbb{E}_\alpha[g] \approx \sum_{j=1}^N g(\alpha_j) w_j$, where $\{\alpha_j\}_{j=1}^N$ are the distinct and real roots of $p_N(\alpha)$, $w_j := \int_\Omega l_j(\alpha) d\mu(\alpha)$ are the weights, and $l_j(\alpha)$ are the Lagrange interpolation polynomials with respect to $\{\alpha_j\}_{j=1}^N$ [9], yielding

$$(3.6) \quad \hat{f}(n) \approx \hat{f}_N(n) := \sum_{j=1}^N f(\alpha_j) p_n(\alpha_j) w_j, \quad n = 0, 1, \dots, N-1.$$

The gPC collocation approximation is defined by

$$(3.7) \quad f_N^{\text{gpc}}(\alpha) := \sum_{n=0}^{N-1} \hat{f}_N(n) p_n(\alpha),$$

where $\{\hat{f}_N(n)\}_{n=0}^{N-1}$ are given by (3.6); see [60].

The spectral accuracy of the gPC approximation in L^2 implies a similar accuracy for the approximation of moments.

Corollary 3.1. *Let f be analytic, and let f_N^{gpc} be its gPC collocation approximation of order N ; see (3.7). Then the moments (2.2) of f can be approximated by the respective moments of f_N^{gpc} with exponential accuracy as $N \rightarrow \infty$.*

Proof. See Appendix A. ■

For a smooth quantity of interest f , this spectral convergence rate is superior to the Monte Carlo's $1/\sqrt{N}$ convergence rate, which explains the popularity of the gPC collocation method.

In [41] we used the gPC approximation for moments and density estimation. Because of its spectral accuracy (Corollary 3.1), the number of sample points that is required for gPC to achieve a certain precision is considerably smaller than for Monte Carlo. To the best of our knowledge, however, there is no convergence result for density estimation using gPC which is analogous to Corollary 3.1.

Algorithm 3.1 can also approximate *nonsmooth* quantities of interest $f(\alpha)$, as long as $u(\cdot; \alpha)$ is smooth; see section 7 and [41]. The choice of the histogram method in step 4 is discussed in section 9.

The evaluation of $\{f(u_N^{\text{gpc}}(\cdot, \tilde{\alpha}_m))\}_{m=1}^M$ in step 3 is computationally cheap, as it amounts to a substitution in a polynomial. Therefore, there is essentially no computational cost for choosing M to be sufficiently high for the histogram method. This algorithm is also nonintrusive, in the sense that it only requires direct simulations of the deterministic system (2.1) with specific α_j values (as opposed to, e.g., Galerkin-type methods [12, 31, 61]). Our choice of the histogram method for density estimation will be explained in section 4.1.

4. Density estimation and spline-based UQ. Despite the prevalence of surrogate models in numerical methods and of density estimation in UQ applications [1, 6, 7, 21, 32, 41, 53], to the best of our knowledge, the adequacy of surrogate models for density estimation has not

Algorithm 3.1. gPC-based estimation [41].

Let $\{\alpha_j, w_j\}_{j=1}^N$ be the points and weights of the Gaussian quadrature rule of order N that correspond to the weight function $c(\alpha)$, and let $\{p_n(\alpha)\}_{n=0}^\infty$ be the respective orthogonal polynomials.

- 1: For $j = 1, \dots, N$, solve (2.1) with $\alpha = \alpha_j$ to obtain $u(t, \mathbf{x}; \alpha_j)$.
- 2: Approximate

$$u(t, x; \alpha) \approx u_N^{\text{gpc}}(t, \mathbf{x}; \alpha),$$

where

$$u_N^{\text{gpc}}(t, \mathbf{x}; \alpha) := \sum_{n=0}^{N-1} \hat{u}_N(t, \mathbf{x}; n) p_n(\alpha)$$

and

$$\hat{u}_N(u, \mathbf{x}; n) = \sum_{j=1}^N p_n(\alpha_j) u(t, \mathbf{x}; \alpha_j) w_j, \quad n = 0, \dots, N - 1.$$

- 3: Approximate $f(\tilde{\alpha}_m) \approx f(u_N^{\text{gpc}}(\cdot, \tilde{\alpha}_m))$ on a sample of $M \gg N$ points $\{\tilde{\alpha}_m\}_{m=1}^M$ which are i.i.d. according to $c(\alpha)$.
- 4: **if goal is moment estimation: then**
- 5: Use the trapezoidal integration rule with $\{f(\tilde{\alpha}_m)\}_{m=1}^M$ to approximate $\mathbb{E}_\alpha[f]$.⁴
- 6: **else if goal is density estimation: then**
- 7: Use the histogram method (3.1) with $\{f(\tilde{\alpha}_m)\}_{m=1}^M$ to estimate the PDF of f .
- 8: **end if**

been addressed in the UQ literature. To study this problem, we first write an explicit relation between a function $f : \Omega \rightarrow \mathbb{R}$ and the PDF that it induces on \mathbb{R} .

Lemma 4.1. *Let f be a real, piecewise monotone, continuously differentiable function on $[a, b]$, where $-\infty \leq a < b \leq \infty$, and let μ be an absolutely continuous probability measure on $[a, b]$, i.e., there is $c \in L^1([a, b])$ such that $d\mu(\alpha) = c(\alpha)d\alpha$. Then*

$$(4.1) \quad p_f(y) = \sum_{f(\alpha_j)=y} \frac{c(\alpha_j)}{|f'(\alpha_j)|},$$

where $p(y)$ is the PDF of f .

Proof. See Appendix C. ■

Because polynomial approximations (e.g., gPC) tend to be oscillatory, they “add” many artificial extremal points. Hence, by Lemma 4.1, the PDFs that they induce might deviate considerably from the exact one. To elucidate this point, in Lemma 4.2 we consider a smooth function f which is approximated by a highly oscillatory function g . In this example, having an upper bound on $\|f - g\|_r$ for some $r \geq 1$ does not yield an upper bound on $\|p_f - p_g\|_q$,

⁴Any standard integration technique could work here, provided sufficient smoothness. If $f(\alpha)$ is smooth, one can approximate $\mathbb{E}_\alpha = \hat{f}(0) \approx \hat{f}_N(0)$; see [59].

where p_f and p_g are the PDFs induced by $f(\alpha)$ and $g(\alpha)$, respectively, and $q \geq 1$, because of the numerous “artificial” extremal points of g .

Lemma 4.2. *Let $\Omega = [0, 1]$ equipped with the Lebesgue measure. Under the above notation, for every $\epsilon > 0$, there exist two functions f and g such that $\|f - g\|_\infty \leq \epsilon$, but $\|p_f - p_g\|_\infty \geq 1/2$.*

Proof. Let $f(\alpha) = \alpha$ and $g(\alpha) = \alpha + \delta \sin((2\delta)^{-1}\alpha)$. By direct differentiation $g'(\alpha) = 1 + 2^{-1} \cos((2\delta)^{-1}\alpha)$ and $f'(\alpha) \equiv 1$. Since f is monotone, and since g is monotone for sufficiently small δ , then by Lemma 4.1 with $c(\alpha) \equiv 1$, and so $p_f(y) = 1/f'(f^{-1}(y)) \equiv 1$ and $p_g(y) = 1/g'(g^{-1}(y))$. Specifically, there exists $y \in \mathbb{R}$ such that $p_g(y) = 1/2$, and so $\|p_f - p_g\|_\infty \geq 1/2$, irrespective of $\|f - g\|_\infty = \delta$, which can be made arbitrarily small. ■

Remark 4.3. A similar argument also shows that $\|f - g\|_r$ does not control $\|p_f - p_g\|_q$ for any $1 \leq q, r \leq \infty$.

To propose a surrogate model for which accurate density estimation is guaranteed, we first note that $f_N^{\text{GPC}}(\alpha)$ is the interpolating polynomial of f of order $N - 1$ at the Gauss quadrature points $\{\alpha_j\}_{j=1}^N$ [8, 27]. This suggests that other interpolants of $f(\alpha)$ can be used in Algorithm 3.1. In what follows, we argue that for our computational tasks, splines provide a better way to approximate $f(\alpha)$ and its associated PDF.

We recall that splines are piecewise polynomials of degree m , with $k < m$ smooth derivatives. Given an interval $\Omega = [\alpha_{\min}, \alpha_{\max}]$ and a grid $\alpha_{\min} = \alpha_1 < \alpha_2 < \dots < \alpha_N = \alpha_{\max}$, the interpolating cubic spline $s_N(\alpha)$ is a C^2 , piecewise-cubic polynomial that interpolates $f(\alpha)$ at $\{\alpha_j\}_{j=1}^N$, endowed with two additional boundary conditions. Three standard choices are (i) the natural cubic spline, for which $\frac{d^2}{d\alpha^2} f_N^{\text{spline}}(\alpha_1) = \frac{d^2}{d\alpha^2} f_N^{\text{spline}}(\alpha_N) = 0$, (ii) the “not-a-knot” spline, for which $\frac{d^3}{d\alpha^3} f_N^{\text{spline}}$ is continuous at α_2 and α_{N-1} , and (iii) the clamped spline, for which $\frac{d}{d\alpha} f_N^{\text{spline}}(\alpha_j) = \frac{d}{d\alpha} f(\alpha_j)$ for $j = 1, N$. Our decision to use splines is motivated by the following reasons:

1. The error of spline interpolation is guaranteed to be “small” for any sample size, in the following sense.

Theorem 4.4 (see [3, 25]). *Let $f \in C^{m+1}([\alpha_{\min}, \alpha_{\max}])$, and let f_N^{spline} be its “not-a-knot,” clamped or natural m th-order spline interpolant. Then*

$$(4.2) \quad \|(f(\alpha) - f_N^{\text{spline}}(\alpha))^{(j)}\|_{L^\infty[\alpha_{\min}, \alpha_{\max}]} \leq C_{\text{spl}}^{(j,m)} \|f^{(m+1)}\|_\infty h_{\max}^{m+1-j}, \quad j = 0, 1, \dots, m-1,$$

where $C_{\text{spl}}^{(j,m)} > 0$ is a universal constant that depends only on the type of boundary condition, m , and j , and $h_{\max} = \max_{1 < j \leq N} |\alpha_j - \alpha_{j-1}|$.

2. Spline interpolation is predominantly local. For further details, see Appendix B.

Thus, although $f_N^{\text{spline}}(\alpha)$ depends on $\{f(\alpha_1), \dots, f(\alpha_N)\}$, it predominantly depends on the few values $f(\alpha_j)$ for which α_j is adjacent to α . Therefore, large derivatives and discontinuities of $f(\alpha)$ may impair the accuracy of $f_N^{\text{spline}}(\alpha)$ only locally.⁵ This is in contrast to gPC (and polynomial interpolation in general), where discontinuities and

⁵For a review of cubic splines that are strictly local, see [4].

Algorithm 4.1. Spline-based estimation.

Let $\Lambda = \{\alpha_1, \dots, \alpha_N\}$ be a uniform grid on $[\alpha_{\min}, \alpha_{\max}]$.

- 1: For each $\alpha_j \in \Lambda$, solve (2.1) with $\alpha = \alpha_j$ to obtain $u(t, \mathbf{x}; \alpha_j)$.
- 2: Approximate $u(t, x; \alpha) \approx u_N^{\text{spline}}(t, \mathbf{x}; \alpha)$, where u_N^{spline} is a cubic spline interpolant on Λ .
- 3: Approximate $f(\tilde{\alpha}_m) \approx f(u_N^{\text{spline}}(\cdot, \tilde{\alpha}_m))$ on a sample of $M \gg N$ points $\{\tilde{\alpha}_m\}_{m=1}^M$ which are i.i.d. according to $c(\alpha)$.
- 4: **if goal is moment estimation: then**
- 5: Use the trapezoidal integration rule with $\{f(\tilde{\alpha}_m)\}_{m=1}^M$ to approximate $\mathbb{E}_\alpha[f]$.
- 6: **else if goal is density estimation: then**
- 7: Use the histogram method (3.1) with $\{f(\tilde{\alpha}_m)\}_{m=1}^M$ to approximate the PDF of f .
- 8: **end if**

large derivatives of f decrease the approximation accuracy across the entire domain.

In addition, splines can be constructed using any choice of sampling points.

In light of these considerations, we propose to replace the gPC interpolant with a spline.

Remark 4.5. See Appendix D for a MATLAB implementation of this algorithm.

Which cubic spline should be used in line 2? If $f'(\alpha_{\min})$ and $f'(\alpha_{\max})$ are known, then one should use the *clamped* cubic spline (or the natural cubic spline if these derivatives are zero). When the boundary derivatives are unknown, however, the “not-a-knot” interpolating cubic spline should be used (as indeed was done in this manuscript). See [4] for further discussion.

Algorithm 4.1 is *identical* to Algorithm 3.1, except for two substantial points:

1. The sampling grid is uniform, rather than the Gauss quadrature grid.⁶
2. The gPC interpolant u_N^{gpc} is replaced by a cubic spline interpolant u_N^{spline} .

Remark 4.6. This method is not to be confused with *spline-smoothing*, in which one approximates the PDF p with splines [15, 55]. Thus, Algorithm 4.1 approximates u with a spline, but the resulting approximation of the PDF p is *not a spline*.

4.1. Accuracy of Algorithm 4.1 for density estimation. The density estimation error of Algorithm 4.1 has two components—the error of the spline approximation (line 3) and that of the histogram method (line 7).⁷

The accuracy of the histogram method in line 7 depends on the number of bins L and on the number of samples M at lines 3 and 7. If the number of bins is chosen to be

$$(4.3) \quad L_{\text{opt}} = K_f M^{-\frac{1}{3}}, \quad K_f = \left(\frac{\|f'\|_2^2 [\max f - \min f]}{6} \right)^{\frac{1}{3}},$$

⁶Algorithm 4.1 can be performed with *any* choice of grid points. For clarity, we present it only with a uniform grid.

⁷In terms of density estimators, this can be explained by the following argument. Denote by p , p_N , and $\hat{p}_{N,M}$ the density of f , f_N and the density estimator of Algorithm 3.1 or 4.1, respectively. Then the approximation error (in any norm) satisfies $\|p - \hat{p}_{N,M}\| \leq \|p - p_N\| + \|p_N - \hat{p}_{N,M}\|$. The second term vanishes as $M \rightarrow \infty$ and L is given by (4.3), in which case the density estimation error is roughly the bias incurred from approximating f by f_N .

the mean squared L^2 error (MISE) of the histogram method decays as $M^{-\frac{2}{3}}$ [58].⁸ Because the computational cost of increasing L and M is negligible, they can be set sufficiently large so that the accuracy of Algorithm 4.1 mainly depends on the difference between the PDFs of f and f_N^{spline} , denoted by p_f and p_{f_N} , respectively. We motivate the choice of the histogram method to estimate the density by four factors:

1. Implementing the histogram method is straightforward and can be done with a few lines of code (see Appendix D).
2. The accuracy of the histogram method can be improved and controlled by varying the number of samples M , with a negligible computational cost.
3. The histogram method can be used even when the quantity of interest f is not smooth.
4. The histogram method can be used for a multidimensional random parameter α .

In principle, we could have used the explicit relation (4.1) to compute the PDF. Because this approach does not have the above advantages, however, the histogram method was chosen.

4.2. Accuracy of spline-based density estimation. In section 4.1 we showed that the accuracy of density estimation of Algorithms 3.1 and 4.1 is determined by the error of approximating the density with that of the surrogate model, and not by the error of the histogram method. By Lemma 4.1, if $f'(\alpha)$ is bounded away from zero, then p is smooth. As noted, however, the gPC polynomial interpolant $f_N^{\text{gpc}}(\alpha)$ tends to be oscillatory, and so it might add artificial external points where $\frac{d}{d\alpha} f_N^{\text{gpc}}(\alpha) = 0$; see, e.g., Figure 2(c). At every such point where $\frac{d}{d\alpha} f_N^{\text{gpc}}(\alpha) = 0$, the PDF approximation becomes unbounded, and so a large error in the PDF estimation occurs. This is seldom the case with the spline interpolant, which due to its local nature (see Lemma B.2) does not produce numerical oscillations throughout its domain Ω . Indeed, the natural cubic spline $f_N^{\text{spline}}(\alpha)$ has the “minimum curvature” property, which implies that it oscillates “very little” about the original function [38]. This notion is made precise by the following result.

Theorem 4.7. *Let $f \in C^{m+1}([\alpha_{\min}, \alpha_{\max}])$ with $|f'(\alpha)| \geq a > 0$, let α be distributed by $c(\alpha)d\alpha$, where $c \in C^1([\alpha_{\min}, \alpha_{\max}])$, and let p_f and p_{f_N} be the PDFs of $f(\alpha)$ and of $f_N = f_N^{\text{spline}}$, its natural, “not-a-knot,” or clamped m th order spline interpolant on a uniform grid of size N . Then, for any $1 \leq q < \infty$*

$$(4.4) \quad \|p_f - p_{f_N}\|_q \leq KN^{-m}, \quad N > \sqrt[m]{\frac{2C_{\text{spl}}^{(1,m)} \|f^{(m+1)}\|_{\infty}}{a}} (\alpha_{\max} - \alpha_{\min}),$$

where $C_{\text{spl}}^{(1,m)}$ is given by Theorem 4.4 and K depends only on $f(\alpha)$, $c(\alpha)$, q , and $|\alpha_{\max} - \alpha_{\min}|$.

Proof. See Appendix E. ■

The proof of Theorem 4.7 only makes use of two properties of spline interpolation: the accurate approximation of the function and its derivative in L^∞ , and the uniform bound on the second derivatives (Theorem 4.4). Therefore, Theorem 4.7 immediately generalizes to a broad family of surrogate models, denoted by $\{g_N\}$.

⁸In practice, f and f' are often unknown, and so K_f needs to be estimated.

Corollary 4.8. Let $f(\alpha)$ and $c(\alpha)$ be as in Theorem 4.7, and let $g_N \in C^1([\alpha_{\min}, \alpha_{\max}])$ be a sequence of approximations of f for which

$$\|f - g_N\|_\infty, \|f' - g'_N\|_\infty \leq KN^{-\tau}, \quad \|g''_N\|_\infty < C_g < \infty,$$

where $\tau > 0$, C_g , and K are independent of N . Then

$$\|p_f - p_{g_N}\|_q \leq \tilde{K}N^{-\tau}$$

for any $1 \leq q < \infty$, where p_f and p_{g_N} are the PDFs of $f(\alpha)$ and $g_N(\alpha)$, respectively, and \tilde{K} is independent of N .

Remark 4.9. If f is only piecewise C^{m+1} , then N^{-m} convergence is guaranteed when the grid points include the discontinuity points of $f(\alpha)$, since the proof can be repeated in each interval on which the function is C^{m+1} in the same way.

Remark 4.10. Although Theorem 4.7 applies only to functions whose derivatives are bounded away from 0, in practice we observe cubic convergence for nonmonotone functions as well (see section 7). Whether Theorem 4.7 generalizes to nonmonotone cases is unclear.

In our numerical simulations (see Figures 2, 4, 8, and 9), we observe that the cubic convergence is often reached well before N satisfies (4.4). We also observe that the density approximation error $\|p_f - p_{f_N}\|_1$ decays at a faster than cubic rate. A possible explanation for this observation is provided by the following.

Lemma 4.11. Assume the conditions of Theorem 4.7 for $m = 3$, and let J_N be the number of times that $\frac{d}{d\alpha}(f(\alpha) - f_N^{\text{spline}}(\alpha))$ changes its sign on $[\alpha_{\min}, \alpha_{\max}]$. If $J_N = O(N^r)$ for $0 \leq r \leq 1$, then $\|p_f - p_{f_N}\|_1 \leq KN^{-4+r}$. Specifically, if J_N is uniformly bounded for all $N \in \mathbb{N}$, then $\|p_f - p_{f_N}\|_1 \leq KN^{-4}$.

Proof. See Appendix F. ■

4.3. Accuracy of moment estimation. While the main focus of this paper is on density estimation using a surrogate model, we also point out two disadvantages of the gPC method for moment estimation:

1. The spectral convergence of the gPC method is attained only *asymptotically* as the number of sample points N becomes sufficiently large. For small or moderate values of N , however, its accuracy may be quite poor, due to insufficient resolution, and the global nature of spectral approximation.
2. The sample points $\{\alpha_j\}_{j=1}^N$ of the gPC method are predetermined by the quadrature rule. Therefore, if one wants to *adaptively* improve the accuracy, one cannot use the samples from the “old” low-resolution grid in the “new” high-accuracy approximation.

Similarly to density estimation, the error of the moment estimation of Algorithm 4.1 comes from both the numerical integration (line 5) and interpolation (line 2). The trapezoidal rule integration error can be made sufficiently small by increasing the number of samples M at line 3, at a negligible computational cost. Moreover, if $c(\alpha) \equiv 1$, the integration over f_N^{spline} can be done exactly.⁹ Hence, the moment estimation error of Algorithm 4.1 is determined by the accuracy of the spline interpolation.

⁹When f is sufficiently smooth and α is uniformly distributed, one can approximate $\mathbb{E}_\alpha[f] \approx \mathbb{E}_\alpha[f_N^{\text{spline}}]$ and compute the right-hand side explicitly (in MATLAB, this can be done using the `fnint` command).

Corollary 4.12. Let $f \in C^4([\alpha_{\min}, \alpha_{\max}])$, let f_N^{spline} be the natural, “not-a-knot,” or clamped cubic spline interpolant of f , and let α be distributed by $c(\alpha)d\alpha$, where $c(\alpha) \geq 0$, and $\int_{\alpha_{\min}}^{\alpha_{\max}} c(\alpha) d\alpha = 1$. Then

$$\left| \mathbb{E}_\alpha[f] - \mathbb{E}_\alpha[f_N^{\text{spline}}] \right| \leq C_{\text{spl}}^{(0)} \|f\|_\infty h_{\max}^4,$$

where $C_{\text{spl}}^{(0)}$ and h_{\max} are defined in Theorem 4.4.

Proof. By Theorem 4.4, $\|f - f_N^{\text{spline}}\|_\infty \leq C_{\text{spl}}^{(0)} \|f^{(4)}\|_\infty h_{\max}^4$. Hence,

$$\begin{aligned} & \left| \int_{\alpha_{\min}}^{\alpha_{\max}} (f(\alpha) - f_N^{\text{spline}}(\alpha)) c(\alpha) d\alpha \right| \\ & \leq \|f - f_N^{\text{spline}}\|_\infty \int_{\alpha_{\min}}^{\alpha_{\max}} c(\alpha) d\alpha \leq \|f - f_N^{\text{spline}}\|_\infty \cdot 1 \leq C_{\text{spl}}^{(0)} \|f^{(4)}\|_\infty h_{\max}^4. \quad \blacksquare \end{aligned}$$

Typically, $C_{\text{spl}}^{(0)} < 1$. For example, for the natural and “not-a-knot” cubic spline, $C_{\text{spl}}^{(0)}$ is equal to $\frac{5}{384}$ and $\frac{1}{25}$, respectively [4, 25]. On a uniform grid, $h_j = \frac{\alpha_{\max} - \alpha_{\min}}{N-1}$ for $1 < j \leq N$, and so $\mathbb{E}_\alpha[f] - \mathbb{E}_\alpha[f_N^{\text{spline}}] = O(N^{-4})$.

As $N \rightarrow \infty$, the polynomial convergence rate of the spline approximation (Corollary 4.12) is outperformed by gPC’s spectral convergence rate (Corollary 3.1). Quite often, however, the spline approximation is more accurate for moderate N values. To see that, note that by (3.3), (3.6), and (3.7), $\mathbb{E}_\alpha[f_N^{\text{gpc}}] = \sum_{j=1}^N f(\alpha_j) w_j$, which is the Gauss quadrature rule. Hence, if $f \in C^{2N}$, then

$$\mathbb{E}_\alpha[f] - \mathbb{E}_\alpha[f_N^{\text{gpc}}] = \frac{f^{(2N)}(\xi)}{k_N^2 (2N)!}, \quad \xi \in (\alpha_{\min}, \alpha_{\max}),$$

where k_N is the leading coefficient of $p_N(\alpha)$ [9]. If for small N , $\|f^{(2N)}\|_\infty$ increases faster than $k_N^2 (2N)!$, the error initially increases with N . In these cases, the exponential convergence is only achieved at large N .¹⁰ Even when gPC does converge exponentially, i.e., $|\mathbb{E}_\alpha[f] - \mathbb{E}_\alpha[f_N^{\text{gpc}}]| \leq K e^{-\gamma N}$, but γ is small, then the error of the spline approximation may be smaller for moderate values of N ; see, e.g., Figure 1(c). To conclude, the accuracy of spline-based moment approximation is guaranteed also with few samples, and not only asymptotically as $N \rightarrow \infty$.

5. Multidimensional noises. To generalize the spline-based density-estimation approach (Algorithm 4.1) to the case where $\boldsymbol{\alpha} \in \Omega = [0, 1]^d$, we use tensor-product splines, which are defined in the following way. Let $m \geq 1$, let $f(\boldsymbol{\alpha}) \in C^{m+1}(\Omega)$, let Λ be the one-dimensional grid $0 = \alpha_1 < \dots < \alpha_n = 1$, and let Λ^d be the respective d -dimensional tensor-product grid. An m th degree tensor-product spline interpolant of f is a function $s(\boldsymbol{\alpha}) \in C^{m-1}(\Omega)$ that interpolates f on Λ^d and reduces to a one-dimensional m th degree spline on every line on Λ^d ; see [44] for a more precise definition.^{11,12} The multidimensional extension of Algorithm 4.1 for density estimation is as follows.

¹⁰For example, if the numerator grows as K^{2N} , the error only decays for $N > K$.

¹¹That is, when $d-1$ coordinates of $\boldsymbol{\alpha}$ are fixed in Λ .

¹² $s(\boldsymbol{\alpha})$ is unique when endowed with sufficiently many boundary conditions; see the discussion on the one-dimensional case in section 4. Theorem 5.1 holds for many possible choices of boundary conditions, including the not-a-knot conditions which we have also used in our simulations.

Algorithm 5.1. Multidimensional spline-based density estimation.

Let $\Lambda^d = \{\alpha_1, \dots, \alpha_N\}^D$ be a tensor-product uniform grid on $[0, 1]^d$.

- 1: For each $\alpha_j \in \Lambda^d$, solve (2.1) with $\alpha = \alpha_j$ to obtain $u(t, \mathbf{x}; \alpha_j)$.
- 2: Approximate $u(t, x; \alpha) \approx u_N^{\text{spline}}(t, \mathbf{x}; \alpha)$, where u_N^{spline} is a tensor-product spline interpolant of order m on Λ^d .
- 3: Approximate $f(\tilde{\alpha}_m) \approx f(u_N^{\text{spline}}(\cdot, \tilde{\alpha}_m))$ on a sample of $M \gg N$ points $\{\tilde{\alpha}_m\}_{m=1}^M$ which are i.i.d. according to $c(\alpha)$.
- 4: Use the histogram method (3.1) with $\{f(\tilde{\alpha}_m)\}_{m=1}^M$ to approximate the PDF of f .

As in the one-dimensional Algorithm 4.1, the analysis of the density-estimation error in Algorithm 5.1 is based on two components:

1. A pointwise error bound for tensor-product spline interpolants, due to Schultz, as in the following.

Theorem 5.1 (see [40, 44]). *Let $\Omega = [0, 1]^d$, $f \in C^{m+1}(\Omega)$, and let $s(\alpha)$ be its m th degree tensor-product spline interpolant. Then for any $\alpha \in \Omega$,*

$$(5.1) \quad |D^j(f - s)| < C_m h^{m+1-j}, \quad j = 0, 1, \dots, m - 1,$$

where D^j is any j th order derivative,¹³ $C_m = C_m(\|D^{m+1}f\|_\infty)$ depends only on the L^∞ norms of the $m + 1$ order derivatives of f , and $h = \max_{1 \leq j < n} |\alpha_{j+1} - \alpha_j|$.

2. A multidimensional generalization of Lemma 4.1,¹⁴ as follows.

Lemma 5.2. *Let $\Omega \subset \mathbb{R}^d$ be a Jordan set, denote by $|\cdot|$ the Euclidean norm in \mathbb{R}^d , let f be piecewise-differentiable with $|\nabla f| \neq 0$ on $\bar{\Omega}$, let α be an absolutely continuous random variable in Ω , i.e., $d\mu(\alpha) = c(\alpha)d\alpha$ for some nonnegative $c \in L^1(\Omega)$, and denote the PDF associated with $f(\alpha)$ by p_f . Then*

$$(5.2) \quad p_f(y) = \frac{1}{\mu(\Omega)} \int_{f^{-1}(y)} \frac{c(\alpha)}{|\nabla f(\alpha)|} d\sigma,$$

where $d\sigma$ is a $(d - 1)$ dimensional surface element on $f^{-1}(y)$.

Proof. See Appendix G. ■

The generalization of Theorem 4.7 to the case of multidimensional random parameter is as follows.

Theorem 5.3. *Let $\Omega = [0, 1]^d$, let $m \geq 1$, let $f \in C^{m+1}(\Omega)$, let s be the m -degree tensor-product spline interpolant of f , let α be uniformly distributed in Ω , and let p_f and p_s be the PDFs of f and s , respectively. If $\kappa_f := \min_\Omega |\nabla f| > 0$, then for sufficiently small h and for any $1 \leq q < \infty$,*

¹³More explicitly, $D^j = \prod_{k=1}^d (\partial_{\alpha_k})^{\ell_k}$, where $\ell_1 + \dots + \ell_d = j$, and each ℓ_k is a nonnegative integer.

¹⁴When $\Omega \subset \mathbb{R}$ is a one-dimensional interval, Lemma 5.2 reduces to Lemma 4.1. Indeed, since $|f'| \neq 0$ on $\bar{\Omega}$, then f is piecewise monotonic, and so $f^{-1}(y)$ consists of a finite number of points. In addition, the surface element $d\sigma$ is a point-mass distribution. Hence, (5.2) reduces to (4.1).

$$(5.3) \quad \|p_f - p_s\|_q \leq Kh^m$$

for some constant $K > 0$, where h is defined in Theorem 5.1.

Proof. See Appendix H. ■

Theorem 5.3 can be extended to any approximation \tilde{f} of f and to any bounded domain $\Omega \subseteq \mathbb{R}^d$, provided that the bound (5.1) holds for $j = 0$ and $j = 1$.

The total number of sample points in the special case where Λ is the uniform one-dimensional grid on $[0, 1]$ is $N = n^d \sim h^{-d}$. Therefore, we have the following.

Corollary 5.4. *Let Λ be the uniform grid on $[0, 1]$. Then under the conditions of Theorem 5.3, for sufficiently large N ,*

$$\|p_f - p_s\|_1 \leq KN^{-\frac{m}{d}}$$

for some constant $K > 0$.

As noted in section 3.1, the L^1 error of the KDE method scales as $N^{-\frac{2}{5}}$ [13]. Therefore, by Corollary 5.4, Algorithm 5.1 outperforms KDEs for dimensions $d \leq \frac{5}{2}m$. Finally, as in the one-dimensional case (Corollary 4.8), the proof of Theorem 5.3 only makes use of two properties of spline interpolation: the L^∞ approximation of the function and of its gradient, and the uniform bound on the second derivatives (Theorem 5.1). Theorem 5.3 therefore generalizes immediately to density estimation using nonspline surrogate models.

Corollary 5.5. *Under the conditions and notation of Theorem 5.3, consider $g_h \in C^1[0, 1]^d$ with uniformly bounded second derivatives such that*

$$\|f - g_h\|_\infty, \|\nabla f - \nabla g_h\|_\infty \leq Kh^{-\tau}$$

for some $\tau > 0$ independent of f and $K = K(f) > 0$. Then $\|p_f - p_{g_h}\|_q \leq \tilde{K}h^{-\tau}$ for any $1 \leq q < \infty$.

6. Simulations. In this section, we compute the density and the moments of the function

$$(6.1) \quad f(\alpha) = \tanh(9\alpha) + \frac{\alpha}{2}, \quad \alpha \in [-1, 1],$$

which is smooth but has a narrow high-derivative region.¹⁵

6.1. Interpolation. With $N = 12$ samples, the spline interpolant f_N^{spline} of (6.1) is nearly indistinguishable from f , whereas the gPC interpolant f_n^{gpc} slightly oscillates “around” f ; see Figure 1(a). Although f_N^{gpc} converges exponentially to f in L^2 (see Figure 1(b)), its L^2 approximation error $\|f - f_N\|_2 = \left(\int_{-1}^1 |f(\alpha) - f_N(\alpha)|^2 d\alpha\right)^{\frac{1}{2}}$ with few samples ($10 \leq N \leq 40$) is larger than that of the spline interpolant by more than an order of magnitude. With sufficiently many samples ($N > 70$), however, the gPC approximation exponential convergence outperforms the spline’s polynomial convergence rate. This example shows that with few samples, the occurrence of a “jump” in f hurts the accuracy of the gPC interpolant. Spline interpolation, on the other hand, is less sensitive to the “jump,” because it “confines” the approximation error induced by the jump to the jump interval (roughly $\alpha \in (-0.1, 0.1)$); see Lemma B.2.

¹⁵The $\frac{\alpha}{2}$ term was added so that $\frac{df}{d\alpha}$ is bounded away from zero, in order to prevent singularities in the PDF; see section 6.3.

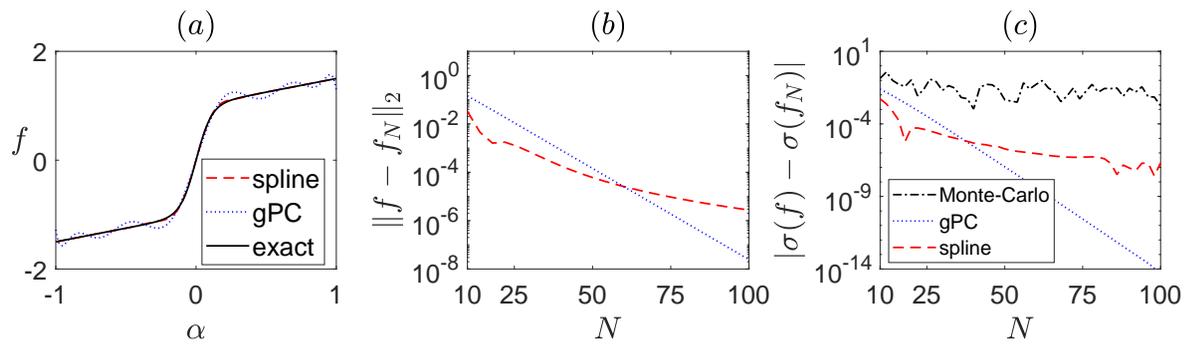


Figure 1. (a) $f(\alpha)$ (solid) (see (6.1)) and its spline interpolant (dashes) are nearly indistinguishable, whereas the gPC interpolant (dots) oscillates “around” f . Both interpolants use $N = 12$ grid points. (b) L^2 error of both interpolants as a function of the number of samples. (c) Error of the standard deviation when it is approximated using the Monte Carlo method (dash-dot), the gPC-based method (dots), and the spline-based method (dashes).

6.2. Moment approximation. The interpolation accuracy is relevant to moment approximation, because a small L^2 error implies a small moment-approximation error (Lemma A.1). For example, Figure 1(c) shows the standard deviation error $|\sigma(f) - \sigma(f_N)|$ (see (6.1)) when α is uniformly distributed in $[-1, 1]$. As expected, the spline-based method (Algorithm 4.1) is more accurate than the gPC-based method (Algorithm 3.1) with few samples, but the gPC is more accurate with sufficiently many samples. A purely statistical approach such as Monte Carlo converges poorly compared to both the spline and the gPC approach, with about 10% error with $N \leq 100$ sample points.

6.3. Density estimation. Consider the PDF induced by $f(\alpha)$ (see (6.1)) when α is uniformly distributed in $[-1, 1]$. The PDF computed by the gPC-based Algorithm 3.1 with $N = 18$ sample points deviates considerably from the exact PDF (see Figure 2(a)), whereas the PDF computed by the spline-based Algorithm 4.1 with $N = 18$ sample points is nearly indistinguishable from the exact PDF (see Figure 2(b)).¹⁶ This is consistent with our discussion in section 4. Indeed, the derivative of the spline interpolant $\frac{d}{d\alpha} f_N^{\text{spline}}$ approximates $f'(\alpha)$ with cubic accuracy, whereas the derivative of the gPC interpolant $\frac{d}{d\alpha} f_N^{\text{gPC}}$ has many artificial extremal points where $\frac{d}{d\alpha} f_N^{\text{gPC}}(\alpha) = 0$, but $\frac{d}{d\alpha} f(\alpha) \neq 0$ (see Figure 2(c)).

The L^1 distance $\|p_f - p_{f_N}\|_1$ between the exact PDF p_f and its approximation p_{f_N} is presented in Figure 2(d). For $10 \leq N \leq 100$ the spline-based approximation is more accurate than the gPC-based one by nearly two orders of magnitude. This is in contrast to moment estimation (see Figure 1(c)), in which the gPC approximation becomes more accurate for $N \geq 40$. Furthermore, we observe numerically that the spline-based method converges even faster than the N^{-3} rate predicted by Theorem 4.7. The KDE approximation has roughly 10% error for $N \leq 100$.¹⁷ Other frequently used distances between

¹⁶The MATLAB code that generates this PDF approximation is given in Appendix D.

¹⁷The poor accuracy of the KDE method is due to the fact that the KDE does not use the “functional information” $\{f_j = f(\alpha_j)\}_{j=1}^N$ but only the set $\{f_j\}_{i=1}^N$.

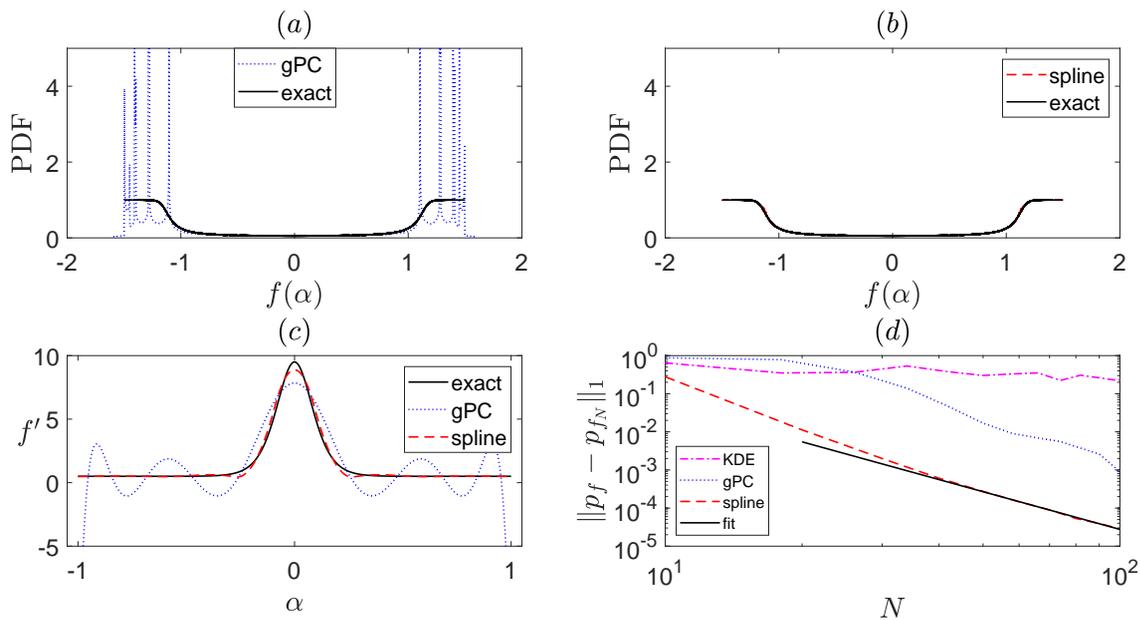


Figure 2. The PDF of $f(\alpha)$ (see (6.1)), where α is uniformly distributed in $[-1, 1]$. (a) exact PDF (solid) and its approximation by the gPC-based Algorithm 3.1 (dots) with $N = 18$ sample points. (b) Same, with the spline-based Algorithm 4.1 (dashes). The two lines are nearly indistinguishable. (c) Derivatives of f (solid), f_N^{spline} (dashes), and f_N^{gPC} (dots). (d) L^1 error of the PDF approximations as a function of the number of sample points, for the KDE (dash-dot), gPC-based approximation (dots), the spline-based approximation (dashes), and its power-law fit $103.2N^{-3.29}$ (solid).

distributions, such as the Hellinger distance $\frac{1}{\sqrt{2}} \|\sqrt{p_f} - \sqrt{p_{f_N}}\|_2$ [29] and the Kullback–Leibler (KL) divergence¹⁸ [28]

$$(6.2) \quad \int_{-\infty}^{\infty} p_f(y) \log \left(\frac{p_f(y)}{p_{f_N}(y)} \right) dy ,$$

produce similar results (data not shown).

6.4. Density estimation of nonsmooth functions.

Let

$$(6.3) \quad g(\alpha) = f(\alpha) \bmod (0.7) ,$$

where f is given by (6.1).¹⁹ Because (6.3) is nonsmooth, with few samples neither the spline nor the gPC interpolant is even remotely close to $g(\alpha)$; see Figure 3. Therefore, to approximate the PDF associated with $g(\alpha)$, we first use Algorithms 3.1 and 4.1 to approximate $f(\alpha) \approx f_N(\alpha)$. Since f is smooth, both approximations are reasonable with few samples;

¹⁸Intuitively, the d_{KL} measures the entropy added, or conversely, the information lost, in approximating p by p_{f_N} .

¹⁹This example is motivated by our study of the nonlinear Schrödinger equation [41], where the cumulative phase $\varphi(t; \alpha) = \arg[\psi(t, 0; \alpha)]$ is smooth, but the quantity of interest, the angle $\varphi \bmod (2\pi)$, is discontinuous. See section 7 for another optics application which motivates this example.

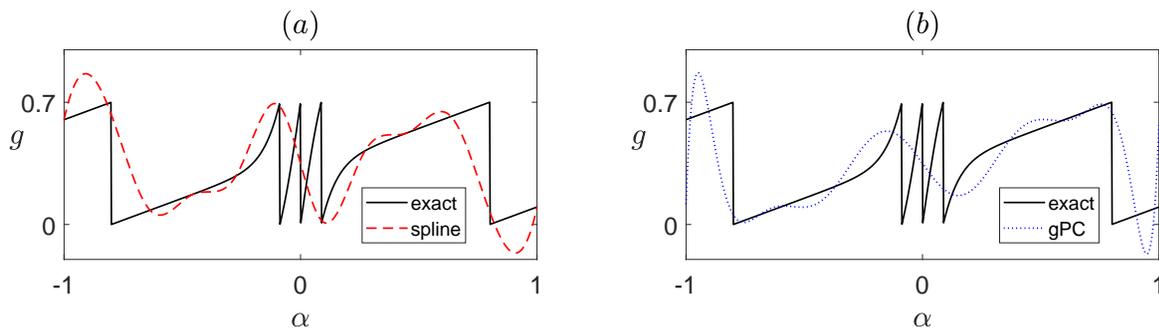


Figure 3. The discontinuous function $g(\alpha)$ (solid line; see (6.3)) and its spline interpolation with $N = 12$ sample points (dashes). (b) Same with the gPC interpolant (dots).

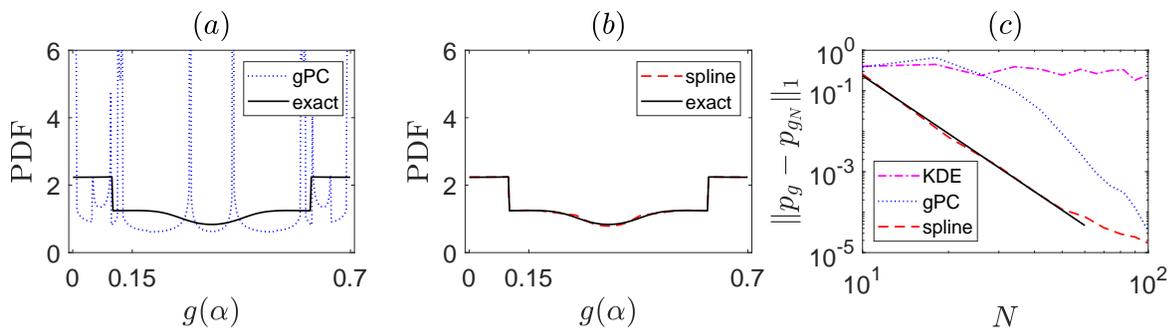


Figure 4. Same as Figure 2 for the discontinuous function $g(\alpha)$; see (6.3). The solid line in subplot (c) is the power-law fit $1.33 \cdot 10^4 N^{-4.75}$ of the spline-based approximate PDF.

see Figure 1. Next, we approximate $g(\alpha_m) \approx f_N(\alpha_m) \bmod (0.7)$ and compute the PDF of g using the histogram method on a high-resolution sampling grid ($M = 2 \cdot 10^6$). We again stress that evaluating f_N is computationally cheap and therefore can be easily done with such a large sample. As in the smooth case (see Figure 2), the PDF approximated by the gPC-based Algorithm 3.1 with $N = 18$ sample points has large deviations and converges poorly (see Figure 4(a)), whereas the PDF approximated by the spline-based Algorithm 4.1 with $N = 18$ sample points is nearly identical to the exact PDF (see Figure 4(b)). Indeed the L^1 error of the spline-based PDF is smaller than that of the gPC-based PDF by at least an order of magnitude, for $20 < N < 50$; see Figure 4(c). Although Theorem 4.7 applies only to C^4 functions, we observe numerically that the convergence rate of the spline-based PDF is faster than N^{-3} . The KDE approximation for the PDF of $g(\alpha)$ is less accurate than that of the spline-based and gPC-based approximations.

6.5. Multidimensional noise. To numerically confirm the error bound of the density estimation (Algorithm 5.1) for $d > 1$, we first consider the two-dimensional function

$$(6.4) \quad f_{2d}(\alpha_1, \alpha_2) = \tanh(6\alpha_1\alpha_2 + \alpha_1/2) + (\alpha_1 + \alpha_2)/3.$$

where α_1 and α_2 are independent and uniformly distributed in $[-1, 1]$. As in the one-dimensional example (see (6.1)), f_{2d} is analytic with high-gradients regions; see Figure 5(a).

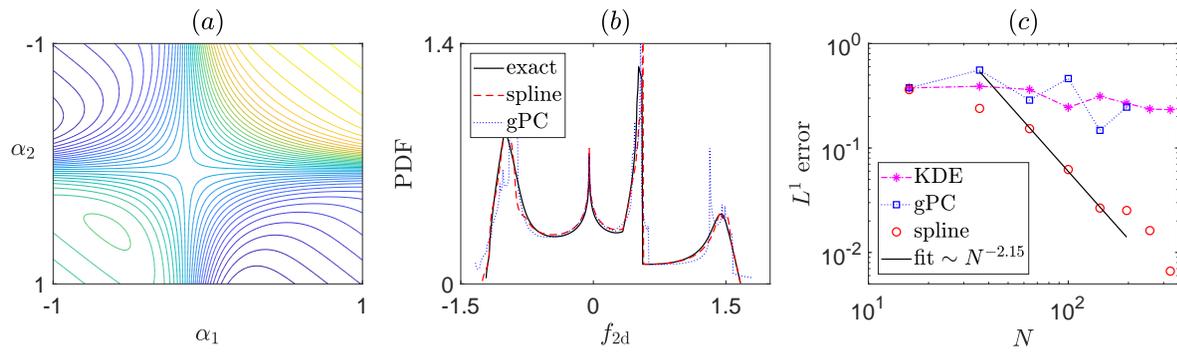


Figure 5. (a) Contours of the function $f_{2d}(\boldsymbol{\alpha})$; see (6.4). (b) The PDF of $f_{2d}(\boldsymbol{\alpha})$ (solid), its approximation by the spline-based Algorithm 4.1 (dashes), and by the gPC-based Algorithm 3.1 (dots). Here $\boldsymbol{\alpha}$ is uniformly distributed in $[-1, 1]^2$, and both approximations use $N = 64$ sample points. (c) L^1 error of the PDF approximations as a function of the number of sample points for the KDE (dash-dots), gPC-based approximation (dots-squares), and spline-based approximation (circles). The solid line is the power-law fit $1208N^{-2.15}$ (solid).

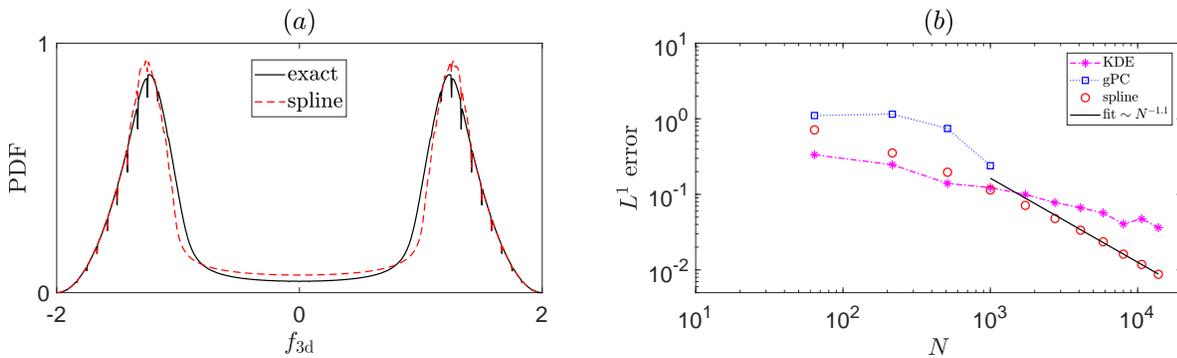


Figure 6. (a) The PDF of $f_{3d}(\boldsymbol{\alpha})$ (see (6.5)), where $\boldsymbol{\alpha}$ is uniformly distributed in $[-1, 1]^3$ (solid) and its approximation by the spline-based Algorithm 4.1 (dashes) with $N = 8^3$ sample points. (b) L^1 error of the PDF approximations as a function of the number of sample points for the KDE (dash-dots), the gPC-based PDF (rectangles), the spline-based PDF (circles), and its power-law fit $354N^{-1.11}$ (solid).

The spline-based PDF approximation with $N = 8^2$ sample points is very close to the exact PDF of $f(\alpha_1, \alpha_2)$, whereas the gPC-based PDF deviates from it substantially (Figure 5(b)). The convergence rate of Algorithm 5.1 with cubic splines is $N^{-2.15}$ (Figure 5(c)), which is consistent with the theoretical $N^{-\frac{3}{2}}$ error bound (Corollary 5.4). The convergence rates of both the KDE and the gPC methods are considerably slower for “small” sample sizes ($N \leq 200$).

Next, consider the three-dimensional function

$$(6.5) \quad f_{3d}(\alpha_1, \alpha_2, \alpha_3) = \tanh(8\alpha_1 + 5\alpha_2 + 10\alpha_3) + (\alpha_1 + \alpha_2 + \alpha_3)/3,$$

where α_1 , α_2 , and α_3 are independent and uniformly distributed in $[-1, 1]$. The spline-based PDF with $N = 10^3$ sample points approximates the exact PDF well (see Figure 6(a)), and its convergence rate is $N^{-1.1}$ (see Figure 6(b)), which is consistent with the theoretical N^{-1} convergence rate (Corollary 5.4). For comparison, the fitted convergence rate of the KDE is $N^{-0.39}$, which is consistent with the theoretical $N^{-\frac{2}{5}}$ rate [13]. Therefore, the spline-based

method is more accurate than the KDE for sufficiently many samples ($N > 10^3$). For smaller values of N (e.g., $N = 216$), however, the KDE achieves a slightly better accuracy than the spline-based method. This can be explained by what is known as the “curse of dimensionality.” Thus, in the three-dimensional tensor-grid spline, $N = 216$ sample points correspond to a mere *six* sample points in each dimension, which leads to insufficient resolution. The KDE method, on the other hand, does not approximate the underlying function f_{3d} and is therefore “indifferent” to the noise dimension. See section 9 for further discussion.

7. Application 1: Nonlinear Schrödinger equation. The one-dimensional coupled nonlinear Schrödinger equation (CNLS)

$$(7.1) \quad i \frac{\partial A_{\pm}(t, x)}{\partial t} + \frac{\partial^2 A_{\pm}}{\partial x^2} + \frac{2}{3} \frac{|A_{\pm}|^2 + 2|A_{\mp}|^2}{1 + \epsilon (|A_{\pm}|^2 + |A_{\mp}|^2)} A_{\pm} = 0,$$

where $0 < \epsilon \ll 1$, $t \geq 0$, and $x \in \mathbb{R}$, describes the propagation of elliptically polarized, ultrashort pulses in optical fibers [2], of elliptically polarized continuous-wave (CW) beams in a bulk medium [36, 45], Stokes and anti-Stokes radiation in Raman amplifiers [39], and rogue water-waves formation at the interaction of crossing seas [1]. We consider (7.1) with an elliptically polarized Gaussian input pulse with a random amplitude [36, 45]

$$(7.2) \quad \begin{pmatrix} A_+ \\ A_- \end{pmatrix} = (1 + 0.1\alpha) \begin{pmatrix} 8 \\ 4 \end{pmatrix} e^{-x^2},$$

where A_+ and A_- are the clockwise and counterclockwise circularly polarized components, respectively. The *on-axis ellipse rotation angle* is defined as

$$(7.3) \quad \theta(t; \alpha) := (\varphi_+(t; \alpha) - \varphi_-(t; \alpha)) \bmod (2\pi),$$

where $\varphi_{\pm}(t; \alpha) := \arg [A_{\pm}(t, 0; \alpha)]$ are the on-axis phases of the components. The distribution of $\theta(t; \alpha)$ indicates to what extent the ellipse rotation angle is “deterministic.”²⁰

Interpolation. For a given sample grid $\{\alpha_j\}_{j=1}^N$, we compute $\theta(t; \alpha_j)$ for each $1 \leq j \leq N$ by solving (7.1)–(7.2) and using (7.3). Figure 7(a) shows the spline and gPC interpolants of $\theta(t = 0.15; \alpha)$ with $N = 64$ points.²¹ While these interpolants seem nearly identical, the spline interpolant is more accurate than the gPC interpolant by more than an order of magnitude (cf. Figures 7(b) and 7(c)). Indeed, the L^2 error of the gPC interpolant (0.17%) is an order of magnitude larger than that of the spline interpolant (0.017%).

Density estimation. The gPC-based approximation with $N = 64$ differs substantially from the exact PDF; see Figure 8(a). In contrast, the spline-based approximated PDF with $N = 64$ sample points is indistinguishable from the exact PDF; see Figure 8(b). Indeed, the KL divergence of the gPC-based approximation (see (6.2)) is *about 16,000 times larger* than that of the spline-based approximation, and the L^1 error is 200 times larger (46% versus 0.2%).

²⁰We solve the CNLS using a fourth-order, compact finite-difference scheme for the spatial discretization and a predictor-corrector Crank–Nicolson scheme for the temporal integration of the semidiscrete problem [17].

²¹Because we have no explicit solution for $\theta(t; \alpha)$, the errors in this section are measured by comparison with $\theta_{513}^{\text{spline}}(0.15, \alpha)$ with $N = 513$ sample points. We verified that $\|\theta_{513}^{\text{spline}}(0.15, \alpha) - \theta_{513}^{\text{gpc}}(0.15, \alpha)\|_2 \approx 5 \cdot 10^{-5}$, which is an order of magnitude smaller than the approximation errors noted in the text.

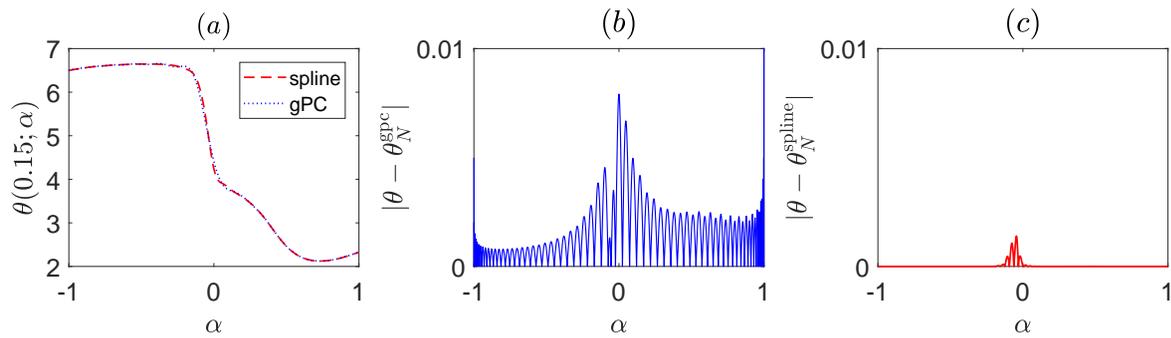


Figure 7. The polarization angle $\theta(t = 0.15; \alpha)$ for solutions of the CNLS (7.2) with $\epsilon = 10^{-5}$ and an elliptically polarized Gaussian initial condition (7.2). (a) Spline interpolation (dashes) and gPC interpolation (dots), with $N = 64$ sample points. The two lines are nearly indistinguishable. (b) Pointwise error of the gPC interpolant. (c) Same for the spline interpolant.

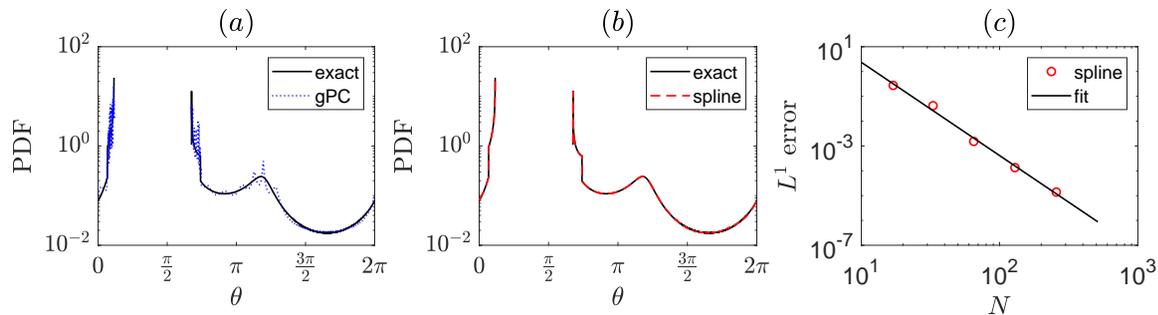


Figure 8. Same settings as in Figure 7. The PDF of $\theta(0.15, \alpha)$, where $\alpha \sim U(-1, 1)$. (a) Exact PDF (solid), and GPC-based approximation using $N = 64$ sample points (dots). (b) Same with the spline-based approximation (dashes). The two lines are indistinguishable. (c) L^1 error of the spline-based PDF as a function of N (circles) and the power-law fit $1.35 \cdot 10^4 N^{-3.76}$ (solid).

With $N = 32$, the spline-based is 32 times more accurate than the gPC-approximated PDF, in terms of KL divergence, and 11 times more accurate in terms of the L^1 error (41% versus 4.5%). The L^1 error of the spline-based PDF decays as $N^{-3.76}$; see Figure 8(c). This results exceed expectations with respect to Theorem 4.7, since $\theta'(0.15; \alpha)$ is not bounded away from 0 (see Figure 7(a)), and so Theorem 4.7 should not, in principle, apply to this case. Since the PDF of $\theta(0.15; \alpha)$ has discontinuities and high derivatives, spline smoothing techniques and KDE methods with smooth kernels were not considered in this case.

Moment approximation. The mean and standard deviation of circular quantities can be defined as [33]²²

$$(7.4) \quad \mathbb{E}_\alpha^{\text{circ}}[\theta(t; \alpha)] = \int_{-1}^1 e^{i\theta(t; \alpha)} d\alpha, \quad \sigma^{\text{circ}}(\theta) = \sqrt{-2 \ln |\mathbb{E}_\alpha^{\text{circ}}[\theta(t; \alpha)]|}.$$

²²To motivate why a different definition for circular moments is needed, consider $y \sim U(-\pi, \pi)$ and $z \sim U(0, 2\pi)$. If we consider y and z as angles, or points on the circle, they are identical. Using the conventional mean definition, however, yields $\mathbb{E}[y] = 0$, but $\mathbb{E}[z] = \pi$.

Table 1

Approximation error of the circular mean and standard deviation (see (7.4)) of $\theta(0.15, \alpha)$ (see (7.3)) with gPC- and spline-based approximations, using N sample points.

	N	gPC error	Spline error	$\frac{\text{gPC error}}{\text{spline error}}$
$\mathbb{E}_\alpha^{\text{circ}}[\theta(0.15; \alpha)]$	32	2.2%	0.54%	4
$\mathbb{E}_\alpha^{\text{circ}}[\theta(0.15; \alpha)]$	64	0.089%	0.006%	14
$\sigma^{\text{circ}}(\theta(0.15; \alpha))$	32	0.64%	0.054%	12
$\sigma^{\text{circ}}(\theta(0.15; \alpha))$	64	0.031%	0.0009%	33

The advantage of splines over gPC with few samples for moments approximation can be seen in Table 1. The approximation of $\mathbb{E}_\alpha^{\text{circ}}[\theta(0.15; \alpha)]$ using the spline approximation with $N = 32$ is 4 times more accurate than that of the gPC; with $N = 64$ it is 14 times more accurate. The approximation of the standard deviation using the spline-based method with $N = 32$ is 12 times more accurate than the gPC; with $N = 64$ it is 33 times more accurate than the gPC-based approximation.

8. Application 2: Inviscid Burgers equation. The inviscid Burgers equation

$$(8.1) \quad u_t(t, x) + \frac{1}{2}(u^2)_x = \frac{1}{2}(\sin^2(x))_x, \quad x \in [0, \pi], \quad t \geq 0,$$

with the initial and boundary conditions $u(0, x) = u_0(x)$ and $u(t, 0) = u(t, \pi) = 0$ models isentropic gas flow in a dual-throat nozzle. Solutions of this equation can develop a static shock wave at a lateral location $x = X_s$ [42]. Following [6], we consider the case in which α is a random variable with a known distribution, $u_0(x) = u_0(x; \alpha)$ is random, and we wish to compute the PDF of X_s using Algorithms 3.1 and 4.1. In general, to do that requires, for each $1 \leq j \leq N$, computing $X_s(\alpha_j)$ by solving (8.1) with α_j . For the special initial condition

$$(8.2a) \quad u_0(x) = \alpha \sin(x);$$

however, the shock location is explicitly given by [6]

$$(8.2b) \quad \alpha = -\cos(X_s).$$

This explicit expression allows us to sample $X_s(\alpha)$ without solving (8.1).

Consider the case where

$$(8.3) \quad \alpha = \begin{cases} \frac{-1 + \sqrt{1 + 4\nu^2}}{2\nu} & \text{if } \nu \neq 0, \\ 0 & \text{if } \nu = 0, \end{cases}$$

and $\nu \sim \mathcal{N}(0, \sigma)$, i.e., it is normally distributed with a zero mean. Because α is not distributed by a classical, standard measure, there is no obvious choice of quadrature points to sample by, nor is there a “natural” orthogonal polynomials basis to expand the solution by. Therefore, the gPC approach cannot be straightforwardly applied.²³ We can, however, apply the gPC

²³Nevertheless, even for nonstandard distributions, the expansion of α by a classical orthogonal-polynomials basis can still converge spectrally, under certain conditions [14].

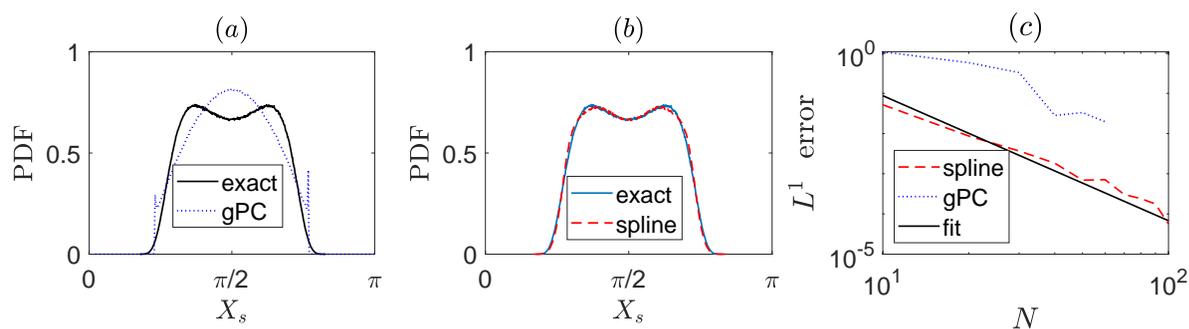


Figure 9. PDF of $X_s(\alpha)$, where $\alpha(\nu)$ is given by (8.3), and $\nu \sim \mathcal{N}(0, 0.6)$. (a) Exact PDF (solid) and gPC-based approximation (dots) with $N = 7$ sample points. (b) Same with the spline-based approximation (dashes). (c) L^1 error of the PDF approximations as a function of the number of sample points, and the power-law fit $112N^{-3.11}$ (solid).

approach to this problem by denoting $X_s(\nu) = X_s(\alpha(\nu))$ and approximating $X_s(\nu)$ using the Hermite polynomials (which are orthogonal with respect to the normal distribution).²⁴ The gPC-based approximated PDF with $N = 7$ sample points differs considerably from the exact PDF; see Figure 9(a). In contrast, the spline-based approximated PDF can be directly applied to $X_s(\alpha)$, and it is nearly indistinguishable from the exact PDF already with $N = 7$ sample points; see Figure 9(b). In general, the spline-based PDF approximation is more accurate than the gPC-based approximation by more than one order of magnitude for $5 < N < 50$; see Figure 9(c). The L^1 error of the spline-based PDF is observed numerically to decay as $N^{-3.11}$, in accordance with Theorem 4.7.

We repeated these simulations for the case with $\alpha \sim B(r, s)$, where $B(r, s)$ is the Beta distribution on $[-1, 1]$.²⁵ The spline-based approximations are nearly identical to the exact PDF, whereas the gPC method was less accurate by an order of magnitude with few samples (results not shown).

9. Discussion. In this paper, we introduced a spline-based method for density and moment estimation. The advantages of this method are as follows:

1. Our m th-order spline-based method approximates the density at a guaranteed convergence rate of $N^{-\frac{m}{d}}$, where N is the sample size and d is the noise dimension. Thus, our method outperforms KDEs for noise dimensions $1 \leq d \leq \frac{5}{2}m$.
2. It provides reasonable approximations for the density and moments using small sample sizes.
3. Its accuracy is relatively unimpaired by the presence of large derivatives.
4. It is nonintrusive, i.e., it is based solely on solving the underlying deterministic model.
5. It is easy to implement.
6. It is applicable with many choices of sample points.
7. It can be applied to nonsmooth quantities of interest.

²⁴Indeed, in [6] the authors use the gPC-Galerkin method with the Hermite polynomials [23, 61].

²⁵The PDF of the Beta distribution on $[0, 1]$ is $p(\alpha) = \frac{\alpha^{r-1}(1-\alpha)^{s-1}\Gamma(r+s)}{\Gamma(r)\Gamma(s)}$.

When $f \in C^{m+1}$, it is tempting to use splines of order $m > 3$ for density estimation, in order to attain faster than cubic convergence rate. If one generalizes Algorithm 4.1 to splines of order m , then, similarly to Theorem 4.7, a convergence of order N^{-m} is guaranteed. Even if f is analytic, however, it is not advisable to take a large m , for two reasons. First, for $s(\alpha)$ to be monotone (and so, by Lemma 4.1 for the PDF to be continuous), N should scale as $\sqrt[m]{\|f^{(m+1)}\|_\infty}$; see (E.1). Therefore, for a large m , high-order convergence might only be attained for very large sample sizes. Second, the density approximation error depends linearly on $\|f^{(m+1)}\|_\infty$ (see Appendix E), and so it might “blow-up” exponentially with m . To conclude, although we do not know whether the optimal spline order is $m = 3$, an arbitrarily high-order spline should not be used.

When approximating a d -dimensional function with a resolution h at each dimension, the total number of samples N scales as h^{-d} . As a result, for a prescribed accuracy, the computational cost grows exponentially with the dimension (the “curse of dimensionality”). In other words, for a given N , the accuracy decays exponentially with the dimension. Indeed, this is consistent with the $N^{-\frac{m}{d}}$ error estimate of the spline-based Algorithm 5.1 (Corollary 5.4). In contrast, the KDE method, which is a standard nonparametric statistical density estimator, converges at a rate of $N^{-\frac{2}{5}}$, regardless of d . Hence, our method will outperform KDE for “low” dimensions ($d < \frac{5}{2}m$) but may become inferior to KDE at higher dimensions.

A popular approach for moment estimation of high-dimensional noise is the use of sparse sampling grids [22, 59]. Recently, a spline approximation based on sparse grids was used in the context of forward uncertainty propagation [54]. Most sparse-grid methods, however, are designed with moment estimation in mind. As we have seen, even in the one-dimensional case (see section 4.1), an accurate moment approximation does not necessarily imply an accurate density estimation. Whether sparse-grids methods can be adapted to density estimation remains an open question. The proof of Theorem 5.3 in Appendix H, however, suggests sufficient conditions by which new approximation methods can be tested for efficient density estimation: (1) The settings should be such that Lemma 5.2 applies, and (2) the approximation method should have pointwise error bounds similar to Theorem 5.1.

In this paper we showed that spline-based density estimation is better than gPC-based density estimation, because it does not produce numerous artificial extremal points (see Lemma 4.1). An interpolating cubic spline, however, might still produce artificial extremal points, though not as much as the gPC polynomial. To absolutely prevent artificial extremal points from being produced, it may be better to use spline interpolants [20] and quasi-interpolants [11] which are *monotonicity-preserving* (i.e., splines which are monotone wherever the sampled data is monotone). Hence, although these methods have the same *order* of error (with respect to h) as spline interpolation, they may provide better approximations for small samples, as they are *guaranteed* not to produce artificial extremal points. We leave it to future research to check whether monotonicity-preserving interpolants provide more accurate PDF approximations than a standard interpolating cubic spline.

As noted throughout the paper, the L^∞ error bounds on the quantity of interest and its gradient are key for the success of our algorithm; see Corollaries 4.8 and 5.5. Since locality plays an important role in the existence of such error bounds for splines, it is natural to explore the use of other local approximations such as NURBS [37, 52] and radial basis

functions [19, 43]. An additional improvement may be achieved by designing surrogate models that are on one hand local but on the other hand supported on an unbounded domain, e.g., Gaussian mixtures [49]. While moment approximation in the case of unbounded input random parameters (e.g., normally or exponentially distributed α) are theoretically well understood, the rigorous study of density estimation in these setting is left for future research.

Appendix A. Proof of Corollary 3.1. We begin with the following lemma.

Lemma A.1. *Let (Ω, μ) be a probability space, denote $\|\cdot\|_p := \|\cdot\|_{L^p(\Omega)}$, and let $f, g \in L^2 \cap L^1$. Then*

$$(A.1a) \quad |\mathbb{E}_\alpha[f] - \mathbb{E}_\alpha[g]| \leq \|f - g\|_2,$$

$$(A.1b) \quad |\text{Var}(f) - \text{Var}(g)| \leq (\sigma(f) + \sigma(g)) \cdot \|f - g\|_2,$$

$$(A.1c) \quad |\sigma(f) - \sigma(g)| \leq \|f - g\|_2.$$

Proof. For all $f, g \in L^2$,

$$|\mathbb{E}_\alpha[f] - \mathbb{E}_\alpha[g]| \leq \int_\Omega |f(\alpha) - g(\alpha)| d\mu(\alpha) = \int_\Omega 1 \cdot |f(\alpha) - g(\alpha)| d\mu(\alpha) \leq \|1\|_2 \cdot \|f - g\|_2 = \|f - g\|_2,$$

where in the second inequality we used the Cauchy–Schwarz inequality. Thus, we proved (A.1a).

For $h \in L^2 \cap L^1$, let $\tilde{h} := h - \mathbb{E}_\alpha[h]$. By definition, $\text{Var}(h) = \|\tilde{h}\|_2^2$ and $\sigma(h) = \|\tilde{h}\|_2$. Hence,

$$(A.2) \quad \begin{aligned} |\text{Var}(f) - \text{Var}(g)| &= |\mathbb{E}_\alpha[\tilde{f}^2 - \tilde{g}^2]| = \left| \int_\Omega (\tilde{f} - \tilde{g})(\tilde{f} + \tilde{g}) d\mu(\alpha) \right| \leq \|\tilde{f} + \tilde{g}\|_2 \cdot \|\tilde{f} - \tilde{g}\|_2 \\ &\leq (\|\tilde{f}\|_2 + \|\tilde{g}\|_2) \cdot \|\tilde{f} - \tilde{g}\|_2 = (\sigma(f) + \sigma(g)) \cdot \|f - g\|_2. \end{aligned}$$

In addition, $\|\tilde{h}\|_2^2 = \text{Var}(h) = \mathbb{E}_\alpha[h^2] - \mathbb{E}_\alpha^2[h] \leq \mathbb{E}_\alpha[h^2] = \|h\|_2^2$, and so $\|\tilde{h}\|_2 \leq \|h\|_2$. Applying this inequality with $h = f - g$ to (A.2) yields (A.1b). Finally, by (A.1b),

$$|\sigma(f) - \sigma(g)| = \left| \frac{\sigma^2(f) - \sigma^2(g)}{\sigma(f) + \sigma(g)} \right| = \frac{|\text{Var}(f) - \text{Var}(g)|}{|\sigma(f) + \sigma(g)|} \leq \frac{\sigma(f) + \sigma(g)}{\sigma(f) + \sigma(g)} \|f - g\|_2 = \|f - g\|_2,$$

which proves (A.1c). ■

In the case of gPC, let $g = f_N^{\text{gpc}}$, the collocation gPC approximation of f ; see (3.7). Since f_N^{gpc} converges exponentially to f in the L^2 norm [59, 26], Lemma A.1 implies that the moments of f_N^{gpc} converge exponentially to the moments of f .

Appendix B. Local properties of spline interpolation. Let us first recall a classical result of Birkhoff and de Boor.

Theorem B.1 (see [5, 10]). *Let $s_i(\alpha)$ be the natural cubic spline that satisfies $s_i(\alpha_k) = \delta_{i,k}$, where $1 \leq i, k \leq N$ and $\alpha_{\min} = \alpha_1 < \alpha_2 < \dots < \alpha_N = \alpha_{\max}$ is given. Then*

$$\max_{\alpha \notin (\alpha_{i-k}, \alpha_{i+k})} |s_i(\alpha)| \leq A 2^{-k}, \quad 1 < i < N,$$

where $A > 0$ is a constant that depends on the global mesh ratio $\frac{\max_{1 < j \leq N} \alpha_j - \alpha_{j-1}}{\min_{1 < k \leq N} \alpha_k - \alpha_{k-1}}$.

Therefore, the natural cubic spline $f_N^{\text{spline}}(\alpha)$ is essentially a local approximation.

Corollary B.2. Denote the natural cubic spline $f_N^{\text{spline}} = f_N^{\text{spline}}(\alpha; f_1, \dots, f_N)$ to emphasize the dependence of the spline interpolation on the sampled values. Then

$$\max_{\alpha \notin (\alpha_{i-k}, \alpha_{i+k})} \left| \frac{\partial f_N^{\text{spline}}(\alpha; f_1, \dots, f_N)}{\partial f_i} \right| \leq A 2^{-k}, \quad 1 < i < N, \quad 1 \leq k \leq N,$$

where $A > 0$ is given by Theorem B.1.

Proof. The function $S(\alpha) = \sum_{i=1}^N f_i s_i(\alpha)$, where $s_i(\alpha)$ are defined in Theorem B.1, is a C^2 cubic spline, which by definition satisfies $S(\alpha_i) = f_i$ and $\frac{d}{d\alpha} S(\alpha_1) = \frac{d}{d\alpha} S(\alpha_N) = 0$. By the uniqueness of the natural cubic spline, $S(\alpha) = f_N^{\text{gpc}}(\alpha)$, so, $\frac{\partial f_N^{\text{spline}}(\alpha; f_1, \dots, f_N)}{\partial f_i} = s_i(\alpha)$. Hence, by Theorem B.1, the corollary is proven. ■

Appendix C. Proof of Lemma 4.1. When f is strictly increasing, its CDF is given by

$$P_f(y) := \int_{\alpha_{\min}}^{f^{-1}(y)} c(\alpha) d\alpha.$$

By the Leibniz rule and the inverse function theorem,

$$p_f(y) = \frac{dP_f(y)}{dy} = c(f^{-1}(y)) (f^{-1})' = c(f^{-1}(y)) \frac{1}{f'(f^{-1}(y))}.$$

Similarly, if f is monotonically decreasing, then $P_f(y) = \int_{f^{-1}(y)}^{\alpha_{\max}} c(\alpha) d\alpha$, and so

$$p_f(y) = -\frac{c(f^{-1}(y))}{f'(f^{-1}(y))}.$$

Note that since $f' < 0$, then $p_f(y) \geq 0$. Finally, if f is piecewise monotonic, we apply this method separately on each subinterval on which it is monotonic and sum up the contributions.

Appendix D. Sample MATLAB code for Algorithm 4.1. The following MATLAB code generates the dashed curve in Figure 2(b):

```

1 alpha_min = -1; alpha_max = 1; N = 18; %sample size
2 f = @(x) tanh(9*x) + x/2;
3 %define the initial sample on the grid [alpha_1, ... ,
  alpha_N]
4 samplingGrid = linspace(alpha_min, alpha_max, N);
5 samples = f(samplingGrid); % step 1
6 %define the refined sample grid [tilde_alpha_1, ...
  tilde_alpha_M]
7 M = 2e6;
8 denseGrid = linspace(alpha_min, alpha_max, M);
9 fN_spline = spline(samplingGrid, samples, denseGrid);
  % steps 2+3

```

```

10     %When f is given explicitly , the optimal number of bins (L)
11     %is given by (14)
12     Cf = 1.69; L =Cf*M^(1/3);
13     %step 4 – histogram of fN on denseGrid, not normalized
14     [ histogram , binsEdges ] = hist( fN_spline , L );
15     binWidth = (max( binsEdges )-min( binsEdges ))/L;
16     %normalize the histogram so that it would be a PDF
17     pdf = histogram / (sum( histogram ) * binWidth );
18     plot( binsEdges , pdf )

```

Appendix E. Proof of Theorem 4.7. Without loss of generality, we can assume that $f'(\alpha) \geq a > 0$. For brevity, denote $s(\alpha) = f_N^{\text{spline}}(\alpha)$, $h = h_{\max}$, and $\|f^{(m+1)}\|_{\infty} = \|f^{(m+1)}\|_{L^{\infty}[\alpha_{\min}, \alpha_{\max}]}$. In general, $s(\alpha)$ can be nonmonotone. By Theorem 4.4, however, $|s'(\alpha) - f'(\alpha)| < C_{\text{spl}}^{(1,m)} \|f^{(m+1)}\|_{\infty} h^m$. Hence

$$(E.1) \quad s'(\alpha) \geq \frac{a}{2} > 0, \quad N > \sqrt[m]{\frac{2C_{\text{spl}}^{(1,m)} \|f^{(m+1)}\|_{\infty}}{a}} (\alpha_{\max} - \alpha_{\min}),$$

and so $s(\alpha)$ is monotonically increasing and invertible for sufficiently large N .²⁶ Because $s(\alpha)$ interpolates $f(\alpha)$, and because both functions are monotone, then $\text{range}(s) = \text{range}(f)$. Since $s, f \in C^1$ and are invertible, by Lemma 4.1

$$(E.2) \quad \|p_f - p_s\|_1 := \int_{f(\alpha_{\min})}^{f(\alpha_{\max})} |p_f(y) - p_s(y)| dy = \int_{f(\alpha_{\min})}^{f(\alpha_{\max})} \left| \frac{c(f^{-1}(y))}{f'(f^{-1}(y))} - \frac{c(s^{-1}(y))}{s'(s^{-1}(y))} \right| dy.$$

Denote $y = f(\alpha)$ and $\alpha_{\star} := \alpha_{\star}(\alpha) = s^{-1}(f(\alpha))$. Then by a change of variable

$$(E.3) \quad \|p_f - p_s\|_1 = \int_{\alpha_{\min}}^{\alpha_{\max}} \left| \frac{c(\alpha)}{f'(\alpha)} - \frac{c(\alpha_{\star})}{s'(\alpha_{\star})} \right| f'(\alpha) d\alpha = \int_{\alpha_{\min}}^{\alpha_{\max}} |s'(\alpha_{\star})c(\alpha) - f'(\alpha)c(\alpha_{\star})| \frac{1}{s'(\alpha_{\star})} d\alpha.$$

For all $\alpha \in [\alpha_{\min}, \alpha_{\max}]$,

$$|s'(\alpha_{\star})c(\alpha) - f'(\alpha)c(\alpha_{\star})| \leq c(\alpha)|s'(\alpha_{\star}) - s'(\alpha)| + c(\alpha)|s'(\alpha) - f'(\alpha)| + f'(\alpha)|c(\alpha) - c(\alpha_{\star})|.$$

Because $s'(\alpha)$ and $c(\alpha)$ are differentiable,

$$(E.4) \quad |s'(\alpha_{\star})c(\alpha) - f'(\alpha)c(\alpha_{\star})| \leq D|\alpha - \alpha_{\star}| + \|c\|_{\infty}|f'(\alpha) - s'(\alpha)|,$$

where $D = [\|c\|_{\infty} \cdot \|s''\|_{\infty} + \|c'\|_{\infty} \cdot \|f\|_{\infty}]$.²⁷ By Lagrange's mean-value theorem, there exists β between α and α_{\star} such that

$$s(\alpha) - s(\alpha_{\star}) = s'(\beta)(\alpha - \alpha_{\star}).$$

²⁶In the numerical example (6.1), this lower bound is roughly $N > 30$.

²⁷By the same argument as (E.1), for a fixed $\epsilon > 0$ there exists a sufficiently large N_0 such that $s''(\alpha) \leq f''(\alpha) + \epsilon$ for all $N > N_0$. Therefore $\max \|s''\|_{\infty} \leq \max \|f''\|_{\infty} + \epsilon$, and so D is independent of N and depends only on $f(\alpha)$, $c(\alpha)$, α_{\min} , and α_{\max} .

On the other hand, since $\alpha_\star = s^{-1}(f(\alpha))$, then $s(\alpha_\star) = f(\alpha)$, and so

$$s(\alpha) - s(\alpha_\star) = s(\alpha) - f(\alpha).$$

Therefore $\alpha - \alpha_\star = \frac{s(\alpha) - f(\alpha)}{s'(\beta)}$. By (E.1), $s'(\beta) \geq \frac{a}{2}$, and by Theorem 4.4, we have $|f(\alpha) - s(\alpha)| \leq C_{\text{spl}}^{(0,m)} \|f^{(m+1)}\|_\infty h^{m+1}$. Hence,

$$|\alpha - \alpha_\star| \leq \frac{2C_{\text{spl}}^{(0,m)} \|f^{(m+1)}\|_\infty h^{m+1}}{a}.$$

By Theorem 4.4, $|f'(\alpha) - s'(\alpha)| \leq C_{\text{spl}}^{(1,m)} \|f^{(m+1)}\|_\infty h^m$. Hence (E.4) reads

$$(E.5) \quad |s'(\alpha_\star)c(\alpha) - f'(\alpha)c(\alpha_\star)| \leq K_1 h^m + K_2 h^{m+1},$$

where $K_1 = C_{\text{spl}}^{(0,m)} \|c\|_\infty \|f^{(m+1)}\|_\infty$ and $K_2 = \frac{2}{a} C_{\text{spl}}^{(0,m)} \|f^{(m+1)}\|_\infty D$. Substituting $\frac{1}{s'(\alpha_\star)} \leq \frac{2}{a}$, see (E.1), and (E.5) in (E.3), for sufficiently large N such that $h = \frac{\alpha_{\max} - \alpha_{\min}}{N-1} < 1$ we have that

$$\|p_f - p_s\|_1 \leq \int_{\alpha_{\min}}^{\alpha_{\max}} \frac{2(K_1 + K_2)}{a} h^m d\alpha = \frac{2(K_1 + K_2)}{a} (\alpha_{\max} - \alpha_{\min}) h^m \leq \frac{K}{N^m},$$

where $K = \frac{2(K_1 + K_2)}{a} (\alpha_{\max} - \alpha_{\min})^{m+1}$.

Similarly, by (E.5), we have that for $1 \leq q < \infty$,

$$\|p_f - p_s\|_q^q \leq \int_{\alpha_{\min}}^{\alpha_{\max}} \left| \frac{2(K_1 + K_2)}{a} h^m \right|^q d\alpha \leq K^q(q) h^{qm}$$

for a suitable $K(q) > 0$, and so $\|p_f - p_s\|_q \leq K(q) h^m \leq K(q) N^{-m}$.

Remark E.1. If $f'(\alpha) = 0$ for some values of α , the approximation p_s is not guaranteed to converge in the L^1 norm. By (E.5), however, we can guarantee a third-order convergence for the pointwise error $p_f(y) - p_s(y)$ for every real number y such that $f'(\alpha)$ does not vanish on $\{\alpha \mid f(\alpha) = y\}$.

Appendix F. Proof of Lemma 4.11. Similarly to the proof of (E.4),

$$|s'(\alpha_\star)c(\alpha) - f'(\alpha)c(\alpha_\star)| \leq D|\alpha - \alpha_\star| + c(\alpha)|f'(\alpha) - s'(\alpha)|.$$

Because $|\alpha - \alpha_\star| \leq K_2 h^4$, then by (E.3),

$$(F.1) \quad \|p_f - p_s\|_1 \leq \frac{2K_2}{a} h^4 + \int_{\alpha_{\min}}^{\alpha_{\max}} |f'(\alpha) - s'(\alpha)| c(\alpha) d\alpha.$$

Since $f'(\alpha) - s'(\alpha)$ is continuous on $[\alpha_{\min}, \alpha_{\max}]$, it vanishes and changes its sign only at $J_N < \infty$ points, denoted by $\alpha_{\min} = \gamma_0 < \gamma_1 < \dots < \gamma_{J_N} = \alpha_{\max}$. Using integration by parts, the last integral reads

$$\begin{aligned} \int_{\alpha_{\min}}^{\alpha_{\max}} |f'(\alpha) - s'(\alpha)| c(\alpha) d\alpha &= \eta \sum_{j=1}^{J-1} (-1)^j \int_{\gamma_j}^{\gamma_{j+1}} (f'(\alpha) - s'(\alpha)) c(\alpha) d\alpha \\ &= \eta \sum_{j=1}^{J-1} (-1)^j \left[c(\gamma_{j+1}) (f(\gamma_{j+1}) - s(\gamma_{j+1})) - c(\gamma_j) (f(\gamma_j) - s(\gamma_j)) \right. \\ &\quad \left. - \int_{\gamma_j}^{\gamma_{j+1}} (f(\alpha) - s(\alpha)) c'(\alpha) d\alpha \right], \end{aligned}$$

where $\eta = \text{sign} [f'(\alpha_{\min}) - s'(\alpha_{\min})]$. By Theorem 4.4,

$$|c(\gamma_j)(f(\gamma_j) - s(\gamma_j))| \leq \|c\|_{\infty} C_{\text{spl}}^{(0)} \|f^{(4)}\|_{\infty} h^4, \quad 1 \leq j \leq J_N,$$

and

$$\left| \int_{\gamma_j}^{\gamma_{j+1}} (f(\alpha) - s(\alpha)) c'(\alpha) d\alpha \right| \leq \|c'\|_{\infty} (\gamma_{j+1} - \gamma_j) C_{\text{spl}}^{(0)} \|f^{(4)}\|_{\infty} h^4, \quad 1 < j \leq J_N.$$

Substituting these bounds in (F.1) yields

$$\|p_f - p_s\|_1 \leq \frac{2K_2}{a} h^4 + K_3 h^4 + K_4 J_N h^4,$$

where $K_3 = \|c'\|_{\infty} |\alpha_{\max} - \alpha_{\min}| C_{\text{spl}}^{(0)} \|f^{(4)}\|_{\infty}$ and $K_4 = 2\|c\|_{\infty} C_{\text{spl}}^{(0)} \|f^{(4)}\|_{\infty}$. In the case of a uniform grid, the first two terms are $O(N^{-4})$, and the last term is $O(N^{-4}J_N)$, which completes the proof.

Appendix G. Proof of Lemma 5.2. For any $y \in \mathbb{R}$, the CDF of f is

$$(G.1) \quad P_f(y) = \text{Prob} \{f(\boldsymbol{\alpha}) \leq y\} = \frac{1}{\mu(\Omega)} \int_{D(y)} d\mu(\boldsymbol{\alpha}) = \frac{1}{\mu(\Omega)} \int_{D(y)} c(\boldsymbol{\alpha}) d\boldsymbol{\alpha},$$

where

$$(G.2) \quad D(y) := \{\boldsymbol{\alpha} \in \Omega \mid f(\boldsymbol{\alpha}) \leq y\}.$$

To compute the PDF $p_f(y) := \frac{d}{dy} P_f(y)$, we recall the *co-area formula*.

Lemma G.1 (see [16]). Let $A \subseteq \mathbb{R}^d$ be a Jordan set, let $u : A \rightarrow \mathbb{R}$ be Lipschitz and piecewise-differentiable such that $u^{-1}(z) \subseteq A$ is a $(d-1)$ dimensional manifold for all $z \in \mathbb{R}$, and let $g \in L^1(A)$. Then

$$(G.3) \quad \int_A g(\boldsymbol{\alpha}) |\nabla u(\boldsymbol{\alpha})| d\boldsymbol{\alpha} = \int_{z \in u(A)} dz \int_{u^{-1}(z)} g(\boldsymbol{\alpha}) d\sigma,$$

where $d\sigma$ is the $(d-1)$ dimensional surface element of $u^{-1}(z)$.

We apply the co-area formula to the right-hand side of (G.1) by substituting $A = D(y)$, $g = \frac{c}{|\nabla f|}$ and $u = f$ in (G.3). The use of (G.3) is justified because of the following:

1. $D(y)$ is bounded, since Ω is bounded. We can therefore show that $D(y)$ is Jordan by proving that $m(\partial D(y)) = 0$, where m is the Lebesgue measure in \mathbb{R}^d . Since $\partial D(y) \subseteq f^{-1}(y) \cup \partial\Omega$, it is sufficient to show that each of these sets is of measure zero. Indeed, Ω is Jordan, and so $m(\partial\Omega) = 0$. In addition, since $|\nabla f| \neq 0$ on $f^{-1}(y)$, by the implicit function theorem $f^{-1}(y)$ is a $(d - 1)$ dimensional manifold, and so $m(f^{-1}(y)) = 0$.
2. f is piecewise-differentiable by the conditions of Lemma 5.2. Furthermore, because f is piecewise-differentiable on a compact set Ω , it is also Lipschitz.
3. Since f is continuous and $|\nabla f| \neq 0$ on $\bar{\Omega}$, then $\frac{1}{|\nabla f|}$ is bounded from above. Therefore, since $c \in L^1$, so is $g = \frac{c}{|\nabla f|}$.

Thus, by Lemma G.1 and (G.1),

$$(G.4) \quad P_f(y) = \frac{1}{\mu(\Omega)} \int_{D(y)} c(\boldsymbol{\alpha}) \, d\boldsymbol{\alpha} = \frac{1}{\mu(\Omega)} \int_{-\infty}^y dz \int_{f^{-1}(z)} \frac{c}{|\nabla f|} \, d\sigma.$$

The outer integral on the right-hand side is over $(-\infty, y)$ since $f(D(y)) \subseteq (-\infty, y)$; see (G.2). Finally, since $p_f(y) = \frac{d}{dy} P_f(y)$, differentiating the last integral using the (one-dimensional) Leibnitz integral rule yields (5.2).

Appendix H. Proof of Theorem 5.3. Since $f \in C^{m+1}(\Omega)$ and Ω is compact, f is also Lipschitz. Hence, Lemma 5.2 can be applied with $m([0, 1]^d) = 1$ and $c(\boldsymbol{\alpha}) \equiv 1$, yielding

$$(H.1) \quad \|p_f - p_s\|_q^q = \int_{-\infty}^{\infty} I^q(y) \, dy, \quad I(y) := \left| \int_{f^{-1}(y)} \frac{1}{|\nabla f|} \, d\sigma - \int_{s^{-1}(y)} \frac{1}{|\nabla s|} \, d\sigma \right|,$$

where σ is the $(d - 1)$ dimensional surface measures induced by the Lebesgue measure.

The outline of the proof is as follows:

1. For a fixed y in the image of $s(\boldsymbol{\alpha})$, we construct a cover $\{A_j(y)\}_{j=1}^d$ of $s^{-1}(y)$.
2. We then construct a set of maps $\phi_j : A_j(y) \rightarrow f^{-1}(y)$, which are characterized in Lemma H.1.
3. We construct a disjoint cover $\tilde{A}_j \subseteq A_j(y)$ for $1 \leq j \leq d$. Lemma H.3 proves that $\{\phi_j(\tilde{A}_j)\}_{j=1}^d$ are mutually disjoint, up to an $O(h^m)$ error, and almost cover $f^{-1}(y)$, up to an $O(h^m)$ error.
4. By an inclusion-exclusion argument and the implicit function theorem, we split the integral of $I(y)$ to d integrals over compact domains in \mathbb{R}^{d-1} .
5. By Theorem 5.1, and similar to the proof of the one-dimensional counterpart (Theorem 4.7), we bound each of the integrals obtained in step 3. Thus, we obtain a pointwise bound on $p_f(y) - p_s(y)$.
6. Finally, we use compactness of Ω and the fact that $f, s \in C^1(\Omega)$ to bound $\|p_f - p_s\|_1$.

Step 1. For brevity, denote by $\partial_{\alpha_j} = \frac{\partial}{\partial \alpha_j}$ the partial derivative along the j th axis for $1 \leq j \leq d$. Fix y , and let $A_j = A_j(y) \subseteq s^{-1}(y)$ be defined by

$$(H.2) \quad A_j := \left\{ \boldsymbol{\alpha} \in s^{-1}(y) \mid |\partial_{\alpha_j} f(\boldsymbol{\alpha})| > \frac{\kappa_f}{d} \right\}, \quad j = 1, \dots, d.$$

Since $|\nabla f| = \sqrt{\sum_{j=1}^d (\partial_{\alpha_j} f)^2} \geq \kappa_f$ on Ω , for every $\alpha \in s^{-1}(y)$ at least one component of $\nabla f(\alpha)$ satisfies $|\partial_{\alpha_j} f| \geq \frac{\kappa_f}{d}$.²⁸ Hence, $\alpha \in A_j(y)$ for some $1 \leq j \leq d$, and so

$$(H.3) \quad s^{-1}(y) = \cup_{j=1}^d A_j(y).$$

Step 2. Next, we prove the existence of the maps $\phi_j : A_j \rightarrow f^{-1}(y)$.

Lemma H.1. *Let $\alpha \in A_j(y)$ and let h be defined as in Theorem 5.1. Then for a sufficiently small $h > 0$, there exists a real number $\delta = \delta(\alpha)$ such that*

1. $\alpha + \delta(\alpha)\hat{e}_j \in f^{-1}(y)$, where \hat{e}_j is the unit vector in the direction of the j th axis;
2. the maps

$$(H.4) \quad \phi_j(\alpha) := \alpha + \delta(\alpha)\hat{e}_j, \quad j = 1, \dots, d,$$

are injective from $A_j = A_j(y)$ to $f^{-1}(y)$;

3. for every $\alpha \in A_j$,

$$(H.5) \quad \delta(\alpha) = O(h^{m+1});$$

4. for every $E \subseteq A_j$,

$$(H.6) \quad |\sigma(E) - \sigma(\phi_j(E))| = O(h^m),$$

where as in (H.1), σ is the $(d-1)$ dimensional surface measure induced by the Lebesgue measure on Ω .

Proof. 1. We prove this for the case where $y > f(\alpha)$ and $\partial_{\alpha_j} f(\alpha) > 0$ on Ω . The proofs for the three other cases are similar. Since $f \in C^{m+1}(\Omega)$ and Ω is compact, all the second derivatives of f are bounded, and so $|\partial_{\alpha_j}^2 f| < M_2 < \infty$ on Ω . Hence, since $\partial_{\alpha_j} f(\alpha) > \frac{\kappa_f}{d}$, there exists a segment $L = L(\alpha) = \{\alpha + \xi\hat{e}_j, |\xi| < \xi_{\max}\}$, where ξ_{\max} depends *only* on M_2 , such that $\partial_{\alpha_j} f > \frac{\kappa_f}{2d}$ on L . Therefore $f(\alpha + \xi_{\max}\hat{e}_j) > f(\alpha) + \frac{\kappa_f}{2d}\xi_{\max}$. By the mean-value theorem, f attains on L all values in $[f(\alpha), f(\alpha) + \frac{\kappa_f}{2d}\xi_{\max}]$. Now, by Theorem 5.1, since $\alpha \in s^{-1}(y)$ and since $y > f(\alpha)$,

$$(H.7) \quad y - f(\alpha) = s(\alpha) - f(\alpha) \leq C_m h^{m+1}.$$

Hence, for h sufficiently small, $y \in [f(\alpha), f(\alpha) + \frac{\kappa_f \xi_{\max}}{2d}]$, and so there exists a point $\alpha + \delta(\alpha)\hat{e}_j \in L$ such that $f(\alpha + \delta(\alpha)\hat{e}_j) = y$.

2. Assume by negation that ϕ_j is *not* injective. Then there exist $\alpha^1, \alpha^2 \in A_j$ such that $\phi_j(\alpha^1) = \phi_j(\alpha^2) = \lambda$. Since ϕ_j only changes the j th coordinate (see (H.4)), we can regard s and f as single-variable functions of the j th coordinate α_j . Since $\phi_j(\alpha^1) = \phi_j(\alpha^2) = \lambda$, from the proof of item 1 in this lemma it follows that $\lambda \in L(\alpha^1) \cap L(\alpha^2)$. Hence, the segment between α^1 and α^2 is contained in $L(\alpha^1) \cup L(\alpha^2)$, where we know that $|\partial_{\alpha_j} f| > \frac{\kappa_f}{2d}$. By Theorem 5.1, this means that if h is sufficiently small, $|\partial_{\alpha_j} s| > 0$ on the segment between α^1 and α^2 . This leads to a contradiction, since on the one hand $\alpha^1, \alpha^2 \in A_j(y) \subseteq s^{-1}(y)$, and so $s(\alpha^1) = s(\alpha^2) = y$, but on the other hand $s(\alpha)$ is strictly monotone on the segment between α^1 and α^2 .

²⁸Since $\kappa_f \leq \sqrt{\sum_{j=1}^d (\partial_{\alpha_j} f)^2} \leq \sqrt{d} \max_{j=1, \dots, d} |\partial_{\alpha_j} f|$, then $\max_{j=1, \dots, d} |\partial_{\alpha_j} f| \geq \frac{\kappa_f}{\sqrt{d}} > \frac{\kappa_f}{d}$.

3. Since $f \in C^2$, and by (H.4),

$$(H.8) \quad \partial_{\alpha_j} f(\phi_j(\boldsymbol{\alpha})) - \partial_{\alpha_j} f(\boldsymbol{\alpha}) = \partial_{\alpha_j} f(\boldsymbol{\alpha} + \delta(\boldsymbol{\alpha})\hat{e}_j) - \partial_{\alpha_j} f(\boldsymbol{\alpha}) = O(\delta(\boldsymbol{\alpha})).$$

In addition, by the Lagrange mean-value theorem, for any $\boldsymbol{\alpha} \in s^{-1}(y)$

$$s(\boldsymbol{\alpha}) - f(\boldsymbol{\alpha}) = y - f(\boldsymbol{\alpha}) = f(\boldsymbol{\alpha} + \delta(\boldsymbol{\alpha})\hat{e}_j) - f(\boldsymbol{\alpha}) = \partial_{\alpha_j} f(\boldsymbol{\alpha} + \zeta\hat{e}_j) \cdot \delta(\boldsymbol{\alpha}), \quad 0 \leq \zeta \leq \delta.$$

Hence, using Theorem 5.1, and since $|\partial_{\alpha_j} f| \geq \frac{\kappa_f}{2d}$ on the segment between $\boldsymbol{\alpha}$ and $\phi_j(\boldsymbol{\alpha})$ (see the proof of item 1 in this lemma), we have that

$$(H.9) \quad |\delta(\boldsymbol{\alpha})| = \left| \frac{s(\boldsymbol{\alpha}) - f(\boldsymbol{\alpha})}{\partial_{\alpha_j} f(\boldsymbol{\alpha} + \zeta\hat{e}_j)} \right| \leq \frac{C_m h^{m+1}}{\frac{\kappa_f}{2d}} = O(h^{m+1}).$$

4. For brevity of notation and without loss of generality, fix $j = d$, and let $E \subseteq A_d$. In this case, $\partial_{\alpha_d} s \neq 0$ on E ,²⁹ and so by the implicit function theorem there exists a function S such that if $s(\alpha_1, \dots, \alpha_d) = y$, then $\alpha_d = S(\alpha_1, \dots, \alpha_{d-1})$. The domain of S is

$$G_E := \{(\alpha_1, \dots, \alpha_{d-1}) \mid \exists \alpha_d \in [0, 1] \text{ s.t. } (\alpha_1, \dots, \alpha_d) \in E\}.$$

In particular, if $(\alpha_1, \dots, \alpha_{d-1}) \in G_E$, then $s(\alpha_1, \dots, \alpha_{d-1}, S(\alpha_1, \dots, \alpha_{d-1})) = y$. Therefore

$$\sigma(E) = \int_E 1 \, d\sigma = \int_{G_E} \sqrt{1 + |\nabla S|^2} \, d\alpha_1 \cdots d\alpha_{d-1}.$$

Furthermore, by the implicit function theorem, $\partial_{\alpha_j} S = -\frac{\partial_{\alpha_j} s}{\partial_{\alpha_d} s}$ for $1 \leq j < d$, and so

$$\sqrt{1 + |\nabla S|^2} = \sqrt{1 + \sum_{j=1}^{d-1} \left(\frac{\partial_{\alpha_j} s}{\partial_{\alpha_d} s}\right)^2} = \frac{1}{|\partial_{\alpha_d} s|} \sqrt{(\partial_{\alpha_d} s)^2 + \sum_{j=1}^{d-1} (\partial_{\alpha_j} s)^2} = \frac{1}{|\partial_{\alpha_d} s|} |\nabla s|.$$

Hence,

$$(H.10a) \quad \sigma(E) = \int_{G_E} \frac{|\nabla s|}{|\partial_{\alpha_d} s|} \, d\alpha_1 \cdots d\alpha_{d-1}.$$

Next, since $|\partial_{\alpha_d} f| \geq \frac{\kappa_f}{2d}$ on $\phi_d(E)$ (see the proof of item 1 in this lemma), we similarly apply the implicit function on $\phi_d(E)$: there exists function $F : G_{\phi_d(E)} \rightarrow \mathbb{R}$ where $G_{\phi_d(E)} \subset \mathbb{R}^{d-1}$ such that $f(\alpha_1, \dots, \alpha_{d-1}, F(\alpha_1, \dots, \alpha_{d-1})) = y$. Hence, since $\phi_d(E) \subseteq f^{-1}(y)$,

$$(H.10b) \quad \sigma(\phi_d(E)) = \int_{\phi_d(E)} 1 \, d\sigma = \int_{G_{\phi_d(E)}} \frac{|\nabla f|}{|\partial_{\alpha_d} f|} \, d\alpha_1 \cdots d\alpha_{d-1}.$$

Next, by item 2 of this lemma, ϕ_d induces a bijection $\varphi_d : G_E \rightarrow G_{\phi_d(E)}$. But, because ϕ_d only alters the α_d coordinate, $\varphi_d = \text{Id}$, and so $G_E = G_{\phi_d(E)}$. Using this equality and (H.10) yields

²⁹ $\partial_{\alpha_d} s \neq 0$ on A_d for sufficiently small h since $|\partial_{\alpha_d} f| \geq \frac{\kappa_f}{d}$ on A_d , and since by Theorem 5.1 $|\partial_{\alpha_d} s - \partial_{\alpha_d} f| = O(h^m)$.

$$\begin{aligned}
 (H.11) \quad |\sigma(E) - \sigma(\phi_d(E))| &= \left| \int_{G_E} \left(\frac{|\nabla f(\phi_d(\boldsymbol{\beta}))|}{|\partial_{\alpha_d} f(\phi_d(\boldsymbol{\beta}))|} - \frac{|\nabla s(\boldsymbol{\beta})|}{|\partial_{\alpha_d} s(\boldsymbol{\beta})|} \right) d\alpha_1 \cdots d\alpha_{d-1} \right| \\
 &\leq \int_{G_E} \frac{||\nabla f(\phi_d(\boldsymbol{\beta}))| \cdot |\partial_{\alpha_d} s(\boldsymbol{\beta})| - |\nabla s(\boldsymbol{\beta})| \cdot |\partial_{\alpha_d} f(\phi_d(\boldsymbol{\beta}))||}{|\partial_{\alpha_d} f(\phi_d(\boldsymbol{\beta}))| \cdot |\partial_{\alpha_d} s(\boldsymbol{\beta})|} d\alpha_1 \cdots d\alpha_{d-1},
 \end{aligned}$$

where for brevity we denote $\boldsymbol{\beta} := (\alpha_1, \dots, \alpha_{d-1}, S(\alpha_1, \dots, \alpha_{d-1})) \in E$ and note that by (H.4)

$$(\alpha_1, \dots, \alpha_{d-1}, F(\alpha_1, \dots, \alpha_{d-1})) = \phi_d(\boldsymbol{\beta}).$$

To bound the right-hand side of (H.11), note that since $|\partial_{\alpha_d} f| > \frac{\kappa_f}{d}$ on E , and since by Theorem 5.1, $|\partial_{\alpha_d} s - \partial_{\alpha_d} f| \leq C_m h^m$, then for a sufficiently small h , $|\partial_{\alpha_d} s| > \frac{\kappa_f}{2d}$ on E . Substituting these bounds in (H.11) yields

$$\begin{aligned}
 (H.12) \quad &|\sigma(E) - \sigma(\phi_d(E))| \\
 &\leq \frac{2d^2}{\kappa_f^2} \int_{G_E} ||\nabla f(\phi_d(\boldsymbol{\beta}))| \cdot |\partial_{\alpha_d} s(\boldsymbol{\beta})| - |\nabla s(\boldsymbol{\beta})| \cdot |\partial_{\alpha_d} f(\phi_d(\boldsymbol{\beta}))|| d\alpha_1 \cdots d\alpha_{d-1}.
 \end{aligned}$$

Therefore, we can rewrite and bound the right-hand-side integrand by

$$\begin{aligned}
 (H.13) \quad &||\nabla f(\phi_d(\boldsymbol{\beta}))| \cdot |\partial_{\alpha_d} s(\boldsymbol{\beta})| - |\nabla s(\boldsymbol{\beta})| \cdot |\partial_{\alpha_d} f(\phi_d(\boldsymbol{\beta}))|| \\
 &\leq |\partial_{\alpha_d} s(\boldsymbol{\beta})| \cdot ||\nabla f(\phi_d(\boldsymbol{\beta}))| - |\nabla f(\boldsymbol{\beta})|| + |\nabla f(\boldsymbol{\beta})| \cdot ||\partial_{\alpha_d} s(\boldsymbol{\beta})| - |\partial_{\alpha_d} s(\phi_d(\boldsymbol{\beta}))|| \\
 &\quad + |\nabla f(\boldsymbol{\beta})| \cdot ||\partial_{\alpha_d} s(\phi_d(\boldsymbol{\beta}))| - |\partial_{\alpha_d} f(\phi_d(\boldsymbol{\beta}))|| + |\partial_{\alpha_d} f(\phi_d(\boldsymbol{\beta}))| \cdot ||\nabla f(\boldsymbol{\beta})| - |\nabla s(\boldsymbol{\beta})||.
 \end{aligned}$$

Since $s, f \in C^2(\Omega)$ and Ω is compact, $\partial_{\alpha_d} s, \partial_{\alpha_d} f$, and ∇f are bounded on Ω . Furthermore, since $s, f \in C^2$, the first and second terms in the right-hand side of (H.13) are $O(\delta)$, and so by (H.5) both of these terms are $O(h^{m+1})$. In addition, by Theorem 5.1 the third and fourth terms on the right-hand side of (H.13) are $O(h^m)$. Hence, the left-hand side of (H.13) is $O(h^m)$, and so finally, (H.12) reads

$$|\sigma(E) - \sigma(\phi_d(E))| \leq \frac{2d^2}{\kappa_f^2} \int_{G_E} K h^m d\alpha_1, \dots, d\alpha_{d-1} \leq \tilde{K} h^m$$

for some constant $\tilde{K} > 0$. ■

We finish this step by noting that Lemma H.1 would still hold if we interchange f and s . Hence, we have the following.

Corollary H.2. *There exist sets $B_j \subseteq f^{-1}(y)$ such that $f^{-1}(y) = \cup_{j=1}^d B_j$ and maps $\tilde{\phi}_j : B_j \rightarrow s^{-1}(y)$ for which items 1–4 of Lemma H.1 hold, interchanging f and s .*

Step 3. Next, we repartition $s^{-1}(y)$ into disjoint sets $\{\tilde{A}_j\}_{j=1}^d$, where $\tilde{A}_j \subseteq A_j$ for every $1 \leq j \leq d$. Let $\tilde{A}_1 := A_1$, and define

$$(H.14) \quad \tilde{A}_j := A_j \setminus \left(\cup_{k=1}^{j-1} \tilde{A}_k \right), \quad 1 < j \leq d.$$

Since by construction $\cup_{j=1}^d \tilde{A}_j = \cup_{j=1}^d A_j$, and since by (H.3) $\cup_{j=1}^d A_j = s^{-1}(y)$, then

$$\cup_{j=1}^d \tilde{A}_j = s^{-1}(y).$$

Hence, since the sets $\{\tilde{A}_j\}_{j=1}^d$ are disjoint, we can rewrite the first component of $I(y)$ (see (H.1)) as

$$(H.15) \quad \int_{s^{-1}(y)} \frac{1}{|\nabla s|} d\sigma = \sum_{j=1}^d \int_{\tilde{A}_j} \frac{1}{|\nabla s|} d\sigma.$$

To prove a counterpart of (H.15) for $\int_{f^{-1}(y)} \frac{1}{|\nabla f|} 1 d\sigma$, we first prove the following lemma.

Lemma H.3. *Let σ be the surface measure on $f^{-1}(y)$, let $\{\tilde{A}_j\}_{j=1}^d$ be defined by (H.15), and let $\{\phi_j\}_{j=1}^d$ be defined by (H.4).*

1. For any $1 \leq k, j \leq d$ with $k \neq j$, then

$$(H.16) \quad \sigma(\phi_j(\tilde{A}_j) \cap \phi_k(\tilde{A}_k)) = O(h^m).$$

- 2.

$$(H.17) \quad \sigma(f^{-1}(y) \setminus \cup_{j=1}^d \phi_j(\tilde{A}_j)) = O(h^m).$$

Proof. 1. Fix the indices $j \neq k$ and denote for brevity $D_{jk} = \phi_j(\tilde{A}_j) \cap \phi_k(\tilde{A}_k)$. Let $\beta \in D_{jk}$. By injectivity of ϕ_j and ϕ_k (see Lemma H.1), there exist unique points $\alpha^{(j)} \in \tilde{A}_j$ and $\alpha^{(k)} \in \tilde{A}_k$ such that $\phi_j(\alpha^{(j)}) = \phi_k(\alpha^{(k)}) = \beta$. By definition (H.4),

$$\beta - \alpha^{(j)} = \delta(\alpha^{(j)})\hat{e}_j, \quad \beta - \alpha^{(k)} = \delta(\alpha^{(k)})\hat{e}_k.$$

Since $\hat{e}_j \perp \hat{e}_k$ and since by (H.5) $\delta(\alpha^{(j)}), \delta(\alpha^{(k)}) = O(h^{m+1})$, then³⁰

$$|\alpha^{(j)} - \alpha^{(k)}| = O(h^{m+1}).$$

Next, denote the geodesic distance on s^{-1} by $|\cdot|_s$. Since $s \in C^1$, then $|\nabla s|$ is bounded from above on Ω and so $|\alpha^{(j)} - \alpha^{(k)}|_s = O(h^{m+1})$ as well. But since the interiors of \tilde{A}_j and \tilde{A}_k are disjoint, then the geodesic path between $\alpha^{(j)}$ and $\alpha^{(k)}$ must pass through a point $\alpha^* \in \partial\tilde{A}_j \cap \partial\tilde{A}_k$. Hence,

$$(H.18) \quad |\alpha^* - \alpha^{(j)}|_s = O(h^{m+1}).$$

Since (H.18) holds for any $\beta \in D_{jk}$ and $\alpha^{(j)} = \phi_j^{-1}(\beta)$, then

$$\phi_j^{-1}(D_{jk}) \subseteq E_{jk}(h) := \left\{ \alpha \in s^{-1}(y) \mid \inf_{\alpha^* \in \partial\tilde{A}_j \cap \partial\tilde{A}_k} |\alpha - \alpha^*|_s \leq Kh^{m+1} \right\}$$

³⁰Geometrically, the points $\alpha^{(j)}$, $\alpha^{(k)}$, and β are the vertices of a right-angle triangle, where both legs are $O(h^{m+1})$. Hence, by the Pythagorean theorem, the length of the hypotenuse is also $O(h^{m+1})$.

for some $K > 0$. It is therefore sufficient to show that $\sigma(E_{jk}(h)) = O(h^m)$ for $0 < h \ll 0$.

By construction, $\partial\tilde{A}_j \cap \partial\tilde{A}_k \subseteq \cup_{j=1}^d \partial A_j$. Since $f \in C^1$, then $\sigma(\cup_{j=1}^d \partial A_j) = 0$ and so by monotonicity of measure $\sigma(\partial\tilde{A}_j \cap \partial\tilde{A}_k) = 0$ as well.³¹ Furthermore $\partial\tilde{A}_j \cap \partial\tilde{A}_k$ is a finite union of smooth subsurface of $s^{-1}(y)$, each of finite $(d-2)$ dimensional surface measure.³² Finally, since $\partial\tilde{A}_j \cap \partial\tilde{A}_k$ is compact in the topology of the smooth $(d-1)$ dimensional manifold $s^{-1}(y)$ (it is bounded and close), and since $E_{jk}(h)$ is of *geodesic* radius Kh^{m+1} from $\partial\tilde{A}_j \cap \partial\tilde{A}_k$, then $\sigma(E_{jk}) = O((h^{m+1})^{(d-1)}) \leq O(h^m)$. Hence,

$$(H.19) \quad \sigma(\phi_j^{-1}(D_{jk})) \leq \sigma(E_{jk}(h)) = O(h^m).$$

In addition, since ϕ_j is injective, $\phi_j(\phi_j^{-1}(D_{jk})) = D_{jk}$. Hence, by taking $E = \phi_j^{-1}(D_{jk})$ in (H.6) yields

$$|\sigma(\phi_j^{-1}(D_{jk})) - \sigma(D_{jk})| = |\sigma(E) - \sigma(\phi_j(E))| \leq O(h^m).$$

Combined with (H.19) this proves that $\sigma(D_{jk}) = O(h^m)$, as required.

2. Since $\cup_{j=1}^d \phi_j(\tilde{A}_j) \subseteq f^{-1}(y)$, then

$$(H.20a) \quad \sigma\left(\cup_{j=1}^d \phi_j(\tilde{A}_j)\right) \leq \sigma(f^{-1}(y)).$$

On the other hand, by item (H.16) and by the inclusion-exclusion argument

$$\begin{aligned} & \sigma\left(\cup_{j=1}^d \phi_j(\tilde{A}_j)\right) \\ &= \sum_{j=1}^d \sigma\left(\phi_j(\tilde{A}_j)\right) - \sum_{j_1, j_2} \sigma\left(\phi_{j_1}(\tilde{A}_{j_1}) \cap \phi_{j_2}(\tilde{A}_{j_2})\right) \\ & \quad + \dots + (-1)^{d+1} \sigma\left(\phi_1(\tilde{A}_1) \cap \dots \cap \phi_d(\tilde{A}_d)\right) \\ &= \sum_{j=1}^d \sigma\left(\phi_j(\tilde{A}_j)\right) + O(h^m) = \sum_{j=1}^d \sigma(\tilde{A}_j) + O(h^m), \end{aligned}$$

where the last equality is due to (H.6). Hence,

$$(H.20b) \quad \sigma\left(\cup_{j=1}^d \phi_j(\tilde{A}_j)\right) = \sum_{j=1}^d \sigma(\tilde{A}_j) + O(h^m) = \sigma(\cup_{j=1}^d \tilde{A}_j) + O(h^m) = \sigma(s^{-1}(y)) + O(h^m),$$

where the second equality follows from the fact that the sets $\{\tilde{A}_j\}_{j=1}^d$ are disjoint, and the third equality follows from $\cup_{j=1}^d \tilde{A}_j = s^{-1}(y)$.

Since the left-hand sides of (H.20a) and (H.20b) are identical, it follows that

³¹For each $1 \leq j \leq d$, the set ∂A_j is the boundary of the smooth manifold A_j , and so it is of measure zero.

³²For example, if $d = 3$, then $\partial\tilde{A}_j \cap \partial\tilde{A}_k$ is a finite set of curves, each with a finite length.

$$(H.21a) \quad \sigma(s^{-1}(y)) + O(h^m) \leq \sigma(f^{-1}(y)).$$

Crucially, since by Corollary H.2, both Lemma H.1 and item 1 of this lemma remain valid if we interchange f and s , we also have that

$$(H.21b) \quad \sigma(f^{-1}(y)) + O(h^m) \leq \sigma(s^{-1}(y)).$$

Combining the two inequalities of (H.21) yields that

$$(H.22) \quad |\sigma(f^{-1}(y)) - \sigma(s^{-1}(y))| = O(h^m).$$

Finally

$$\begin{aligned} \sigma\left(f^{-1}(y) \setminus \cup_{j=1}^d \phi_j(\tilde{A}_j)\right) &= \sigma(f^{-1}(y)) - \sigma\left(\cup_{j=1}^d \phi_j(\tilde{A}_j)\right) \\ &\leq |\sigma(f^{-1}(y)) - \sigma(s^{-1}(y))| + O(h^m) = O(h^m), \end{aligned}$$

where the inequality in the first line is due to (H.20b), and the last equality is due to (H.22). ■

Step 4. By (H.17), and since $\frac{1}{|\nabla f|} \leq \frac{1}{\kappa_f}$, then

$$\int_{f^{-1}(y)} \frac{1}{|\nabla f|} d\sigma = \int_{\cup_{j=1}^d \phi_j(\tilde{A}_j)} \frac{1}{|\nabla f|} d\sigma + O(h^m).$$

Hence, by an inclusion-exclusion argument,

$$(H.23) \quad \begin{aligned} \int_{f^{-1}(y)} \frac{1}{|\nabla f|} d\sigma &= O(h^m) + \sum_{j=1}^d \int_{\phi_j(\tilde{A}_j)} \frac{1}{|\nabla f|} d\sigma \\ &\quad - \sum_{\substack{j_1 < j_2 \\ j_1=1}}^d \int_{\phi_{j_1}(\tilde{A}_{j_1}) \cap \phi_{j_2}(\tilde{A}_{j_2})} \frac{1}{|\nabla f|} d\sigma + \dots + (-1)^{d-1} \int_{\phi_1(\tilde{A}_1) \cap \dots \cap \phi_d(\tilde{A}_d)} \frac{1}{|\nabla f|} d\sigma. \end{aligned}$$

But, by (H.16), we can reduce all of the higher-order terms to yield

$$(H.24) \quad \int_{f^{-1}(y)} \frac{1}{|\nabla f|} d\sigma = \sum_{j=1}^d \int_{\phi_j(\tilde{A}_j)} \frac{1}{|\nabla f|} d\sigma + O(h^m).$$

Hence, substituting (H.15) and (H.24) into (H.1) yields

$$(H.25) \quad I(y) \leq \sum_{j=1}^d \left| \int_{\phi_j(\tilde{A}_j)} \frac{1}{|\nabla f|} d\sigma - \int_{\tilde{A}_j} \frac{1}{|\nabla s|} d\sigma \right| + O(h^m).$$

Step 5. By (H.25), in order to show that $I(y) = O(h^m)$, it is sufficient to prove that

$$(H.26) \quad I_j(y) := \left| \int_{\phi_j(\tilde{A}_j)} \frac{1}{|\nabla f|} d\sigma - \int_{\tilde{A}_j} \frac{1}{|\nabla s|} d\sigma \right| = O(h^m), \quad 1 \leq j \leq d.$$

This proof is similar to that of item 4 in Lemma H.1. For ease of notation, we assume without loss of generality that $j = d$. In this case, $\partial_{\alpha_d} s \neq 0$ on \tilde{A}_j ,³³ and so by the implicit function theorem there exists a function S such that if $s(\alpha_1, \dots, \alpha_d) = y$, then $\alpha_d = S(\alpha_1, \dots, \alpha_{d-1})$. The domain of S is

$$G_{\tilde{A}_d} := \{(\alpha_1, \dots, \alpha_{d-1}) \mid \exists \alpha_d \in [0, 1] \text{ s.t. } (\alpha_1, \dots, \alpha_d) \in \tilde{A}_d\}.$$

In particular, if $(\alpha_1, \dots, \alpha_{d-1}) \in G_{\tilde{A}_d}$, then $s(\alpha_1, \dots, \alpha_{d-1}, S(\alpha_1, \dots, \alpha_{d-1})) = y$. Therefore

$$\int_{\tilde{A}_d} \frac{1}{|\nabla s|} d\sigma = \int_{G_{\tilde{A}_d}} \frac{1}{|\nabla s(\alpha_1, \dots, \alpha_{d-1}, S(\alpha_1, \dots, \alpha_{d-1}))|} \sqrt{1 + |\nabla S|^2} d\alpha_1 \cdots d\alpha_{d-1}.$$

Furthermore, by the implicit function theorem, $\partial_{\alpha_j} S = -\frac{\partial_{\alpha_j} s}{\partial_{\alpha_d} s}$ for $1 \leq j < d$, and so

$$\sqrt{1 + |\nabla S|^2} = \sqrt{1 + \sum_{j=1}^{d-1} \left(\frac{\partial_{\alpha_j} s}{\partial_{\alpha_d} s}\right)^2} = \frac{1}{|\partial_{\alpha_d} s|} \sqrt{(\partial_{\alpha_d} s)^2 + \sum_{j=1}^{d-1} (\partial_{\alpha_j} s)^2} = \frac{1}{|\partial_{\alpha_d} s|} |\nabla s|.$$

Hence,

$$(H.27a) \quad \int_{\tilde{A}_d} \frac{1}{|\nabla s|} d\sigma = \int_{G_{\tilde{A}_d}} \frac{1}{|\partial_{\alpha_d} s|} d\alpha_1 \cdots d\alpha_{d-1}.$$

Similarly, since $|\partial_{\alpha_d} f| \geq \frac{\kappa_f}{2d} > 0$ on $\phi_j(\tilde{A}_j)$, applying the implicit function theorem to f yields a function $F: G_{\phi_d(\tilde{A}_d)} \rightarrow \mathbb{R}$ where $G_{\phi_d(\tilde{A}_d)} \subset \mathbb{R}^{d-1}$, such that $f(\alpha_1, \dots, \alpha_{d-1}, F(\alpha_1, \dots, \alpha_{d-1})) = y$, and

$$(H.27b) \quad \int_{\phi_d(\tilde{A}_d)} \frac{1}{|\nabla f|} d\sigma = \int_{G_{\phi_d(\tilde{A}_d)}} \frac{1}{|\partial_{\alpha_d} f|} d\alpha'_1 \cdots d\alpha'_{d-1}.$$

Next, by item 2 of Lemma H.1, ϕ_d induces a surjective map $\varphi_d : G_{\tilde{A}_d} \rightarrow G_{\phi_d(\tilde{A}_d)}$. But, because ϕ_d only alters the α_d coordinate, $\varphi_d = \text{Id}$, and so $G_{\tilde{A}_d} = G_{\phi_d(\tilde{A}_d)}$. Substituting this equality and (H.27) into (H.26) yields

$$(H.28) \quad I_d(y) \leq \left| \int_{G_{\tilde{A}_d}} \left(\frac{1}{|\partial_{\alpha_d} f|} - \frac{1}{|\partial_{\alpha_d} s|} \right) d\alpha'_1 \cdots d\alpha'_{d-1} \right| \leq \int_{G_{\tilde{A}_d}} \frac{|\partial_{\alpha_d} f - \partial_{\alpha_d} s|}{|\partial_{\alpha_d} f| \cdot |\partial_{\alpha_d} s|} d\alpha'_1 \cdots d\alpha'_{d-1}.$$

Bounding (H.28) is similar to its one-dimensional counterpart in Appendix E. Since $|\partial_{\alpha_d} f| > \frac{\kappa_f}{d}$, and since by Theorem 5.1, $|\partial_{\alpha_d} s - \partial_{\alpha_d} f| \leq C_m h^m$, then for a sufficiently small h , $|\partial_{\alpha_d} s| > \frac{\kappa_f}{2d}$ on $\phi_d(\tilde{A}_d)$. Substituting these bounds in (H.28) yields

³³As before, this follows for sufficiently small h from the fact that $|\partial_{\alpha_d} f| \geq \frac{\kappa_f}{d}$ on $A_d(y)$ and from Theorem 5.1.

$$I_d(y) \leq \frac{2d^2}{\kappa_f^2} \int_{G_{\tilde{A}_d}} |\partial_{\alpha_d} s(\alpha_1, \dots, \alpha_{d-1}, S(\alpha_1, \dots, \alpha_{d-1})) - \partial_{\alpha_d} f(\alpha_1, \dots, \alpha_{d-1}, F(\alpha_1, \dots, \alpha_{d-1}))| d\alpha_1 \cdots d\alpha_{d-1}.$$

Next, if we denote $\beta = (\alpha_1, \dots, \alpha_{d-1}, S(\alpha_1, \dots, \alpha_{d-1}))$, then by (H.4)

$$\phi_d(\beta) = (\alpha_1, \dots, \alpha_{d-1}, F(\alpha_1, \dots, \alpha_{d-1})).$$

Therefore, we can rewrite and bound the left-hand-side integrand by

$$(H.29) \quad |\partial_{\alpha_d} s(\beta) - \partial_{\alpha_d} f(\phi_d(\beta))| \leq |\partial_{\alpha_d} s(\beta) - \partial_{\alpha_d} f(\beta)| + |\partial_{\alpha_d} f(\beta) - \partial_{\alpha_d} f(\phi_d(\beta))|.$$

This bound is very similar to its one-dimensional counterpart in (E.3). The first term on the right-hand side of (H.29) is $O(h^m)$; see Theorem 5.1. In addition, since $f \in C^2$, the second term in the right-hand side of (H.29) reads

$$|\partial_{\alpha_d} f(\beta) - \partial_{\alpha_d} f(\phi_d(\beta))| \leq M_2 |\beta - \phi_d(\beta)| = M_2 |\delta(\beta)| = O(h^{m+1}),$$

where, as before, $M_2 = \max_{\Omega} |\partial_{\alpha_d}^2 f|$ and the last equality is due to (H.7). Applying these bounds to (H.29) yields

$$(H.30) \quad I_d(y) \leq \frac{2d^2}{\kappa_f^2} \tilde{K} h^m \int_{G_{\tilde{A}_d}} d\alpha_1 \cdots d\alpha_{d-1} = K h^m$$

for some constants $\tilde{K}, K > 0$. Moreover, since (H.30) holds for $I_j(y)$ for all indices $1 \leq j \leq d$, then by (H.25)

$$I(y) \leq \sum_{j=1}^d I_j(y) + O(h^m) \leq dK h^m + O(h^m).$$

Step 6. Although $\|p_f - p_s\|_1 = \int_{-\infty}^{\infty} I(y) dy$, since Ω is compact and s and f are continuous,

$$Q_1 \leq s(\alpha), f(\alpha) \leq Q_2,$$

and so $I(y) = 0$ for $y \notin [Q_1, Q_2]$. Hence, by (H.30)

$$\|p_f - p_s\|_q = \left(\int_{-\infty}^{\infty} I^q(y) dy \right)^{\frac{1}{q}} = \left(\int_{Q_1}^{Q_2} I^q(y) dy \right)^{\frac{1}{q}} \leq (K^q h^{qm} (Q_2 - Q_1))^{\frac{1}{q}} \leq K (Q_2 - Q_1)^{\frac{1}{q}} h^m.$$

Acknowledgments. The authors thank Y. Harness, B. Brill, R. Kats, F. Abramovich, and D. Levin for useful comments and conversations.

REFERENCES

[1] M. J. ABLOWITZ AND T. P. HORIKIS, *Interacting nonlinear wave envelopes and rogue wave formation in deep water*, Phys. Fluids, 27 (2015), 012107.
 [2] G. P. AGRAWAL, *Nonlinear Fiber Optics*, 5th ed., Academic Press, New York, 2012.

- [3] R. K. BEATSON, *On the convergence of some cubic spline interpolation schemes*, SIAM J. Numer. Anal., 23 (1986), pp. 903–912.
- [4] R. K. BEATSON AND E. CHACKO, *Which cubic spline should one use?* SIAM J. Sci. Stat. Comput., 13 (1992), pp. 1009–1024.
- [5] G. BIRKHOFF AND C. DE BOOR, *Error bounds for spline interpolation*, J. Math. Mech., 13 (1964), pp. 827–836.
- [6] Q. Y. CHEN, D. GOTTLIEB, AND J. S. HESTHAVEN, *Uncertainty analysis for the steady-state flows in a dual throat nozzle*, J. Comput. Phys., 204 (2005), pp. 378–398.
- [7] I. COLOMBO, F. NOBILE, G. PORTA, A. SCOTTI, AND L. TAMELLINI, *Uncertainty quantification of geochemical and mechanical compaction in layered sedimentary basins*, Comput. Methods Appl. Mech. Engrg., 328 (2018), pp. 122–146.
- [8] P. G. CONSTANTINE, M. S. ELDRED, AND E. T. PHIPPS, *Sparse pseudospectral approximation method*, Comput. Methods Appl. Mech. Engrg., 229 (2012), pp. 1–12.
- [9] P. J. DAVIS AND P. RABINOWITZ, *Numerical Integration*, Academic Press, New York, 1975.
- [10] C. DE BOOR, *On cubic spline functions that vanish on all knots*, Adv. Math., 20 (1976), pp. 1–17.
- [11] C. DE BOOR, *A Practical Guide to Splines*, Springer, New York, 1978.
- [12] B. J. DEBUSSCHERE, H. N. NAJM, P. P. PÉBAY, O. M. KNIO, R. GHANEM, AND O. LE MAITRE, *Numerical challenges in the use of polynomial chaos representations for stochastic processes*, SIAM J. Sci. Comput., 26 (2004), pp. 698–719.
- [13] L. DEVROYE AND L. GYÖFRI, *Nonparametric Density Estimation—The L_1 View*, Wiley, New York, 1985.
- [14] A. DITKOWSKI AND R. KATS, *On spectral approximations with nonstandard weight functions and their implementations to generalized chaos expansions*, J. Sci. Comput., 79 (2019), pp. 1981–2005.
- [15] R. L. EUBANK, *Nonparametric Regression and Spline Smoothing*, Marcel Dekker, New York, 1999.
- [16] L. C. EVANS AND R. F. GARIEPY, *Measure Theory and Fine Properties of Functions*, CRC Press, Boca Raton, FL, 1991.
- [17] G. FIBICH, *The Nonlinear Schrödinger Equation*, Springer, New York, 2015.
- [18] R. J. FIELD AND R. M. NOYES, *Oscillations in chemical systems. IV. Limit cycle behavior in a model of a real chemical reaction*, J. Chem. Phys., 60 (1974), pp. 1877–1884.
- [19] B. FORNBERG AND N. FLYER, *A Primer on Radial Basis Functions with Applications to the Geosciences*, CBMS-NSE Regional Conf. Ser. in Appl. Math. 87, SIAM, Philadelphia, 2015.
- [20] F. N. FRITSCH AND R. E. CARLSON, *Monotone piecewise cubic interpolation*, SIAM J. Numer. Anal., 17 (1980), pp. 238–246.
- [21] B. GANAPATHYSUBRAMANIAN AND N. ZABARAS, *Sparse grid collocation schemes for stochastic natural convection problems*, J. Comput. Phys., 225 (2007), pp. 652–685.
- [22] R. GHANEM, D. HIGDON, AND H. OWHADI, *Handbook of Uncertainty Quantification*, Springer, New York, 2017.
- [23] R. GHANEM AND P. D. SPANOS, *Stochastic Finite Elements: A Spectral Approach*, Springer, New York, 1991.
- [24] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, John Hopkins University, Baltimore, MD, 2012.
- [25] C. A. HALL AND W. W. MEYER, *Optimal error bounds for cubic spline interpolation*, J. Approx. Theory, 16 (1976), pp. 105–122.
- [26] J. S. HESTHAVEN, S. GOTTLIEB, AND D. GOTTLIEB, *Spectral Methods for Time-Dependent Problem*, Cambridge Monogr. Appl. Comput. Math. 21, Cambridge University Press, Cambridge, UK, 2007.
- [27] S. C. HSU AND R. BHATTACHARYA, *Design of stochastic collocation based linear parameter varying quadratic regulator*, in Proceedings of the American Control Conference, IEEE, 2017, pp. 2375–2380.
- [28] S. KULLBACK AND R. A. LEIBLER, *On information and sufficiency*, Ann. Math. Stat., 22 (1951), pp. 79–86.
- [29] L. LE CAM AND G. L. YANG, *Asymptotics in Statistics: Some Basic Concepts*, Springer, New York, 2012.
- [30] O. LE MAÎTRE AND O. M. KNIO, *Spectral Methods for Uncertainty Quantification: With Applications to Computational Fluid Dynamics*, Springer, New York, 2010.
- [31] O. P. LE MAÎTRE, O. M. KNIO, H. N. NAJM, AND R. GHANEM, *Uncertainty propagation using Wiener-Haar expansions*, J. Comput. Phys., 197 (2004), pp. 28–57.

- [32] O. P. LE MAÎTRE, L. MATHELIN, O. M. KNIO, AND M. Y. HUSSAINI, *Asynchronous time integration for polynomial chaos expansion of uncertain periodic dynamics*, *Discrete Contin. Dyn. Syst.*, 28 (2010), pp. 199–226.
- [33] K. V. MARDIA AND P. E. JUPP, *Directional Statistics*, John Wiley & Sons, Chichester, UK, 2009.
- [34] H. N. NAJM, *Uncertainty quantification and polynomial chaos techniques in computational fluid dynamics*, *Annu. Rev. Fluid Mech.*, 41 (2009), pp. 35–52.
- [35] A. O’HAGAN, *Polynomial chaos: A tutorial and critique from a statistician’s perspective*, *SIAM/ASA J. Uncertain. Quantif.*, 20 (2013), pp. 1–20.
- [36] G. PATWARDHAN, X. GAO, A. SAGIV, A. DUTT, J. GINSBERG, A. DITKOWSKI, G. FIBICH, AND A. L. GAETA, *Loss of polarization of elliptically polarized collapsing beams*, *Phys. Rev. A*, 99 (2019), pp. 033824.
- [37] L. PIEGL AND W. TILLER, *The NURBS Book*, Springer, New York, 2012.
- [38] P. M. PRENTER, *Splines and Variational Methods*, Courier, New York, 2008.
- [39] S. RANDOUX, N. DALLOZ, AND P. SURET, *Intracavity changes in the field statistics of Raman fiber lasers*, *Opt. Lett.*, 36 (2011), pp. 790–792.
- [40] J. R. RICE, *Multivariate piecewise polynomial approximation*, in *Multivariate Approximation*, D.G. Handscorn, ed., Academic Press, New York, 1978.
- [41] A. SAGIV, A. DITKOWSKI, AND G. FIBICH, *Loss of phase and universality of stochastic interactions between laser beams*, *Opt. Exp.*, 25 (2017), pp. 24387–24399.
- [42] M. D. SALAS, S. ABARBANEL, AND D. GOTTLIEB, *Multiple steady states for characteristic initial value problems*, *Appl. Numer. Math.*, 2 (1986), pp. 193–210.
- [43] R. SCHABACK, *Error estimates and condition numbers for radial basis function interpolation*, *Adv. Comput. Math.* 3 (1995), pp. 251–264.
- [44] M. H. SCHULTZ, *L^∞ -Multivariate approximation theory*, *SIAM J. Numer. Anal.*, 6 (1969), pp. 161–183.
- [45] A. H. SHEINFUX, E. SCHLEIFER, J. PAPEER, G. FIBICH, B. ILAN, AND A. ZIGLER, *Measuring the stability of polarization orientation in high intensity laser filaments in air*, *Appl. Phys. Lett.*, 101 (2012), 201105.
- [46] G. STEFANO, *The stochastic finite element method: past, present and future*, *Comput. Methods Appl. Mech. Engrg.*, 198 (2009), pp. 1031–1051.
- [47] B. SUDRET AND A. DER KIUREGHIAN, *Stochastic Finite Element Methods and Reliability: A State-of-the-Art Report*, Department of Civil and Environmental Engineering, University of California Berkeley, Berkeley, CA, 2000.
- [48] G. SZEGO, *Orthogonal Polynomials*, Amer. Math. Soc. Colloq. Publ. 23, AMS, Providence, RI, 1939.
- [49] G. TEREJANU, P. SINGLA, T. SINGH, AND P. D. SCOTT, *Uncertainty propagation for nonlinear dynamic systems using Gaussian mixture models*, *J. Guid. Cont. Dyn.*, 31 (2008), pp. 1623–1633.
- [50] L. N. TREFETHEN, *Approximation Theory and Approximation Practice*, SIAM, Philadelphia, 2013.
- [51] A. B. TSYBAKOV, *Introduction to Nonparametric Estimation*, Springer, New York, 2009.
- [52] C. J. TURNER AND R. H. CRAWFORD, *N -dimensional nonuniform rational B -splines for metamodeling*, *J. Comput. Inf. Sci. Engrg.*, 9 (2009), 031002.
- [53] S. ULLMANN AND J. LANG, *POD-Galerkin modeling and sparse-grid collocation for a natural convection problem with stochastic boundary conditions*, in *Sparse Grids and Applications*, Munich, Lect. Notes Comput. Sci. Eng. 97, Springer, New York, 2012, pp. 295–315.
- [54] Y. VAN HALDER, B. SANDERSE, AND B. KOREN, *An Adaptive Minimum Spanning Tree Multi-Element Method for Uncertainty Quantification of Smooth and Discontinuous Responses*, <https://arxiv.org/abs/1803.06833>, 2018.
- [55] G. WAHBA, *Spline Models for Observational Data*, CBMS-NSF Regional Conf. Ser. in Appl. Math. 59, SIAM, Philadelphia, 1990.
- [56] X. WAN AND G. E. KARNIADAKIS, *An adaptive multi-element generalized polynomial chaos method for stochastic differential equations*, *J. Comput. Phys.*, 209 (2005), pp. 617–642.
- [57] H. WANG AND S. XIANG, *On the convergence rates of Legendre approximation*, *Math. Comp.*, 81 (2012), pp. 861–877.
- [58] L. WASSERMAN, *All of Statistics: A Concise Course in Statistical Inference*, Springer, New York, 2004.
- [59] D. XIU, *Numerical Methods for Stochastic Computations: A Spectral Method Approach*, Princeton University Press, Princeton, NJ, 2010.

- [60] D. XIU AND J. S. HESTHAVEN, *High-order collocation methods for differential equations with random inputs*, SIAM J. Sci. Comput., 27 (2005), pp. 1118–1139.
- [61] D. XIU AND G. E. KARNIADAKIS, *The Wiener–Askey polynomial chaos for stochastic differential equations*, SIAM J. Sci. Comput., 24 (2002), pp. 619–644.