

## FULLY SYMMETRIC KERNEL QUADRATURE\*

TONI KARVONEN<sup>†</sup> AND SIMO SÄRKKÄ<sup>†</sup>

**Abstract.** Kernel quadratures and other kernel-based approximation methods typically suffer from prohibitive cubic time and quadratic space complexity in the number of function evaluations. The problem arises because a system of linear equations needs to be solved. In this article we show that the weights of a kernel quadrature rule can be computed efficiently and exactly for up to tens of millions of nodes if the kernel, integration domain, and measure are fully symmetric and the node set is a union of fully symmetric sets. This is based on the observations that in such a setting there are only as many distinct weights as there are fully symmetric sets and that these weights can be solved from a linear system of equations constructed out of row sums of certain submatrices of the full kernel matrix. We present several numerical examples that show feasibility, both for a large number of nodes and in high dimensions, of the developed fully symmetric kernel quadrature rules. Most prominent of the fully symmetric kernel quadrature rules we propose are those that use sparse grids.

**Key words.** numerical integration, kernel quadrature, Bayesian quadrature, reproducing kernel Hilbert spaces, fully symmetric sets, sparse grids

**AMS subject classifications.** 46E22, 47B32, 60G15, 65C05, 65C50, 65D30, 65D32

**DOI.** 10.1137/17M1121779

**1. Introduction.** Let  $\Omega$  be a subset of  $\mathbb{R}^d$ ,  $\mu$  a measure on  $\Omega$ , and  $f: \Omega \rightarrow \mathbb{R}$  a function that is integrable with respect to  $\mu$ . Computation of the integral  $\mu(f) := \int_{\Omega} f \, d\mu$  is a recurring problem in applied mathematics and statistics. In most cases, this integral has no readily available analytical form, and one must resort to a quadrature rule (or, occasionally, a cubature rule if  $d > 1$ ) for its approximation. A quadrature rule  $Q$  is a linear functional of the form

$$Q(f) := \sum_{i=1}^n w_i f(\mathbf{x}_i) \approx \int_{\Omega} f \, d\mu,$$

where  $\mathbf{x}_i \in \mathbb{R}^d$  are the nodes and  $w_i \in \mathbb{R}$  are the weights. The nodes and weights are often chosen so that the quadrature approximation is exact whenever the integrand is a low-degree polynomial [12, 11]—such methods are called classical or polynomial quadrature rules in this article (we reserve the term Gaussian for rules that use  $n$  nodes to integrate polynomials up to degree  $2n - 1$  exactly). Another possibility is to use Monte Carlo or quasi Monte Carlo methods [8].

Here we study *kernel quadrature rules* that are, for arbitrary fixed nodes, optimal in the reproducing kernel Hilbert space (RKHS) induced by a user-specified positive-definite kernel. In this setting, optimality is measured in terms of the worst-case error (or, equivalently, the average-case error [55, 45]). Kernel quadrature rules go back at least to the work of Larkin [32, 33] in the 1970s. Lately, these rules have been the subject of renewed interest because they can be used for numerical integration on scattered data

---

\*Submitted to the journal's Methods and Algorithms for Scientific Computing section March 20, 2017; accepted for publication (in revised form) January 3, 2018; published electronically March 1, 2018.

<http://www.siam.org/journals/sisc/40-2/M112177.html>

**Funding:** This work was supported by Aalto ELEC Doctoral School as well as Academy of Finland projects 266940 and 273475.

<sup>†</sup>Department of Electrical Engineering and Automation, Aalto University, Espoo, Finland (toni.karvonen@aalto.fi, simo.sarkka@aalto.fi).

sets [3, 57], and they carry a probabilistic interpretation as posterior means for Gaussian processes assigned to the integrand [48]. The probabilistic interpretation, equivalent to the RKHS formulation we use, is interesting because it allows for statistical modeling of error in numerical integration and is one of main motivators behind this article. The above topics, including the probabilistic interpretation, are reviewed in section 2.

An advantage of kernel quadrature rules is that the nodes are *not* prescribed, as opposed to polynomial quadrature rules (polynomial rules with arbitrary nodes could probably be developed along the lines of de Boor and Ron interpolation [13, 39]). The flexibility comes with the price of having to solve the weights from a linear system of  $n$  equations, a task of cubic time and quadratic space complexity. It would not be practical to tabulate the weights beforehand as they depend on the kernel, integration domain, and measure. Partially due to the computational cost, only low numbers of nodes have been used in kernel quadrature, and much of the literature is concerned with optimal selection of the nodes. See [32, 48, 49, 37, 58] for some optimal node configurations, [46] for an algorithm to generate such nodes in one dimension, and [5, 4] for other nonoptimal alternatives. Efficient computation of the optimal nodes in higher dimensions is an open problem and not the topic of this article. Instead, we want to find nodes for which the weights can be computed easily and fast.

There is not much work on extending applicability of kernel quadrature to integration problems where it is necessary to use a large number of nodes due to high dimensionality or high level of desired accuracy. O'Hagan [48, 49] proposed some computationally beneficial product grid (number of nodes grows quickly in dimension and when the grid is refined) and simplex (only  $d + 1$  nodes) designs that are too inflexible to be of much use in many situations. The most exciting work is by Oettershagen [46], who has recently shown that the standard approach to sparse grid quadrature can be used to achieve quadratic time complexity. Furthermore, several fast and approximate kernel-based methods have been developed in the scattered data approximation, statistics, and machine learning literature (a compendium can be found in, e.g., [6, Supplement C]). However, the accuracy of quadrature rules is often strongly dependent on the weights having been computed exactly, and approximate weights also give rise to some philosophical objections if they are to be used for statistical modeling of error of an integral approximation.

In this article we show that if a certain structure is imposed on the node set, then the kernel quadrature weights can be computed *exactly* and in a very simple manner. Our approach is based on *fully symmetric sets* [18, 19] which are point sets that can be obtained from a given vector through permutations and sign changes of its coordinates. In section 3 we show that some symmetricity assumptions on  $\Omega$  and  $\mu$  (see Assumption 3.4) lead, for a large class of kernels that includes all isotropic kernels, to tractable computation of the weights if the node set is a union of fully symmetric sets. Depending on the dimension, the weights can be computed for sets of this type that contain up to tens of millions of nodes. The crucial observation under our assumptions is that there are only as many distinct weights as there are fully symmetric sets making up the node set. The *fully symmetric kernel quadratures* we construct exhibit the following advantageous properties:

- The algorithm (see section 3.4) for exact computation of the weights is exceedingly simple and easy to implement.
- If there are  $J$  fully symmetric sets containing  $n$  nodes in total, only  $Jn$  kernel evaluations are needed. In all situations we can envision,  $J$  is at most a few hundred, while  $n$  can, as mentioned, go up to millions (section 5.4 contains an example where  $J = 832$  and  $n = 15,005,761$ ). The weights are solved from

- a system of  $J$  linear equations, and only a  $J \times J$  matrix needs to be stored.
- The node selection scheme remains quite flexible, and the number of nodes does not grow too fast with the dimensions (see (3.2) and Table 3.1), as happens when, for example, full Cartesian grids are used. The smallest nontrivial fully symmetric sets contain  $2d$  points.

In section 4 we discuss a number of possible ways of selecting the fully symmetric sets. Sparse grids [7], popular in polynomial-based high-dimensional quadrature, are maybe the most obvious and promising choice. For kernel quadratures that use Clenshaw–Curtis sparse grids [41] we also provide some theoretical convergence guarantees in Theorem 4.2. Even though kernel quadrature rules on sparse grids can be efficiently constructed without the use of fully symmetric sets [46], our approach appears to be computationally competitive. In any case, sparse grids serve as a straightforward node selection scheme for showcasing that our algorithm indeed works.

The fast weight algorithm for computing the weights is not easily extended to fitting of the kernel parameters that often have considerable effect on accuracy of the integral approximation. Our experiments show that fully symmetric kernel quadratures are feasible, but we have to either resort to ad hoc solutions for fitting the kernel parameters or know them beforehand. Development of efficient fitting procedures is left for future research. This is discussed in section 5.1.

Finally, it is worth remarking that this article is not the first instance of fully symmetric sets being used in conjunction with kernel quadrature. Arguably the simplest nontrivial fully symmetric kernel quadrature rule (this rule appears briefly in section 4.3) has seen use in approximate filtering of nonlinear systems [58, 50, 51], but without an efficient weight computation algorithm.

**2. Kernel quadrature.** This section reviews the basics of quadrature rules in reproducing kernel Hilbert spaces. We also briefly discuss connections to probabilistic modeling of numerical algorithms. See [32, 6, 46] for proofs and additional references. Standard references on reproducing kernel Hilbert spaces are [1, 2].

**2.1. Quadrature in reproducing kernel Hilbert spaces.** A kernel  $k: \Omega \times \Omega \rightarrow \mathbb{R}$  is said to be positive-definite if the  $n \times n$  kernel Gram matrix  $[\mathbf{K}]_{ij} := k(\mathbf{x}_i, \mathbf{x}_j)$  is positive-definite for every  $n \geq 0$  and any distinct  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \Omega$ . Every continuous positive-definite kernel defines a unique reproducing kernel Hilbert space  $\mathcal{H}$  of functions  $f: \Omega \rightarrow \mathbb{R}$  through the properties (i)  $k(\cdot, \mathbf{x}) \in \mathcal{H}$  for every  $\mathbf{x} \in \Omega$ , and (ii) pointwise evaluations of any  $f \in \mathcal{H}$  can be represented in terms of inner product with the kernel as  $\langle f, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} = f(\mathbf{x})$ . The latter of these is called the reproducing property. The integral operator  $\mu$  and the quadrature rule  $Q$  are bounded linear functionals on  $\mathcal{H}$  under the nonrestrictive assumption  $\int_{\Omega} \sqrt{k(\mathbf{x}, \mathbf{x})} d\mu(\mathbf{x}) < \infty$ . The worst-case error (WCE)  $e(Q)$  of a quadrature rule  $Q$  is defined in terms of the dual norm

$$(2.1) \quad e(Q) := \|\mu - Q\|_{\mathcal{H}^*} = \sup_{\|f\|_{\mathcal{H}} \leq 1} |\mu(f) - Q(f)|,$$

which can be also written as  $e(Q) = \|\mu[k(\cdot, \mathbf{x})] - Q[k(\cdot, \mathbf{x})]\|_{\mathcal{H}}$ . Why this is a reasonable measure of error of the quadrature rule is apparent after an application of the reproducing property and the Cauchy–Schwarz inequality:

$$|\mu(f) - Q(f)| = |\langle f, \mu[k(\cdot, \mathbf{x})] - Q[k(\cdot, \mathbf{x})] \rangle_{\mathcal{H}}| \leq e(Q) \|f\|_{\mathcal{H}}$$

for  $f \in \mathcal{H}$ . That is, if the integrand belongs to the RKHS, convergence in the usual sense of diminishing integration error is implied by convergence to zero of the WCE.

The relationship between the kernel and its induced RKHS is further discussed in section 4.4, where we also provide two convergence theorems for the WCE.

The quadrature rule that, for arbitrary fixed distinct nodes  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , minimizes the WCE (2.1) is called the kernel quadrature rule and denoted by  $Q_k$ . This rule is unique and the optimal weights  $\mathbf{w} = (w_1, \dots, w_n)$  can be solved from

$$(2.2) \quad \underbrace{\begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}}_{=\mathbf{K}} \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = \underbrace{\begin{pmatrix} k_\mu(\mathbf{x}_1) \\ \vdots \\ k_\mu(\mathbf{x}_n) \end{pmatrix}}_{=k_\mu(\mathcal{X})},$$

where the kernel Gram matrix  $\mathbf{K}$  is positive-definite—and hence nonsingular—and  $k_\mu(\mathbf{x}) := \int_\Omega k(\mathbf{x}, \mathbf{x}') d\mu(\mathbf{x}')$  is the kernel mean, an object of much independent interest [38]. The kernel quadrature rule and its WCE are

$$(2.3) \quad Q_k(f) = \sum_{i=1}^n [\mathbf{K}^{-1} k_\mu(\mathcal{X})]_i f(\mathbf{x}_i) = \mathbf{y}^\top \mathbf{K}^{-1} k_\mu(\mathcal{X}),$$

$$e(Q_k)^2 = \int_\Omega \int_\Omega k(\mathbf{x}, \mathbf{x}') d\mu(\mathbf{x}) d\mu(\mathbf{x}') - k_\mu(\mathcal{X})^\top \mathbf{K}^{-1} k_\mu(\mathcal{X}) = \mu(k_\mu) - Q_k(k_\mu),$$

where  $\mathbf{y} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ .

An *optimal kernel quadrature rule* minimizes the WCE also over the node set (of fixed cardinality). Such rules cannot be constructed efficiently at the moment in dimensions larger than one. We discuss structurally constrained versions in section 4.3.

**2.2. Probabilistic interpretation.** The probabilistic interpretation of kernel quadrature as *Bayesian quadrature* is a part of the emergent field of *probabilistic numerical computing* [14, 49, 25, 9], the origins of which can be traced back at least to the work of Larkin [33]. This interpretation is a major motivator behind the present article.

In Bayesian quadrature, the integrand  $f$  is typically modeled as a Gaussian process [47, 53] (prompting the alternative term *Gaussian process quadrature*) with the covariance kernel  $k$ . With the node locations  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and function evaluations  $f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)$  considered the “data”  $\mathcal{D}$ , the posterior  $f | \mathcal{D}$  is a Gaussian process with the mean and covariance

$$\mathbb{E}[f(\mathbf{x}) | \mathcal{D}] = \mathbf{y}^\top \mathbf{K}^{-1} k(\mathcal{X}, \mathbf{x}),$$

$$\mathbb{C}[f(\mathbf{x}), f(\mathbf{x}') | \mathcal{D}] = k(\mathbf{x}, \mathbf{x}') - k(\mathbf{x}, \mathcal{X})^\top \mathbf{K}^{-1} k(\mathbf{x}, \mathcal{X}),$$

where  $[k(\mathbf{x}, \mathcal{X})]_i = k(\mathbf{x}, \mathbf{x}_i)$ . Because  $\mu$  is a linear operator, this induces the Gaussian posterior distribution  $\mu(f) | \mathcal{D}$  on the integral with the mean  $\mathbb{E}[\mu(f) | \mathcal{D}]$  and variance  $\mathbb{V}[\mu(f) | \mathcal{D}]$  that turn out to be precisely  $Q_k(f)$  and  $e(Q_k)^2$  from the preceding section. The WCE can therefore be interpreted as a measure of numerical uncertainty over the integral approximation and then exploited in, for instance, uncertainty quantification and allocation of limited computational resources in computational pipelines [9]. Clear expositions of this probabilistic viewpoint to numerical integration are [48, 37, 6], and the methodology is quite popular in machine learning (see, e.g., [52, 24]).

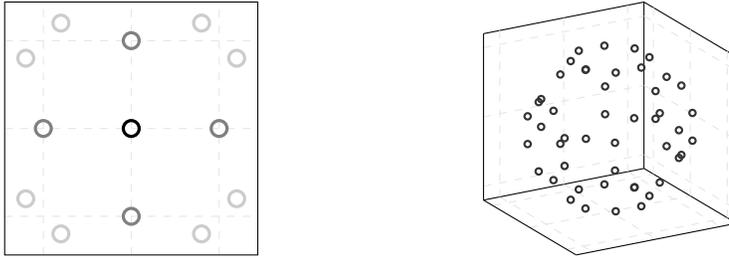


FIG. 3.1. Examples of fully symmetric sets in two and three dimensions. Left: the fully symmetric sets  $[0, 0]$ ,  $[1, 0]$ , and  $[1.2, 0.8]$  in  $\mathbb{R}^2$ . Right: the fully symmetric set  $[1, 0.5, 0.2]$  that consists of 48 elements in  $\mathbb{R}^3$ .

**3. Fully symmetric kernel quadrature.** This is the main section of the article. We introduce fully symmetric sets and their connection to multivariate quadrature rules, and we prove our main result, Theorem 3.6, on the computational benefits of doing kernel quadrature with node sets that are unions of fully symmetric sets.

**3.1. Fully symmetric sets.** A fully symmetric set is a point set in  $\mathbb{R}^d$  that is obtained from a given vector through permutations and sign changes of its coordinates. Let  $\Pi_d$  be the set of all permutations  $\mathbf{q} = (q_1, \dots, q_d)$  of the integers  $1, \dots, d$ , and let  $S_d$  be the set of all vectors of the form  $\mathbf{s} = (s_1, \dots, s_d)$  with each  $s_i$  either 1 or  $-1$ . Then, given  $d$  nonnegative scalars  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_d)$ , called *generators*, the point set

$$(3.1) \quad [\boldsymbol{\lambda}] = [\lambda_1, \dots, \lambda_d] := \bigcup_{\mathbf{q} \in \Pi_d} \bigcup_{\mathbf{s} \in S_d} \{(s_1 \lambda_{q_1}, \dots, s_d \lambda_{q_d})\} \subset \mathbb{R}^d$$

is the fully symmetric set generated by the *generator vector*  $\boldsymbol{\lambda}$ . With  $m$  the number of nonzero generators,  $m_0$  the number of zero generators (i.e.,  $m = d - m_0$ ), and  $m_1, \dots, m_l$  multiplicities of distinct nonzero generators so that  $\sum_{i=1}^l m_i = m$ , cardinality of the fully symmetric set (3.1) is

$$(3.2) \quad \#[\lambda_1, \dots, \lambda_d] = \frac{2^m d!}{m_0! \cdots m_l!}.$$

See Table 3.1 for cardinalities of a number of fully symmetric sets.

An alternative way of writing (3.1) is via *permutation matrices* as  $[\boldsymbol{\lambda}] = \bigcup_{\mathbf{P}} \mathbf{P}\boldsymbol{\lambda}$ , where the union is over all  $d \times d$  permutation and sign change matrices  $\mathbf{P}$ . These are matrices that have on each row and column exactly one element that is either 1 or  $-1$  and the rest are zero. Any element of a fully symmetric set can be obtained from any other via linear transformation by an appropriate permutation matrix. How (3.1) works and what the resulting point sets look like is illustrated in two and three dimensions in Example 3.2 and Figure 3.1. Note that all elements of a fully symmetric set are equidistant from the origin, which is to say that if  $\mathbf{x} \in [\lambda_1, \dots, \lambda_d]$ , then  $\|\mathbf{x}\|^2 = \|\boldsymbol{\lambda}\|^2 = \lambda_1^2 + \dots + \lambda_d^2$ . We also need the concept of a fully symmetric function.

**DEFINITION 3.1.** A function  $f: \Omega \rightarrow \mathbb{R}$  is fully symmetric if it is constant in every fully symmetric set. That is, with  $\boldsymbol{\lambda}$  any generator vector, it holds that  $f(\mathbf{x}) = f(\mathbf{x}')$  for any  $\mathbf{x}, \mathbf{x}' \in \Omega \cap [\boldsymbol{\lambda}]$ . Alternatively,  $f(\mathbf{P}\mathbf{x}) = f(\mathbf{x})$  for any  $\mathbf{x} \in \Omega$  and any permutation and sign change matrix  $\mathbf{P}$  such that  $\mathbf{P}\mathbf{x} \in \Omega$ .

TABLE 3.1

Cardinalities, as computed from (3.2), of fully symmetric sets generated by  $m = 1, \dots, 9$  distinct nonzero generators for dimensions  $d = 2, \dots, 9$ .

		Dimension							
$m$	2	3	4	5	6	7	8	9	
1	4	6	8	10	12	14	16	18	
2	8	24	48	80	120	168	224	288	
3		48	192	480	960	1,680	2,688	4,032	
4			384	1,920	5,760	13,440	26,880	48,384	
5				3,840	23,040	80,640	215,040	483,840	
6					46,080	322,560	1,290,240	3,870,720	
7						645,120	5,160,960	23,224,320	
8							10,321,920	92,897,280	
9								185,794,560	

Example 3.2. In  $\mathbb{R}^3$ , the nonzero and distinct generators  $\lambda_1$  and  $\lambda_2$  generate the fully symmetric set

$$[\lambda_1, \lambda_2, 0] = \{(\lambda_1, \lambda_2, 0), (-\lambda_1, \lambda_2, 0), (\lambda_1, -\lambda_2, 0), (-\lambda_1, -\lambda_2, 0), (\lambda_2, \lambda_1, 0), (-\lambda_2, \lambda_1, 0), (\lambda_2, -\lambda_1, 0), (-\lambda_2, -\lambda_1, 0), (0, \lambda_1, \lambda_2), (0, -\lambda_1, \lambda_2), (0, \lambda_1, -\lambda_2), (0, -\lambda_1, -\lambda_2), (0, \lambda_2, \lambda_1), (0, -\lambda_2, \lambda_1), (0, \lambda_2, -\lambda_1), (0, -\lambda_2, -\lambda_1), (\lambda_1, 0, \lambda_2), (-\lambda_1, 0, \lambda_2), (\lambda_1, 0, -\lambda_2), (-\lambda_1, 0, -\lambda_2), (\lambda_2, 0, \lambda_1), (-\lambda_2, 0, \lambda_1), (\lambda_2, 0, -\lambda_1), (-\lambda_2, 0, -\lambda_1)\},$$

which has  $2^2 \times 3! / (1! \times 1! \times 1!) = 24$  elements. In terms of permutation matrices, the element  $(-\lambda_1, 0, \lambda_2)$  is

$$\begin{pmatrix} -\lambda_1 \\ 0 \\ \lambda_2 \end{pmatrix} = \begin{pmatrix} -1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ 0 \end{pmatrix}.$$

The method we have used to generate fully symmetric sets out of user-specified generator vectors is detailed in Algorithm 1 in section 3.4. There are many other possibilities; we do not claim that the one presented is the optimal implementation.

**3.2. Fully symmetric quadrature rules.** The notation  $f[\boldsymbol{\lambda}] = f[\lambda_1, \dots, \lambda_d]$  stands for the sum of evaluations of  $f$  at the points of the fully symmetric set:

$$f[\boldsymbol{\lambda}] := \sum_{\mathbf{x} \in [\boldsymbol{\lambda}]} f(\mathbf{x}).$$

A fully symmetric quadrature rule is a quadrature rule of the form

$$Q(f) = \sum_{\boldsymbol{\lambda} \in \Lambda} w_{\boldsymbol{\lambda}} f[\boldsymbol{\lambda}] = \sum_{\boldsymbol{\lambda} \in \Lambda} w_{\boldsymbol{\lambda}} \sum_{\mathbf{x} \in [\boldsymbol{\lambda}]} f(\mathbf{x}),$$

where  $\Lambda$  is a given finite collection of distinct generator vectors  $\boldsymbol{\lambda}$ . Such a rule uses only  $\#\Lambda$  distinct weights, each corresponding to often a very large number of nodes. In Theorem 3.6 we establish conditions under which a kernel quadrature rule is fully symmetric. This will yield significant computational savings because only  $\#\Lambda$  (instead of  $n = \sum_{\boldsymbol{\lambda} \in \Lambda} \#[\boldsymbol{\lambda}]$ ) distinct weights need to be computed.

Fully symmetric quadrature rules are prominent among classical polynomial quadrature rules, with work on them going back to [35, 36]. To the best of our knowledge, the most general and efficient constructions have been given by Genz [18] for the uniform distribution on a square and by Genz and Keister [19] for Gaussians on the whole real space (a case studied also in [34]). See, for example, the review [11] for more examples and discussion on the highly related invariant theory. To achieve high algebraic order of precision, the classical fully symmetric quadrature rules rely on symmetry of the underlying measure and advantageous properties of polynomials when integrated with respect to such measures. In contrast to the kernel quadrature rules we are about to construct, the aforementioned rules do not permit free selection of the fully symmetric sets that are to be used.

Many of the popular sparse grid rules are also fully symmetric [44, 43]. We exploit this useful fact in section 4.2 for construction of sparse grid kernel quadrature rules whose weights can be computed efficiently.

**3.3. Fully symmetric kernels.** We can now introduce the class of kernels that this article is concerned with as well as the necessary assumptions on the integration domain and measure.

**DEFINITION 3.3.** *Suppose that  $\mathbf{P}\mathbf{x} \in \Omega$  for any  $\mathbf{x} \in \Omega$  and any permutation and sign change matrix  $\mathbf{P}$ . A kernel  $k$  is fully symmetric if  $k(\mathbf{P}\mathbf{x}, \mathbf{P}\mathbf{x}') = k(\mathbf{x}, \mathbf{x}')$  for any such matrix  $\mathbf{P}$  and any  $\mathbf{x}, \mathbf{x}' \in \Omega$ .*

This class of kernels includes (i) isotropic kernels, (ii) products of an isotropic kernel  $k_1$  of the form  $k(\mathbf{x}, \mathbf{x}') = \prod_{i=1}^d k_1(|x_i - x'_i|)$ , and (iii) sums of an isotropic kernel  $k_1$  of the form  $k(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^d k_1(|x_i - x'_i|)$ . Some polynomial kernels<sup>1</sup> of the form  $k(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^p P_i(\mathbf{x})P_i(\mathbf{x}')$  for suitable multivariate polynomials  $P_i$  are also fully symmetric. For example, the selection  $P_1 \equiv 1$  and  $P_i = x_{i-1}^2$  for  $i = 2, \dots, p = d + 1$  results in a fully symmetric kernel. See [58, 30] for some results on how quadrature rules for such kernels are related to classical quadrature rules.

*Assumption 3.4.* We assume the following:

- (i) The integration domain  $\Omega \subset \mathbb{R}^d$  is invariant under permutations and sign changes of coordinates of its elements. That is,  $\Omega = \mathbf{P}\Omega = \{\mathbf{P}\omega : \omega \in \Omega\}$  for any permutation and sign change matrix  $\mathbf{P}$ .
- (ii) The measure  $\mu$  is fully symmetric in the sense that its density  $f_\mu$  (with respect to the Lebesgue measure) is a fully symmetric function.
- (iii) The kernel  $k$  is positive-definite and fully symmetric.

This assumption holds, for example, for  $\Omega = [-1, 1]^d$  equipped with the uniform measure and  $\Omega = \mathbb{R}^d$  equipped with the Gaussian measure, as well as for many other cases of interest. The numerical examples in section 5 are for these two cases.

**LEMMA 3.5.** *The kernel mean  $k_\mu$  is fully symmetric under Assumption 3.4.*

*Proof.* Let  $\mathbf{x}$  and  $\mathbf{x}'$  be elements of the same fully symmetric set. That is,  $\mathbf{x}' = \mathbf{P}\mathbf{x}$  for some permutation and sign change matrix  $\mathbf{P}$ . A change of variables yields

$$\int_{\Omega} k(\mathbf{x}', \mathbf{z})f_\mu(\mathbf{z}) \, d\mathbf{z} = \int_{\Omega} |\det \mathbf{P}| k(\mathbf{P}\mathbf{x}, \mathbf{P}\mathbf{z})f_\mu(\mathbf{P}\mathbf{z}) \, d\mathbf{z} = \int_{\Omega} k(\mathbf{x}, \mathbf{z})f_\mu(\mathbf{z}) \, d\mathbf{z},$$

where we have used the fact that  $|\det \mathbf{P}| = 1$ . That is,  $k_\mu(\mathbf{x}') = k_\mu(\mathbf{x})$ . □

<sup>1</sup>Strictly speaking, these kernels are not positive-definite, as the kernel matrix does not remain nonsingular for any number of distinct points. See also Remark 3.7.

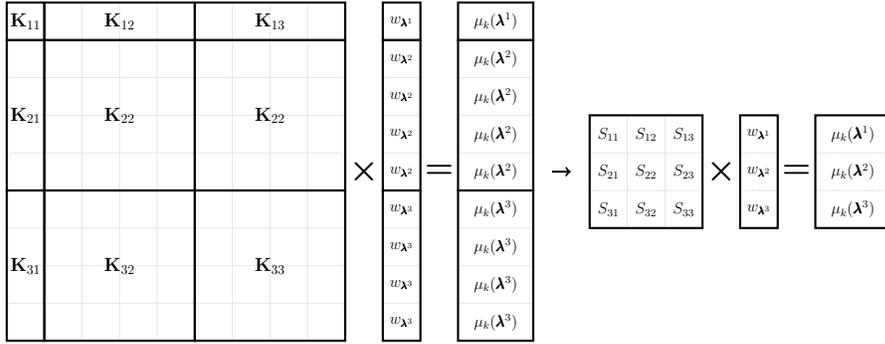


FIG. 3.2. Illustration of Theorem 3.6 for a node set that is a union of three fully symmetric sets: one containing one element and two containing four elements. All row sums of the matrices  $K_{ij}$ , defined in (3.4), are equal to  $S_{ij}$ .

**3.4. Fully symmetric kernel quadrature.** Let  $[\lambda^1], \dots, [\lambda^J]$  be distinct fully symmetric sets generated by  $\lambda^1, \dots, \lambda^J \in \mathbb{R}^d$ . If the node set  $\mathcal{X}$  is the union  $\mathcal{X} = \cup_{j=1}^J [\lambda^j]$  of these fully symmetric sets, then a kernel quadrature rule using this node set is a fully symmetric quadrature rule in the sense of section 3.2. Furthermore, its  $J$  distinct weights can be computed extremely efficiently when compared to naively solving the linear system (2.2) of  $\#\mathcal{X} = n$  equations. This is formalized in the following theorem. Figure 3.2 illustrates the simplified weight computation process in the case of a node set that is a union of three fully symmetric sets.

**THEOREM 3.6.** *Suppose that  $\Omega$ ,  $\mu$ , and  $k$  satisfy Assumption 3.4. If the node set  $\mathcal{X}$  is a union of  $J$  distinct fully symmetric sets  $[\lambda^1], \dots, [\lambda^J]$ , then the kernel quadrature rule  $Q_k$  is fully symmetric:*

$$Q_k(f) = \sum_{j=1}^J w_{\lambda^j} f[\lambda^j].$$

Furthermore, the  $J$  weights  $w_{\lambda^1}, \dots, w_{\lambda^J}$  corresponding to the fully symmetric sets can be solved from the nonsingular linear system of  $J$  equations

$$(3.3) \quad \begin{pmatrix} S_{11} & \cdots & S_{1J} \\ \vdots & \ddots & \vdots \\ S_{J1} & \cdots & S_{JJ} \end{pmatrix} \begin{pmatrix} w_{\lambda^1} \\ \vdots \\ w_{\lambda^J} \end{pmatrix} = \begin{pmatrix} k_\mu(\lambda^1) \\ \vdots \\ k_\mu(\lambda^J) \end{pmatrix},$$

where

$$S_{ij} = \sum_{\mathbf{x} \in [\lambda^j]} k(\mathbf{x}^i, \mathbf{x}) \quad \text{for any } \mathbf{x}^i \in [\lambda^i].$$

*Proof.* Let  $\mathcal{X} = \cup_{j=1}^J [\lambda^j]$  be ordered such that all elements of a single fully symmetric set appear consecutively, and the fully symmetric sets themselves are in ascending order in terms of their index  $j$ .

We denote  $n^i = \#[\lambda^i]$  and enumerate each fully symmetric set as  $[\lambda^i] = \{\mathbf{x}_1^i, \dots, \mathbf{x}_{n^i}^i\}$ . By Lemma 3.5, the kernel mean is fully symmetric and, consequently, the kernel mean vector  $k_\mu(\mathcal{X}) \in \mathbb{R}^n$ ,  $n = n^1 + \dots + n^J$ , is

$$k_\mu(\mathcal{X}) = (k_\mu([\lambda^1]), \dots, k_\mu([\lambda^J])),$$

where  $k_\mu([\lambda^j]) = (k_\mu(\lambda^j), \dots, k_\mu(\lambda^j)) \in \mathbb{R}^{n^j}$ . That is,  $k_\mu(\mathcal{X})$  contains only  $J$  distinct elements that occur in blocks of  $n^j$ . Consider then the kernel matrix  $\mathbf{K}$  that can be partitioned into  $J^2$  submatrices  $\mathbf{K}_{ij}$  of dimensions  $n^i \times n^j$ , each containing all the kernel evaluations  $k(\mathbf{x}^i, \mathbf{x}^j)$  for  $\mathbf{x}^i \in [\lambda^i]$  and  $\mathbf{x}^j \in [\lambda^j]$ :

$$(3.4) \quad \mathbf{K} = \begin{pmatrix} \mathbf{K}_{11} & \cdots & \mathbf{K}_{1J} \\ \vdots & \ddots & \vdots \\ \mathbf{K}_{J1} & \cdots & \mathbf{K}_{JJ} \end{pmatrix}, \text{ where } \mathbf{K}_{ij} = \begin{pmatrix} k(\mathbf{x}_1^i, \mathbf{x}_1^j) & \cdots & k(\mathbf{x}_1^i, \mathbf{x}_{n^j}^j) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_{n^i}^i, \mathbf{x}_1^j) & \cdots & k(\mathbf{x}_{n^i}^i, \mathbf{x}_{n^j}^j) \end{pmatrix}.$$

Any row of any submatrix  $\mathbf{K}_{ij}$  can be obtained from any other of its rows by a permutation of elements of the row. To confirm this, consider any distinct rows  $p, p' \leq n^i$  of  $\mathbf{K}_{ij}$ . There exists a permutation and sign change matrix  $\mathbf{P}$  such that  $\mathbf{x}_{p'}^i = \mathbf{P}\mathbf{x}_p^i$  because fully symmetric sets are closed under such transformations. Note that  $\mathbf{P}$  is nonsingular and its inverse  $\mathbf{P}^{-1}$  is also a permutation and sign change matrix. Since the kernel is fully symmetric, for any  $l \leq n^j$  we have

$$k(\mathbf{x}_p^i, \mathbf{x}_l^j) = k(\mathbf{P}^{-1}\mathbf{x}_p^i, \mathbf{P}^{-1}\mathbf{x}_l^j) = k(\mathbf{x}_{p'}^i, \mathbf{P}^{-1}\mathbf{x}_l^j),$$

where  $\mathbf{P}^{-1}\mathbf{x}_l^j \in [\lambda^j]$ . This means that for every  $l$  there is an element on the row  $p'$  that equals the  $l$ th element of the  $p$ th row. That is, the rows are permutations of each other. Consequently, the row sums  $S_{ij} := \sum_{\mathbf{x} \in [\lambda^j]} k(\mathbf{x}^i, \mathbf{x})$  of  $\mathbf{K}_{ij}$  do not depend on  $\mathbf{x}^i \in [\lambda^i]$ .

Consider the  $J \times J$  matrix  $[\mathbf{S}]_{ij} = S_{ij}$  composed of the row sums defined above. Then the equation  $\mathbf{S}\mathbf{a} = \mathbf{b}$  for some vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^J$  implies that

$$\begin{pmatrix} \mathbf{K}_{11} & \cdots & \mathbf{K}_{1J} \\ \vdots & \ddots & \vdots \\ \mathbf{K}_{J1} & \cdots & \mathbf{K}_{JJ} \end{pmatrix} \begin{pmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_J \end{pmatrix} = \begin{pmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_J \end{pmatrix},$$

where  $\mathbf{a}_i = (a_i, \dots, a_i) \in \mathbb{R}^{n^i}$  and  $\mathbf{b}_i = (b_i, \dots, b_i) \in \mathbb{R}^{n^i}$ , because

$$b_i = \sum_{j=1}^J a_j S_{ij} = \sum_{j=1}^J a_j \sum_{l=1}^{n^j} [\mathbf{K}_{ij}]_{pl} = \sum_{j=1}^J a_j \sum_{l=1}^{n^j} k(\mathbf{x}_p^i, \mathbf{x}_l^j)$$

for every  $p \leq n^i$ . The matrix  $\mathbf{S}$  is nonsingular, for if it were singular there would exist a nonzero vector  $\mathbf{a} \in \mathbb{R}^J$  such that  $\mathbf{S}\mathbf{a} = \mathbf{0}$ . But by the above argument this would imply that  $\mathbf{K}$  is singular, which is not the case, because the kernel  $k$  is positive-definite. All this implies that if  $(w_{\lambda^1}, \dots, w_{\lambda^J})$  is the unique solution to the linear system of equations (3.3), then

$$\mathbf{w} = (\mathbf{w}_{\lambda^1}, \dots, \mathbf{w}_{\lambda^J}) \in \mathbb{R}^n, \quad \text{where } \mathbf{w}_{\lambda^j} = (w_{\lambda^j}, \dots, w_{\lambda^j}) \in \mathbb{R}^{n^j},$$

must be the solution to  $\mathbf{K}\mathbf{w} = k_\mu(\mathcal{X})$ . That is, weights for nodes in each fully symmetric set are equal and the kernel quadrature rule is fully symmetric. This concludes the proof.  $\square$

*Remark 3.7.* Theorem 3.6 also applies to kernels whose kernel matrix is positive-definite only for every collection of  $m \leq p$  distinct points for some  $p > 0$  if the total number  $n$  of nodes does not exceed  $p$ . The polynomial kernels briefly mentioned in section 3.3 are examples of such kernels.

The full algorithm for fully symmetric kernel quadrature is presented in high-level pseudocode in Algorithm 1 below. We expect that  $J$ , the number of fully symmetric sets, is rarely more than a few hundred (the example in section 5.4 has  $J = 832$ , but this results in  $n \approx 15,000,000$ ), so solving the weights from the linear system (3.3) is not a computational bottleneck. Instead, it is usually the  $Jn$  kernel evaluations that take the most time.

---

**Algorithm 1.** Fully symmetric kernel quadrature.

---

*Construct the fully symmetric sets*

1. Select  $J$  distinct generator vectors  $\boldsymbol{\lambda}^j \in \mathbb{R}^d$  with nonnegative elements.

**For each**  $j = 1, \dots, J$  construct the fully symmetric set  $[\boldsymbol{\lambda}^j]$ :

2. Sort  $\boldsymbol{\lambda}^j$  in descending order.
3. Identify the unique nonzero elements  $\mathbf{u} \in \mathbb{R}^{d_{\mathbf{u}}}$ ,  $d_{\mathbf{u}} \leq d$ , of  $\boldsymbol{\lambda}^j$  and their multiplicities  $\mathbf{m} \in \mathbb{N}^{d_{\mathbf{u}}}$ . Denote  $\Sigma_{\mathbf{m}} = \sum_{l=1}^{d_{\mathbf{u}}} m_l$ . That is,

$$\boldsymbol{\lambda}^j = (\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_{d_{\mathbf{u}}}, \mathbf{0}_{(d-\Sigma_{\mathbf{m}}) \times 1}) \in \mathbb{R}^d,$$

where  $\tilde{\mathbf{u}}_l = (u_l, \dots, u_l) \in \mathbb{R}^{m_l}$  for  $l = 1, \dots, d_{\mathbf{u}}$ .

4. Construct all  $d_{\mathbf{a}}$  possible vectors  $\mathbf{a}^i \in \mathbb{N}^{d_{\mathbf{u}}}$  such that  $a_l^i \leq m_l$  for each  $l = 1, \dots, d_{\mathbf{u}}$ .
5. Set  $[\boldsymbol{\lambda}^j] = \emptyset$ .

**For each**  $i = 1, \dots, d_{\mathbf{a}}$ :

6. Construct the vector

$$\boldsymbol{\lambda}_i^j = (\tilde{\mathbf{u}}_1^s, \dots, \tilde{\mathbf{u}}_{d_{\mathbf{u}}}^s, \mathbf{0}_{(d-\Sigma_{\mathbf{m}}) \times 1}) \in \mathbb{R}^d,$$

where, for each  $l = 1, \dots, d_{\mathbf{u}}$ , the first  $a_l^i$  elements of  $\tilde{\mathbf{u}}_l^s \in \mathbb{R}^{m_l}$  are  $-u_l$  and the rest are  $u_l$ . The vector  $\boldsymbol{\lambda}_i^j$  essentially corresponds to one possible sign combination  $(s_1 \lambda_1^j, \dots, s_d \lambda_d^j)$  in (3.1).

7. Compute the collection  $U$  of all unique permutations of  $\boldsymbol{\lambda}_i^j$  and append it to the fully symmetric set:  $[\boldsymbol{\lambda}^j] = [\boldsymbol{\lambda}^j] \cup U$ .

*Compute the kernel quadrature weights*

8. Construct an empty matrix  $\mathbf{S} \in \mathbb{R}^{J \times J}$ .

**For each**  $(i, j) \in \{1, \dots, J\}^2$ :

9. Select any  $\mathbf{x}^i \in [\boldsymbol{\lambda}^i]$  and set  $[\mathbf{S}]_{ij} = \sum_{\mathbf{x} \in [\boldsymbol{\lambda}^j]} k(\mathbf{x}^i, \mathbf{x})$ .

10. Solve the  $J$  distinct weights  $\mathbf{w}_{\boldsymbol{\lambda}} = (w_{\boldsymbol{\lambda}^1}, \dots, w_{\boldsymbol{\lambda}^J})$  from the linear system of equations  $\mathbf{S} \mathbf{w}_{\boldsymbol{\lambda}} = \mathbf{b}$ , where  $b_l = k_{\mu}(\boldsymbol{\lambda}^l)$ .

*Compute the quadrature approximation*

11. Compute  $Q_k(f) \approx \int_{\Omega} f \, d\mu$  as

$$Q_k(f) = \sum_{j=1}^J w_{\boldsymbol{\lambda}^j} f[\boldsymbol{\lambda}^j] = \sum_{j=1}^J w_{\boldsymbol{\lambda}^j} \sum_{\mathbf{x} \in [\boldsymbol{\lambda}^j]} f(\mathbf{x}).$$


---

**4. Selection of the fully symmetric sets.** This section presents three different approaches for constructing the node set as a union of fully symmetric sets. Of these the sparse grids of section 4.2 are the most promising alternative. We also discuss convergence properties of some of the kernel quadrature rules we construct in section 4.4. We expect there to exist many other competitive schemes, as one of the main advantages

of fully symmetric kernel quadrature is that there are no restrictions in selecting the generator vectors.

**4.1. Random generators.** Arguably, the simplest approach, both conceptually and algorithmically, is to draw a number of generator vectors randomly from the underlying distribution. However, unless additional constraints are enforced, all the generators will be distinct and nonzero, resulting in unrealistic numbers of integrand evaluations needed if  $d \geq 6$ , as seen from Table 3.1. One could heuristically set some generators to zero to reduce the number of nodes, but it is not entirely clear how this should be done. In any case, for at least  $d < 6$ , the random generator approach seems realistic. We call this method the *fully symmetric kernel Monte Carlo method* (FSKMC).

Theorem 4.1 provides theoretical convergence guarantees, and section 5.3 demonstrates that FSKMC can also work in practice. Nevertheless, the method does not seem very promising, as it comes across that a large number of random generator vectors—and thus an even larger number of nodes—is required to capture the underlying distribution.

This approach bears some similarity to the stochastic radial and spherical integration rules developed in [20, 21]. These rules are less flexible due to the usual constraints of integrating low-degree polynomials exactly and are more involved in their implementation.

**4.2. Sparse grids.** An iterated quadrature rule of degree  $m$  based on a regular Cartesian product grid requires  $m^d$  nodes—a number that quickly becomes impractically large. Sparse grids that originate in the work of Smolyak [56] are “sparsified” product sets widely used in numerical integration [41, 23, 43, 44]. See also the general survey by Bungartz and Griebel [7] and [26] for a wealth of financial applications. Recently, Oettershagen [46] has shown that the standard approach to sparse grids is also applicable to fast computation of the weights of kernel quadrature rules. This approach is different from ours, which is based on identifying the fully symmetric sets a sparse grid is a union of, and is specific to sparse grids. Other sparse grid-based kernel methods appear in [17, 22, 15, 59]. The construction of sparse grids that we present in this section is not the most general possible, since we work in the fully symmetric framework. More general constructions are contained in some of the aforementioned references. We assume that  $\Omega = [-a, a]^d$  for a possibly infinite  $a > 0$ .

Let  $X^1 = \{0\}$  and  $X^i \subset X^{i+1} \subset [-a, a]$  for  $i > 1$  be finite, nested, and symmetric (i.e., if  $x \in X^i$ , then  $-x \in X^i$ ) point sets. Then the *sparse grid* of level  $q \geq 1$  is the set

$$H(q, d) := \bigcup_{|\alpha|=d+q} (X^{\alpha_1} \times \dots \times X^{\alpha_d}),$$

where  $\alpha \in \mathbb{N}^d$  is a  $d$ -dimensional multi-index with the elements  $\alpha_i = \alpha(i)$  and  $|\alpha| = \alpha_1 + \dots + \alpha_d$ . Note that the largest  $X^i$  that is needed for a sparse grid of level  $q$  is  $X^{q+1}$ . As the basis sets  $X^i$  are nested and symmetric, it is fairly easy to see that  $H(q, d)$  is the union of fully symmetric sets and can be explicitly written thus:

$$\begin{aligned} H(q, d) &= \bigcup_{\substack{|\alpha|=d+q \\ \alpha_i \geq \alpha_{i+1}}} \bigcup_{\mathbf{q} \in \Pi_d} (X^{\alpha(q_1)} \times \dots \times X^{\alpha(q_d)}) \\ &= \bigcup_{\substack{|\alpha|=d+q \\ \alpha_i \geq \alpha_{i+1}}} \bigcup_{\mathbf{q} \in \Pi_d} \bigcup_{\mathbf{s} \in S_d} \bigcup_{\substack{\boldsymbol{\lambda} \in \mathbb{R}^d \\ \lambda_j \in X^{\alpha(q_j)} \\ \lambda_j \geq 0}} \{(s_1 \lambda_1, \dots, s_d \lambda_d)\} \end{aligned}$$

$$= \bigcup_{\substack{|\alpha|=d+q \\ \alpha_i \geq \alpha_{i+1}}} \{[\lambda_1, \dots, \lambda_d] : \lambda_j \in X^{\alpha_j} \text{ and } \lambda_j \geq 0 \text{ for } j = 1, \dots, d\},$$

where the restriction  $\alpha_1 \geq \alpha_{i+1}$  eliminates a large number of permutations that would be otherwise duplicated when generating fully symmetric sets. That is, Theorem 3.6 applies to sparse grids. We call the resulting kernel quadrature rules the *sparse grid kernel quadrature rules*.

We are left with selection of the nested point sets  $X^i$ . In polynomial-based sparse grid quadrature rules these sets come coupled with univariate quadrature rules whose weights are used to construct the final sparse grid weights, but we are under no such restrictions. An obvious idea for selecting  $X^i$  would be to sequentially minimize the worst-case error in one dimension—provided that the kernel is one for which this makes sense, for example, any of the three examples of fully symmetric kernels given in section 3.3. Discussion on different sequential kernel quadrature methods (usually known as *sequential Bayesian quadratures*) can be found in, for example, [10, 27, 24, 5]. A different approach appears in [46].

However, owing to difficulties in setting the kernel length-scale, we do not employ this selection scheme. Instead, we use (i) the Clenshaw–Curtis point sets, rather standard in sparse grid literature, for  $\Omega = [-1, 1]^d$  and the uniform measure; and (ii) nested sets formed out of Gauss–Hermite nodes in the Gaussian case:

- (i) For  $i > 1$  and  $m_i = 2^{i-1} + 1$ , the nested Clenshaw–Curtis sets are

$$X^i = \{x_1^i, \dots, x_{m_i}^i\} \quad \text{with} \quad x_j^i = -\cos\left(\frac{\pi(j-1)}{m_i-1}\right) \in [-1, 1].$$

The points  $x_j^i$  are the roots and extrema of Chebyshev polynomials. The corresponding sparse grid kernel quadrature rule is called *Clenshaw–Curtis sparse grid kernel quadrature* (CCSGKQ). Numerical results for this kernel quadrature are given in section 5.4, and convergence for sufficiently smooth functions is the topic of Theorem 4.2.

- (ii) In the *Gauss–Hermite sparse grid kernel quadrature* (GHSGKQ) we use the classical Gauss–Hermite nodes that are the roots of the Hermite polynomials

$$H_p(x) = (-1)^p \exp(x^2/2) \frac{d^p}{dx^p} \exp(-x^2/2).$$

Given a level  $q$ , we generate the  $2q + 1$  symmetric roots of  $H_{2q+1}$  and for  $i = 1, \dots, q + 1$  select

$$X^i = \text{the } 2i - 1 \text{ smallest roots by absolute value.}$$

The number of nodes, in terms of the level  $q$ , grows significantly slower than with Clenshaw–Curtis sparse grids. A numerical experiment involving a financial problem is given in section 5.5. As is usual for quadrature rules on the whole of  $\mathbb{R}^d$ , there are no theoretical convergence guarantees for GHSGKQ. Because  $H(q, d)$  is not a subset of  $H(q + 1, d)$  for the Gauss–Hermite points, these grids are only suitable for cases where the number of nodes that can be used is determined beforehand based on, for example, the computational budget available. Note that these grids are completely different from several other sparse grids in the literature that use nodes of Gaussian quadrature rules.

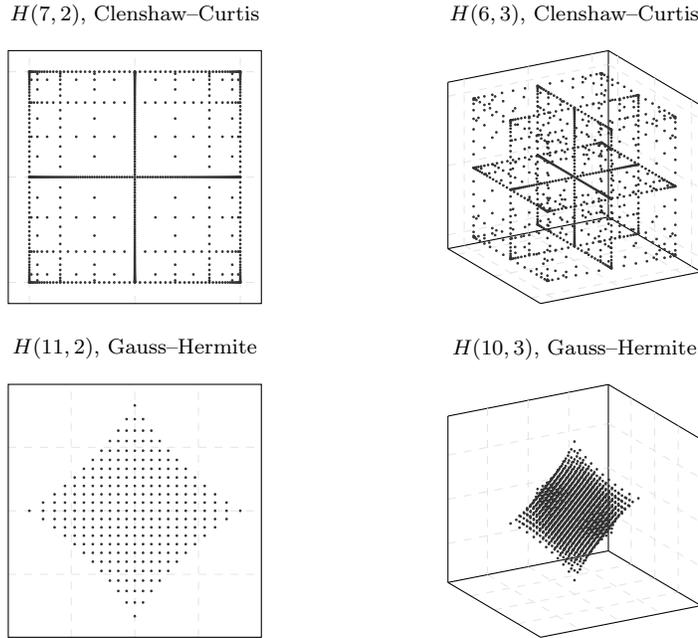


FIG. 4.1. Examples of sparse grids with Clenshaw–Curtis and Gauss–Hermite nodes. The numbers of nodes in the sparse grids are 705 (upper left), 1,073 (upper right), 265 (lower left), and 1,561 (lower right). Compare these to the cardinalities of the corresponding full grids that are  $129^2 = 16,641$ ,  $65^3 = 274,625$ ,  $23^2 = 529$ , and  $21^3 = 9,261$ .

Four sparse grids based on these two point sequences are depicted in Figure 4.1. There is a large array of other possibilities available in the literature. For example, Gerstner and Griebel [23] use Gauss–Patterson nodes, and Genz and Keister [19] have developed a nested version of the Gauss–Hermite rule. The rule (ii) can also be trivially extended for other integration domains and measures if a different sequence of orthogonal polynomials is used (e.g., Legendre or Chebyshev polynomials on  $[-1, 1]$ ).

**4.3. Worst-case error minimization with respect to the generators.** The third methodology for choosing the fully symmetric sets is that of principled worst-case error (WCE) minimization. Suppose that one, based on, for example, the number of nodes desired, fixes a number of generators of a fully symmetric kernel quadrature rule to zero or sets some equality constraints. Then the WCE  $e(Q_k)$  of the kernel quadrature rule can be minimized with respect to the generator vectors obeying these constraints. Especially in higher dimensions, this is a task vastly simpler than trying to minimize the error over a node set of unconstrained geometry. Optimal kernel quadrature rules under certain structural constraints have been previously experimented with, at least by O’Hagan [48, 49].

As a simplistic and somewhat arbitrary example, suppose that one desires a good fully symmetric kernel quadrature rule having about 80 nodes in  $\Omega \subset \mathbb{R}^3$ . A rule of the form

$$(4.1) \quad Q_k^\lambda(f) = w_1 f[0, 0, 0] + w_2 f[\lambda_1, \lambda_2, \lambda_3] + w_3 f[\lambda_3, \lambda_3, \lambda_3] + w_4 f[\lambda_4, \lambda_4, \lambda_5]$$

has  $1 + 48 + 8 + 24 = 81$  nodes if the generators  $\lambda = (\lambda_1, \dots, \lambda_5)$  are distinct and nonzero. Optimal generators  $\lambda^* = (\lambda_1^*, \lambda_2^*, \lambda_3^*, \lambda_4^*, \lambda_5^*)$  in the sense of minimal WCE

would then be

$$\lambda^* = \arg \min_{\lambda \in \Omega, \lambda_i > 0} e(Q_k^\lambda).$$

That is,  $Q_k^{\lambda^*}$  has the smallest WCE among all rules of the form (4.1). The optimal generators cannot in general be solved analytically, nor does this minimization problem appear to be convex.

In some very simple cases minimization is trivial. Consider rules of the form

$$Q_k^\lambda(f) = w_1 f[0, \dots, 0] + w_2 f[\lambda, 0, \dots, 0]$$

in  $\Omega \subset \mathbb{R}^d$ . If the kernel is Gaussian with length-scale  $\ell$  and  $\mu$  the standard Gaussian measure (see section 5.2), the task of finding the optimal generator  $\lambda^*$  reduces to

$$\lambda^* = \arg \min_{\lambda > 0} e(Q_k^\lambda) = \arg \max_{\lambda > 0} \left[ \frac{e^{-\lambda^2/2(1+\ell^2)} - e^{-\lambda^2/2\ell^2}}{1 - e^{-\lambda^2/\ell^2}} \right].$$

That is, the optimal generator is dimension-independent and can be easily computed. However, with increasing dimension and constant length-scale, this results in a negative weight for the origin, which often impairs numerical stability. It is somewhat questionable if such a rule is actually “good” (removing the origin of course yields positive weights).

We do not attempt to construct efficient fully symmetric rules using the technique described above in this article. The topic, alongside node selection for sparse grids via sequential WCE minimization briefly discussed in section 4.2, is left for future research.

**4.4. Convergence analysis.** In this section we provide convergence theorems for the fully symmetric kernel Monte Carlo and the Clenshaw–Curtis sparse grid kernel quadrature. The theorems are straightforward corollaries of some well-known results in the literature. For stating the results, we need to introduce the following three standard function classes. We assume that  $\Omega = [-a, a]^d$  and  $a < \infty$  in this section. The general principle on the convergence results for kernel quadrature is that the rates obtained are at least as good as those for any other method using the same nodes if the integrand belongs to the RKHS induced by the kernel. This should be quite clear from the definition of kernel quadrature.

With  $\alpha \in \mathbb{N}^d$  a multi-index and  $f: \Omega \rightarrow \mathbb{R}$  a sufficiently smooth function, the derivative operator is  $D^\alpha f = \partial^{|\alpha|} f / \partial x_1^{\alpha_1} \dots \partial x_d^{\alpha_d}$ . By  $\alpha \leq r$  we mean that  $\alpha_1, \dots, \alpha_d \leq r$ . The function classes we need are as follows:

- (i) The Sobolev space  $W_2^r$  is a Hilbert space defined as

$$W_2^r := \{f \in L^2(\mu) : D^\alpha f \in L^2(\mu) \text{ exists for all } |\alpha| \leq r\},$$

with the norm  $\|f\|_{W_2^r} = \sum_{|\alpha| \leq r} \|D^\alpha f\|_{L^2(\mu)}$ .

- (ii) The class  $C^r$  is the class of functions that have bounded derivatives:

$$C^r := \{f: \Omega \rightarrow \mathbb{R} : \|D^\alpha f\|_\infty < \infty \text{ for all } |\alpha| \leq r\}.$$

This space is equipped with the norm  $\|f\|_{C^r} = \max\{\|D^\alpha f\|_\infty : |\alpha| \leq r\}$ .

- (iii) The class  $F^r$  is the class of functions that have bounded mixed derivatives:

$$F^r := \{f: \Omega \rightarrow \mathbb{R} : \|D^\alpha f\|_\infty < \infty \text{ for all } \alpha \leq r\}.$$

This space is equipped with the norm  $\|f\|_{F^r} = \max\{\|D^\alpha f\|_\infty : \alpha \leq r\}$ .

For relations of the above function classes to RKHSs induced by different kernels, see, for example, [60, Chapter 10]. Two norms  $\|\cdot\|_1$  and  $\|\cdot\|_2$  of an arbitrary vector space are *norm-equivalent* if there are positive constants  $C_1$  and  $C_2$  such that  $C_1 \|x\|_1 \leq \|x\|_2 \leq C_2 \|x\|_1$  for all elements  $x$  of the vector space. Recall from section 2 that  $\mathcal{H}$  is the RKHS induced by the kernel  $k$ . Then  $\mathcal{H}$  is norm-equivalent to the Sobolev space  $W_2^r$  if  $r > d/2$  and the Fourier transform  $\mathcal{F}(\omega)$  of the kernel  $k$  decays at the rate  $(1 + \|\omega\|^2)^{-r}$ . This holds if the kernel is, for example, of the Matérn class with  $\nu \geq r - 1/2$ . In a similar manner,  $\mathcal{H}$  is norm-equivalent to  $F^r$  if the kernel is a product of one-dimensional Matérn kernels.

The following convergence theorems are simple consequences of results available in the literature. In these theorems  $Q_{k,n}$  stands for an  $n$ -point kernel quadrature rule with its type specified by a superscript. Extensions to the misspecified setting analyzed in [28, 29] may be possible.

**THEOREM 4.1** (convergence of FSKMC). *Let  $\Omega = [-a, a]^d$  with  $a < \infty$ . If  $\mathcal{H}$  is norm-equivalent to the Sobolev space  $W_2^r$  with  $r > d/2$ , then*

$$(4.2) \quad \mathbb{E}[e(Q_{k,n}^{\text{FSKMC}})] = \mathcal{O}(n^{-r/d+\varepsilon})$$

for any  $\varepsilon > 0$ . The expectation above is with respect to the joint distribution of the random generator vectors.

*Proof.* The WCE is decreasing in the number of nodes, so we know that  $e(Q_{k,n}^{\text{FSKMC}}) \leq e(Q_{k,n}^{\text{GEN}})$ , where the rule  $Q_{k,n}^{\text{GEN}}$  uses only the  $J$  generator vectors as its nodes. The rate (4.2) is realized by the regular kernel Monte Carlo [6, Theorem 1] under the norm-equivalence assumption. Therefore,

$$\mathbb{E}[e(Q_{k,n}^{\text{FSKMC}})] \leq \mathbb{E}[e(Q_{k,n}^{\text{GEN}})] = \mathcal{O}(J^{-r/d+\varepsilon}) = \mathcal{O}(n^{-r/d+\varepsilon})$$

because there is a dimension-dependent upper bound  $2^d d!$  (see (3.2)) for the number of nodes one fully symmetric set can contain.  $\square$

It is clear that the above rate is extremely crude with respect to  $d$  as the dimension also enters through the multiplicative factor  $2^d d!$ . It is likely that this factor can be eliminated or diminished with more careful analysis.

**THEOREM 4.2** (convergence of CCSGKQ). *Let  $\Omega = [-a, a]^d$  with  $a < \infty$ , and let  $\mu$  be the uniform measure on  $\Omega$ . If  $\mathcal{H}$  is norm-equivalent to  $C^r$ , then*

$$(4.3) \quad e(Q_{k,n}^{\text{CCSGKQ}}) = \mathcal{O}(n^{-r/d}(\log n)^{(d-1)(r/d+1)}).$$

If  $\mathcal{H}$  is norm-equivalent to  $F^r$ , then

$$(4.4) \quad e(Q_{k,n}^{\text{CCSGKQ}}) = \mathcal{O}(n^{-r}(\log n)^{(d-1)(r+1)}).$$

*Proof.* The rates (4.3) and (4.4) hold for the standard Clenshaw–Curtis sparse grid quadrature [41, 42] if the WCE (2.1) is over the unit balls of  $C^r$  and  $F^r$ , respectively. Because kernel quadrature rules have minimal WCEs in the induced RKHS among all quadrature rules with fixed nodes, the convergence rates follow from the assumptions of norm-equivalence.  $\square$

**5. Numerical examples and computational aspects.** This section contains three numerical examples for the fully symmetric kernel Monte Carlo, the Clenshaw–Curtis sparse grid kernel quadrature, and the (modified) Gauss–Hermite

sparse grid kernel quadrature, as well as discussion on some computational aspects. The examples and algorithms, implemented in MATLAB, are available at <https://github.com/tskarvone/fskq>. Numerous classical sparse grid quadrature methods for MATLAB are implemented in the Sparse Grid Interpolation Toolbox [31]. Parts of our code make use of this toolbox.

We emphasize that the examples are not meant to demonstrate superiority of fully symmetric kernel quadrature to other numerical integration methods. Comparisons to other methods are merely to show that fully symmetric kernel quadratures can achieve roughly comparable accuracy. Rather, we aim to show that fully symmetric sets make it possible to apply kernel quadrature rules to large-scale and high-dimensional situations that have been beyond the scope of these quadrature rules before.

**5.1. Choosing the length-scale parameter.** Accuracy of any approximation based on a stationary kernel is heavily dependent on the length-scale parameter  $\ell > 0$  whose effect is via  $k_\ell(\mathbf{x} - \mathbf{x}') = k((\mathbf{x} - \mathbf{x}')/\ell)$ ; see, for example, the Gaussian kernel (5.1). Choosing in some sense the best value of this parameter efficiently is an important topic of research both in scattered data approximation literature [54, 16] and statistics and machine learning [53, Chapter 5]. See also [6, section 4.1] for discussion in the context of kernel quadrature.

Unfortunately, we have not been able to come up with a way to exploit the fully symmetric structure of the node set in any of the existing parameter fitting methods, such as marginal likelihood maximization or cross-validation. Consequently, in large-scale applications that go beyond the limits of naive methods based on inverting the kernel matrix, one has to resort to ad hoc techniques to fit the length-scale. In the examples below we either use few enough nodes that naive computations are possible (section 5.3), integrate a function whose length-scale is known beforehand (section 5.4), or set the length-scale somewhat heuristically (section 5.5). We recognize that the lack of a principled method for choosing the length-scale is a significant shortcoming and hope to fix this in the future.

When the length-scale is changed, the interpretation of the WCE as an indicator of accuracy of the quadrature rule is confounded because the RKHS norm  $\|\cdot\|_{\mathcal{H}}$  depends on the length-scale. This occurs in sections 5.3 and 5.5. Nevertheless, if one follows the paradigm presented in section 2.2 the WCE still carries a meaningful probabilistic interpretation as the integral posterior standard deviation (STD). As such, it is plotted in all the examples. However, one should not draw too many conclusions from these plots, as we have not made any effort to fit the kernel scale parameter (i.e., the constant multiplier of the kernel).

**5.2. Closed-form kernel means.** In kernel quadrature, one needs to be able to evaluate the kernel mean  $k_\mu(\mathbf{x}_i) = \int_{\Omega} k(\mathbf{x}_i, \mathbf{x}) d\mu(\mathbf{x})$  at the nodes  $\mathbf{x}_i$ . A number of kernel-measure pairs that yield tractable kernel means are tabulated in [6]. It is also possible to evaluate the kernel mean numerically [57]. In fact, when fully symmetric sets are used, numerical evaluation may be quite viable, as the kernel mean needs to be evaluated only at each generator vector instead of each node.

All our examples use the standard Gaussian kernel

$$(5.1) \quad k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right)$$

with length-scale  $\ell > 0$  and unit scale parameter (this parameter only affects the magnitude of the WCE). The integration domain  $\Omega$  and measure  $\mu$  are either (i) the whole of  $\mathbb{R}^d$  and the standard Gaussian measure with the density

$\varphi(\mathbf{x}) = (2\pi)^{-d/2} \exp(-\|\mathbf{x}\|^2/2)$ , or (ii) the hypercube  $[-1, 1]^d$  and the (normalizing) uniform measure. In the former case the kernel mean and its integral (needed for computing the WCE using (2.3)) are

$$k_\mu(\mathbf{x}) = \int_{\mathbb{R}^d} k(\mathbf{x}, \mathbf{x}')\varphi(\mathbf{x}') d\mathbf{x}' = \left(\frac{\ell^2}{1 + \ell^2}\right)^{d/2} \exp\left(-\frac{\|\mathbf{x}\|^2}{2(1 + \ell^2)}\right),$$

$$\mu(k_\mu) = \int_{\mathbb{R}^d} k_\mu(\mathbf{x})\varphi(\mathbf{x}) d\mathbf{x} = \left(\frac{\ell^2}{2 + \ell^2}\right)^{d/2}$$

and in the latter

$$k_\mu(\mathbf{x}) = 2^{-d} \int_{[-1,1]^d} k(\mathbf{x}, \mathbf{x}') d\mathbf{x}' = \left(\frac{\pi\ell^2}{8}\right)^{d/2} \prod_{i=1}^d \left[ \operatorname{erf}\left(\frac{x_i + 1}{\ell\sqrt{2}}\right) - \operatorname{erf}\left(\frac{x_i - 1}{\ell\sqrt{2}}\right) \right],$$

$$\mu(k_\mu) = 2^{-d} \int_{[-1,1]^d} k_\mu(\mathbf{x}) d\mathbf{x} = \left(\frac{\pi\ell^2}{8}\right)^{d/2} \left(\sqrt{\frac{2\ell^2}{\pi}} (e^{-2/\ell^2} - 1) + 2 \operatorname{erf}(\sqrt{2}/\ell)\right)^d,$$

where  $\operatorname{erf}(x) = \pi^{-1/2} \int_{-x}^x \exp(-t^2) dt$  is the standard error function.

**5.3. Example 1: Random generators.** Our first example is just a proof of concept to demonstrate that the fully symmetric kernel Monte Carlo (FSKMC) from section 4.1 indeed works (though not necessarily that well). We try numerically integrating the nonradial function

$$(5.2) \quad f(\mathbf{x}) = \exp\left(\sin(5\|\mathbf{x}\|)^2 - (x_1^2 + 0.5x_2^2 + 2x_3^4)\right)$$

over  $\mathbb{R}^3$  and with respect to the standard normal distribution. Results for the kernel Monte Carlo (KMC) [52, 6], where the nodes to be used in kernel quadrature are drawn randomly, and FSKMC are presented in Figure 5.1.

For both KMC and FSKMC, the kernel length-scale was fit by the method of maximum likelihood (see [53, Chapter 5]) using the Monte Carlo samples of KMC. We have also experimented with fitting the FSKMC length-scale using the randomly generated fully symmetric sets, but in this case the fitted length-scale was markedly larger and the integral approximations much worse.

It is clear that FSKMC fares worse. However, FSKMC has a tremendous advantage in computational scalability in the number of nodes. In general, when the number of nodes exceeds some tens of thousands, kernel quadrature methods based on naively solving the weights from the linear system (2.2), such as KMC, become infeasible due to the cubic time and quadratic space complexity. In contrast, fully symmetric kernel quadratures such as FSKMC remain feasible: only  $Jn$  (recall that  $J \leq n$  is the number of fully symmetric sets) kernel evaluations and solving a linear system of  $J$  equations—as opposed to  $n^2$  and  $n$ , respectively, of naive methods—are required. For instance, using FSKMC with 1,000 fully symmetric sets (i.e., 48,000 nodes) in this example would require 48,000,000 kernel evaluations and solving a linear system of 1,000 equations, neither of which is a computational challenge, while the KMC weights for 48,000 nodes cannot be computed on a standard computer.

The difference becomes even more pronounced in higher dimensions where fully symmetric sets contain significantly more points. The next two examples demonstrate the superior scalability of fully symmetric kernel quadratures.

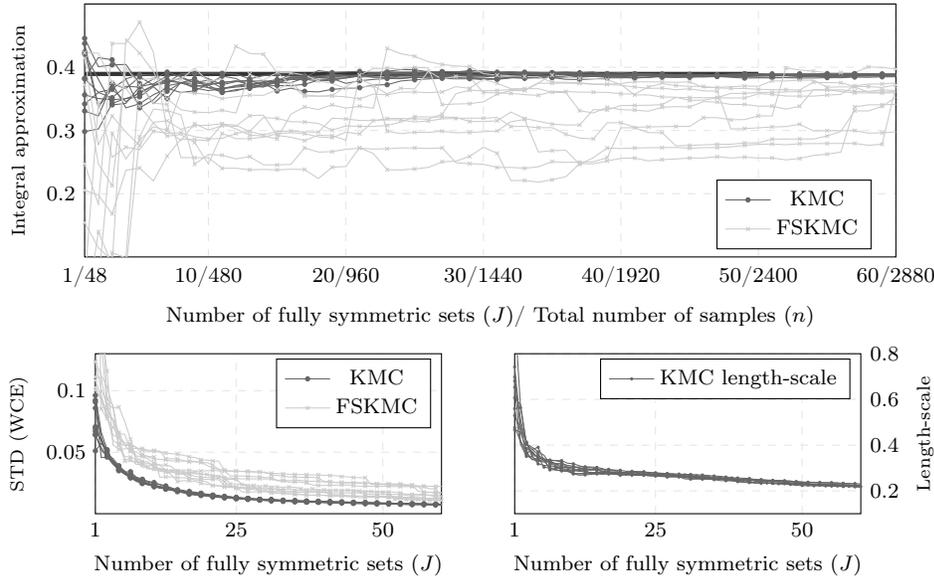


FIG. 5.1. Numerical integration of the function (5.2) with respect to the standard normal distribution. The upper figure shows integral approximations by the kernel Monte Carlo (KMC) and the fully symmetric kernel Monte Carlo (FSKMC) as a function of the number  $J$  of fully symmetric sets and the total number  $n$  of Monte Carlo samples for ten realizations. The lower figures display the worst-case errors (standard deviations) and the length-scales fitted. Each fully symmetric set contains 48 nodes (see Table 3.1). The underlying black line is the value of the integral that is approximately 0.389. The generator vectors for FSKMC have been generated independently of the KMC samples. Both methods use the same length-scale that has been fit using the Monte Carlo samples of KMC.

**5.4. Example 2: A priori known length-scale.** This simple example demonstrates that sparse grid kernel quadrature based on fully symmetric sets is numerically stable, is consistent, and works well for an extremely large number of nodes.

We work in the domain  $\Omega = [-1, 1]^d$ ,  $d = 11$ , equipped with the normalizing uniform measure. The integrand is

$$(5.3) \quad f(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_f\|^2}{2\ell_f^2}\right)$$

with  $\ell_f = 0.8$  and  $\mathbf{x}_f$  a vector of 11 evenly spaced points on the interval  $[0.2, 0.5]$  (with the end points included). The integral we seek to approximate is

$$2^{-d} \int_{[-1,1]^d} f(\mathbf{x}) \, d\mathbf{x} = \left(\frac{\pi\ell_f^2}{8}\right)^{d/2} \prod_{i=1}^d \left[ \operatorname{erf}\left(\frac{x_{f,i} + 1}{\ell_f\sqrt{2}}\right) - \operatorname{erf}\left(\frac{x_{f,i} - 1}{\ell_f\sqrt{2}}\right) \right] \approx 0.0392.$$

We use the Gaussian kernel with  $\ell = \ell_f = 0.8$  and the Clenshaw–Curtis sparse grid kernel quadrature (CCSGKQ). Results for the relative error  $|\mu(f) - Q_k(f)|/\mu(f)$  and the kernel WCE (or standard deviation) are shown in Figure 5.2 for the levels  $q = 1, \dots, 9$ , the last of them corresponding to the total of 15,005,761 nodes. Table 5.1 contains a breakdown of the computational times required. We also display results for KMC using up to 12,000 nodes. For this many nodes, the time taken by CCSGKQ is negligible, while the KMC is noticeably slowing down.

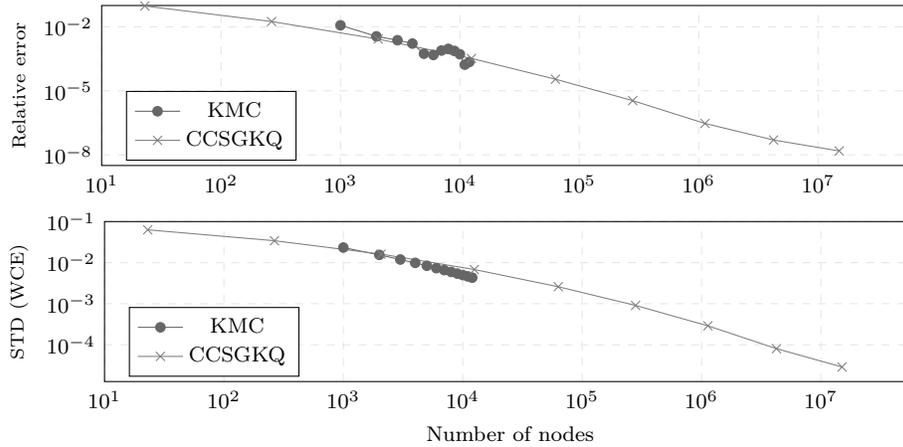


FIG. 5.2. Relative error  $|\mu(f) - Q_k(f)|/\mu(f)$  (upper) and the worst-case error (lower; standard deviation) for integration of the function (5.3) on the 11-dimensional hypercube with the kernel Monte Carlo quadrature (KMC) and the Clenshaw–Curtis sparse grid kernel quadrature (CCSGKQ). The number of nodes for the KMC varied from 1,000 to 12,000 (with increments of 1,000), and CCSGKQ used levels from 1 to 9. These corresponded to 23, 265, 2,069, 12,497, 63,097, 280,017, 1,129,569, 4,236,673, and 15,005,761 nodes and 2, 4, 8, 17, 36, 79, 172, 379, and 832 fully symmetric sets.

TABLE 5.1

Computational times in seconds for KMC (left) and CCSGKQ (right) in Example 2. The columns indicate the time taken by kernel evaluations (kernel), computing the weights from a linear system of equations (weights), and constructing the fully symmetric sets (FSS). The MATLAB code was run on a laptop with an Intel Core i5-6300 2.40 GHz processor and 8 GB of RAM.

Computational times (seconds) for KMC / CCSGKQ						
$n$	Kernel	Weights	$J / n$	Kernel	Weights	FSS
1k	0.08	0.01	2 / 23	< 0.01	< 0.001	0.07
2k	0.15	0.07	4 / 265	< 0.01	< 0.001	0.07
3k	0.32	0.20	8 / 2k	< 0.01	< 0.001	0.04
4k	0.54	0.47	17 / 12k	0.02	< 0.001	0.04
5k	0.79	0.85	36 / 63k	0.14	< 0.001	0.10
6k	1.17	1.42	79 / 280k	1.26	0.003	0.27
7k	1.49	2.21	172 / 1.1m	10.91	0.004	0.82
8k	1.92	3.20	379 / 4.2m	90.41	0.004	2.63
9k	2.43	4.47	832 / 15m	760.00	0.072	8.61
10k	2.98	6.03				
11k	3.62	8.21				
12k	4.45	10.61				

It is seen that for a similar number of nodes the two methods are roughly equivalent. When the level, and consequently the number of nodes, increases, CCSGKQ becomes more accurate, which shows that the nodes are selected well enough and that the weights are being computed correctly. For the highest levels 8 and 9, the sparse grids consisted of 379 and 832 fully symmetric sets or 4,236,673 and 15,005,761 nodes, resulting in  $379 \times 4,236,673 = 1,605,699,067$  and  $832 \times 15,005,761 = 12,484,793,152$  kernel evaluations needed to compute the weights. It is not possible to compute the KMC weights for this many nodes.

**5.5. Example 3: Zero coupon bonds.** This example demonstrates that fully symmetric kernel quadrature rules are also feasible in high dimensions. We use the toy

example from [40] (see also [26, section 6.1]) that is concerned with pricing zero coupon bonds through simulation of a discretized stochastic differential equation model. The model is convenient for our purposes, as there is a closed-form solution that serves as a baseline and finer discretizations correspond to higher integration dimensions.

Consider the stochastic differential equation (going by the name the Vasicek model)

$$dr(t) = \kappa(\theta - r(t)) dt + \sigma dW(t),$$

where  $W(t)$  is the standard Brownian motion and  $\kappa$ ,  $\theta$ , and  $\sigma$  are positive parameters. We want to solve this SDE at the time  $t = T$ . The Euler–Maruyama discretization with the uniform step size  $\Delta t = T/d$  is

$$r_k = r_{k-1} + \kappa(\theta - r_{k-1})\Delta t + \sigma x_k, \quad k = 1, \dots, d,$$

where  $x_k \sim \mathcal{N}(0, \Delta t)$  are independent and  $r_0$  is a free parameter. The quantity we are interested in is the Gaussian integral

$$\begin{aligned} (5.4) \quad P(0, T) &:= \mathbb{E} \left[ \exp \left( -\Delta t \sum_{k=0}^{d-1} r_k \right) \right] \\ &= \exp(-\Delta t r_0) \int_{\mathbb{R}^{d-1}} \exp \left( -\Delta t f(\sqrt{\Delta t} \mathbf{x}) \right) \varphi(\mathbf{x}) d\mathbf{x}, \end{aligned}$$

where  $f(\mathbf{x}) = \sum_{k=1}^{d-1} r_k$ . This integral admits a closed-form solution

$$(5.5) \quad P(0, T) = \exp \left( -\frac{(\gamma + \beta_d r_0) T}{d} \right)$$

with  $\beta_k = \sum_{i=1}^k (1 - \kappa \Delta t)^{j-1}$  and  $\gamma = \sum_{k=1}^{d-1} (\beta_k \kappa \theta \Delta t - (\beta_k \sigma \Delta t)^2 / 2)$ . As can be seen, the number  $d$  of discretization steps controls the integration dimension, which is  $d - 1$ .

In the integration experiment, we set

$$\kappa = 0.1817303, \quad \theta = 0.0825398957, \quad \sigma = 0.0125901, \quad r_0 = 0.021673, \quad T = 5.$$

These values are equal to those used in [40, 26]. We consider numerical integration of (5.4) for  $d = 10, \dots, 300$  using the Gauss–Hermite sparse grid kernel quadrature (GHSGKQ) with  $q = 2$ . We use the Gaussian kernel with the somewhat heuristic choice  $\ell = d$  of the length-scale. The central node (i.e., the origin) tended to have a fairly large negative weight, so it was removed to improve numerical stability. Results for the relative error and the WCE (standard deviation) are depicted in Figure 5.3. For comparison, we have also included a Monte Carlo estimate (KMC is feasible only for dimensions somewhat less than 100, so it was excluded).

The results show that GHSGKQ is able to maintain an accuracy that is generally better than that of the standard MC. This indicates that fully symmetric kernel quadratures have potential also in very high dimensions. Note that the dimension-adaptive methods used in [26] would be more accurate in this example. It is probable that fully symmetric kernel quadratures could be combined with these methods.

**6. Conclusions and discussion.** We introduced fully symmetric kernel quadrature rules and showed that their weights can be computed exactly with a very simple algorithm under some assumptions on the integration domain and measure and the kernel. We also proposed using sparse grids in conjunction with this algorithm and

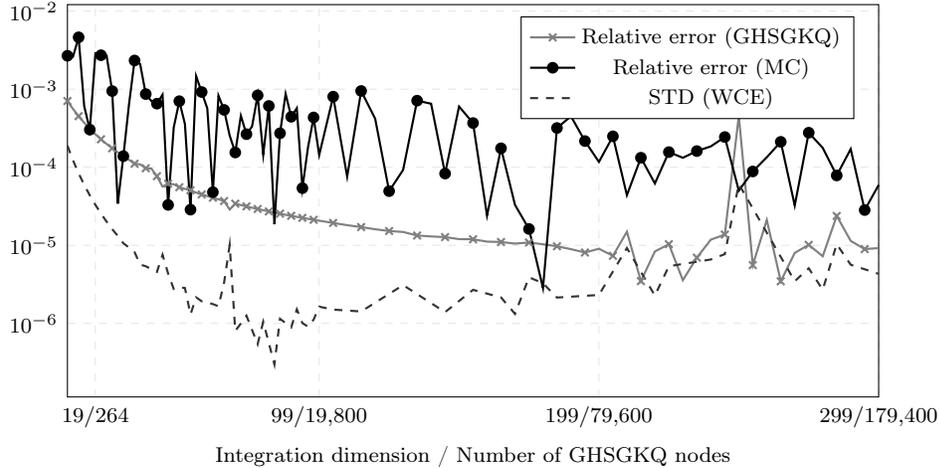


FIG. 5.3. *Relative error and the WCE (standard deviation) of the Gauss–Hermite sparse grid kernel quadrature (GHSKGKQ) for the zero coupon bond setup of section 5.5. Value of the integral (5.4), as computed from (5.5), is between 0.81 and 0.815 for all dimensions. Also depicted is an integral estimate by the standard Monte Carlo using the same number of points as GHSKGKQ in each dimension.*

provided some simple theoretical convergence analysis for Clenshaw–Curtis sparse grids. In the schemes presented, the nodes can be selected in a comparatively flexible manner. Three numerical experiments demonstrated that the approach is sound and can cope with both a very large number of nodes and high-dimensional domains.

Even with the tremendous computational simplifications provided by the fully symmetric sets, kernel quadrature rules remain computationally more demanding than most classical quadrature rules. In the end, the decision on which method to use is highly dependent on the computational complexity of evaluating the integrand. Extremes where the rules we have developed are not necessarily useful are easy to identify: (i) In section 5.4 it is clearly absurd that an integrand as cheap to evaluate as the kernel is evaluated 15 million times, while 12 billion kernel evaluations are used to compute the weights, (ii) whereas when the integrand, being, for example, a complex computer simulation, is very expensive, the computational overhead from nonsymmetric and likely more accurate kernel quadrature rules is going to be negligible. Consequently, we believe that the method presented in this article is best suited for “moderately” expensive integrands in the case when high accuracy is required or probabilistic modeling of uncertainty in the integral estimate desired. This is of course somewhat ambiguous. Precise (and useful) analysis is complicated by, among other things, the facts that we do not know how the accuracy of a fully symmetric kernel quadrature rule compares to that of a nonsymmetric one (e.g., Theorem 4.1 is only about rates, not the associated constant coefficients) and that it is difficult to account for the value one places on the uncertainty measure—nor is it easy to decide how much one should value this measure to begin with.

Besides what was discussed in the preceding paragraph, there are a number of topics that could be pursued in the future:

- Developing principled methods for choosing the kernel length-scale for large-scale problems.
- Proper probabilistic approach to large-scale integration problems. We anticipate that much can be gained in pursuing this direction.

- As discussed in sections 4.2 and 4.3, there is much room for improvement via optimization of the fully symmetric sets.
- Rows of the submatrices  $\mathbf{K}_{ij}$  in (3.4) typically contain several nondistinct elements. Minor computational improvements might be possible.

**Acknowledgments.** We thank François-Xavier Briol, Jon Cockayne, Chris Oates, Jens Oettershagen, and Filip Tronarp for discussion and constructive comments. Suggestions by the anonymous reviewers helped to improve many parts of the article.

## REFERENCES

- [1] N. ARONSZAJN, *Theory of reproducing kernels*, Trans. Amer. Math. Soc., 68 (1950), pp. 337–404.
- [2] A. BERLINET AND C. THOMAS-AGNAN, *Reproducing Kernel Hilbert Spaces in Probability and Statistics*, Springer, 2004.
- [3] A. YU. BEZHAEV, *Cubature formulae on scattered meshes*, Soviet J. Numer. Anal. Math. Modelling, 6 (1991), pp. 95–106.
- [4] F.-X. BRIOL, C. J. OATES, J. COCKAYNE, W. Y. CHEN, AND M. GIROLAMI, *On the sampling problem for kernel quadrature*, in 34th International Conference on Machine Learning, Proceedings of Machine Learning Research 70, 2017, pp. 586–595.
- [5] F.-X. BRIOL, C. J. OATES, M. GIROLAMI, AND M. A. OSBORNE, *Frank-Wolfe Bayesian quadrature: Probabilistic integration with theoretical guarantees*, in Advances in Neural Information Processing Systems, Vol. 25, 2015, pp. 1162–1170.
- [6] F.-X. BRIOL, C. J. OATES, M. GIROLAMI, M. A. OSBORNE, AND D. SEJDINOVIC, *Probabilistic Integration: A Role for Statisticians in Numerical Analysis?*, preprint, <https://arxiv.org/abs/1512.00933v5>, 2016.
- [7] H.-J. BUNGARTZ AND M. GRIEBEL, *Sparse grids*, Acta Numer., 13 (2004), pp. 147–269.
- [8] R. E. CAFLISCH, *Monte Carlo and quasi-Monte Carlo methods*, Acta Numer., 7 (1998), pp. 1–49.
- [9] J. COCKAYNE, C. OATES, T. SULLIVAN, AND M. GIROLAMI, *Bayesian Probabilistic Numerical Methods*, preprint, <https://arxiv.org/abs/1702.03673>, 2017.
- [10] T. D. COOK AND M. K. CLAYTON, *Sequential Bayesian Quadrature*, tech. report, Department of Statistics, University of Wisconsin, 1998.
- [11] R. COOLS, *Constructing cubature formulae: The science behind the art*, Acta Numer., 6 (1997), pp. 1–54.
- [12] P. J. DAVIS AND P. RABINOWITZ, *Methods of Numerical Integration*, 2nd ed., Academic Press, 1984.
- [13] C. DE BOOR AND A. RON, *On multivariate polynomial interpolation*, Constr. Approx., 6 (1990), pp. 287–302.
- [14] P. DIACONIS, *Bayesian numerical analysis*, in Statistical Decision Theory and Related Topics IV, Vol. 1, Springer-Verlag, New York, 1988, pp. 163–175.
- [15] Z. DONG, E. H. GEORGIOULIS, J. LEVESLEY, AND F. USTA, *Fast Multilevel Sparse Gaussian Kernels for High-Dimensional Approximation and Integration*, preprint, <https://arxiv.org/abs/1501.03296>, 2015.
- [16] G. E. FASSHAUER AND J. G. ZHANG, *On choosing “optimal” shape parameters for RBF approximation*, Numer. Algorithms, 45 (2007), pp. 345–368.
- [17] J. GARCKE, *A dimension adaptive sparse grid combination technique for machine learning*, ANZIAM J., 48 (2006), pp. 725–740.
- [18] A. GENZ, *Fully symmetric interpolatory rules for multiple integrals*, SIAM J. Numer. Anal., 23 (1986), pp. 1273–1283, <https://doi.org/10.1137/0723086>.
- [19] A. GENZ AND B. D. KEISTER, *Fully symmetric interpolatory rules for multiple integrals over infinite regions with Gaussian weight*, J. Comput. Appl. Math., 71 (1996), pp. 299–309.
- [20] A. GENZ AND J. MONAHAN, *Stochastic integration rules for infinite regions*, SIAM J. Sci. Comput., 19 (1998), pp. 426–439, <https://doi.org/10.1137/S1064827595286803>.
- [21] A. GENZ AND J. MONAHAN, *A stochastic algorithm for high-dimensional integrals over unbounded regions with Gaussian weight*, J. Comput. Appl. Math., 112 (1999), pp. 71–81.
- [22] E. H. GEORGIOULIS, J. LEVESLEY, AND F. SUBHAN, *Multilevel sparse kernel-based interpolation*, SIAM J. Sci. Comput., 35 (2013), pp. A815–A831, <https://doi.org/10.1137/110859610>.
- [23] T. GERSTNER AND M. GRIEBEL, *Numerical integration using sparse grids*, Numer. Algorithms, 18 (1998), pp. 209–232.
- [24] T. GUNTER, M. A. OSBORNE, R. GARNETT, P. HENNIG, AND S. J. ROBERTS, *Sampling for inference in probabilistic models with fast Bayesian quadrature*, in Advances in Neural

- Information Processing Systems, Vol. 24, 2014, pp. 2789–2797.
- [25] P. HENNIG, M. A. OSBORNE, AND M. GIROLAMI, *Probabilistic numerics and uncertainty in computations*, Proc. Roy. Soc. London A Math. Phys. Engrg. Sci., 471 (2015), 20150142.
- [26] M. HOLTZ, *Sparse Grid Quadrature in High Dimensions with Applications in Finance and Insurance*, Lecture Notes in Comput. Sci. Engrg. 77, Springer, 2011.
- [27] F. HUSZÁR AND D. DUVENAUD, *Optimally-weighted herding is Bayesian quadrature*, in 28th Conference on Uncertainty in Artificial Intelligence, 2012, pp. 377–385.
- [28] M. KANAGAWA, B. K. SRIPERUMBUDUR, AND K. FUKUMIZU, *Convergence guarantees for kernel-based quadrature rules in misspecified settings*, in Advances in Neural Information Processing Systems, Vol. 29, 2016, pp. 3288–3296.
- [29] M. KANAGAWA, B. K. SRIPERUMBUDUR, AND K. FUKUMIZU, *Convergence Analysis of Deterministic Kernel-Based Quadrature Rules in Misspecified Settings*, preprint, <https://arxiv.org/abs/1709.00147>, 2017.
- [30] T. KARVONEN AND S. SÄRKKÄ, *Classical quadrature rules via Gaussian processes*, in 27th IEEE International Workshop on Machine Learning for Signal Processing, 2017.
- [31] A. KLIMKE AND B. WOHLMUTH, *Algorithm 847: Spinterp: Piecewise multilinear hierarchical sparse grid interpolation in MATLAB*, ACM Trans. Math. Software, 31 (2005), pp. 561–579.
- [32] F. M. LARKIN, *Optimal approximation in Hilbert spaces with reproducing kernel functions*, Math. Comp., 24 (1970), pp. 911–921.
- [33] F. M. LARKIN, *Gaussian measure in Hilbert space and applications in numerical analysis*, Rocky Mountain J. Math., 2 (1972), pp. 379–422.
- [34] J. LU AND D. L. DARMOFAL, *Higher-dimensional integration with Gaussian weight for applications in probabilistic design*, SIAM J. Sci. Comput., 26 (2004), pp. 613–624, <https://doi.org/10.1137/S1064827503426863>.
- [35] J. N. LYNESS, *Symmetric integration rules for hypercubes I. Error coefficients*, Math. Comp., 19 (1965), pp. 260–276.
- [36] J. MCNAMEE AND F. STENGER, *Construction of fully symmetric numerical integration formulas*, Numer. Math., 10 (1967), pp. 327–344.
- [37] T. MINKA, *Deriving Quadrature Rules from Gaussian Processes*, tech. report, Statistics Department, Carnegie Mellon University, 2000.
- [38] K. MUANDET, K. FUKUMIZU, B. SRIPERUMBUDUR, AND B. SCHÖLKOPF, *Kernel mean embedding of distributions: A review and beyond*, Found. Trends Mach. Learning, 10 (2017), pp. 1–141.
- [39] A. NARAYAN AND D. XIU, *Stochastic collocation methods on unstructured grids in high dimensions via interpolation*, SIAM J. Sci. Comput., 34 (2012), pp. A1729–A1752, <https://doi.org/10.1137/110854059>.
- [40] S. NINOMIYA AND S. TEZUKA, *Toward real-time pricing of complex financial derivatives*, Appl. Math. Finance, 3 (1996), pp. 1–20.
- [41] E. NOVAK AND K. RITTER, *High dimensional integration of smooth functions over cubes*, Numer. Math., 75 (1996), pp. 79–97.
- [42] E. NOVAK AND K. RITTER, *The curse of dimension and a universal method for numerical integration*, in Multivariate Approximation and Splines, Birkhäuser, Basel, 1997, pp. 177–187.
- [43] E. NOVAK AND K. RITTER, *Simple cubature formulas with high polynomial exactness*, Constr. Approx., 15 (1999), pp. 499–522.
- [44] E. NOVAK, K. RITTER, R. SCHMITT, AND A. STEINBAUER, *On an interpolatory method for high dimensional integration*, J. Comput. Appl. Math., 112 (1999), pp. 215–228.
- [45] E. NOVAK AND H. WOŹNIAKOWSKI, *Tractability of Multivariate Problems, Volume II: Standard Information for Functionals*, EMS Tracts Math. 12, European Mathematical Society, 2010.
- [46] J. OETTERSCHAGEN, *Construction of Optimal Cubature Algorithms with Applications to Econometrics and Uncertainty Quantification*, Ph.D. thesis, Institut für Numerische Simulation, Universität Bonn, 2017.
- [47] A. O’HAGAN, *Curve fitting and optimal design for prediction*, J. Roy. Statist. Soc. Ser. B, 40 (1978), pp. 1–42.
- [48] A. O’HAGAN, *Bayes–Hermite quadrature*, J. Statist. Plann. Inference, 29 (1991), pp. 245–260.
- [49] A. O’HAGAN, *Some Bayesian numerical analysis*, Bayesian Statist., 4 (1992), pp. 345–363.
- [50] J. PRÜHER AND O. STRAKA, *Gaussian process quadrature moment transform*, IEEE Trans. Automat. Control, 2017, <https://doi.org/10.1109/TAC.2017.2774444>.
- [51] J. PRÜHER, F. TRONARP, T. KARVONEN, S. SÄRKKÄ, AND O. STRAKA, *Student-t process quadratures for filtering of non-linear systems with heavy-tailed noise*, in 20th International Conference on Information Fusion, 2017.
- [52] C. E. RASMUSSEN AND Z. GHAHRAMANI, *Bayesian Monte Carlo*, in Advances in Neural Information Processing Systems, Vol. 15, 2002, pp. 505–512.

- [53] C. E. RASMUSSEN AND C. K. I. WILLIAMS, *Gaussian Processes for Machine Learning*, Adaptive Computation and Machine Learning, MIT Press, 2006.
- [54] S. RIPPA, *An algorithm for selecting a good value for the parameter  $c$  in radial basis function interpolation*, *Adv. Comput. Math.*, 11 (1999), pp. 193–210.
- [55] K. RITTER, *Average-Case Analysis of Numerical Problems*, Lecture Notes in Math. 1733, Springer, 2000.
- [56] S. A. SMOLYAK, *Quadrature and interpolation formulas for tensor products of certain classes of functions*, *Dokl. Akad. Nauk SSSR*, 4 (1963), pp. 240–243.
- [57] A. SOMMARIVA AND M. VIANELLO, *Numerical cubature on scattered data by radial basis functions*, *Computing*, 76 (2006), pp. 295–310.
- [58] S. SÄRKKÄ, J. HARTIKAINEN, L. SVENSSON, AND F. SANDBLOM, *On the relation between Gaussian process quadratures and sigma-point methods*, *J. Adv. Inform. Fusion*, 11 (2016), pp. 31–46.
- [59] F. USTA AND J. LEVESLEY, *Multilevel quasi-interpolation on a sparse grid with the Gaussian*, *Numer. Algorithms*, <https://doi.org/10.1007/s11075-017-0340-y>.
- [60] H. WENDLAND, *Scattered Data Approximation*, Cambridge Monogr. Appl. Comput. Math. 17, Cambridge University Press, 2010.