# Optimal point sets for quasi–Monte Carlo integration of bivariate periodic functions with bounded mixed derivatives

Aicke Hinrichs and Jens Oettershagen

**Abstract** We investigate quasi-Monte Carlo (QMC) integration of bivariate periodic functions with dominating mixed smoothness of order one. While there exist several QMC constructions which asymptotically yield the optimal rate of convergence of $\mathscr{O}(N^{-1}\log(N)^{\frac{1}{2}})$, it is yet unknown which point set is optimal in the sense that it is a global minimizer of the worst case integration error. We will present a computer-assisted proof by exhaustion that the Fibonacci lattice is the unique minimizer of the QMC worst case error in periodic $H^1_{\mathrm{mix}}$ for small Fibonacci numbers $N$. Moreover, we investigate the situation for point sets whose cardinality $N$ is not a Fibonacci number. It turns out that for $N = 1, 2, 3, 5, 7, 8, 12, 13$ the optimal point sets are integration lattices.

## 1 Introduction

Quasi-Monte Carlo (QMC) rules are equal-weight quadrature rules which can be used to approximate integrals defined on the $d$-dimensional unit cube $[0,1)^d$

$$\int_{[0,1)^d} f(\boldsymbol{x})\,\mathrm{d}\boldsymbol{x} \approx \frac{1}{N}\sum_{i=1}^{N} f(\boldsymbol{x}_i),$$

where $\mathscr{P}_N = \{\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N\}$ are deterministically chosen quadrature points in $[0,1)^d$. The integration error for a specific function $f$ is given as

Aicke Hinrichs
Institut für Analysis, Johannes-Kepler-Universität Linz, Altenberger Straße 69, 4040 Linz, Austria
e-mail: aicke.hinrichs@uni-rostock.de

Jens Oettershagen
Institute for Numerical Simulation, Wegelerstraße 6, 53115 Bonn, Germany
e-mail: oettershagen@ins.uni-bonn.de

$$\left| \int_{[0,1)^d} f(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - \frac{1}{N} \sum_{i=1}^{N} f(\boldsymbol{x}_i) \right|.$$

To study the behavior of this error as $N$ increases for $f$ from a Banach space $(\mathscr{H}, \|\cdot\|)$ one considers the worst case error

$$\mathrm{wce}(\mathscr{H}, \mathscr{P}_N) = \sup_{\substack{f \in \mathscr{H} \\ \|f\| \leq 1}} \left| \int_{[0,1)^d} f(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} - \frac{1}{N} \sum_{i=1}^{N} f(\boldsymbol{x}_i) \right|.$$

Particularly nice examples of such function spaces are reproducing kernel Hilbert spaces [1]. Here, we will consider the reproducing kernel Hilbert space $H_{\mathrm{mix}}^1$ of 1-periodic functions with mixed smoothness. Details on these spaces are given in Section 2. The reproducing kernel is a tensor product kernel of the form

$$K_{d,\gamma}(\boldsymbol{x}, \boldsymbol{y}) = \prod_{j=1}^{d} K_{1,\gamma}(x_j, y_j) \text{ for } \boldsymbol{x} = (x_1, \dots, x_d), \boldsymbol{y} = (y_1, \dots, y_d) \in [0,1)^d$$

with $K_{1,\gamma}(x, y) = 1 + \gamma k(|x - y|)$ and $k(t) = \frac{1}{2}(t^2 - t + \frac{1}{6})$ and a parameter $\gamma > 0$. It turns out that minimizing the worst case error $\mathrm{wce}(H_{\mathrm{mix}}^1, \mathscr{P}_N)$ among all $N$-point sets $\mathscr{P}_N = \{\boldsymbol{x}_1, \dots, \boldsymbol{x}_N\}$ with respect to the Hilbert space norm corresponding to the kernel $K_{d,\gamma}$ is equivalent to minimizing the double sum

$$G_{\gamma}(\boldsymbol{x}_1, \dots, \boldsymbol{x}_N) = \sum_{i,j=1}^{N} K_{d,\gamma}(\boldsymbol{x}_i, \boldsymbol{x}_j).$$

There is a general connection between the discrepancy of a point set and the worst case error of integration. Details can be found in [11, Chapter 9]. In our case, the relevant notion is the $L_2$-norm of the periodic discrepancy. We describe the connection in detail in Section 2.3.

There are many results on the rate of convergence of worst case errors and of the optimal discrepancies for $N \to \infty$, see e.g. [10, 11], but results on the optimal point configurations for fixed $N$ and $d > 1$ are scarce. For discrepancies, we are only aware of [21], where the point configurations minimizing the standard $L_\infty$-star-discrepancy for $d = 2$ and $N = 1, 2, \dots, 6$ are determined, [14], where for $N = 1$ the point minimizing the standard $L_\infty$- and $L_2$-star discrepancy for $d \geq 1$ is found, and [6], where this is extended to $N = 2$.

It is the aim of this paper to provide a method which for $d = 2$ and $N > 2$ yields the optimal points for the periodic $L_2$-discrepancy and worst case error in $H_{\mathrm{mix}}^1$. Our approach is based on a decomposition of the global optimization problem into exponentially many local ones which each possess unique solutions that can be approximated efficiently by a nonlinear block Gauß-Seidel method. Moreover, we use the symmetries of the two-dimensional torus to significantly reduce the number of local problems that have to be considered.

It turns out that in the case that $N$ is a (small) Fibonacci number, the Fibonacci lattice yields the optimal point configuration. It is common wisdom, see e.g. [3, 8, 15, 16], that the Fibonacci lattice provides a very good point set for integrating periodic functions. Now our results support the conjecture that they are actually the best points.

These results may suggest that the optimal point configurations are integration lattices or at least lattice point sets. This seems to be true for some numbers $N$ of points, for example for Fibonacci numbers, but not always. However, it can be shown that integration lattices are always *local* minima of $\mathrm{wce}(H_{\mathrm{mix}}^1, \mathscr{P}_N)$. Moreover, our numerical results also suggest that for small $\gamma$ the optimal points are always *close* to a lattice point set, i.e. $N$-point sets of the form

$$\left\{\left(\frac{i}{N}, \frac{\sigma(i)}{N}\right) : i = 0, \ldots, N-1\right\},$$

where $\sigma$ is a permutation of $\{0, 1, \ldots, N-1\}$.

The remainder of this article is organized as follows: In Section 2 we recall Sobolev spaces with bounded mixed derivatives, the notion of the worst case integration error in reproducing kernel Hilbert spaces and the connection to periodic discrepancy. In Section 3 we discuss necessary and sufficient conditions for optimal point sets and derive lower bounds of the worst case error on certain local patches of the whole $[0, 1)^{2N}$. In Section 4 we compute candidates for optimal point sets up to machine precision. Using arbitrary precision rational arithmetic we prove that they are indeed near the global minimum which also turns out to be unique up to torus-symmetries. For certain point numbers the global minima are integration lattices as is the case if $N$ is a Fibonacci number. We close with some remarks in Section 5.

## 2 Quasi–Monte Carlo Integration in $H_{\mathrm{mix}}^1(\mathbb{T}^2)$

### 2.1 Sobolev Spaces of Periodic Functions

We consider univariate 1-periodic functions $f : \mathbb{R} \to \mathbb{R}$ which are given by their values on the torus $\mathbb{T} = [0, 1)$. For $k \in \mathbb{Z}$, the $k$-th Fourier coefficient of a function $f \in L_2(\mathbb{T})$ is given by $\hat{f}_k = \int_0^1 f(x) \exp(2\pi \mathrm{i} kx)\, \mathrm{d}x$. The definition

$$\|f\|_{H^{1,\gamma}}^2 = \hat{f}_0^2 + \gamma \sum_{k \in \mathbb{Z}} |2\pi k|^2 \hat{f}_k^2 = \left(\int_{\mathbb{T}} f(x)\, \mathrm{d}x\right)^2 + \gamma \int_{\mathbb{T}} f'(x)^2\, \mathrm{d}x \tag{1}$$

for a function $f$ in the univariate Sobolev space $H^1(\mathbb{T}) = W^{1,2}(\mathbb{T}) \subset L_2(\mathbb{T})$ of functions with first weak derivatives bounded in $L_2$ gives a Hilbert space norm $\|f\|_{H^{1,\gamma}}$ on $H^1(\mathbb{T})$ depending on the parameter $\gamma > 0$. The corresponding inner product is given by

$$(f,g)_{H^{1,\gamma}(\mathbb{T})} = \left( \int_0^1 f(x) \, \mathrm{d}x \right) \left( \int_0^1 g(x) \, \mathrm{d}x \right) + \gamma \int_0^1 f'(x) g'(x) \, \mathrm{d}x.$$

We denote the Hilbert space $H^1(\mathbb{T})$ equipped with this inner product by $H^{1,\gamma}(\mathbb{T})$.

Since $H^{1,\gamma}(\mathbb{T})$ is continuously embedded in $C^0(\mathbb{T})$ it is a reproducing kernel Hilbert space (RKHS), see [1], with a symmetric and positive definite kernel $K_{1,\gamma} : \mathbb{T} \times \mathbb{T} \to \mathbb{R}$, given by [20]

$$\begin{aligned} K_{1,\gamma}(x,y) :=& 1 + \gamma \sum_{k \in \mathbb{Z} \setminus \{0\}} |2\pi k|^{-2} \exp(2\pi i k(x-y)) \\ =& 1 + \gamma k(|x-y|), \end{aligned} \tag{2}$$

where $k(t) = \frac{1}{2}(t^2 - t + \frac{1}{6})$ is the Bernoulli polynomial of degree two divided by two.

This kernel has the property that it reproduces point evaluations in $H^1$, i.e. $f(x) = (f(\cdot), K(\cdot, x))_{H^{1,\gamma}}$ for all $f \in H^1$. The reproducing kernel of the tensor product space $H_{\mathrm{mix}}^{1,\gamma}(\mathbb{T}^2) := H^1(\mathbb{T}) \otimes H^1(\mathbb{T}) \subset C(\mathbb{T}^2)$ is the product of the univariate kernels, i.e.

$$\begin{aligned} K_{2,\gamma}(\boldsymbol{x}, \boldsymbol{y}) =& K_{1,\gamma}(x_1, y_1) \cdot K_{1,\gamma}(x_2, y_2) \\ =& 1 + \gamma k(|x_1 - y_1|) + \gamma k(|x_2 - y_2|) + \gamma^2 k(|x_1 - y_1|) k(|x_2 - y_2|). \end{aligned} \tag{3}$$

### 2.2 Quasi–Monte Carlo Cubature

A linear cubature algorithm $Q_N(f) := \frac{1}{N} \sum_{i=1}^N f(\boldsymbol{x}_i)$ with uniform weights $\frac{1}{N}$ on a point set $\mathscr{P}_N = \{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\}$ is called a QMC cubature rule. Well-known examples for point sets used in such quadrature methods are digital nets, see e.g. [4, 10], and lattice rules [15]. A two-dimensional integration lattice is a set of $N$ points given as

$$\left\{ \left( \frac{i}{N}, \frac{ig}{N} \mod 1 \right) : i = 0, \ldots, N-1 \right\}$$

for some $g \in \{1, \ldots, N-1\}$ coprime to $N$. A special case of such a rank-1 lattice rule is the so called Fibonacci lattice that only exists for $N$ being a Fibonacci number $F_n$ and is given by the generating vector $(1, g) = (1, F_{n-1})$, where $F_n$ denotes the $n$-th Fibonacci number. It is well known that the Fibonacci lattices yield the optimal rate of convergence in certain spaces of periodic functions.

In the setting of a reproducing kernel Hilbert space with kernel $K$ on a general domain $D$, the worst case error of the QMC-rule $Q_N$ can be computed as

$$\mathrm{wce}(\mathscr{H}, \mathscr{P}_N)^2 = \int_D \int_D K(\boldsymbol{x}, \boldsymbol{y}) \, \mathrm{d}\boldsymbol{x} \, \mathrm{d}\boldsymbol{y} - \frac{2}{N} \sum_{i=1}^N \int_D K(\boldsymbol{x}_i, \boldsymbol{y}) \, \mathrm{d}y + \frac{1}{N^2} \sum_{i,j=1}^N K(\boldsymbol{x}_i, \boldsymbol{x}_j),$$

which is the norm of the error functional, see e.g. [4, 11]. For the kernel $K_{2,\gamma}$ we obtain

$$\mathrm{wce}(H_{\mathrm{mix}}^{1;\gamma}(\mathbb{T}^2), \mathscr{P}_N)^2 = -1 + \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} K_{2,\gamma}(\boldsymbol{x}_i, \boldsymbol{x}_j).$$

There is a close connection between the worst case error of integration in $\mathrm{wce}(H_{\mathrm{mix}}^{1;\gamma}(\mathbb{T}^2), \mathscr{P}_N)$ for the case $\gamma = 6$ and periodic $L_2$-discrepancy, which we will describe in the following.

## 2.3 Periodic Discrepancy

The periodic $L_2$-discrepancy is measured with respect to periodic boxes. In dimension $d = 1$, periodic intervals $I(x, y)$ for $x, y \in [0, 1)$ are given by

$$I(x, y) = [x, y) \text{ if } x \leq y \qquad \text{and} \qquad I(x, y) = [x, 1) \cup [0, y) \text{ if } x > y.$$

In dimension $d > 1$, the periodic boxes $B(\boldsymbol{x}, \boldsymbol{y})$ for $\boldsymbol{x} = (x_1, \dots, x_d)$ and $\boldsymbol{y} = (y_1, \dots, y_d) \in [0, 1)^d$ are products of the one-dimensional intervals, i.e.

$$B(\boldsymbol{x}, \boldsymbol{y}) = I(x_1, y_1) \times \cdots \times I(x_d, y_d).$$

The discrepancy of a set $\mathscr{P}_N = \{\boldsymbol{x}_1, \dots, \boldsymbol{x}_N\} \subset [0, 1)^d$ with respect to such a periodic box $B = B(\boldsymbol{x}, \boldsymbol{y})$ is the deviation of the relative number of points of $\mathscr{P}_N$ in $B$ from the volume of $B$

$$D(\mathscr{P}_N, B) = \frac{\#\mathscr{P}_N \cap B}{N} - \mathrm{vol}(B).$$

Finally, the periodic $L_2$-discrepancy of $\mathscr{P}_N$ is the $L_2$-norm of the discrepancy function taken over all periodic boxes $B = B(\boldsymbol{x}, \boldsymbol{y})$, i.e.

$$D_2(\mathscr{P}_N) = \left( \int_{[0,1)^d} \int_{[0,1)^d} D(\mathscr{P}_N, B(\boldsymbol{x}, \boldsymbol{y}))^2 \, \mathrm{d}\boldsymbol{y} \, \mathrm{d}\boldsymbol{x} \right)^{1/2}.$$

It turns out, see [11, page 43] that the periodic $L_2$-discrepancy can be computed as

$$D_2(\mathscr{P}_N)^2 = -3^{-d} + \frac{1}{N^2} \sum_{\boldsymbol{x}, \boldsymbol{y} \in \mathscr{P}_N} \tilde{K}_d(\boldsymbol{x}, \boldsymbol{y})$$

$$= 3^{-d} \mathrm{wce}(H_{\mathrm{mix}}^{1,6}(\mathbb{T}^d), \mathscr{P}_N)^2,$$

where $\tilde{K}_d$ is the tensor product of $d$ kernels $\tilde{K}_1(x, y) = |x - y|^2 - |x - y| + \frac{1}{2}$. So minimizing the periodic $L_2$-discrepancy is equivalent to minimizing the worst case error in $H_{\mathrm{mix}}^{1;\gamma}$ for $\gamma = 6$. Let us also remark that the periodic $L_2$-discrepancy is (up to a factor) sometimes also called diaphony. This terminology was introduced in [22].

# 3 Optimal Cubature Points

In this section we deal with (local) optimality conditions for a set of two-dimensional points $\mathscr{P}_N \equiv (\boldsymbol{x}, \boldsymbol{y}) \subset \mathbb{T}^2$, where $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{T}^N$ denote the vectors of the first and second components of the points, respectively.

## 3.1 Optimization Problem

We want to minimize the squared worst case error

$$
\begin{aligned}
\mathrm{wce}(H_{\mathrm{mix}}^{1,\gamma}(\mathbb{T}^2), \mathscr{P}_N)^2 &= -1 + \frac{1}{N^2} \sum_{i,j=0}^{N-1} K_{1,\gamma}(x_i, x_j) K_{1,\gamma}(y_i, y_j) \\
&= -1 + \frac{1}{N^2} \sum_{i,j=0}^{N-1} \left(1 + \gamma k(|x_i - x_j|) + \gamma k(|y_i - y_j|) + \gamma^2 k(|x_i - x_j|)k(|y_i - y_j|)\right) \\
&= \frac{\gamma}{N^2} \sum_{i,j=0}^{N-1} \left(k(|x_i - x_j|) + k(|y_i - y_j|) + \gamma k(|x_i - x_j|)k(|y_i - y_j|)\right) \\
&= \frac{\gamma(2k(0) + \gamma k(0)^2)}{N} \\
&\quad + \frac{2\gamma}{N^2} \sum_{i=0}^{N-2} \sum_{j=i+1}^{N-1} \left(k(|x_i - x_j|) + k(|y_i - y_j|) + \gamma k(|x_i - x_j|)k(|y_i - y_j|)\right)
\end{aligned}
$$

Thus, minimizing $\mathrm{wce}(H_{\mathrm{mix}}^{1,\gamma}(\mathbb{T}^2), \mathscr{P}_N)^2$ is equivalent to minimizing either

$$
F_\gamma(\boldsymbol{x}, \boldsymbol{y}) := \sum_{i=0}^{N-2} \sum_{j=i+1}^{N-1} \left(k(|x_i - x_j|) + k(|y_i - y_j|) + \gamma k(|x_i - x_j|)k(|y_i - y_j|)\right) \quad (4)
$$

or

$$
G_\gamma(\boldsymbol{x}, \boldsymbol{y}) := \sum_{i,j=0}^{N-1} (1 + \gamma k(|x_i - x_j|))(1 + \gamma k(|y_i - y_j|)). \quad (5)
$$

For theoretical considerations we will sometimes use $G_\gamma$, while for the numerical implementation we will use $F_\gamma$ as objective function, since it has less summands.

Let $\tau, \sigma \in S_N$ be two permutations of $\{0, 1, \ldots, N-1\}$. Define the sets

$$
D_{\tau,\sigma} = \left\{ \boldsymbol{x} \in [0,1)^N, \boldsymbol{y} \in [0,1)^N : \begin{array}{l} x_{\tau(0)} \le x_{\tau(1)} \le \cdots \le x_{\tau(N-1)} \\ y_{\sigma(0)} \le y_{\sigma(1)} \le \cdots \le y_{\sigma(N-1)} \end{array} \right\} \quad (6)
$$

on which all points maintain the same order in both components and hence it holds $|x_i - x_j| = s_{i,j}(x_i - x_j)$ for $s_{i,j} \in \{-1, 1\}$. It follows that the restriction of $F_\gamma$ to $D_{\tau,\sigma}$,

i.e. $F_\gamma(\pmb{x},\pmb{y})_{|D_{\tau,\sigma}}$, is a polynomial of degree 4 in $(\pmb{x},\pmb{y})$. Moreover, $F_{\gamma|D_{\tau,\sigma}}$ is convex for sufficiently small $\gamma$.

**Proposition 1.** $F_\gamma(\pmb{x},\pmb{y})_{|D_{\tau,\sigma}}$ and $G_\gamma(\pmb{x},\pmb{y})_{|D_{\tau,\sigma}}$ are convex if $\gamma \in [0,6]$.

*Proof.* It is enough to prove the claim for

$$G_\gamma(\pmb{x},\pmb{y}) = \sum_{i,j=0}^{N-1} (1+\gamma k(|x_i - x_j|))(1+\gamma k(|y_i - y_j|)).$$

Since the sum of convex functions is convex and since $f(x-y)$ is convex if $f$ is, it is enough to show that $f(s,t) = \big(1+\gamma k(s)\big)\big(1+\gamma k(t)\big)$ is convex for $s,t \in [0,1]$. To this end, we show that the Hesse matrix $\mathscr{H}(f)$ is positive definite if $0 \le \gamma < 6$. First, $f_{ss} = \gamma\big(1+\gamma k(t)\big)$ is positive if $\gamma < 24$. Hence is is enough to check that the determinant of $\mathscr{H}(f)$ is positive, which is equivalent to the inequality

$$\big(1+\gamma k(s)\big)\big(1+\gamma k(t)\big) > \gamma^2 \left(s - \frac{1}{2}\right)^2 \left(t - \frac{1}{2}\right)^2.$$

So it remains to see that

$$1+\gamma k(s) = 1 + \frac{\gamma}{2}\left(s^2 - s + \frac{1}{6}\right) > \gamma\left(s - \frac{1}{2}\right)^2.$$

But this is elementary to check for $0 \le \gamma < 6$ and $s \in [0,1]$. In the case $\gamma = 6$ the determinant of $\mathscr{H}(f) = 0$ and some additional argument is necessary which we omit here. $\qquad\square$

Since
$$[0,1)^N \times [0,1)^N = \bigcup_{(\tau,\sigma)\in S_N \times S_N} D_{\tau,\sigma},$$

one can obtain the global minimum of $F_\gamma$ on $[0,1)^N \times [0,1)^N$ by computing $\arg\min_{(\pmb{x},\pmb{y})\in D_{\tau,\sigma}} F_\gamma(\pmb{x},\pmb{y})$ for all $(\tau,\sigma) \in S_N \times S_N$ and choose the global minimum as the smallest of all the local ones.


### *3.2 Using the Torus Symmetries*

We now want to analyze how symmetries of the two dimensional torus $\mathbb{T}^2$ allow to reduce the number of regions $D_{\tau,\sigma}$ for which the optimization problem has to be solved.

The symmetries of the torus $\mathbb{T}^2$ which do not change the worst case error for the considered classes of periodic functions are generated by

1. Shifts in the first coordinate $x \mapsto x + c \mod 1$ and shifts in the second coordinate $y \mapsto y + c \mod 1$.

2. Reflection of the first coordinate $x \mapsto 1 - x$ and reflection of the second coordinate $y \mapsto 1 - y$.
3. Interchanging the first coordinate $x$ and the second coordinate $y$.
4. The points are indistinguishable, hence relabeling the points does not change the worst case error.

Applying finite compositions of these symmetries to all the points in the point set $\mathscr{P}_N = \{(x_0, y_0), \ldots, (x_{N-1}, y_{N-1})\}$ leads to an equivalent point set with the same worst case integration error. This shows that the group of symmetries $G$ acting on the pairs $(\tau, \sigma)$ indexing $D_{\tau,\sigma}$ generated by the following operations

1. replacing $\tau$ or $\sigma$ by a shifted permutation: $\tau \mapsto (\tau(0) + k \mod N, \ldots, \tau(N-1) + k \mod N)$ or $\sigma \mapsto (\sigma(0) + k \mod N, \ldots, \sigma(N-1) + k \mod N)$
2. replacing $\tau$ or $\sigma$ by its flipped permutation: $\tau \mapsto (\tau(N-1), \tau(N-2), \ldots, \tau(1), \tau(0))$ or $\sigma \mapsto (\sigma(N-1), \sigma(N-2), \ldots, \sigma(1), \sigma(0))$
3. interchanging $\sigma$ and $\tau$: $(\tau, \sigma) \mapsto (\sigma, \tau)$
4. applying a permutation $\pi \in S_N$ to both $\tau$ and $\sigma$ : $(\tau, \sigma) \mapsto (\pi\tau, \pi\sigma)$

lead to equivalent optimization problems. So let us call the pairs $(\tau, \sigma)$ and $(\tau', \sigma')$ in $S_N \times S_N$ equivalent if they are in the same orbit with respect to the action of $G$. In this case we write $(\tau, \sigma) \sim (\tau', \sigma')$.

Using the torus symmetries 1. and 4. it can always be arranged that $\tau = \mathrm{id}$ and $\sigma(0) = 0$, which together with fixing the point $(x_0, y_0) = (0, 0)$ leads to the sets

$$D_\sigma = \left\{ \boldsymbol{x} \in [0,1)^N, \boldsymbol{y} \in [0,1)^N : \begin{array}{l} 0 = x_0 \le x_1 \le \ldots \le x_{N-1} \\ 0 = y_0 \le y_{\sigma(1)} \le \cdots \le y_{\sigma(N-1)} \end{array} \right\}, \qquad (7)$$

where $\sigma \in S_{N-1}$ denotes a permutation of $\{1, 2, \ldots, N-1\}$.

But there are many more symmetries and it would be algorithmically desirable to cycle through exactly one representative of each equivalence class without ever touching the other equivalent $\sigma$. This seems to be difficult to implement, hence we settled for a little less which still reduces the amount of permutations to be handled considerably.

To this end, let us define the symmetrized metric

$$d(i, j) = \min\{|i - j|, N - |i - j|\} \qquad \text{for} \qquad 0 \le i, j \le N - 1 \qquad (8)$$

and the following subset of $S_N$.

**Definition 1.** The set of *semi-canonical* permutations $\mathfrak{C}_N \subset S_N$ consists of permutations $\sigma$ which fulfill

(i) $\sigma(0) = 0$
(ii) $d(\sigma(1), \sigma(2)) \le d(0, \sigma(N-1))$
(iii) $\sigma(1) = \min\{d(\sigma(i), \sigma(i+1)) \mid i = 0, 1, \ldots, N-1\}$
(iv) $\sigma$ is lexicographically smaller than $\sigma^{-1}$.

Here we identify $\sigma(N)$ with $0 = \sigma(0)$.

This means that $\sigma$ is *semi-canonical* if the distance between $0 = \sigma(0)$ and $\sigma(1)$ is minimal among all distances between $\sigma(i)$ and $\sigma(i+1)$, which can be arranged by a shift. Moreover, the distance between $\sigma(1)$ and $\sigma(2)$ is at most as large as the distance between $\sigma(0)$ and $\sigma(N-1)$, which can be arranged by a reflection and a shift if it is not the case. Hence we have obtained the following lemma.

**Lemma 1.** *For any permutation $\sigma \in S_N$ with $\sigma(0) = 0$ there exists a semi-canonical $\sigma'$ such that the sets $D_\sigma$ and $D_{\sigma'}$ are equivalent up to torus symmetry.*

Thus we need to consider only semi-canonical $\sigma$ which is easy to do algorithmically.

*Remark 1.* If $\sigma \in S_N$ is semi-canonical, it holds $\sigma(1) \leq N/2$.

Another main advantage in considering our objective function only in domains $D_\sigma$ is that it is not only convex but strictly convex here. This is due to the fact that we fix $(x_0, y_0) = (0, 0)$.

**Proposition 2.** $F_\gamma(\boldsymbol{x}, \boldsymbol{y})_{|D_\sigma}$ *and* $G_\gamma(\boldsymbol{x}, \boldsymbol{y})_{|D_\sigma}$ *are strictly convex if* $\gamma \in [0, 6]$.

*Proof.* Again it is enough to prove the claim for

$$G_\gamma(\boldsymbol{x}, \boldsymbol{y}) = \sum_{i,j=0}^{N-1} (1 + \gamma k(|x_i - x_j|))(1 + \gamma k(|y_i - y_j|)).$$

Now we use that the sum of a convex and a strictly convex function is again strictly convex. Hence it is enough to show that the function

$$
\begin{aligned}
f(x_1, \ldots, x_{N-1}, y_1, \ldots, y_{N-1}) &= \sum_{i=1}^{N-1} (1 + \gamma k(|x_i - x_0|))(1 + \gamma k(|y_i - y_0|)) \\
&= \sum_{i=1}^{N-1} (1 + \gamma k(x_i))(1 + \gamma k(y_i))
\end{aligned}
$$

is strictly convex on $[0, 1]^{N-1} \times [0, 1]^{N-1}$. In the proof of Proposition 1 it was actually shown that $f_i(x_i, y_i) = (1 + \gamma k(x_i))(1 + \gamma k(y_i))$ is strictly convex for $(x_i, y_i) \in [0, 1]^2$ for each fixed $i = 1, \ldots, N-1$. Hence the strict convexity of $f$ follows from the following easily verified lemma. $\qquad\square$

**Lemma 2.** *Let* $f_i : D_i \to \mathbb{R}, i = 1, \ldots, m$ *be strictly convex functions on the convex domains* $D_i \in \mathbb{R}^{d_i}$. *Then the function*

$$f : D = D_1 \times \cdots \times D_m \to \mathbb{R}, (z_1, \ldots, z_m) \mapsto \sum_{i=1}^{m} f_i(z_i)$$

*is strictly convex.*

Hence we have indeed a unique point in each $D_\sigma$ where the minimum of $F_\gamma$ is attained.

### 3.3 Minimizing $F_\gamma$ on $D_\sigma$

Our strategy will be to compute the local minimum of $F_\gamma$ on each region $D_\sigma \subset [0,1)^N \times [0,1)^N$ for all semi-canonical permutations $\sigma \in \mathfrak{C}_N \subset S_N$ and determine the global minimum by choosing the smallest of all the local ones.

This gives for each $\sigma \in \mathfrak{C}_N$ the constrained optimization problem

$$\min_{(\boldsymbol{x},\boldsymbol{y}) \in D_\sigma} F_\gamma(\boldsymbol{x},\boldsymbol{y}) \quad \text{subject to } v_i(\boldsymbol{x}) \geq 0 \text{ and } w_i(\boldsymbol{y}) \geq 0 \text{ for all } i = 1,\ldots,N-1, \quad (9)$$

where the inequality constraints are linear and given by

$$v_i(\boldsymbol{x}) = x_i - x_{i-1} \quad \text{and} \quad w_i(\boldsymbol{y}) = y_{\sigma(i)} - y_{\sigma(i-1)} \quad \text{for } i = 1,\ldots,N-1. \quad (10)$$

In order to use the necessary (and due to local strict convexity also sufficient) conditions for local minima

$$\frac{\partial}{\partial x_k} F_\gamma(\boldsymbol{x},\boldsymbol{y}) = 0 \quad \text{and} \quad \frac{\partial}{\partial y_k} F_\gamma(\boldsymbol{x},\boldsymbol{y}) = 0 \quad \text{for } k = 1,\ldots,N-1$$

for $(\boldsymbol{x},\boldsymbol{y}) \in D_\sigma$ we need to evaluate the partial derivatives of $F_\gamma$.

**Proposition 3.** *For a given permutation $\sigma \in \mathfrak{C}_N$ the partial derivative of $F_{\gamma|D_\sigma}$ with respect to the second component $\boldsymbol{y}$ is given by*

$$\frac{\partial}{\partial y_k} F_\gamma(\boldsymbol{x},\boldsymbol{y})_{|D_\sigma} = y_k \left( \sum_{\substack{i=0 \\ i \neq k}}^{N-1} c_{i,k} \right) - \sum_{\substack{i=0 \\ i \neq k}}^{N-1} c_{i,k} y_i + \frac{1}{2} \left( \sum_{i=0}^{k-1} c_{i,k} s_{i,k} - \sum_{j=k+1}^{N-1} c_{k,j} s_{k,j} \right),$$
$$(11)$$

*where $s_{i,j} = sgn(y_i - y_j)$ and $c_{i,j} := 1 + \gamma k(|x_i - x_j|) = c_{j,i}$.*

*Interchanging $\boldsymbol{x}$ and $\boldsymbol{y}$ the same result holds for the partial derivatives with respect to $\boldsymbol{x}$ with the obvious modification to $c_{i,j}$ and the simplification that $s_{i,j} = -1$.*

*The second order derivatives with respect to $\boldsymbol{y}$ are given by*

$$\frac{\partial^2}{\partial y_k \partial y_j} F(\boldsymbol{x},\boldsymbol{y})_{|D_\sigma} = \begin{cases} \sum_{i=0}^{k-1} c_{i,k} + \sum_{i=k+1}^{N-1} c_{i,k} & \text{for } j = k \\ -c_{k,j} & \text{for } j \neq k \end{cases}, \quad k,j \in \{1,\ldots,N-1\}$$
$$(12)$$

*Again, the analogue for $\frac{\partial^2}{\partial x_k \partial x_j} F(\boldsymbol{x},\boldsymbol{y})_{|D_\sigma}$ is obtained with the obvious modification $c_{i,j} = 1 + \gamma k(|y_i - y_j|)$.*

*Proof.* We prove the claim for the partial derivative with respect to $\boldsymbol{y}$:

$$\frac{\partial}{\partial y_k} F_\gamma(\boldsymbol{x}, \boldsymbol{y}) = \sum_{i=0}^{N-2} \sum_{j=i+1}^{N-1} \frac{\partial}{\partial y_k} k(|y_i - y_j|) \underbrace{(1 + \gamma k(|x_i - x_j|))}_{=:c_{i,j}} + \frac{\partial}{\partial y_k} k(|x_i - x_j|)$$

$$= \sum_{i=0}^{N-2} \sum_{j=i+1}^{N-1} c_{i,j} \frac{\partial}{\partial y_k} k(|y_i - y_j|)$$

$$= \sum_{i=0}^{N-2} \sum_{j=i+1}^{N-1} c_{i,j} \, k'(s_{i,j}(y_i - y_j)) \cdot \begin{cases} s_{i,j} & \text{for } i = k \\ -s_{i,j} & \text{for } j = k \\ 0 & \text{else} \end{cases}$$

$$= \sum_{j=k+1}^{N-1} c_{k,j} s_{k,j} \left( s_{k,j}(y_k - y_j) - \frac{1}{2} \right) - \sum_{i=0}^{k-1} c_{i,k} s_{i,k} \left( s_{i,k}(y_i - y_k) - \frac{1}{2} \right)$$

$$= y_k \left( \sum_{\substack{i=0 \\ i \neq k}}^{N-1} c_{i,k} \right) - \sum_{\substack{i=0 \\ i \neq k}}^{N-1} c_{i,k} y_i + \frac{1}{2} \left( \sum_{i=0}^{k-1} c_{i,k} s_{i,k} - \sum_{j=k+1}^{N-1} c_{k,j} s_{k,j} \right).$$

From this we immediately get the second derivative (12).          □

### 3.4 Lower Bounds of $F_\gamma$ on $D_\sigma$

Until now we are capable of approximating local minima of $F_\gamma$ on a given $D_\sigma$. If this is done for all $\sigma \in \mathfrak{C}_N$ we can obtain a candidate for a global minimum, but due to the finite precision of floating point arithmetic one can never be sure to be close to the actual global minimum. However, it is also possible to compute a lower bound for the optimal point set for each $D_\sigma$ using Wolfe-duality for constrained optimization. It is known [12] that for a convex problem with linear inequality constraints like (9) the Lagrangian

$$\mathcal{L}_F(\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{\lambda}, \boldsymbol{\mu}) := F(\boldsymbol{x}, \boldsymbol{y}) - \boldsymbol{\lambda}^T \boldsymbol{v}(\boldsymbol{x}) - \boldsymbol{\mu}^T \boldsymbol{w}(\boldsymbol{y}) \tag{13}$$

$$= F(\boldsymbol{x}, \boldsymbol{y}) - \sum_{i=1}^{N-1} (\lambda_i v_i(\boldsymbol{x}) + \mu_i w_i(\boldsymbol{y})) \tag{14}$$

gives a lower bound on $F$, i.e.

$$\min_{(\boldsymbol{x}, \boldsymbol{y}) \in D_\sigma} F(\boldsymbol{x}, \boldsymbol{y}) \geq \mathcal{L}_F(\tilde{\boldsymbol{x}}, \tilde{\boldsymbol{y}}, \boldsymbol{\lambda}, \boldsymbol{\mu})$$

for all $(\tilde{\boldsymbol{x}}, \tilde{\boldsymbol{y}}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ that fulfill the constraint

$$\nabla_{(\boldsymbol{x}, \boldsymbol{y})} \mathcal{L}_F(\tilde{\boldsymbol{x}}, \tilde{\boldsymbol{y}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = 0 \quad \text{and} \quad \boldsymbol{\lambda}, \boldsymbol{\mu} \geq 0 \text{ (component-wise)}. \tag{15}$$

Here, $\nabla_{(\boldsymbol{x}, \boldsymbol{y})} = (\nabla_{\boldsymbol{x}}, \nabla_{\boldsymbol{y}})$, where $\nabla_{\boldsymbol{x}}$ denotes the gradient of a function with respect to the variables in $\boldsymbol{x}$. Hence it is our goal to find for each $D_\sigma$ such an admissible point

$(\tilde{\boldsymbol{x}}, \tilde{\boldsymbol{y}}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ which yields a lower bound that is larger than some given candidate for the global minimum. If the relevant computations are carried out in infinite precision rational number arithmetic these bounds are mathematically reliable.

In order to accomplish this we first have to compute the Lagrangian of (9). To this end, let $\mathbf{P}_\sigma \in \{-1,0,1\}^{(N-1)\times(N-1)}$ denote the permutation matrix corresponding to $\sigma \in S_{N-1}$ and

$$
\boldsymbol{B} := \begin{pmatrix} 1 & -1 & 0 & \dots & 0 & 0 \\ 0 & 1 & -1 & \dots & 0 & 0 \\ \vdots & & & \ddots & & \vdots \\ 0 & & \dots & 0 & 1 & -1 \\ 0 & & \dots & & 0 & 1 \end{pmatrix} \in \mathbb{R}^{(N-1)\times(N-1)}. \tag{16}
$$

Then the partial derivatives of $\mathscr{L}_F$ with respect to $\boldsymbol{x}$ and $\boldsymbol{y}$ are given by

$$
\nabla_{\boldsymbol{x}}\mathscr{L}_F(\boldsymbol{x},\boldsymbol{y},\boldsymbol{\lambda},\boldsymbol{\mu}) = \nabla_{\boldsymbol{x}}F(\boldsymbol{x},\boldsymbol{y}) - \begin{pmatrix} \lambda_1 - \lambda_2 \\ \vdots \\ \lambda_{N-2} - \lambda_{N-1} \\ \lambda_{N-1} \end{pmatrix} = \nabla_{\boldsymbol{x}}F(\boldsymbol{x},\boldsymbol{y}) - \boldsymbol{B}\boldsymbol{\lambda} \tag{17}
$$

and

$$
\nabla_{\boldsymbol{y}}\mathscr{L}_F(\boldsymbol{x},\boldsymbol{y},\boldsymbol{\lambda},\boldsymbol{\mu}) = \nabla_{\boldsymbol{y}}F(\boldsymbol{x},\boldsymbol{y}) - \begin{pmatrix} \mu_{\sigma(1)} - \mu_{\sigma(2)} \\ \vdots \\ \mu_{\sigma(N-2)} - \mu_{\sigma(N-1)} \\ \mu_{\sigma(N-1)} \end{pmatrix} = \nabla_{\boldsymbol{y}}F(\boldsymbol{x},\boldsymbol{y}) - \boldsymbol{B}\mathbf{P}_\sigma\boldsymbol{\mu}. \tag{18}
$$

This leads to the following theorem.

**Theorem 1.** *For $\sigma \in \mathfrak{C}_N$ and $\delta > 0$ let the point $(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) \in D_\sigma$ fulfill*

$$
\frac{\partial}{\partial x_k}F(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) = \delta \quad \text{and} \quad \frac{\partial}{\partial y_k}F(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) = \delta \quad \text{for } k = 1, \dots, N-1. \tag{19}
$$

*Then*

$$
F(\boldsymbol{x},\boldsymbol{y}) \geq F(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) - \delta \sum_{i=1}^{N-1} \left( (N-i)\cdot v_i(\tilde{\boldsymbol{x}}_\sigma) + \sigma(N-i)w_i(\tilde{\boldsymbol{y}}_\sigma) \right) \tag{20}
$$

$$
> F(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) - \delta N^2 \tag{21}
$$

*holds for all $(\boldsymbol{x},\boldsymbol{y}) \in D_\sigma$.*

*Proof.* Choosing

$$
\boldsymbol{\lambda} = \boldsymbol{B}^{-1}\nabla_{\boldsymbol{x}}F(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) \quad \text{and} \quad \boldsymbol{\mu} = \mathbf{P}_\sigma^{-1}\boldsymbol{B}^{-1}\nabla_{\boldsymbol{y}}F(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) \tag{22}
$$

yields

$$\nabla_{\boldsymbol{x}} F(\tilde{\boldsymbol{x}}, \tilde{\boldsymbol{y}}) = \boldsymbol{B}\boldsymbol{\lambda} \quad \text{and} \quad \nabla_{\boldsymbol{y}} F(\tilde{\boldsymbol{x}}, \tilde{\boldsymbol{y}}) = \boldsymbol{B}\mathbf{P}_{\sigma}\boldsymbol{\mu}. \tag{23}$$

A short computation shows that the inverse of $\boldsymbol{B}$ from (16) is given by

$$\boldsymbol{B}^{-1} := \begin{pmatrix} 1 & 1 & \dots & 1 \\ 0 & 1 & \dots & 1 \\ \vdots & 0 & \ddots & \vdots \\ 0 & \dots & 0 & 1 \end{pmatrix} \in \mathbb{R}^{(N-1)\times(N-1)},$$

which yields $\boldsymbol{y}, \boldsymbol{\lambda} > 0$ and hence by Wolfe duality gives (20). The second inequality (21) then follows from noting that both $|v_i(\boldsymbol{x})|$ and $|w_i(\boldsymbol{y})|$ are bounded by 1 and $2\sum_{i=1}^{N-1} \sigma(N-i) = 2\sum_{i=1}^{N-1} i = (N-1)(N-2) < N^2$. $\qquad \square$

Now, suppose we had some candidate $(\boldsymbol{x}^*, \boldsymbol{y}^*) \in D_{\sigma^*}$ for an optimal point set. If we can find for all other $\sigma \in \mathfrak{C}_N$ points $(\tilde{\boldsymbol{x}}_{\sigma}, \tilde{\boldsymbol{y}}_{\sigma})$ that fulfills (19) and

$$F(\tilde{\boldsymbol{x}}_{\sigma}, \tilde{\boldsymbol{y}}_{\sigma}) - \delta N^2 \geq F_{\gamma}(\boldsymbol{x}^*, \boldsymbol{y}^*)$$

for some $\delta > 0$, we can be sure that $D_{\sigma^*}$ is (up to torus symmetry) the unique domain $D_{\sigma}$ that contains the globally optimal point set.

## 4 Numerical Investigation of Optimal point sets

In this section we numerically obtain optimal point sets with respect to the worst case error in $H_{\mathrm{mix}}^1$. Moreover, we present a proof by exhaustion that these point sets are indeed approximations to the unique (modulo torus symmetry) minimizers of $F_{\gamma}$. Since integration lattices are local minima, if the $D_{\sigma}$ containing the global minimizer corresponds to an integration lattice, this integration lattice is the exact global minimizer.

### 4.1 Numerical Minimization with Alternating Directions

In order to obtain the global minimum $(\boldsymbol{x}^*, \boldsymbol{y}^*)$ of $F_{\gamma}$ we are going to compute

$$\sigma^* := \arg\min_{\sigma \in \mathfrak{C}_N} \min_{(\boldsymbol{x}, \boldsymbol{y}) \in D_{\sigma}} F_{\gamma}(\boldsymbol{x}, \boldsymbol{y}), \tag{24}$$

where the inner minimum has a unique solution due to Proposition 2. Moreover, since $D_{\sigma}$ is a convex domain we know that the local minimum of $F_{\gamma}(\boldsymbol{x}, \boldsymbol{y})_{|D_{\sigma}}$ is not on the boundary. Hence we can restrict our search for optimal point sets to the interior of $D_{\sigma}$, where $F_{\gamma}$ is differentiable.

Instead of directly employing a local optimization technique, we will make use of the special structure of $F_{\gamma}$. While $F_{\gamma}(\boldsymbol{x}, \boldsymbol{y})_{|D_{\sigma}}$ is a polynomial of degree four, the

---

**Algorithm 1:** Alternating minimization algorithm. For off-set $\delta = 0$ it finds local minima of $F_\gamma$. For $\delta > 0$ it obtains feasible points used by Algorithm 2.

---

**Given:** Permutation $\sigma \in \mathfrak{C}_N$, tolerance $\varepsilon > 0$ and off-set $\delta \geq 0$.
**Initialize:**

1. $\boldsymbol{x}^{(0)} := (0, \frac{1}{N}, \ldots, \frac{N-1}{N})$ and $\boldsymbol{y}^{(0)} = (0, \frac{\sigma(1)}{N}, \ldots, \frac{\sigma(N-1)}{N})$.
2. $k := 0$.

**repeat**

    1. compute $\boldsymbol{H_x} := \left( \partial_{x_i} \partial_{x_j} F_\gamma(\boldsymbol{x}^{(k)}, \boldsymbol{y}^{(k)}) \right)_{i,j=1}^N$ and $\nabla_{\boldsymbol{x}} = \left( \partial_{x_i} F_\gamma(\boldsymbol{x}^{(k)}, \boldsymbol{y}^{(k)}) \right)_{i=1}^N$ by (12) and (11).

    2. Update $\boldsymbol{x}^{(k+1)} := \boldsymbol{H_x}^{-1} (\nabla_{\boldsymbol{x}} + \delta \boldsymbol{1})$ via Cholesky factorization.

    3. compute $\boldsymbol{H_y} := \left( \partial_{y_i} \partial_{y_j} F_\gamma(\boldsymbol{x}^{(k+1)}, \boldsymbol{y}^{(k)}) \right)_{i,j=1}^N$ and $\nabla_{\boldsymbol{y}} = \left( \partial_{y_i} F_\gamma(\boldsymbol{x}^{(k+1)}, \boldsymbol{y}^{(k)}) \right)_{i=1}^N$.

    4. Update $\boldsymbol{y}^{(k+1)} := \boldsymbol{H_y}^{-1} (\nabla_{\boldsymbol{y}} + \delta \boldsymbol{1})$ via Cholesky factorization.

    5. $k := k+1$.

**until** $\sqrt{\|\nabla_{\boldsymbol{x}}\|^2 + \|\nabla_{\boldsymbol{y}}\|^2} < \varepsilon$;
**Output**: point set $(\boldsymbol{x}, \boldsymbol{y}) \in D_\sigma$ with $\nabla_{\boldsymbol{x}} F_\gamma(\boldsymbol{x}, \boldsymbol{y}) \approx \delta \boldsymbol{1}$ and $\nabla_{\boldsymbol{y}} F_\gamma(\boldsymbol{x}, \boldsymbol{y}) \approx \delta \boldsymbol{1}$.

---

functions

$$\boldsymbol{x} \mapsto F_\gamma(\boldsymbol{x}, \boldsymbol{y}_0)_{|D_\sigma} \quad \text{and} \quad \boldsymbol{y} \mapsto F_\gamma(\boldsymbol{x}_0, \boldsymbol{y})_{|D_\sigma}, \tag{25}$$

where one coordinate direction is fixed, are quadratic polynomials, which have unique minima in $D_\sigma$. We are going to use this property within an alternating minimization approach. This means, that the objective function $F$ is not minimized along all coordinate directions simultaneously, but with respect to certain successively alternating blocks of coordinates. If these blocks have size one this method is usually referred to as *coordinate descent* [7] or nonlinear Gauß-Seidel method [5]. It is successfully employed in various applications, like e.g. expectation maximization or tensor approximation [9, 19].

In our case we will alternate between minimizing $F_\gamma(\boldsymbol{x}, \boldsymbol{y})$ along the first coordinate block $\boldsymbol{x} \in (0,1)^{N-1}$ and the second one $\boldsymbol{y} \in (0,1)^{N-1}$, which can be done exactly due to the quadratic polynomial property of the partial objectives (25). The method is outlined in Algorithm 1, which for threshold-parameter $\delta = 0$ approximates the local minimum of $F_\gamma$ on $D_\sigma$. For $\delta > 0$ it obtains feasible points that fulfill (19), i.e. $\nabla_{(\boldsymbol{x},\boldsymbol{y})} F_\gamma = (\delta, \ldots, \delta) = \delta \boldsymbol{1}$. Linear convergence of the alternating optimization method for strictly convex functions was for example proven in [13, 2].

## 4.2 Obtaining Lower Bounds

By now we are able to obtain a point set $(\boldsymbol{x}^*, \boldsymbol{y}^*) \in D_{\sigma^*}$ as a candidate for a global minimum of $F_\gamma$ by finding local minima on each $D_\sigma, \sigma \in \mathfrak{C}_N$. On first sight we can not be sure that we chose the right $\sigma^*$, because the value of $\min_{(\boldsymbol{x},\boldsymbol{y}) \in D_\sigma} F_\gamma(\boldsymbol{x}, \boldsymbol{y})$ can only be computed numerically.

---

**Algorithm 2:** Computation of lower bound on $D_\sigma$.

---

**Given:** Optimal point candidate $\mathscr{P}_N := (\boldsymbol{x}^*, \boldsymbol{y}^*) \in D_\sigma$ with $\sigma \in \mathfrak{C}_N$, tolerance $\varepsilon > 0$ and off-set $\theta \geq 0$.
**Initialize:**

1. Compute $\theta_N := F_\gamma(\boldsymbol{x}^*, \boldsymbol{y}^*)$ (in APR arithmetic).
2. $\Xi_N := \emptyset$.

**for all** $\sigma \in \mathfrak{C}_N$ **do**

    1. Find $(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) \in D_\sigma$ s.t. $\nabla_{(\boldsymbol{x},\boldsymbol{y})} F_\gamma(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) \approx \delta \boldsymbol{1}$ by Algorithm 1.
    2. Compute $\boldsymbol{\lambda} := \boldsymbol{B}^{-1} \nabla_{\boldsymbol{x}} F(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma)$ and $\boldsymbol{\mu} := \boldsymbol{P}_\sigma^{-1} \boldsymbol{B}^{-1} \nabla_{\boldsymbol{y}} F(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma)$ (in APR arithmetic).
    3. Verify $\boldsymbol{\lambda}, \boldsymbol{\mu} > 0$.
    4. Evaluate $\beta_\sigma := \mathscr{L}_{F_\gamma}(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma, \boldsymbol{\lambda}, \boldsymbol{\mu})$ (in APR arithmetic).
    5. **If** ( $\beta_\sigma \leq \theta_N$ ) $\Xi_N := \Xi_N \cup \sigma$.

**end**
**Output**: Set $\Xi$ of permutations $\sigma$ in which $D_\sigma$ contained a lower bound smaller than $\theta_N$.

---

On the other hand, Theorem 1 allows to compute lower bounds for all the other domains $D_\sigma$ with $\sigma \in \mathfrak{C}_N$. If we were able to obtain for each $\sigma$ a point $(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma)$, such that

$$\min_{(\boldsymbol{x},\boldsymbol{y}) \in D_{\sigma^*}} F_\gamma(\boldsymbol{x}, \boldsymbol{y}) \approx \theta_N := F_\gamma(\boldsymbol{x}^*, \boldsymbol{y}^*) < \mathscr{L}_F(\tilde{\boldsymbol{x}}_\sigma, \tilde{\boldsymbol{y}}_\sigma) - 2N^2 \delta \leq F_\gamma(\boldsymbol{x}, \boldsymbol{y}),$$

we could be sure that the global optimum is indeed located in $D_{\sigma^*}$ and $(\boldsymbol{x}^*, \boldsymbol{y}^*)$ is a good approximation to it. Of course certain computations can not be done in standard double floating point arithmetic. Instead we use arbitrary precision rational number (APR) arithmetic from the GNU Multiprecision library GMP from http://www.gmplib.org. Compared to standard floating point arithmetic in double precision this is very expensive, but it has only to be used at certain parts of the algorithm. The resulting procedure is outlined in Algorithm 2, where we marked those parts which require APR arithmetic.

### 4.3 Results

In Figures 1 and 2 the optimal point sets for $N = 2, \ldots, 16$ and both $\gamma = 1$ and $\gamma = 6$ are plotted. It can be seen that they are close to lattice point sets, which justifies using them as start points in Algorithm 1. The distance to lattice points seems to be small if $\gamma$ is small.

In Table 1 we list the permutations $\sigma$ for which $D_\sigma$ contains an optimal set of cubature points. In the second column the total number of semi-canonical permutations $\mathfrak{C}_N$ that had to be considered is shown. It grows approximately like $\frac{1}{2}(N-2)!$. Moreover, we computed the minimal worst case error and periodic $L_2$-discrepancies.

| $N$ | $\lvert \mathfrak{C}_N \rvert$ | $\mathrm{wce}(H^{1,1}_{\mathrm{mix}}, \mathscr{P}^*_N)$ | $D_2(\mathscr{P}^*_N)$ | $\sigma^*$ | Lattice |
|---|---|---|---|---|---|
| **1** | 0 | 0.416667 | 0.372678 | (0) | ✓ |
| **2** | 1 | 0.214492 | 0.212459 | (0 1) | ✓ |
| **3** | 1 | 0.146109 | 0.153826 | (0 1 2) | ✓ |
| 4 | 2 | 0.111307 | 0.121181 | (0 1 3 2) | |
| **5** | 5 | 0.0892064 | 0.0980249 | (0 2 4 1 3) | ✓ |
| 6 | 13 | 0.0752924 | 0.0850795 | (0 2 4 1 5 3) | |
| 7 | 57 | 0.0650941 | 0.0749072 | (0 2 4 6 1 3 5), (0 3 6 2 5 1 4) | ✓ |
| **8** | 282 | 0.056846 | 0.0651562 | (0 3 6 1 4 7 2 5) | ✓ |
| 9 | 1,862 | 0.0512711 | 0.0601654 | (0 2 6 3 8 5 1 7 4), (0 2 7 4 1 6 3 8 5) | |
| 10 | 14,076 | 0.0461857 | 0.054473 | (0 3 7 1 4 9 6 2 8 5) | |
| 11 | 124,995 | 0.0422449 | 0.050152 | (0 3 8 1 6 10 4 7 2 9 5), (0 3 9 5 1 7 10 4 8 2 6) | |
| 12 | 1,227,562 | 0.0370732 | 0.0456259 | (0 5 10 3 8 1 6 11 4 9 2 7) | ✓ |
| **13** | 13,481,042 | 0.0355885 | 0.0421763 | (0 5 10 2 7 12 4 9 1 6 11 3 8) | ✓ |
| 14 | 160,456,465 | 0.0333232 | 0.0400524 | (0 5 10 2 8 13 4 11 6 1 9 3 12 7), (0 5 10 3 12 7 1 9 4 13 6 11 2 8) | |
| 15 | 2,086,626,584 | 0.0312562 | 0.0379055 | (0 4 9 13 6 1 11 3 8 14 5 10 2 12 7), (0 5 11 2 7 14 9 3 12 6 1 10 4 13 8), (0 5 11 2 8 13 4 10 1 6 14 9 3 12 7), (0 5 11 2 8 13 6 1 10 4 14 7 12 3 9) | |
| 16 | 29,067,602,676 | 0.0294507 | 0.0359673 | (0 3 11 5 14 9 1 7 12 4 15 10 2 6 13 8), (0 3 11 6 13 1 9 4 15 7 12 2 10 5 14 8) | |

**Table 1** List of semi-canonical permutations $\sigma$, such that $D_\sigma$ contains an optimal set of cubature points for $N = 1, \ldots, 16$.

In some cases we found more than one semi-canonical permutation $\sigma$ for which $D_\sigma$ contained a point set which yields the optimal worst case error. Nevertheless, they represent equivalent permutations. In the following list, the torus symmetries used to show the equivalency of the permutations are given. All operations are modulo 1.

- $N = 7$: $(x, y) \mapsto (1 - y, x)$
- $N = 9$: $(x, y) \mapsto (y - 2/9, x - 1/9)$
- $N = 11$: $(x, y) \mapsto (y + 5/11, x - 4/11)$
- $N = 14$: $(x, y) \mapsto (x - 4/14, y + 6/14)$
- $N = 15$: $(x, y) \mapsto (y + 3/15, x + 2/15), (y - 2/15, 12/15 - x), (y - 6/15, 4/15 - x)$
- $N = 16$: $(x, y) \mapsto (1/16 - x, 3/16 - y)$

In all the examined cases $N \in \{2, \ldots, 16\}$ Algorithm 2 produced sets $\Xi_N$ which contained exactly the permutations that were previously obtained by Algorithm 1 and are listed in Table 1. Thus we can be sure, that the respective $D_\sigma$ contained minimizers of $F_\gamma$, which on each $D_\sigma$ are unique. Hence we know that our numerical approximation of the minimum is close to the true global minimum, which (modulo torus symmetries) is unique. In the cases $N = 1, 2, 3, 5, 7, 8, 12, 13$ the obtained global minima are integration lattices.

## 5 Conclusion

In the present paper we computed optimal point sets for quasi–Monte Carlo cubature of bivariate periodic functions with mixed smoothness of order one by decomposing the required global optimization problem into approximately $(N-2)!/2$ local ones. Moreover, we computed lower bounds for each local problem using arbitrary precision rational number arithmetic. Thereby we obtained that our approximation of the global minimum is in fact close to the real solution.

In the special case of $N$ being a Fibonacci number our approach showed that for $N \in \{1,2,3,5,8,13\}$ the Fibonacci lattice is the unique global minimizer of the worst case integration error in $H^1_{\mathrm{mix}}$. We strongly conjecture that this is true for all Fibonacci numbers. Also in the cases $N = 7,12$, the global minimizer is the obtained integration lattice.

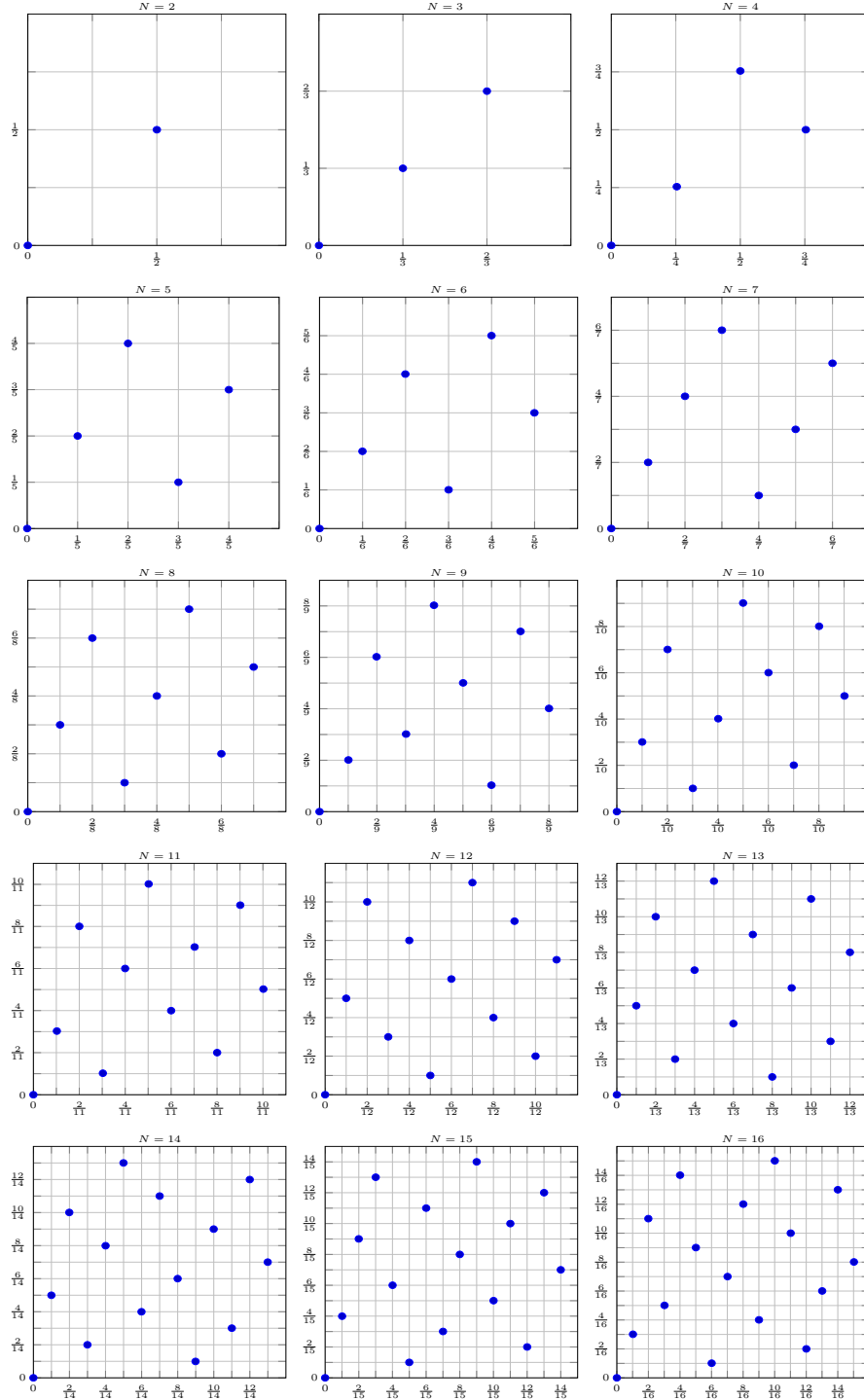In the future we are planning to prove that optimal points are close to lattice points. Moreover, we will investigate $H^r_{\mathrm{mix}}$, i.e. Sobolev spaces with dominating mixed smoothness of order $r \geq 2$ and other suitable kernels and discrepancies.
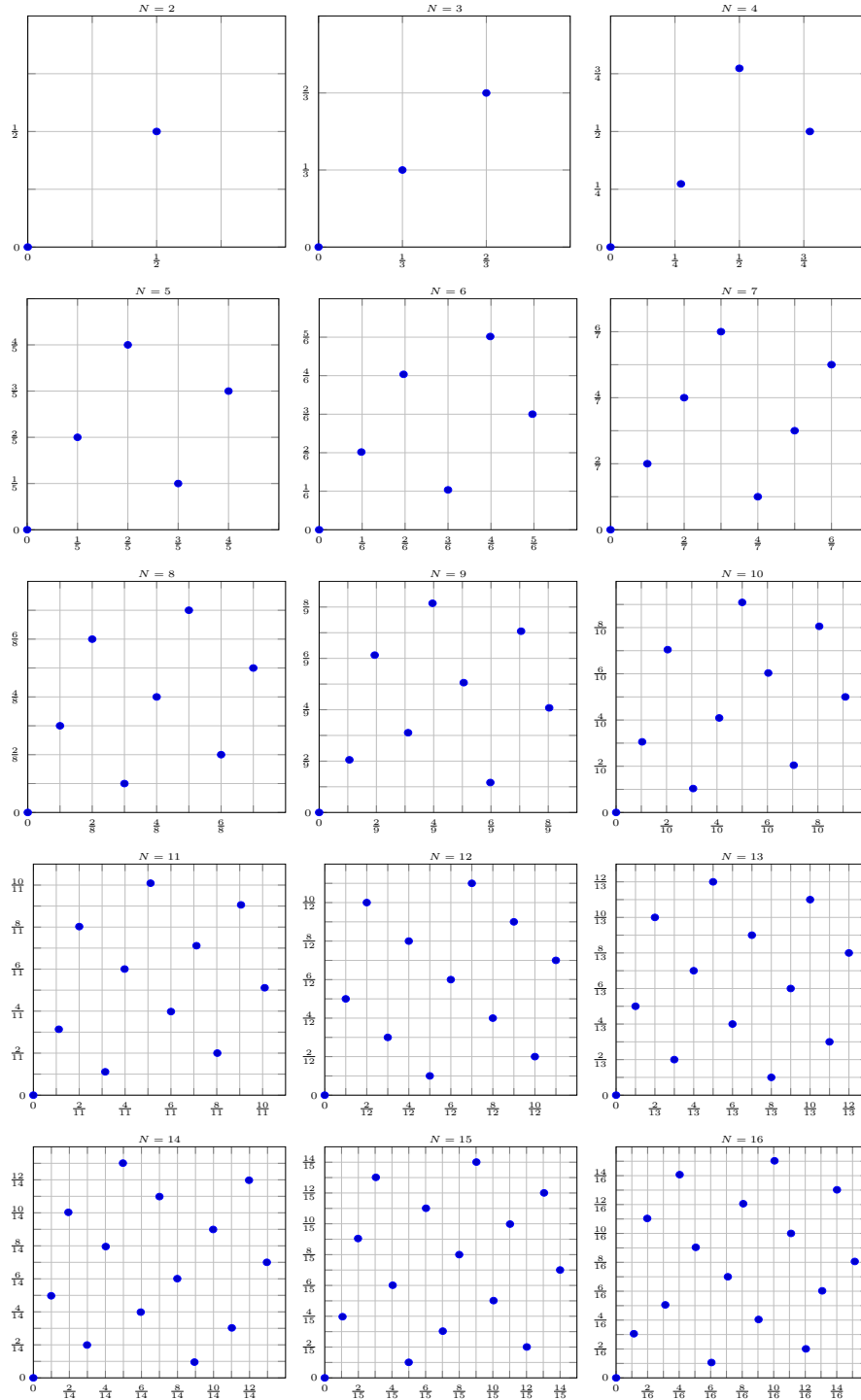
## References

1. N. Aronszajn: Theory of Reproducing Kernels. Transactions of the American Mathematical Society 68:1950, 337–404.
2. J. C. Bezdek, R. J. Hathaway, R. E. Howard, C. A. Wilson, M. P. Windham: Local convergence analysis of a grouped variable version of coordinate descent. J. of Optimization Theory and Applications 54(3):1987, 471–477.
3. D. Bilyk, V. N. Temlyakov, R. Yu: Fibonacci sets and symmetrization in discrepancy theory. J. of Complexity 28:2012, 18–36.
4. J. Dick, F. Pillichshammer: *Digital Nets and Sequences: Discrepancy Theory and Quasi–Monte Carlo Integration*, Cambridge University Press, 2010.
5. L. Grippo, M. Sciandrone: On the convergence of the block nonlinear GaußSeidel method under convex constraints. Operations Research Letters 26(3):2000, 127–136.
6. G. Larcher, F. Pillichshammer: A note on optimal point distributions in $[0,1)^s$. J. of Computational and Applied Mathematics 206:2007, 977–985.
7. Z. Q. Luo, P. Tseng: On the convergence of the coordinate descent method for convex differentiable minimization. J. of Optimization Theory and Applications 72(1):1992, 7–35
8. H. Niederreiter, I. H. Sloan: Integration of nonperiodic functions of two variables by Fibonacci lattice rules. J. of Computational and Applied Mathematics 51:1994, 57–70.
9. G.J. McLachlan and T. Krishnan. The EM Algorithm and Extensions. Wiley series in probability and statistics. John Wiley & Sons, 1997.
10. H. Niederreiter: *Quasi-Monte Carlo Methods and Pseudo-Random Numbers*, Society for Industrial and Applied Mathematics. 1987.
11. E. Novak, H. Woźniakowski: *Tractability of Multivariate Problems. Volume II: Standard Information for Functionals.* European Mathematical Society Publishing House, Zürich, 2010.
12. J. Nocedal, S.J. Wright: *Numerical Optimization*, 2nd edition. Springer, 2006.

13. J. M. Ortega, W. C. Rheinboldt: *Iterative Solution of Nonlinear Equations in Several Variables*, Society for Industrial and Applied Mathematics, 1987.
14. T. Pillards, B. Vandewoestyne, R. Cools: Minimizing the $L_2$ and $L$ star discrepancies of a single point in the unit hypercube. J. of Computational and Applied Mathematics 197:2006, 282–285.
15. I. H. Sloan, S. Joe: *Lattice Methods for Multiple Integration.* Oxford University Press, New York and Oxford, 1994.
16. V. T. Sós, S. K. Zaremba: The mean-square discrepancies of some two-dimensional lattices. Studia Scientiarum Mathematicarum Hungarica 14:1982, 255–271.
17. V. N. Temlyakov: Error estimates for Fibonacci quadrature formulae for classes of functions. Trudy Mat. Inst. Steklov 200:1991, 327–335.
18. T. Ullrich, D. Zung: Lower bounds for the integration error for multivariate functions with mixed smoothness and optimal Fibonacci cubature for functions on the square. Math. Nachr. 288(7):2015, 743–762.
19. A. Uschmajew: Local convergence of the alternating least squares algorithm for canonical tensor approximation. SIAM Journal on Matrix Analysis and Applications 33(2):2012, 639–652.
20. G. Wahba: Smoothing noisy data with spline functions. Numerische Mathematik 24(5):1975, 383–393.
21. B.E. White: On optimal extreme-discrepancy point sets in the square. Numerische Mathematik 27: 1977, 157–164.
22. P. Zinterhof: Über einige Abschätzungen bei der Approximation von Funktionen mit Gleichverteilungsmethoden. Österreich. Akad. Wiss. Math.-Naturwiss. Kl. S.-B. II 185:1976, 121–132.

**Fig. 1** Optimal point sets for $N = 2, \ldots, 16$ and $\gamma = 1$.

**Fig. 2** Optimal point sets for $N = 2, \ldots, 16$ and $\gamma = 6$.