

ON THE LIPSCHITZ CONTINUITY OF THE SPHERICAL CAP DISCREPANCY AROUND GENERIC POINT SETS

HOLGER HEITSCH¹, RENÉ HENRION¹

¹Weierstrass Institute for Applied Analysis and Stochastics, Berlin, Germany.

ABSTRACT. The spherical cap discrepancy is a prominent measure of uniformity for sets on the d -dimensional sphere. It is particularly important for estimating the integration error for certain classes of functions on the sphere. Building on a recently proven explicit formula for the spherical discrepancy, we show as a main result of this paper that this discrepancy is Lipschitz continuous in a neighbourhood of so-called generic point sets (as they are typical outcomes of Monte-Carlo sampling). This property may have some impact (both algorithmically and theoretically for deriving necessary optimality conditions) on optimal quantization, i.e., on finding point sets of fixed size on the sphere having minimum spherical discrepancy.

Communicated by Friedrich Pillichshammer

1. Introduction

Point sets uniformly located on the classical or higher dimensional sphere are of much interest in many disciplines of mathematics. As examples we refer to point cloud interpolation in computer vision [12] or to optimization problems with chance constraints using the so-called spherical-radial decomposition of elliptically distributed (e.g., Gaussian) random vectors [14]. Uniformity of point sets on the sphere can be characterized by various criteria, e.g., the sum of pairwise distances (which should be large) or by its Coulomb energy (which should be small). If the focus is on estimating the integration error when replacing a spherical integral of a function by an average function value on the spherical point set, then the so-called spherical cap discrepancy is a natural measure of goodness for the uniformity of this point set [1],[3],[6]. Contrary to the criteria mentioned above, the spherical cap discrepancy (being defined as a supremum of infinitely many local discrepancies) is originally not endowed with an explicit formula which could be used for its numerical computation or for its minimization as a function of the point set. This did not harm theoretical investigations in the context of the construction of low discrepancy sequences but it became obstructive in numerical experiments. A possible

remedy consisted in reducing the supremum to a maximum over finitely many local discrepancies (e.g. [[1], p.1005]), but, of course, this provides just a lower bound which might deviate considerably from the true value [[7], p.13]. A certainly more precise algorithmic approximation was provided in [2], but still it was not based on an exact formula and moreover restricted to the classical two-dimensional sphere. In [7], a precise enumerative formula for the spherical cap discrepancy was derived, which reduced the supremum over an infinite family of local discrepancies to a finite maximum of fully explicit and numerically easy to compute expressions. Not surprisingly, this formula suffers from a poor complexity. Nonetheless, it could be used for calibration purposes for moderate sizes of the point set and small dimensions of the sphere (in [7], sets with 2000 points in the two-dimensional sphere to 100 points in the five-dimensional sphere were considered). For a practical application of this formula in image analysis, we refer to [12].

It turns out that, apart from its numerical use, the mentioned formula maybe of interest in characterizing the spherical cap discrepancy as a function of the point set. This observation is based on the fact that the finitely many expressions whose maximum constitutes the spherical cap discrepancy are fully explicit functions of the point set. This allows us, beyond proving the continuity of the spherical cap discrepancy by elementary arguments, to verify even its Lipschitz continuity around so-called generic point sets. The latter refers to point sets on the sphere for which each selection of cardinality not larger than the space dimension is linearly independent. Such point sets are typical outcomes of Monte-Carlo (but not of Quasi Monte-Carlo) sampling. The main argument for proving Lipschitz continuity relies on the fact that, locally around a generic point set, the spherical cap discrepancy can be represented as a continuous selection of $C1$ -functions (see [13]). Moreover, we are able to provide explicitly computable Lipschitz constants. This might be of interest in the application of global optimization methods for minimizing the spherical cap discrepancy (optimal quantization) for a fixed sample size. Note that low discrepancy sequences whose design on the sphere is an active field of research have nice asymptotic properties but do not guarantee optimality for a fixed sample size. Apart from algorithmic relevance, the proven Lipschitz continuity paves a way for establishing necessary optimality conditions in optimal quantization on the sphere by means of the Clarke subdifferential [4].

The paper is organized as follows: Section 2 briefly introduces some basic concepts, presents some simple preliminary results needed later on and proves the continuity of the spherical cap discrepancy. In Section 3, a representation formula for the spherical cap discrepancy as a maximum of finitely many (explicit) functions around generic point sets is proven. In section 4, an extended cap discrepancy is introduced and its Lipschitz continuity around generic point sets is verified. As a trivial consequence, the same property for the original discrepancy is derived as the main result of the paper. Section 5 briefly describes how the previous results could applied in order to derive necessary conditions for optimal quantization with respect to the spherical cap discrepancy.

2. Basic concepts and continuity of the spherical cap discrepancy

We start by defining the following family of subsets of \mathbb{R}^d :

$$H(w, t) := \{x \in \mathbb{R}^d \mid \langle w, x \rangle \geq t\} (w \in \mathbb{R}^d, t \in \mathbb{R}).$$

If $w \neq 0$, then $H(w, t)$ represents a closed half space in \mathbb{R}^d , otherwise it coincides with either \mathbb{R}^d or the empty set depending on whether $t \leq 0$ or $t > 0$. With each of these sets, we associate its so-called cap measure on the sphere:

$$\mu^{cap}(w, t) := \sigma(\mathbb{S}^{d-1} \cap H(w, t)) \quad (\sigma = \text{law of uniform distribution on } \mathbb{S}^{d-1}),$$

where \mathbb{S}^{d-1} refers to the $(d-1)$ -dimensional Euclidean unit sphere in \mathbb{R}^d . We assume in the following that $d \geq 2$.

For a matrix $X = (x^{(1)}, \dots, x^{(N)})$ of order (d, N) with $N \geq 1$ representing a set of points $\{x^{(1)}, \dots, x^{(N)}\} \subseteq \mathbb{S}^{d-1}$, the empirical measure induced from this point set assigns to the set $\mathbb{S}^{d-1} \cap H(w, t)$ its empirical probability

$$\begin{aligned} \mu^{emp}(X, w, t) &:= N^{-1} \cdot \#\{i \in \{1, \dots, N\} \mid x^{(i)} \in \mathbb{S}^{d-1} \cap H(w, t)\} \\ &= N^{-1} \cdot \#\{i \in \{1, \dots, N\} \mid x^{(i)} \in H(w, t)\}. \end{aligned}$$

As a side remark we note that the following relation is immediate from the definition:

$$\mu^{emp}(X, w, t) + \mu^{emp}(X, -w, -t) = 1 + N^{-1} \cdot \#\{i \in \{1, \dots, N\} \mid \langle w, x^{(i)} \rangle = t\}. \quad (1)$$

In order to measure the uniformity of a point set on the sphere, one might compare the deviation between its cap measure and empirical measure on all sets $\mathbb{S}^{d-1} \cap H(w, t)$:

$$\Delta^0(X) := \sup_{w \in \mathbb{R}^d, t \in \mathbb{R}} |\mu^{emp}(X, w, t) - \mu^{cap}(w, t)| \quad (X \in (\mathbb{S}^{d-1})^N) \quad (2)$$

Clearly, the smaller Δ^0 , the better both measures coincide on the chosen family of sets. Such quantities are called discrepancies. If one restricts the family of sets $H(w, t)$ to those with $(w, t) \in \mathbb{S}^{d-1} \times [-1, 1]$, then one obtains the so-called spherical cap discrepancy (e.g., [3])

$$\Delta(X) := \sup_{w \in \mathbb{S}^{d-1}, t \in [-1, 1]} |\mu^{emp}(X, w, t) - \mu^{cap}(w, t)| \quad (X \in (\mathbb{S}^{d-1})^N). \quad (3)$$

Observe, that for $(w, t) \in \mathbb{S}^{d-1} \times [-1, 1]$, the sets $H(w, t)$ represent closed half spaces with normal vector w and height t . Their intersections $\mathbb{S}^{d-1} \cap H(w, t)$, on which the empirical measure and the uniform distribution are compared, are nonempty and called spherical caps. Some authors define the spherical cap discrepancy by using open half spaces instead, i.e., by imposing the strict inequality $\langle w, x \rangle > t$ in the definition of $H(w, t)$ (e.g., [5]). One could formally refer to this alternative definition as a discrepancy $\Delta^1(X)$. It is easy to see that all these three discrepancy definitions coincide, i.e., $\Delta(X) = \Delta^0(X) = \Delta^1(X)$. We provide a proof in Proposition A.1 of the appendix for the reader's convenience. We shall base this paper on the representation (3), but occasionally, the equality with (2) may turn out to be useful.

If $w \in \mathbb{S}^{d-1}$, then the cap measure does not depend on w and we simply write $\mu^{cap}(t) := \mu^{cap}(w, t)$. In this case, the following explicit formula is well known (e.g., [9]):

$$\mu^{cap}(t) = \begin{cases} C_d \int_0^{\arccos(t)} \sin^{d-2}(\tau) d\tau, & \text{if } 0 \leq t \leq 1, \\ 1 - C_d \int_0^{\arccos(-t)} \sin^{d-2}(\tau) d\tau, & \text{if } -1 \leq t < 0, \end{cases} \quad (4)$$

where

$$C_d := \frac{1}{\int_0^\pi \sin^{d-2}(\tau) d\tau} \quad (5)$$

is some normalizing constant. It follows immediately from (4) that μ^{cap} is continuous and that

$$\mu^{cap}(t) = 1 - \mu^{cap}(-t) \quad \forall t \in [-1, 1]. \quad (6)$$

Therefore, we shall work from now on with the following form of (3):

$$\Delta(X) := \sup_{w \in \mathbb{S}^{d-1}, t \in [-1, 1]} |\mu^{emp}(X, w, t) - \mu^{cap}(t)| \quad (X \in (\mathbb{S}^{d-1})^N). \quad (7)$$

We collect three properties of the spherical cap discrepancy that are direct consequences of the definition (7). We observe first, that the supremum in (7) is actually a maximum and that the spherical cap realizing this maximum contains at least one element of the given point set on its relative boundary:

PROPOSITION 2.1. ([7], Proposition 1 & 2). *Let $X \in (\mathbb{S}^{d-1})^N$ be given. Then, there are $w^* \in \mathbb{S}^{d-1}$ and $t^* \in [-1, 1]$ such that*

$$\Delta(X) = |\mu^{emp}(X, w^*, t^*) - \mu^{cap}(t^*)|.$$

Moreover, there exists some $i \in \{1, \dots, N\}$ with $\langle w^, x^{(i)} \rangle = t^*$.*

Secondly, we state a general lower bound for $\Delta(X)$ that depends on the space dimension and the number of points, but not on the position of the points on the sphere.

PROPOSITION 2.2. *Let $\kappa := \min\{d, N\}$. One has that $\Delta(X) \geq \kappa(2N)^{-1} > 0$ for all $X \in (\mathbb{S}^{d-1})^N$.*

Proof. Choose some $w \in \mathbb{S}^{d-1}$ such that $\langle w, x^{(1)} - x^{(j)} \rangle = 0$ for all $j = 2, \dots, \kappa$ and put $t := \langle w, x^{(1)} \rangle$. Then, $|t| \leq 1$, and we have that $\langle w, x^{(i)} \rangle = t$ for $i = 1, \dots, \kappa$. Therefore, owing to (1) and (6),

$$\begin{aligned} 2\Delta(X) &\geq |\mu^{emp}(X, w, t) - \mu^{cap}(t)| + |\mu^{emp}(X, -w, -t) - \mu^{cap}(-t)| \\ &\geq |\mu^{emp}(X, w, t) + \mu^{emp}(X, -w, -t) - \mu^{cap}(t) - \mu^{cap}(-t)| \\ &= |1 + N^{-1} \cdot \#\{i \in \{1, \dots, N\} \mid \langle w, x^{(i)} \rangle = t\} - 1| \\ &= N^{-1} \cdot \#\{i \in \{1, \dots, N\} \mid \langle w, x^{(i)} \rangle = t\} \geq N^{-1}\kappa. \end{aligned}$$

A further property we want to adapt from [7] is a slightly stronger version of [[7], Corollary 1]. We observe that the empirical measure is always strictly greater than the cap measure for any (w^*, t^*) realizing the spherical cap discrepancy.

PROPOSITION 2.3. *For (w^*, t^*) realizing $\Delta(X)$ in Proposition 2.1 it holds that $\mu^{emp}(X, w^*, t^*) > \mu^{cap}(t^*)$.*

Proof. By assumption, we have that $\Delta(X) = |\mu^{emp}(X, w^*, t^*) - \mu^{cap}(t^*)|$. From [[7], Corollary 1] we already know that $\mu^{emp}(X, w^*, t^*) \geq \mu^{cap}(t^*)$. Now, the equality $\mu^{emp}(X, w^*, t^*) = \mu^{cap}(t^*)$ would imply $\Delta(X) = 0$, a contradiction with Proposition 2.2.

As a consequence, we end up at a yet different representation of the spherical cap discrepancy, which allows us to get rid of absolute values:

COROLLARY 2.4. *One has that*

$$\Delta(X) = \sup_{w \in \mathbb{S}^{d-1}, t \in [-1, 1]} \mu^{emp}(X, w, t) - \mu^{cap}(t) \quad \forall X \in (\mathbb{S}^{d-1})^N.$$

Proof. Clearly, the relation ' \geq ' in the claimed equality holds true by (7). On the other hand, by Proposition 2.3, there exists $(w^*, t^*) \in \mathbb{S}^{d-1} \times [-1, 1]$ such that $\Delta(X) = \mu^{emp}(X, w^*, t^*) - \mu^{cap}(t^*)$. Hence, the reverse relation ' \leq ' holds also true in the claimed equality.

Throughout the paper, we understand the sphere \mathbb{S}^{d-1} as a metric space inheriting its metric from the Euclidean norm in \mathbb{R}^d . Next, we are going to prove that the spherical cap discrepancy is continuous.

THEOREM 2.5. *The function $\Delta: (\mathbb{S}^{d-1})^N \rightarrow \mathbb{R}$ is continuous.*

Proof. We show first that Δ is lower semicontinuous. Fix some arbitrary $X = (x^{(1)}, \dots, x^{(N)}) \in (\mathbb{S}^{d-1})^N$ and $\varepsilon > 0$. According to Proposition 2.1 and Proposition 2.3, there exist $w^* \in \mathbb{S}^{d-1}$ and $t^* \in [-1, 1]$ such that

$$\Delta(X) = \mu^{emp}(X, w^*, t^*) - \mu^{cap}(t^*).$$

We claim that $t^* > -1$. Indeed, if $t^* = -1$, then $\mu^{emp}(X, w^*, t^*) = \mu^{cap}(t^*) = 1$, whence the contradiction $\Delta(X) = 0$ with Proposition 2.3.

Define $I := \{i \in \{1, \dots, N\} \mid x^{(i)} \in H(w^*, t^*)\}$. Clearly, we find $c > 0$ such that

$$\begin{aligned} t^* - c &\geq -1; \langle w^*, x^{(i)} \rangle > t^* - c \quad \forall i \in I; \langle w^*, x^{(i)} \rangle < t^* - c \quad \forall i \in I^c; \\ |\mu^{cap}(t^*) - \mu^{cap}(t^* - c)| &< \varepsilon. \end{aligned}$$

By continuity, there exists $\delta > 0$ such that for all $\tilde{X} \in (\mathbb{S}^{d-1})^N$ with $\|\tilde{X} - X\| < \delta$ ($\|\cdot\|$ denoting the Euclidean norm) it holds that

$$\langle w^*, \tilde{x}^{(i)} \rangle > t^* - c \quad \forall i \in I; \langle w^*, \tilde{x}^{(i)} \rangle < t^* - c \quad \forall i \in I^c.$$

Consequently, $\mu^{emp}(\tilde{X}, w^*, t^* - c) = \mu^{emp}(X, w^*, t^*)$ for all such \tilde{X} . Hence, for all $\tilde{X} \in (\mathbb{S}^{d-1})^N$ with $\|\tilde{X} - X\| < \delta$,

$$\begin{aligned} \Delta(\tilde{X}) &\geq |\mu^{emp}(\tilde{X}, w^*, t^* - c) - \mu^{cap}(t^* - c)| \\ &= |\mu^{emp}(X, w^*, t^*) - \mu^{cap}(t^*) + \mu^{cap}(t^*) - \mu^{cap}(t^* - c)| > \Delta(X) - \varepsilon. \end{aligned}$$

Since $\varepsilon > 0$ was arbitrary, this shows the lower semicontinuity of Δ at X .

As for the upper semicontinuity, assume that Δ fails to be upper semicontinuous at some $X \in (\mathbb{S}^{d-1})^N$. Then there exist some $c > 0$ as well as a sequence $X_n \in (\mathbb{S}^{d-1})^N$ with $X_n \rightarrow X$ and $\Delta(X_n) > \Delta(X) + c$. Let (w_n^*, t_n^*) be a sequence that realizes the cap discrepancies $\Delta(X_n)$. Due to Proposition 2.3 we have that

$$\Delta(X_n) = \mu^{emp}(X_n, w_n^*, t_n^*) - \mu^{cap}(t_n^*) \quad \forall n \in \mathbb{N}.$$

Since $(w_n^*, t_n^*) \in \mathbb{S}^{d-1} \times [-1, 1]$, by passing to a subsequence, we may assume that $(w_n^*, t_n^*) \rightarrow (\tilde{w}, \tilde{t}) \in \mathbb{S}^{d-1} \times [-1, 1]$. Altogether, $(X_n, w_n^*, t_n^*) \rightarrow (X, \tilde{w}, \tilde{t})$. With the index set $\tilde{I} := \{i \in \{1, \dots, N\} \mid x^{(i)} \in H(\tilde{w}, \tilde{t})\}$, one has that $\langle \tilde{w}, x^{(i)} \rangle < \tilde{t}$ for all $i \in \tilde{I}^c$. By continuity, there is some n_0 such that $\langle w_n^*, x_n^{(i)} \rangle < t_n^*$ for all $n \geq n_0$ and $i \in \tilde{I}^c$. This entails that $\mu^{emp}(X_n, w_n^*, t_n^*) \leq \mu^{emp}(X, \tilde{w}, \tilde{t})$ for $n \geq n_0$. Moreover, by continuity of μ^{cap} , we have $|\mu^{cap}(\tilde{t}) - \mu^{cap}(t_n^*)| \leq c$ for sufficient large n . Consequently, there exists some $n_1 \in \mathbb{N}$ such that for all $n \geq n_1$

$$\begin{aligned} \Delta(X_n) &= \mu^{emp}(X_n, w_n^*, t_n^*) - \mu^{cap}(t_n^*) \\ &\leq \mu^{emp}(X, \tilde{w}, \tilde{t}) - \mu^{cap}(\tilde{t}) + \mu^{cap}(\tilde{t}) - \mu^{cap}(t_n^*) \\ &\leq |\mu^{emp}(X, \tilde{w}, \tilde{t}) - \mu^{cap}(\tilde{t})| + |\mu^{cap}(\tilde{t}) - \mu^{cap}(t_n^*)| \leq \Delta(X) + c, \end{aligned}$$

a contradiction to the previously established inequality $\Delta(X_n) > \Delta(X) + c$.

A consequence of the continuity property is the existence of an optimal quantization with respect to the spherical cap discrepancy for any fixed number of points on the unit sphere.

COROLLARY 2.6. *For each $N \geq 1$, there exists a point set $X_* = (x_*^{(1)}, \dots, x_*^{(N)})$ with $X_* \in (\mathbb{S}^{d-1})^N$ realizing the minimal spherical cap discrepancy, i.e.,*

$$\Delta(X_*) = \inf_{X \in (\mathbb{S}^{d-1})^N} \Delta(X)$$

3. Generic point sets and a representation formula for the spherical cap discrepancy

Our ultimate goal in this paper is to prove the local Lipschitz continuity of the spherical cap discrepancy. While it is not clear at this point, whether a general Lipschitz result holds true in general, we will be able to derive it for the class of generic point sets, which would be the typical outcomes of Monte-Carlo sampling on the sphere.

DEFINITION 3.1. A point set $X = (x^{(1)}, \dots, x^{(N)}) \in (\mathbb{R}^d)^N$ is called generic if for any index set $I \subseteq \{1, \dots, N\}$ with $\#I \leq d$ the selection $\{x^{(i)} \mid i \in I\}$ is linear independent in \mathbb{R}^d .

Clearly, all point sets close enough to some generic point set are generic themselves, which allows for the following proposition.

PROPOSITION 3.2. *If $\bar{X} \in (\mathbb{R}^d)^N$ is generic, then there exists a neighborhood \mathcal{O} of \bar{X} such that X is generic for each $X \in \mathcal{O}$.*

DEFINITION 3.3. Define the family of index sets

$$\Phi := \{I \subseteq \{1, \dots, N\} \mid 1 \leq \#I \leq d\}. \quad (8)$$

For generic $X = (x^{(1)}, \dots, x^{(N)}) \in (\mathbb{R}^d)^N$ and $I \in \Phi$, let X_I be the matrix whose columns are the $x^{(i)}$ ($i \in I$). Put

$$\mathbf{1} := (1, \dots, 1)^\top \in \mathbb{R}^{\#I}, t_I := (\mathbf{1}^\top (X_I^\top X_I)^{-1} \mathbf{1})^{-1/2}, w_I := t_I X_I (X_I^\top X_I)^{-1} \mathbf{1} \quad (9)$$

which are well-defined by the assumed genericity of X .

PROPOSITION 3.4. *If $X \in (\mathbb{R}^d)^N$ is generic, then $t_I > 0$, $w_I \in (\mathbb{S}^{d-1})$ and $X_I^\top w_I = t_I \mathbf{1}$ for all $I \in \Phi$. If, moreover, $X \in (\mathbb{S}^{d-1})^N$, then $0 < t_I \leq 1$ for all $I \in \Phi$.*

Proof. The first assertion is evident from (9). If $X \in (\mathbb{S}^{d-1})^N$, then the first assertion implies the second one: Fix some arbitrary $I \in \Phi$ and some arbitrary $j \in I$ and obtain

$$t_I = \langle x^{(j)}, w_I \rangle \leq \|w_I\| = 1$$

Next, we shall prove a representation formula for the spherical cap discrepancy of generic point sets which follows from and simplifies the enumerative formula for general point sets proven in [[7], Theorem 1].

Theorem 3.5. *Let $X = (x^{(1)}, \dots, x^{(N)}) \in (\mathbb{S}^{d-1})^N$ be generic. Then, with the notation from Definition 3.3, the spherical cap discrepancy may be represented as*

$$\Delta(X) = \max_{I \in \Phi} \max \{ \mu^{\text{emp}}(X, w_I, t_I) - \mu^{\text{cap}}(t_I), \mu^{\text{emp}}(X, -w_I, -t_I) - \mu^{\text{cap}}(-t_I) \}. \quad (10)$$

Proof. For some $I \in \Phi$, denote by \tilde{X}_I the extension $\tilde{X}_I = \begin{pmatrix} X_I \\ -\mathbf{1}^\top \end{pmatrix}$ of the matrix X_I . From the enumeration formula in [[7], Theorem 1] we know that the cap discrepancy is represented as a maximum of local discrepancies associated with index subsets contained in Φ . Let $I^* \in \Phi$ some index set realizing this maximum. Then, according to [[7], Theorem 1], we have that $\text{rank } \tilde{X}_{I^*} = \#I^*$, $\gamma := \mathbf{1}^\top (\tilde{X}_{I^*}^\top \tilde{X}_{I^*})^{-1} \mathbf{1} \in (0, 1]$ and

$$\Delta(X) = \max \{ |\mu^{\text{emp}}(X, w^*, t^*) - \mu^{\text{cap}}(t^*)|, |\mu^{\text{emp}}(X, -w^*, -t^*) - \mu^{\text{cap}}(-t^*)| \}, \quad (11)$$

where $t^* := \left(\frac{1-\gamma}{\gamma} \right)^{1/2} \geq 0$ and

$$w^* := \frac{1+(t^*)^2}{t^*} X_{I^*} (\tilde{X}_{I^*}^\top \tilde{X}_{I^*})^{-1} \mathbf{1} \text{ if } t^* > 0; w^* \in \text{Ker } X_{I^*}^\top \cap \mathbb{S}^{d-1} \text{ if } t^* = 0. \quad (12)$$

As noted in [[7], Theorem 1], the choice of w^* in the second case of (12) is arbitrary. Then, by virtue of Proposition 2.3, regardless of whether the first or the second term in (11) is dominating,

$$\Delta(X) = \max \{ \mu^{\text{emp}}(X, w^*, t^*) - \mu^{\text{cap}}(t^*), \mu^{\text{emp}}(X, -w^*, -t^*) - \mu^{\text{cap}}(-t^*) \}. \quad (13)$$

To proceed, put

$$v := -(1 + (t^*)^2)(\tilde{X}_{I^*}^\top \tilde{X}_{I^*})^{-1} \mathbf{1}.$$

From here, we get the two relations

$$\mathbf{1}^\top v = -(1 + (t^*)^2)\gamma = -\left(1 + \frac{1-\gamma}{\gamma}\right)\gamma = -1; \tilde{X}_{I^*}^\top \tilde{X}_{I^*} v = -(1 + (t^*)^2)\mathbf{1}.$$

Along with the definition of \tilde{X}_{I^*} as an extended matrix, this yields that

$$-(1 + (t^*)^2)\mathbf{1} = \tilde{X}_{I^*}^\top \tilde{X}_{I^*} v = (X_{I^*}^\top X_{I^*} + \mathbf{1}\mathbf{1}^\top)v = X_{I^*}^\top X_{I^*} v - \mathbf{1}.$$

Therefore, it holds $-(t^*)^2 \mathbf{1} = X_{I^*}^\top X_{I^*} v$. Since $X_{I^*}^\top X_{I^*}$ is regular by genericity of X , one gets that

$$v = -(t^*)^2 (X_{I^*}^\top X_{I^*})^{-1} \mathbf{1} \text{ and } 1 = -\mathbf{1}^\top v = (t^*)^2 \mathbf{1}^\top (X_{I^*}^\top X_{I^*})^{-1} \mathbf{1}. \quad (14)$$

In particular, it must be $t^* > 0$ and we observe that

$$t^* = \left(\mathbf{1}^\top (X_{I^*}^\top X_{I^*})^{-1} \mathbf{1} \right)^{-1/2}.$$

Furthermore, thanks to $t^* > 0$, on the other hand, by (12) and (14) one arrives at

$$w^* = \frac{1+(t^*)^2}{t^*} X_{I^*} (\tilde{X}_{I^*}^\top \tilde{X}_{I^*})^{-1} \mathbf{1} = -\frac{1}{t^*} X_{I^*} v = t^* X_{I^*} (X_{I^*}^\top X_{I^*})^{-1} \mathbf{1}.$$

Altogether, we conclude that $(w^*, t^*) = (w_{I^*}, t_{I^*})$ for w_{I^*}, t_{I^*} defined in (9). Combining this with (13), we get that

$$\Delta(X) = \max\{\mu^{emp}(X, w_{I^*}, t_{I^*}) - \mu^{cap}(t_{I^*}), \mu^{emp}(X, -w_{I^*}, -t_{I^*}) - \mu^{cap}(-t_{I^*})\}.$$

Moreover, because $I^* \in \Phi$, it even holds that

$$\begin{aligned} \Delta(X) &\leq \max_{I \in \Phi} \max\{\mu^{emp}(X, w_I, t_I) - \mu^{cap}(t_I), \mu^{emp}(X, -w_I, -t_I) - \mu^{cap}(-t_I)\} \\ &\leq \Delta(X), \end{aligned}$$

where the last inequality relies on (7) and on the fact that $w_I \in \mathbb{S}^{d-1}$ and $t_I \in [-1, 1]$ for all $I \in \Phi$ by Proposition 3.4. This proves (10).

We may slightly improve the representation formula (10) by excluding singletons from the index family Φ in Theorem 3.5.

PROPOSITION 3.6. *Let $N \geq 2$. Then the assertion of Theorem 3.5 remains valid if replacing the family of index sets Φ in (8) by the (smaller) family of index sets*

$$\bar{\Phi} := \{I \subseteq \{1, \dots, N\} \mid 2 \leq \#I \leq d\} \quad (15)$$

Proof. Since $\bar{\Phi} \subseteq \Phi$, it is sufficient to show that there always exists some $\bar{I}^* \in \bar{\Phi}$ realizing the cap discrepancy $\Delta(X)$ in (10). Assuming to the contrary that

$$\Delta(X) > \max_{I \in \Phi} \max\{\mu^{emp}(X, w_I, t_I) - \mu^{cap}(t_I), \mu^{emp}(X, -w_I, -t_I) - \mu^{cap}(-t_I)\}, \quad (16)$$

$\Delta(X)$ must be realized by some $I^* \in \Phi \setminus \bar{\Phi}$. This implies that I^* is a singleton, i.e., $I^* = \{\ell\}$ for some $\ell \in \{1, \dots, N\}$. Then, by (9) we have $t_{I^*} = 1$ and $w_{I^*} = x^{(\ell)}$, which by $\|x^{(i)}\| = 1$ for $i = 1, \dots, N$ implies that

$$x^{(i)} \in H(w_{I^*}^*, t_{I^*}^*) \Leftrightarrow \langle x^{(i)}, x^{(\ell)} \rangle = 1 \Leftrightarrow x^{(i)} = x^{(\ell)} \quad (i = 1, \dots, N).$$

On the other hand, by genericity of X , we know that $x^{(i)} \neq x^{(\ell)}$ for $i \neq \ell$. Consequently, $\mu^{emp}(X, w_{I^*}, t_{I^*}) = N^{-1}$ and $\mu^{emp}(X, -w_{I^*}, -t_{I^*}) = 1$ due to (1). Moreover, $\mu^{cap}(t_{I^*}) = \mu^{cap}(1) = 0$ and $\mu^{cap}(-t_{I^*}) = \mu^{cap}(-1) = 1$. Thus,

$$\Delta(X) = \max\{\mu^{emp}(X, w_{I^*}, t_{I^*}) - \mu^{cap}(t_{I^*}), \mu^{emp}(X, -w_{I^*}, -t_{I^*}) - \mu^{cap}(-t_{I^*})\} = N^{-1}. \quad (17)$$

Consider $\bar{I} := \{1, 2\} \in \bar{\Phi}$ and define $t_{\bar{I}}, w_{\bar{I}}$ as in (9). For

$$\bar{\Delta} := \max\{\mu^{emp}(X, w_{\bar{I}}, t_{\bar{I}}) - \mu^{cap}(t_{\bar{I}}), \mu^{emp}(X, -w_{\bar{I}}, -t_{\bar{I}}) - \mu^{cap}(-t_{\bar{I}})\}$$

it holds that (similarly to the proof of Proposition 2.2)

$$\begin{aligned} 2\bar{\Delta} &\geq \mu^{emp}(X, w_{\bar{I}}, t_{\bar{I}}) - \mu^{cap}(t_{\bar{I}}) + \mu^{emp}(X, -w_{\bar{I}}, -t_{\bar{I}}) - \mu^{cap}(-t_{\bar{I}}) \\ &= 1 + N^{-1} \cdot \#\{i \in \{1, \dots, N\} \mid \langle w_{\bar{I}}, x^{(i)} \rangle = t_{\bar{I}}\} - 1 \geq 2N^{-1} \end{aligned}$$

From (9), it follows that $X_{\bar{I}}^T w_{\bar{I}} = t_{\bar{I}} \mathbf{1}$, and so, $\langle w_{\bar{I}}, x^{(i)} \rangle = t_{\bar{I}}$ for $i = 1, 2$. Therefore, $2\bar{\Delta} \geq 2N^{-1}$. On the other hand, $\Delta(X) > \bar{\Delta}$ by (16). This yields the contradiction $\Delta(X) > N^{-1}$ with (17).

At the end of this section, we prove a lemma connected with Theorem 3.5 and the quantities defined in (9) which will be of later use.

LEMMA 3.7. *Let $X \in (\mathbb{R}^d)^N$ be generic and $I \in \Phi$ with $\#I < d$. If there exists some index $j \in \{1, \dots, N\} \setminus I$ such that $\langle w_I, x^{(j)} \rangle = t_I$, then for $J := I \cup \{j\}$ it holds that $t_J = t_I$ and $w_J = w_I$.*

Proof. By assumption and by definition of w_I we obtain for $y := x^{(j)}$ that

$$t_I = \langle w_I, y \rangle = t_I \mathbf{1}^T (X_I^T X_I)^{-1} X_I^T y.$$

Hence, with $t_I > 0$ (see Proposition 3.4), we observe that

$$\mathbf{1}^T (X_I^T X_I)^{-1} X_I^T y = y^T X_I (X_I^T X_I)^{-1} \mathbf{1} = 1 \quad (18)$$

We first show that $t_J = t_I$: The genericity of X ensures that $X_J^T X_J$ is regular and that

$$\frac{1}{t_J^2} = \mathbf{1}^T (X_J^T X_J)^{-1} \mathbf{1} = (\mathbf{1}^T \mid 1) \underbrace{\begin{pmatrix} X_I^T X_I & X_I^T y \\ y^T X_I & \|y\|^2 \end{pmatrix}^{-1}}_{=: Z} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Using the Schur complement $S := \|y\|^2 - y^\top X_I (X_I^\top X_I)^{-1} X_I^\top y \neq 0$ of $X_I^\top X_I$, a well-known formula for the inverse of partitioned matrices, yields

$$Z = \begin{pmatrix} (X_I^\top X_I)^{-1} + \frac{1}{S} (X_I^\top X_I)^{-1} X_I^\top y y^\top X_I (X_I^\top X_I)^{-1} & -\frac{1}{S} (X_I^\top X_I)^{-1} X_I^\top y \\ -\frac{1}{S} y^\top X_I (X_I^\top X_I)^{-1} & \frac{1}{S} \end{pmatrix}$$

and together with (18) and the definition of w_I in (9) we have that

$$\begin{aligned} \frac{1}{t_j^2} &= \mathbf{1}^\top (X_I^\top X_I)^{-1} \mathbf{1} + \frac{1}{S} \mathbf{1}^\top (X_I^\top X_I)^{-1} X_I^\top y y^\top X_I (X_I^\top X_I)^{-1} \mathbf{1} \\ &\quad - \frac{2}{S} \mathbf{1}^\top (X_I^\top X_I)^{-1} X_I^\top y + \frac{1}{S} = \frac{1}{t_I^2}. \end{aligned}$$

Thus, $t_j = t_I$. Now we show that also $w_j = w_I$: To this end, referring to (9) and taking into account (18), we compute

$$\begin{aligned} \langle w_j, w_I \rangle &= t_j t_I (\mathbf{1}^\top \mid 1) \begin{pmatrix} X_I^\top X_I & X_I^\top y \\ y^\top X_I & \|y\|^2 \end{pmatrix}^{-1} \begin{pmatrix} X_I^\top \\ y^\top \end{pmatrix} X_I (X_I^\top X_I)^{-1} \mathbf{1} \\ &= t_j t_I (\mathbf{1}^\top \mid 1) \begin{pmatrix} X_I^\top X_I & X_I^\top y \\ y^\top X_I & \|y\|^2 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \frac{t_I t_j}{t_j^2} = 1. \end{aligned}$$

Since $w_I, w_j \in \mathbb{S}^{d-1}$ by Proposition 3.4, we conclude that $w_j = w_I$.

4. Local Lipschitz continuity of the spherical cap discrepancy at generic point sets

In this section, we are going to prove the main result of this paper, namely the local Lipschitz continuity of the spherical cap discrepancy Δ around generic point sets. The main argument would aim at representing Δ as a continuous selection of C^1 -functions. The Lipschitz continuity would allow one to calculate the Clarke subdifferential of Δ and to exploit it in the derivation of necessary optimality conditions for minimizing Δ as a function of the point set (optimal quantization). A technical difficulty arising in this context is the fact that both, the argument of deriving Lipschitz continuity for continuous selections of C^1 -functions and the definition of Clarke's subdifferential are tied to a structure of normed linear spaces, whereas Δ is defined on the sphere. For this reason, we introduce a generalized cap discrepancy Λ that extends the spherical cap discrepancy Δ to arbitrary point sets in the Euclidean space $(\mathbb{R}^d)^N$ in a neighborhood of a given generic point set on the unit sphere. The idea is to prove the local Lipschitz continuity of Λ first and then to get as an immediate corollary the same property for the genuine spherical cap discrepancy Δ which is the restriction of Λ to the sphere around generic point sets.

4.1. Definition and continuity of the generalized cap discrepancy

In order to define the generalized cap discrepancy Λ mentioned above, one could be tempted to directly extend the definition (7) of Δ to arbitrary Euclidean point sets. For deriving the desired Lipschitz property, however, it is beneficial to restrict considerations to generic

point sets and to take the representation formula (10) in Theorem 3.5 as a basis for defining Λ . From now on, we shall assume that $d \geq 3$ which is no substantial restriction because uniformity of point sets on a circle is trivial.

We start by introducing an extended cap measure $\mu^{cap}: \mathbb{R} \rightarrow \mathbb{R}$ (for dimension $d \geq 3$) in a way that it is continuously differentiable on \mathbb{R} and coincides with the original cap measure μ^{cap} from (4) on $[-1, 1]$ (which is continuously differentiable on $(-1, 1)$). This is achieved by the following definition:

$$\mu^{cap}(t) := \begin{cases} \mu^{cap}(t), & t \in [-1, 1] \\ -\frac{1}{2}t + \frac{1}{2}, & |t| > 1, d = 3 \\ 0, & t > 1, d \geq 4 \\ 1, & t < -1, d \geq 4 \end{cases} \quad (19)$$

Indeed, it is easily seen from (4) that $(\mu^{cap})'(-1) = (\mu^{cap})'(1) = 0$, whenever $d \geq 4$. Hence the constant continuation by the respective function values $\mu^{cap}(1) = 0, \mu^{cap}(-1) = 1$ yields a continuously differentiable extension in this case. The special case $d = 3$ cannot be treated in the same way because one easily sees that $\mu^{cap}(t) = -t/2 + 1/2$ for all $t \in [-1, 1]$, so that the derivatives do not vanish at -1 and 1 , respectively. We may therefore simply keep the definition of the function globally in order to end up at a continuously differentiable extension. Note also, that in the case of $d = 2$ (which we excluded), there exists no continuously differentiable extension of μ^{cap} because its derivative converges to $-\infty$ with the argument t converging to ± 1 .

It is easy to show that for all $d \geq 3$ we may extend relation (6) to

$$\mu^{cap}(-t) = 1 - \mu^{cap}(t) \quad \forall t \in \mathbb{R}. \quad (20)$$

Similarly, to the cap measure, we may extend the empirical measure to arbitrary point sets by putting for all $X \in (\mathbb{R}^d)^N$ and $(w, t) \in \mathbb{R}^{d+1}$:

$$\mu^{Emp}(X, w, t) := N^{-1} \cdot \#\{i \in \{1, \dots, N\} \mid x^{(i)} \in H(w, t)\}. \quad (21)$$

Clearly, for all normalized point sets $X \in (\mathbb{S}^{d-1})^N$ it holds that

$$\mu^{Emp}(X, w, t) = \mu^{emp}(X, w, t) \quad \forall (w, t) \in \mathbb{R}^{d+1}. \quad (22)$$

For the following definition, we make reference to the quantities t_I, w_I defined in (9) for $I \in \Phi$ with Φ as in (8). Note that, in the previous section, all results were formulated for a fixed (generic) point set X . Therefore, for notational convenience, we did not emphasize the dependence of t_I, w_I on X . In this section, however, we will investigate Lipschitz continuity of the spherical cap discrepancy based on the representation formula (10). Since now the point set will become a true variable, we will rather use the notations $t_I(X), w_I(X)$ in the definitions (9) in order to stress the dependence on X . It is obvious that t_I, w_I are continuous mappings on the set of generic point sets X .

DEFINITION 4.1. For generic point sets $X \in (\mathbb{R}^d)^N$, we define the generalized cap discrepancy

$$\Lambda(X) := \max_{I \in \Phi} \max \{ \mu^{Emp}(X, w_I(X), t_I(X)) - \mu^{Cap}(t_I(X)), \mu^{Emp}(X, -w_I(X), -t_I(X)) - \mu^{Cap}(-t_I(X)) \}. \quad (23)$$

Thanks to Proposition 3.2, we make the following observation:

REMARK 1. If $\bar{X} \in (\mathbb{R}^d)^N$ is generic, then there exists a neighborhood \mathcal{O} of \bar{X} such that Λ is defined on \mathcal{O} and has the representation (23) for all $X \in \mathcal{O}$.

By Proposition 3.4, it follows for generic $X \in (\mathbb{S}^{d-1})^N$, that $|t_I| \leq 1$ for all $I \in \Phi$. This entails that $\mu^{Cap}(\pm t_I) = \mu^{Cap}(\pm t_I)$ for all $I \in \Phi$. Moreover, by (22), one also has in this case that $\mu^{Emp}(X, \pm w_I, \pm t_I) = \mu^{Emp}(X, \pm w_I, \pm t_I)$ for all $I \in \Phi$. Now, (23) and Theorem 3.5 entail that our generalized cap discrepancy reduces to the original spherical cap discrepancy for generic point sets on the sphere:

COROLLARY 4.2. For generic $X \in (\mathbb{S}^{d-1})^N$ one has that $\Lambda(X) = \Delta(X)$.

The first basic ingredient for proving the local Lipschitz continuity of Λ around a generic point set is the continuity itself at such point. Adding to this property later that Λ is a selection of \mathcal{C}^1 -functions, we will arrive at the desired Lipschitz result. Given the already proven continuity of the genuine discrepancy Δ at arbitrary point sets (Theorem 2.5), the following result shows the continuity of the generalized cap discrepancy Λ at generic point sets.

PROPOSITION 4.3. Let $\bar{X} \in (\mathbb{R}^d)^N$ be generic and \mathcal{O} some open neighborhood of \bar{X} such that all $X \in \mathcal{O}$ are generic too (see Proposition 3.2). Then, $\Lambda: \mathcal{O} \rightarrow \mathbb{R}$ is continuous.

Proof. Of course, it is sufficient to prove continuity of Λ at the arbitrarily fixed generic point set \bar{X} which entails continuity on the whole neighbourhood \mathcal{O} mentioned in the statement of Proposition 4.3. We shall show first the lower and later the upper semicontinuity of Λ at \bar{X} , thus proving continuity itself.

Let $I^* \in \Phi$ be some index set realizing the maximum in (23), so that $\Lambda(\bar{X}) = \mu^{Emp}(\bar{X}, w^*, t^*) - \mu^{Cap}(t^*)$ for some $(w^*, t^*) \in \pm\{(w_{I^*}(X), t_{I^*}(X))\}$. We fix an arbitrary $\varepsilon > 0$. Now, by Lemma A.2 proven in the appendix, we can find some $\delta > 0$, small enough such that $\mathbb{B}_\delta(\bar{X}) \subseteq \mathcal{O}$, and in such a way that choosing an arbitrary $X \in \mathbb{B}_\delta(\bar{X})$, we find J, w and t with $J \in \Phi, (w, t) \in \pm\{(w_J(X), t_J(X))\}$ satisfying

$$\mu^{Emp}(X, w, t) = \mu^{Emp}(\bar{X}, w^*, t^*) \text{ and } \mu^{Cap}(t) < \mu^{Cap}(t^*) + \varepsilon.$$

In particular, by (23), then

$$\begin{aligned} \Lambda(X) &\geq \mu^{Emp}(X, w, t) - \mu^{Cap}(t) \\ &= \mu^{Emp}(\bar{X}, w^*, t^*) - \mu^{Cap}(t^*) + \mu^{Cap}(t^*) - \mu^{Cap}(t) > \Lambda(\bar{X}) - \varepsilon. \end{aligned}$$

This means that Λ is lower semicontinuous at \bar{X} . In order to show that Λ is also upper semicontinuous at \bar{X} , we assume to the contrary that there exist some $c > 0$ as well as a sequence $X_n \rightarrow \bar{X}$ such that

$$\Lambda(X_n) > \Lambda(\bar{X}) + c \quad \forall n \in \mathbb{N}. \quad (24)$$

For each $n \in \mathbb{N}$, choose $I_n^* \in \Phi$ and $(w_n^*(X_n), t_n^*(X_n)) \in \pm \{(w_{I_n^*}(X_n), t_{I_n^*}(X_n))\}$ such that $\Lambda(X_n)$ is realized, i.e.,

$$\Lambda(X_n) = \mu^{Emp}(X_n, w_n^*(X_n), t_n^*(X_n)) - \mu^{Cap}(t_n^*(X_n)).$$

Since Φ is a finite set, there exists some $\emptyset \neq I^* \subseteq \{1, \dots, N\}$ such that, upon passing to a subsequence, $I_n^* = I^*$ for all $n \in \mathbb{N}$. Once more, by passing to a subsequence, we may assume that for all $n \in \mathbb{N}$ either

- a) $(w_n^*(X_n), t_n^*(X_n)) = (w_{I^*}(X_n), t_{I^*}(X_n))$ or
- b) $(w_n^*(X_n), t_n^*(X_n)) = (-w_{I^*}(X_n), -t_{I^*}(X_n))$.

We consider just case a) here (the second case being completely analogous). By continuity of w_{I^*} and t_{I^*} , we have that $w_n^*(X_n) \rightarrow w_{I^*}(\bar{X})$ and $t_n^*(X_n) \rightarrow t_{I^*}(\bar{X})$ as $n \rightarrow \infty$. From the definition in (21) it follows easily for continuity reasons that the empirical measure at some triple (X, w, t) is always larger than or equal to the empirical measure of triples (X', w', t') in a sufficiently small neighborhood of (X, w, t) . Accordingly,

$$\mu^{Emp}(X_n, w_n^*(X_n), t_n^*(X_n)) \leq \mu^{Emp}(\bar{X}, w_{I^*}(\bar{X}), t_{I^*}(\bar{X}))$$

for n large enough. Moreover, the continuity of the cap measure implies that

$$|\mu^{Cap}(t_{I^*}(\bar{X})) - \mu^{Cap}(t_n^*(X_n))| \leq c$$

for sufficient large n . Consequently, there exists some $n_0 \in \mathbb{N}$ with

$$\begin{aligned} \Lambda(X_n) &= \mu^{Emp}(X_n, w_n^*(X_n), t_n^*(X_n)) - \mu^{Cap}(t_n^*(X_n)) \\ &\leq \mu^{Emp}(\bar{X}, w_{I^*}(\bar{X}), t_{I^*}(\bar{X})) - \mu^{Cap}(t_{I^*}(\bar{X})) \\ &\quad + \mu^{Cap}(t_{I^*}(\bar{X})) - \mu^{Cap}(t_n^*(X_n)) \leq \Lambda(\bar{X}) + c \end{aligned}$$

for all $n \geq n_0$, which is a contradiction to inequality (24). Hence, Λ is also upper semicontinuous at \bar{X} .

4.2. Local Lipschitz continuity of the generalized cap discrepancy at generic point sets

Now we turn to the Lipschitz continuity of the generalized cap discrepancy locally around a generic point set \bar{X} . As before, we denote by \mathcal{O} an open neighborhood of \bar{X} of generic point sets. According to (23), we have the representation

$$\Lambda(X) = \max_{I \in \Phi} \max \left\{ \varphi_I^{(1)}(X), \varphi_I^{(2)}(X) \right\} \quad \forall X \in \mathcal{O}, \quad (25)$$

Where

$$\begin{aligned} \varphi_I^{(1)}(X) &:= \mu^{Emp}(X, w_I(X), t_I(X)) - \mu^{Cap}(t_I(X)) \\ \varphi_I^{(2)}(X) &:= \mu^{Emp}(X, -w_I(X), -t_I(X)) - \mu^{Cap}(-t_I(X)) \end{aligned} \quad (26)$$

As a preparatory step, we prove the following Lemma:

LEMMA 4.4. *Let $\bar{X} \in (\mathbb{R}^d)^N$ be generic and \mathcal{O} some open neighborhood of \bar{X} such that all $X \in \mathcal{O}$ are generic too. Then, there exists a neighborhood $\mathcal{U} \subseteq \mathcal{O}$ of \bar{X} such that for all $I \in \Phi$ and all $X \in \mathcal{U}$ there holds:*

$$\begin{aligned} \Lambda(\bar{X}) &= \varphi_I^{(1)}(\bar{X}), \Lambda(X) = \varphi_I^{(1)}(X) \\ &\Rightarrow \mu^{Emp}(X, w_I(X), t_I(X)) = \mu^{Emp}(\bar{X}, w_I(\bar{X}), t_I(\bar{X})) \\ \Lambda(\bar{X}) &= \varphi_I^{(2)}(\bar{X}), \Lambda(X) = \varphi_I^{(2)}(X) \\ &\Rightarrow \mu^{Emp}(X, -w_I(X), -t_I(X)) = \mu^{Emp}(\bar{X}, -w_I(\bar{X}), -t_I(\bar{X})) \end{aligned}$$

Proof. Without loss of generality, we prove just the first implication and assume it would not hold true. Then, there exist sequences $X_n \rightarrow \bar{X}$ and $I_n \in \Phi$ such that

$$\begin{aligned} \Lambda(\bar{X}) &= \varphi_{I_n}^{(1)}(\bar{X}), \Lambda(X_n) = \varphi_{I_n}^{(1)}(X_n) \\ \left| \mu^{Emp}(\bar{X}, w_{I_n}(\bar{X}), t_{I_n}(\bar{X})) - \mu^{Emp}(X_n, w_{I_n}(X_n), t_{I_n}(X_n)) \right| &\geq \frac{1}{N} \end{aligned}$$

In the last inequality, we used the fact that the values of μ^{Emp} are multiples of $\frac{1}{N}$. Moreover, by continuity on \mathcal{O} of $\mu^{Cap} \circ t_I$ for all $I \in \Phi$, we have that, for n large enough,

$$\left| \mu^{Cap}(t_{I_n}(\bar{X})) - \mu^{Cap}(t_{I_n}(X_n)) \right| \leq \frac{1}{2N}$$

whenever U is small enough. Consequently, for n large enough, we arrive at the following contradiction with the continuity of Λ shown in Proposition 4.3:

$$|\Lambda(\bar{X}) - \Lambda(X_n)| = \left| \varphi_{I_n}^{(1)}(\bar{X}) - \varphi_{I_n}^{(1)}(X_n) \right| \geq \frac{1}{2N}$$

A natural idea to show the local Lipschitz continuity around generic points of the maximum function Λ in (25) would rely on checking the continuous differentiability or at least local Lipschitz continuity of the elementary functions $\varphi_I^{(1)}, \varphi_I^{(2)}$. This, however, does not apply because these functions fail to be even continuous as a consequence of the discontinuity of

μ^{Emp} . The fact is illustrated for a numerical example in Figure 1. Here, a generic set \bar{X} of four points in \mathbb{S}^2 is considered and subjected to one-parametric variation (shifting one of the four points while keeping the others fixed). The variation parameter zero corresponds to the nominal point set \bar{X} . The figure plots the local discrepancies $\Delta^*(I) := \max \{ \varphi_I^{(1)}(\bar{X}), \varphi_I^{(2)}(\bar{X}) \}$ given by the functions defined in (26). According to (23), their maximum over all index sets $I \in \Phi$ yields the (global) discrepancy Λ . As can be seen, this maximum Λ is continuous as it should be according to Proposition 4.3. However, all elementary functions (local discrepancies) being active for the maximum at the nominal point set \bar{X} exhibit jumps at that same point set. Still, the maximum function Λ is apparently not only continuous but even Lipschitz continuous. To show this rigorously, we shall represent Λ as a selection (not a maximum though!) of finitely many smooth functions. It is well known that continuous selections of smooth (or more generally: locally Lipschitzian) functions are locally Lipschitzian. The desired selection cannot be made among the original elementary functions $\varphi_I^{(1)}, \varphi_I^{(2)}$ due to their discontinuity. We therefore define smooth modifications of these functions by locally fixing μ^{Emp} around the nominal point set \bar{X} :

$$\begin{aligned} \tilde{\varphi}_I^{(1)}(X) &:= \mu^{Emp}(\bar{X}, w_I(\bar{X}), t_I(\bar{X})) - \mu^{Cap}(t_I(X)) \\ \tilde{\varphi}_I^{(2)}(X) &:= \mu^{Emp}(\bar{X}, -w_I(\bar{X}), -t_I(\bar{X})) - \mu^{Cap}(-t_I(X)) \end{aligned} \tag{27}$$

Clearly, the desired smoothness of the $\tilde{\varphi}_I^{(1)}, \tilde{\varphi}_I^{(2)}$ will follow from the continuous differentiability of the functions $\beta_I := \mu^{Cap} \circ t_I$ for $I \in \Phi$ around some arbitrary generic point set X .

HOLGER HEITSCH — RENÉ HENRION

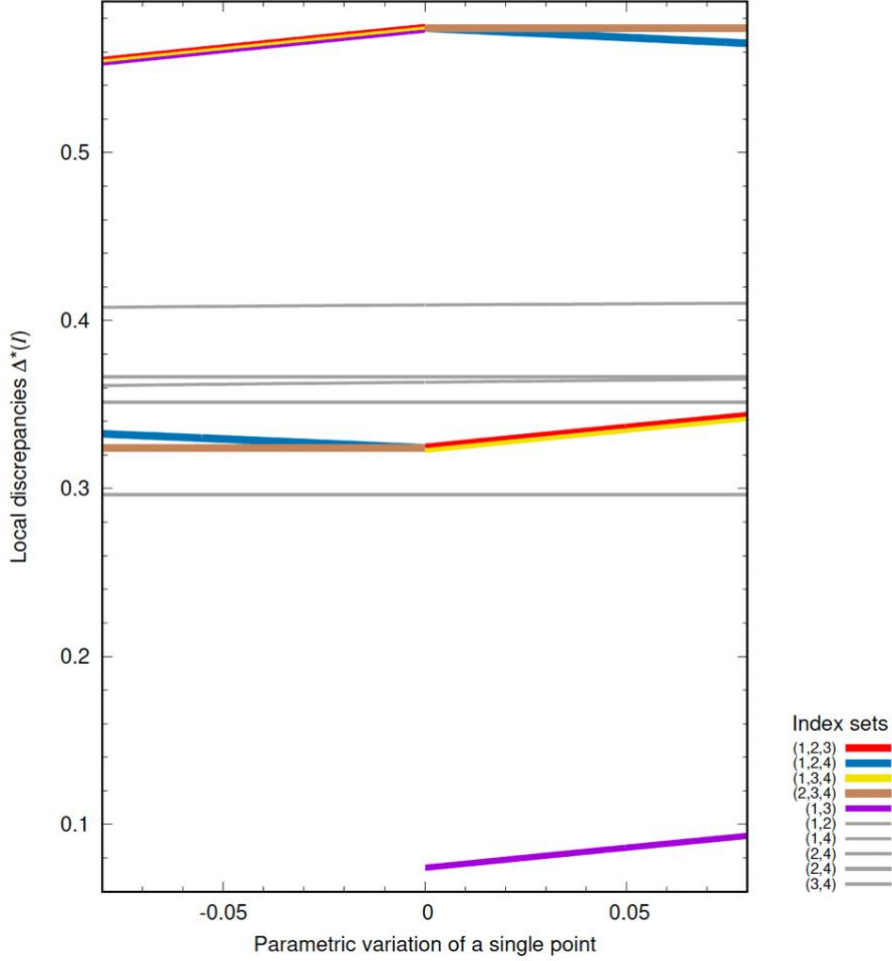


FIGURE 1. The cap discrepancy Δ as a (continuous) maximum of discontinuous functions. The picture shows the behavior of the local discrepancies $\Delta^*(I)$ - whose maximum is Δ (see text) - when parametrically varying a single point out of four points from a given generic point set $\bar{X} \in (\mathbb{S}^2)^4$. The variation parameter zero corresponds to the nominal point set \bar{X} . Local discrepancies realizing the maximum (i.e., Δ) at the nominal set \bar{X} are highlighted by different colors.

LEMMA 4.5. *Let $\bar{X} \in (\mathbb{R}^d)^N$ be generic and \mathcal{O} a neighbourhood of \bar{X} such that all $X \in \mathcal{O}$ are generic too. Then, for each $I \in \Phi$, the function β_I is continuously differentiable on \mathcal{O} with the following partial gradients w. r. t. $x^{(l)}$, ($l \in 1, \dots, N$):*

$$\nabla_{x^{(l)}} \beta_I(X) = \begin{cases} 0, & \text{if } l \notin I \text{ or } |t_l(X)| \geq 1, d \geq 4, \\ -\frac{1}{2} t_l^2(X) c_I^{\tau(l)} w_I(X), & \text{if } l \in I, d = 3, \\ \rho_I c_I^{\tau(l)} w_I(X), & \text{if } l \in I, |t_l(X)| < 1, d \geq 4, \end{cases} \quad \forall X \in \mathcal{O}. \quad (28)$$

Here, with C_d from (5),

$$\rho_I := -C_d t_I^2(X) (1 - t_I^2(X))^{\frac{d-3}{2}} \text{ and } c_I^j := \sum_{i=1}^{\#I} (X_I^T X_I)_{i,j}^{-1} \quad (j = 1, \dots, \#I).$$

Moreover, for $l \in I$, the index $\tau(l)$ refers to the rank of l in the index set I , i.e., if $I = \{\kappa 1, \dots, \kappa_{\#I}\}$, then $l = \kappa_{\tau(l)}$.

Proof. Consider some arbitrary $X \in \mathcal{O}$, whence X is generic. Let $I \in \Phi$ be arbitrary too. We assume that $I = \{\kappa 1, \dots, \kappa_{\#I}\} \subseteq \{1, \dots, N\}$. We want to derive first the function

$$\alpha(X) := \mathbf{1}^T (X_I^T X_I)^{-1} \mathbf{1}. \quad (29)$$

From well-known rules of matrix differential calculus (see, e.g., [10]) one obtains with $M(X) := X_I^T X_I$ that

$$\frac{\partial \alpha}{\partial X_{k,l}}(X) = -2 \sum_{q=1}^{\#I} \left(\sum_{i=1}^{\#I} [M(X)]_{i,\tau(l)}^{-1} \right) \left(\sum_{j=1}^{\#I} [M(X)]_{j,q}^{-1} \right) X_{k,\kappa_q}$$

with $\tau(l)$ as introduced in the statement of this lemma (for a detailed argumentation, we refer to the preprint version of this paper [8]). By definition of $M(X)$ and of the coefficients c_I^j introduced in the statement of this Lemma, we obtain that

$$\nabla_{x^{(l)}} \alpha(X) = -2 c_I^{\tau(l)} \sum_{q=1}^{\#I} c_I^q x^{(\kappa_q)}. \quad (30)$$

Next, we observe that, for all $q = 1, \dots, \#I$,

$$[(X_I^T X_I)^{-1} \mathbf{1}]_q = \sum_{i=1}^{\#I} (X_I^T X_I)_{q,i}^{-1} = \sum_{i=1}^{\#I} (X_I^T X_I)_{i,q}^{-1} = c_I^q.$$

Consequently, by definition (9),

$$w_I(X) = t_I(X) X_I (X_I^T X_I)^{-1} \mathbf{1} = t_I(X) \sum_{q=1}^{\#I} c_I^q x^{(\kappa_q)}.$$

Thanks to (30), this entails that

$$\nabla_{x^{(l)}} \alpha(X) = -2 \frac{c_I^{\tau(l)}}{t_I(X)} w_I(X) \quad (31)$$

We observe next that the function μ^{cap} defined in (4) is continuously differentiable for $d \geq 3$, $t \in (0,1)$ with

$$[\mu^{cap}]'(t) = -C_d (1 - t^2)^{\frac{d-3}{2}}.$$

Along with (19) and the explanations below this equation, this yields that μ^{cap} is continuously differentiable with

$$[\mu^{cap}]'(t) = \begin{cases} 0, & \text{if } |t_I(X)| \geq 1, d \geq 4, \\ -\frac{1}{2}, & \text{if } d = 3, \\ -C_d (1 - t^2)^{\frac{d-3}{2}}, & \text{if } |t_I(X)| < 1. \end{cases} \quad (32)$$

Moreover, the function $t_I = \alpha^{-1/2}$ defined in (9) and (29) is continuously differentiable in the generic point set X because α was shown so in (30). Therefore, the function $\beta = \mu^{cap} \circ t_I$ is continuously differentiable in X with

$$\begin{aligned} \frac{\partial \beta_I}{\partial X_{k,l}}(X) &= [\mu^{cap}]'(t_I(X)) \cdot \frac{\partial t_I}{\partial X_{k,l}}(X) \\ &= -\frac{1}{2} [\mu^{cap}]'(t_I(X)) \cdot [\alpha(X)]^{-3/2} \frac{\partial \alpha}{\partial X_{k,l}}(X), \end{aligned}$$

whence, along with (31)

$$\begin{aligned} \nabla_{x^{(l)}} \beta_I(X) &= -\frac{1}{2} [\mu^{cap}]'(t_I(X)) \cdot [t_I(X)]^3 \nabla_{x^{(l)}} \alpha(X) \\ &= [\mu^{cap}]'(t_I(X)) \cdot [t_I(X)]^2 c_I^{\tau(l)} w_I(X) \end{aligned} \quad (33)$$

Now, lines two and three in (32), yield the corresponding lines in (28). Clearly, the outcomes of (28) depend continuously on X thanks to the continuity of t_I, w_I . This also proves the continuous differentiability of β_I on \mathcal{O} .

COROLLARY 4.6. *For each $I \in \Phi$, the functions $\tilde{\varphi}_I^{(1)}(\cdot), \tilde{\varphi}_I^{(2)}(\cdot)$ defined in (27) are continuously differentiable on $\mathcal{O}(\bar{X})$ with*

$$\nabla \tilde{\varphi}_I^{(1)}(X) = -\nabla \beta_I(X) \quad \text{and} \quad \nabla \tilde{\varphi}_I^{(2)}(X) = \nabla \beta_I(X) \quad \forall X \in \mathcal{O}(\bar{X}).$$

Proof. The first formula above is evident from the definition of $\tilde{\varphi}_I^{(1)}$ in (27). Similarly, the definition of $\tilde{\varphi}_I^{(2)}$ yields that

$$\nabla \tilde{\varphi}_I^{(2)}(X) = [\mu^{cap}]'(-t_I(X)) \nabla t_I(X) = [\mu^{cap}]'(t_I(X)) \nabla t_I(X) = \nabla \beta_I(X),$$

where the second equation follows from (32).

We shall prove now that, locally around generic point sets, Λ is a selection of the continuously differentiable functions $\tilde{\varphi}_I^{(1)}, \tilde{\varphi}_I^{(2)}$.

PROPOSITION 4.7. *Let $\bar{X} \in (\mathbb{R}^d)^N$ be generic and \mathcal{O} some open neighborhood of \bar{X} such that all $X \in \mathcal{O}$ are generic too. Then, there exists a neighborhood $\mathcal{V} \subseteq \mathcal{O}$ of \bar{X} such that for all $X \in \mathcal{V}$ there are $I \in \Phi$ and $s \in \{1, 2\}$ with $\Lambda(X) = \tilde{\varphi}_I^{(s)}(X)$.*

Proof. Let $\mathcal{U} \subseteq \mathcal{O}$ be the neighborhood of \bar{X} from Lemma 4.4 and define the set of active indices as

$$\mathcal{A}(X) := \{(I, s) \in \Phi \times \{1, 2\} \mid \Lambda(X) = \tilde{\varphi}_I^{(s)}(X)\} \quad \forall X \in \mathcal{O}. \quad (34)$$

(see (25)). We claim that there exists a neighborhood $\mathcal{V} \subseteq \mathcal{U}$ of \bar{X} such that

$$\mathcal{A}(X) \subseteq \mathcal{A}(\bar{X}) \quad \forall X \in \mathcal{V}. \quad (35)$$

If this wasn't the case, we could find some sequences $X_n \in (\mathbb{R}^d)^N$ and $(I_n, s_n) \in \Phi \times \{1, 2\}$ such that

$$\Lambda(X_n) = \varphi_{I_n}^{(s_n)}(X_n), \quad \Lambda(\bar{X}) > \varphi_{I_n}^{(s_n)}(\bar{X}) \quad \forall n \in \mathbb{N} \quad \text{and} \quad X_n \rightarrow \bar{X}.$$

Moreover, by passing to a subsequence, we may find some $\bar{I} \in \Phi$ and $\bar{s} \in \{1, 2\}$ such that

$$\Lambda(X_n) = \varphi_{\bar{I}}^{(\bar{s})}(X_n) \quad \forall n \in \mathbb{N} \quad \text{and} \quad \Lambda(\bar{X}) > \varphi_{\bar{I}}^{(\bar{s})}(\bar{X}).$$

Because Λ is continuous at \bar{X} by Proposition 4.3, there is some $n_0 \in \mathbb{N}$ such that

$$\varphi_{\bar{I}}^{(\bar{s})}(X_n) > \varphi_{\bar{I}}^{(\bar{s})}(\bar{X}) + \frac{c}{2} \quad \forall n \geq n_0, \quad \text{where } c := \Lambda(\bar{X}) - \varphi_{\bar{I}}^{(\bar{s})}(\bar{X}) > 0.$$

Next, we use an argument already employed in the proof of Proposition 4.3, namely that the definition in (21) easily implies for continuity reasons that the empirical measure at some triple (X, w, t) is always larger than or equal to the empirical measure of triples (X', w', t') in a sufficiently small neighborhood of (X, w, t) . Assuming, without loss of generality that $\bar{s} = 1$ (the argument being exactly the same for $\bar{s} = 2$), we therefore get for $n \geq n_0$ that

$$\begin{aligned} \mu^{Emp}(\bar{X}, w_{\bar{I}}(\bar{X}), t_{\bar{I}}(\bar{X})) - \mu^{Cap}(t_{\bar{I}}(\bar{X})) + \frac{c}{2} &= \varphi_{\bar{I}}^{(1)}(\bar{X}) + \frac{c}{2} \\ &< \varphi_{\bar{I}}^{(1)}(X_n) = \mu^{Emp}(X_n, w_{\bar{I}}(X_n), t_{\bar{I}}(X_n)) - \mu^{Cap}(t_{\bar{I}}(X_n)) \\ &\leq \mu^{Emp}(\bar{X}, w_{\bar{I}}(\bar{X}), t_{\bar{I}}(\bar{X})) - \mu^{Cap}(t_{\bar{I}}(X_n)) \end{aligned}$$

Passing to the limit on the right-hand side and exploiting the continuity of $\mu^{Cap} \circ t_{\bar{I}}$, we arrive at the contradiction

$$-\mu^{Cap}(t_{\bar{I}}(\bar{X})) + \frac{c}{2} \leq -\mu^{Cap}(t_{\bar{I}}(\bar{X}))$$

which proves (35).

Now, fix an arbitrary $X \in \mathcal{V}$ and choose $I \in \Phi$ and $s \in \{1, 2\}$ such that $\Lambda(X) = \varphi_I^{(s)}(X)$. Then, $(I, s) \in \mathcal{A}(X) \subseteq \mathcal{A}(\bar{X})$ by (35), whence $\Lambda(\bar{X}) = \varphi_I^{(s)}(\bar{X})$. Since $\mathcal{V} \subseteq \mathcal{U}$, Lemma 4.4 yields that $\varphi_I^{(s)}(X) = \bar{\varphi}_I^{(s)}(X)$. Thus, $\Lambda(X) = \bar{\varphi}_I^{(s)}(X)$, as was to be shown.

We are now in a position to formulate the main result of this paper:

THEOREM 4.8. *Let $\bar{X} = (\bar{x}^{(1)}, \dots, \bar{x}^{(N)}) \in (\mathbb{R}^d)^N$ be generic. Then, there exists some neighborhood \mathcal{U} of \bar{X} such that X is generic for all $X \in \mathcal{U}$ and Λ is Lipschitz continuous on \mathcal{U} . In other words, there exists some $L > 0$ such that*

$$|\Lambda(X_1) - \Lambda(X_2)| \leq L \|X_1 - X_2\| \quad \forall X_1, X_2 \in \mathcal{U}$$

Proof. By Propositions 4.3 and 4.7, there exists a neighborhood \mathcal{U} of \bar{X} such that Λ is continuous and a selection of finitely many continuously differentiable functions on \mathcal{U} (which means that Λ is piecewise differentiable in the terminology of Scholtes [[13], page 91]). In particular, Λ is a continuous selection of Lipschitz functions on \mathcal{U} , hence Λ is Lipschitz continuous on \mathcal{U} itself [[13], Proposition 4.1.2.].

As an immediate consequence, we get the desired local Lipschitz continuity of the (original) spherical cap discrepancy Δ around generic points on the sphere:

COROLLARY 4.9. *Let $\bar{X} = (\bar{x}^{(1)}, \dots, \bar{x}^{(N)}) \in (\mathbb{S}^{d-1})^N$ be generic. Then, there exists some neighborhood \mathcal{U} of \bar{X} such that X is generic for all $X \in \mathcal{U} \cap (\mathbb{S}^{d-1})^N$ and Δ is Lipschitz continuous on $\mathcal{U} \cap (\mathbb{S}^{d-1})^N$. In other words, there exists some $L > 0$ such that*

$$|\Delta(X_1) - \Delta(X_2)| \leq L \|X_1 - X_2\| \quad \forall X_1, X_2 \in \mathcal{U} \cap (\mathbb{S}^{d-1})^N.$$

Proof. This follows immediately from Theorem 4.8 and Corollary 4.2.

It is noteworthy that the Lipschitz constant L in Theorem 4.8 (which is the same as in Corollary 4.9) can be explicitly estimated from the data by using the formulae in Lemma 4.5. Indeed, as a consequence of Proposition 4.7 and of [[13], Proposition 4.1.2], we obtain that the Lipschitz constant L of Δ on \mathcal{U} can be represented by the Lipschitz constants $L_I^{(s)}$ of the continuously differentiable functions $\tilde{\varphi}_I^{(s)}$ as

$$L = \max_{(I,s) \in \Phi \times \{1,2\}} L_I^{(s)} \quad (36)$$

Clearly, the $L_I^{(s)}$ can be chosen greater than but arbitrarily close to the Euclidean norms $\|\nabla \tilde{\varphi}_I^{(s)}(\bar{X})\|$ by shrinking the neighbourhood \mathcal{U} . By Corollary 4.6 and Lemma 4.5, a rough upper estimate of the $L_I^{(s)}$ would be (for $d \geq 4$)

$$C_d \max_{I \in \Phi, l \in I} |c_I^{\tau(l)}|$$

(a finer estimate would incorporate the expressions $t_I^2(X)$).

REMARK 2. We make the following observations:

- (1) The best possible Lipschitz constant L in Corollary 4.9 is in general strictly smaller than the best possible Lipschitz constant in Theorem 4.8 because the set of arguments for Δ is restricted to spherical (not arbitrary) point sets.
- (2) The estimate of the Lipschitz constant in (36) can be improved by restricting the maximum to the active indices at \bar{X} as defined in (34):

$$\tilde{L} := \max_{(I,s) \in \mathcal{A}(\bar{X})} L_I^{(s)} \leq L \quad (37)$$

which would keep being a Lipschitz constant for Δ and Δ , respectively because, thanks to (35), nonactive indices at \bar{X} remain nonactive locally and, hence, do not contribute to Δ and Δ .

The following example illustrates the computation of a Lipschitz constant:

EXAMPLE. We consider a set of four points on the classical sphere \mathbb{S}^2 which represent the vertices of a regular tetrahedron:

$$\bar{X} = \begin{pmatrix} \sqrt{\frac{2}{3}} & -\sqrt{\frac{2}{3}} & 0 & 0 \\ 0 & 0 & \sqrt{\frac{2}{3}} & -\sqrt{\frac{2}{3}} \\ \frac{-1}{\sqrt{3}} & \frac{-1}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \end{pmatrix}.$$

We are going to estimate the Lipschitz constants $L_I^{(s)}$ in the improved estimate (37). Numerical computation of the functions $\varphi_I^{(s)}(\bar{X})$ in (26) yields that $\varphi_I^{(1)}(\bar{X}) = 5/12 \approx 0.417$, $\varphi_I^{(2)}(\bar{X}) = 1/3 \approx 0.333$ for all $I \in \Phi$ with $\#I = 3$, and $\varphi_I^{(1)}(\bar{X}) = 1/\sqrt{12} \approx 0.289$, $\varphi_I^{(2)}(\bar{X}) = 1/2 - 1/\sqrt{12} \approx 0.211$ for all $I \in \Phi$ with $\#I = 2$. Hence, the set of active indices introduced in (34) equals $\mathcal{A}(\bar{X}) = (I, 1) | \#I = 3$. We approximate the corresponding elementary Lipschitz constants $L_I^{(1)}$ by the norms of the gradients $\|\nabla \tilde{\varphi}_I^{(1)}(\bar{X})\| = \|\nabla \beta_I(\bar{X})\|$ (see Corollary 4.6). By Lemma 4.5 and Proposition 3.4, the norms of the nonzero components of $\nabla \beta_I(\bar{X})$ calculate as

$$\|\nabla_{x^{(l)}} \beta_I(\bar{X})\| = \left\| -\frac{1}{2} t_I^2(\bar{X}) c_I^{\tau(l)} w_I(\bar{X}) \right\| = \frac{1}{2} t_I^2(\bar{X}) c_I^{\tau(l)} \quad (l \in I)$$

Since, by complete symmetry of a tetrahedron, the occurring quantities on the right-hand side are all the same for all subsets having cardinality 3, we may restrict our considerations to the representative index set $I := \{1, 2, 3\}$. From the definition of \bar{X} , we calculate

$$(\bar{X}_I^T \bar{X}_I)^{-1} = \begin{pmatrix} \frac{3}{2} & \frac{3}{4} & \frac{3}{4} \\ \frac{3}{4} & \frac{3}{2} & \frac{3}{4} \\ \frac{3}{4} & \frac{3}{4} & \frac{3}{2} \end{pmatrix}$$

Then, the sum of elements equals $\mathbf{1}^T (\bar{X}_I^T \bar{X}_I)^{-1} \mathbf{1} = 9$, so that according to (9), $t_I^2(\bar{X}) = 1/9$. Moreover, by definition in Lemma 4.5,

$$c_I^{\tau(1)} = c_I^1 = \sum_{i=1}^3 (X_I^T X_I)^{-1}_{i,1} = 3/2 + 3/4 + 3/4 = 3$$

Hence,

$$\|\nabla_{x^{(1)}} \beta_I(\bar{X})\| = \frac{1}{2} t_I^2(\bar{X}) c_I^{\tau(1)} = \frac{1}{6}.$$

Similarly, $c_I^{\tau(2)} = c_I^{\tau(3)} = 3$ leading to (recall that $4 \notin I$, whence the corresponding partial gradient vanishes by Lemma 3.3)

$$\|\nabla_{x^{(1)}} \beta_I(\bar{X})\| = \|\nabla_{x^{(2)}} \beta_I(\bar{X})\| = \|\nabla_{x^{(3)}} \beta_I(\bar{X})\| = \frac{1}{6}; \quad \|\nabla_{x^{(4)}} \beta_I(\bar{X})\| = 0$$

Therefore,

$$\|\nabla \beta_I(\bar{X})\| = \left(\sum_{i=1}^4 \|\nabla_{x^{(i)}} \beta_I(\bar{X})\|^2 \right)^{1/2} = 1/\sqrt{12} \approx 0.289.$$

Summarizing, upon shrinking the neighbourhood \mathcal{U} in Theorem 4.8, the Lipschitz constant L for Λ can be chosen larger than, but arbitrarily close to 0.289. If one is rather interested in the Lipschitz constant for the original spherical cap discrepancy Δ as in Corollary 4.9, then, one would have to compute the norms of the partial gradients $\nabla_{x^{(l)}} \beta_l(\bar{X})$ projected onto the tangent space of the sphere at the points $x^{(l)}$. This is easily done in the present example and yields the Lipschitz constant 0.272. The gap with the previous Lipschitz constant is explained in Remark 2 item (1).

5. Optimal quantization and necessary optimality conditions

Finding an optimal point set on the sphere minimizing the spherical cap discrepancy amounts to the optimization problem

$$\text{minimize } \Delta(X) \tag{38}$$

where $\Delta: (\mathbb{S}^{d-1})^N \rightarrow [0,1]$ is the spherical cap discrepancy introduced in (3). This problem is also referred to as *optimal quantization* and has to be distinguished from the construction of low discrepancy sequences because the cardinality N of the point set X is fixed. Our aim is to establish necessary optimality conditions a point set has to satisfy in order to be optimal. Note that there is no hope for optimality conditions which are sufficient at the same time due to the lack of convexity of Λ . We will restrict ourselves here to generic point sets. The degenerate case seems to be more delicate to handle and is left for future research.

While (38) is a free (without constraints) optimization problem, the domain of the objective function is a manifold. Standard optimization problems, however, are usually defined on normed spaces subjected to possible further constraints in order to conveniently derive nonsmooth optimality conditions by using tools from generalized differentiation such as the subdifferentials in the sense of Clarke [4] or Mordukhovich [11]. For this reason it is beneficial to equivalently rewrite problem (38) as an optimization problem in the Euclidean space with the additional constraint that the arguments belong to the sphere componentwise:

$$\text{minimize } \Lambda(X) \text{ subject to } X \in (\mathbb{S}^{d-1})^N, X \text{ generic.} \tag{39}$$

The restriction to generic X is necessary because Λ is defined for such point sets only. While the genericity constraint cannot be conveniently described as a classical (in-)equality constraint, it is an open property. This means, that if we are interested in checking whether some generic point set \bar{X} satisfies certain necessary optimality conditions, then we don't have to care about this constraint, because we know it persists to hold in an open neighbourhood \mathcal{O} of \bar{X} and, thus, has no impact on the necessary optimality condition at all. Now, the equivalence of (39) with (38) around some generic $X \in (\mathbb{S}^{d-1})^N$ is evident from Corollary 4.2. We represent the normalization constraint on $X = (x^{(1)}, \dots, x^{(N)})$ as the set of smooth equalities

$$\|x^{(l)}\|^2 = 1 \forall l \in \{1, \dots, N\}.$$

Then the derivative with respect to X of the l -th constraint function equals the matrix

$$2(0 \cdots 0 | x^{(l)} | 0 \cdots 0)$$

Clearly, all these derivatives are linearly independent due to $x^{(l)} \neq 0$. Now, the local Lipschitz continuity of Λ and the continuous differentiability of the constraint functions implies that a generic point set \bar{X} being a (local) solution of the optimal quantization problem (39) has to satisfy the inclusion

$$(\lambda_1 \bar{x}^{(1)} | \dots | \lambda_N \bar{x}^{(N)}) \in \partial^c \Lambda(\bar{X}) \quad (40)$$

for certain multipliers $\lambda_1, \dots, \lambda_N \in \mathbb{R}$, where $\partial^c \Lambda$ refers to the Clarke subdifferential of Λ [[4], p.235-236]. In order to work with such an abstract condition, one has to make the Clarke subdifferential more explicit: From [[13], Proposition 4.3.1.], we know that at generic \bar{X} the identity

$$\partial^c \Lambda(\bar{X}) = \text{conv} \left\{ \nabla \tilde{\varphi}_I^{(s)}(\bar{X}) \mid (I, s) \in \mathcal{A}^*(\bar{X}) \subseteq \Phi \times \{1, 2\} \right\}$$

holds true, where $\mathcal{A}^*(\bar{X})$ refers to the so-called set of essentially active indices (see [[13], p. 92]) and 'conv' refers to the convex hull. Rather than providing a precise definition of the difficult to handle index set $\mathcal{A}^*(\bar{X})$ here, we just recall from its definition in [[13], p. 92], that it is always a subset of the set of active indices $\mathcal{A}(\bar{X})$ defined in (34)

$$\begin{aligned} \mathcal{A}^*(\bar{X}) &\subseteq \left\{ (I, s) \in \Phi \times \{1, 2\} \mid \Lambda(\bar{X}) = \tilde{\varphi}_I^{(s)}(\bar{X}) \right\} \\ &= \left\{ (I, s) \in \Phi \times \{1, 2\} \mid \Lambda(\bar{X}) = \varphi_I^{(s)}(\bar{X}) \right\} = \mathcal{A}(\bar{X}). \end{aligned}$$

Consequently, we arrive at an explicit upper estimate of $\partial^c \Lambda(\bar{X})$ just in terms of active gradients:

$$\partial^c \Lambda(\bar{X}) \subseteq \text{conv} \left\{ \nabla \tilde{\varphi}_I^{(s)}(\bar{X}) \mid (I, s) \in \mathcal{A}(\bar{X}) \right\}.$$

This upper estimate can now be clearly used to establish a weakened but explicit necessary optimality condition as follows: A generic point set \bar{X} being a (local) solution of the optimal quantization problem (39) has to satisfy the inclusion

$$(\lambda_1 \bar{x}^{(1)} | \dots | \lambda_N \bar{x}^{(N)}) \in \text{conv} \left\{ \nabla \tilde{\varphi}_I^{(s)}(\bar{X}) \mid (I, s) \in \mathcal{A}(\bar{X}) \right\}$$

for certain multipliers $\lambda_1, \dots, \lambda_N \in \mathbb{R}$. Resolving for the convex hull, we may extend this statement to: If a generic point set \bar{X} is a (local) solution of the optimal quantization problem (39), then there exist multipliers $\lambda_1, \dots, \lambda_N \in \mathbb{R}$ and $\gamma_{(I,s)} \geq 0$ for $(I, s) \in \mathcal{A}(\bar{X})$ with

$$\lambda_l \bar{x}^{(l)} = \sum_{(I,s) \in \mathcal{A}(\bar{X})} \gamma_{(I,s)} \nabla_{x^{(l)}} \tilde{\varphi}_I^{(s)}(\bar{X}) \quad \text{and} \quad \sum_{(I,s) \in \mathcal{A}(\bar{X})} \gamma_{(I,s)} = 1,$$

$l = 1, \dots, N$. Taking into account that $\nabla_{x^{(l)}} \tilde{\varphi}_I^{(s)}(\bar{X}) = (-1)^s \nabla_{x^{(l)}} \beta_I^{(s)}(\bar{X})$ for $(I, s) \in \Phi \times \{1, 2\}$ and $l = 1, \dots, N$ by Corollary 4.6, and that $\nabla_{x^{(l)}} \beta_I^{(s)}(\bar{X}) = 0$ if $l \notin I$ by Lemma 4.5, we may further rewrite this last relation as

$$\lambda_l \bar{x}^{(l)} = \sum_{\{(I,s) \in \mathcal{A}(\bar{X}) \mid l \in I\}} \gamma_{(I,s)} (-1)^s \nabla_{x^{(l)}} \beta_I^{(s)}(\bar{X}) \quad \text{and} \quad \sum_{(I,s) \in \mathcal{A}(\bar{X})} \gamma_{(I,s)} = 1,$$

$l = 1, \dots, N$. This relation is now fully explicit thanks to the explicit gradient formulae in Lemma 4.5 and it can be used to figure out potential candidates for local minima of the spherical

cap discrepancy (by verifying the necessary optimality conditions) or to exclude certain generic point sets as local or global minima (by showing that the necessary optimality conditions cannot hold). We shall not pursue this concrete application of optimality conditions here and rather leave this to future research.

6. Conclusions

We have proven the Lipschitz continuity of the spherical cap discrepancy around generic point sets on the sphere. Of course, it would be desirable to prove or disprove the Lipschitz continuity on the whole sphere. It seems that we will not be able to show the positive result using the approach taken here (via the representation formula (23)). On the other hand, a counter example isn't easy to construct either. We therefore strongly believe that the following conjecture holds true (note that local Lipschitz continuity around arbitrary point sets implies the global Lipschitz continuity by compactness of the sphere):

CONJECTURE. *The spherical cap discrepancy $\Delta: \mathbb{S}^{(d-1)} \rightarrow [0,1]$ is Lipschitz continuous.*

Apart from proving this conjecture, future research will be devoted to the concrete application of the necessary optimality conditions derived in Section 5 and to a numerical solution of the optimal quantization problem exploiting Lipschitz continuity of the spherical cap discrepancy

7. Appendix A.

PROPOSITION A.1. *For the discrepancies presented in the introduction it holds that $\Delta(X) = \Delta^0(X) = \Delta^1(X)$ for all $X = (x^{(1)}, \dots, x^{(N)}) \in (\mathbb{S}^{d-1})^N$.*

Proof. In order to show that $\Delta(X) = \Delta^0(X)$, it is evidently sufficient to verify the relation

$$|\mu^{emp}(X, w, t) - \mu^{cap}(w, t)| \leq \Delta(X) \forall (w, t) \in \mathbb{R}^d \times \mathbb{R}. \quad (41)$$

Let $(w, t) \in \mathbb{R}^d \times \mathbb{R}$ be arbitrary and assume first that $w = 0$ and $t \leq 0$. Then, $H(w, t) = \mathbb{R}^d$ and, hence, $\mu^{emp}(X, w, t) = \mu^{cap}(w, t) = 1$ and (41) follows trivially. Similarly, if $w = 0$ and $t > 0$, then $H(w, t) = \emptyset$ and, hence, $\mu^{emp}(X, w, t) = \mu^{cap}(w, t) = 0$, so that (41) follows again. Next, let $w \neq 0$ and $t < -\|w\|$. Then, $x \in H(w, t)$ for all $x \in \mathbb{S}^{d-1}$ because of

$$\langle w, x \rangle \geq -\|w\| > t \quad \forall x \in \mathbb{S}^{d-1}$$

which implies $\mu^{emp}(X, w, t) = \mu^{cap}(w, t) = 1$. Similarly, if $t > \|w\|$, then $x \notin H(w, t)$ for all $x \in \mathbb{S}^{d-1}$ because any such x satisfies the relation

$$\langle w, x \rangle \leq \|w\| < t$$

so we have $\mu^{emp}(X, w, t) = \mu^{cap}(w, t) = 0$. In both cases, (41) follows trivially as before. It remains to consider the case that $w \neq 0$ and $|t| \leq \|w\|$. Then, $H(w, t) = H(w^*, t^*)$ for $w^* := w/\|w\| \in \mathbb{S}^{d-1}$ and $t^* := t/\|w\| \in [-1, 1]$. Accordingly, and by virtue of (3),

$$|\mu^{emp}(X, w, t) - \mu^{cap}(w, t)| = |\mu^{emp}(X, w^*, t^*) - \mu^{cap}(w^*, t^*)| \leq \Delta(X).$$

It remains to show that

$$\Delta(X) = \Delta^1(X) := \sup_{w \in \mathbb{S}^{d-1}, t \in [-1, 1]} |\mu_0^{emp}(X, w, t) - \mu_0^{cap}(w, t)|,$$

where, with $H_0(w, t) := \{x \in \mathbb{R}^d \mid \langle w, x \rangle > t\}$, $\mu_0^{cap}(w, t) := \sigma(\mathbb{S}^{d-1} \cap H_0(w, t))$ and

$$\mu_0^{emp}(X, w, t) := N^{-1} \cdot \#\{i \in \{1, \dots, N\} \mid x^{(i)} \in H_0(w, t)\}$$

for all $w \in \mathbb{S}^{d-1}$ and $t \in [-1, 1]$. We immediately check from the definitions, that for arbitrary $(w, t) \in \mathbb{S}^{d-1} \times [-1, 1]$ one has that

$$\mu_0^{emp}(X, w, t) = 1 - \mu^{emp}(X, -w, -t).$$

Moreover, by (6), for arbitrary $(w, t) \in \mathbb{S}^{d-1} \times [-1, 1]$ one has that

$$\mu_0^{cap}(w, t) = \mu^{cap}(w, t) = 1 - \mu^{cap}(-w, -t)$$

Now, the claimed equality $\Delta(X) = \Delta^1(X)$ follows readily from the identity

$$|\mu_0^{emp}(X, w, t) - \mu_0^{cap}(w, t)| = |\mu^{emp}(X, -w, -t) - \mu^{cap}(-w, -t)|$$

for all $(w, t) \in \mathbb{S}^{d-1} \times [-1, 1]$.

LEMMA A.2. *Let $\bar{X} \in (\mathbb{R}^d)^N$ be generic. Let $I^* \in \Phi$ be some index set realizing the maximum in (23), so that $\Lambda(\bar{X}) = \mu^{Emp}(\bar{X}, w^*, t^*) - \mu^{Cap}(t^*)$ for some $(w^*, t^*) \in \pm\{(w_{I^*}(\bar{X}), t_{I^*}(\bar{X}))\}$. Then, the following holds true:*

$$\begin{aligned} \forall \varepsilon > 0 \exists \delta > 0 \forall X \in \mathbb{B}_\delta(\bar{X}) \exists J \in \Phi \exists (w, t) \in \pm\{(w_J(X), t_J(X))\}: \\ \mu^{Emp}(X, w, t) = \mu^{Emp}(\bar{X}, w^*, t^*), \mu^{Cap}(t) < \mu^{Cap}(t^*) + \varepsilon. \end{aligned} \quad (42)$$

Proof. We assume from the very beginning that the δ to be found in (42) is small enough to satisfy the inclusion $\mathbb{B}_\delta(\bar{X}) \subseteq \mathcal{O}$ from Remark 1, so that all X from this ball are generic. We introduce the index sets

$$I_0 := \{i \in \{1, \dots, N\} \mid \langle w^*, \bar{x}^{(i)} \rangle = t^*\}, I_1 := \{i \in \{1, \dots, N\} \mid \langle w^*, \bar{x}^{(i)} \rangle > t^*\}$$

Proposition 3.4 ensures that $\bar{X}_{I^*}^T w_{I^*}(\bar{X}) = t_{I^*}(\bar{X}) \mathbf{1}$, whence $\bar{X}_{I^*}^T w^* = t^* \mathbf{1}$ due to $(w^*, t^*) \in \pm\{(w_{I^*}(\bar{X}), t_{I^*}(\bar{X}))\}$. It follows that $I^* \subseteq I_0$.

Case 1: $I^* = I_0$. Without loss of generality, we may also assume that $(w^*, t^*) = (w_{I^*}(\bar{X}), t_{I^*}(\bar{X}))$ (the opposite case following by absolutely analogous arguments). Then, by definition,

$$\mu^{Emp}(\bar{X}, w_{I^*}(\bar{X}), t_{I^*}(\bar{X})) = N^{-1}(\#I_0 + \#I_1)$$

For arbitrarily given $\varepsilon > 0$ we choose $\delta > 0$ such that $\mathbb{B}_\delta(\bar{X}) \subseteq \mathcal{O}$ for the open neighborhood from the statement of Proposition 4.3 (i.e. all $X \in \mathbb{B}_\delta(\bar{X})$ are generic). Moreover, $\delta > 0$ is chosen small enough to satisfy (by continuity of the mappings t_{I^*}, w_{I^*})

$$\langle w_{I^*}(X), x^{(i)} \rangle > t_{I^*}(X) \quad \forall i \in I_1 \quad \langle w_{I^*}(X), x^{(i)} \rangle < t_{I^*}(X) \quad \forall i \in (I_0 \cup I_1)^c$$

for all $X \in \mathbb{B}_\delta(\bar{X})$ and all $i \in \{1, \dots, N\}$. Moreover, by Proposition 3.4, we have $X_{I^*}^T w_{I^*}(X) = t_{I^*}(X) \mathbf{1}$ for all such X , hence $\langle w_{I^*}(X), x^{(i)} \rangle = t_{I^*}(X)$ for all $i \in I^* = I_0$. Altogether, this implies that $\langle w_{I^*}(X), x^{(i)} \rangle \geq t_{I^*}(X)$ if and only if $i \in I_0 \cup I_1$. Therefore,

$$\mu^{\text{Emp}}(X, w_{I^*}(X), t_{I^*}(X)) = N^{-1}(\#I_0 + \#I_1) = \mu^{\text{Emp}}(\bar{X}, w_{I^*}(\bar{X}), t_{I^*}(\bar{X}))$$

for all $X \in \mathbb{B}_\delta(\bar{X})$. Finally, by continuity of μ^{Cap} , we may further shrink $\delta > 0$ such that $\mu^{\text{Cap}}(t_{I^*}(X)) < \mu^{\text{Cap}}(t_{I^*}(\bar{X})) + \varepsilon$ for all $X \in \mathbb{B}_\delta(\bar{X})$. Thus, we verify (42) by the (constant) selection $J := I^*, w := w_{I^*}(\bar{X}), t := t_{I^*}(\bar{X})$ for each $X \in \mathbb{B}_\delta(\bar{X})$.

Case 2: $I^* \subsetneq I_0$. First, we observe that we may assume that $\#I^* = d$. Indeed, if $\#I^* < d$, then we select some $j \in I_0 \setminus I^*$. From $(w^*, t^*) \in \pm\{(w_{I^*}(\bar{X}), t_{I^*}(\bar{X}))\}$ and by definition of I_0 , we derive that $\langle w_{I^*}(\bar{X}), \bar{x}^{(j)} \rangle = t_{I^*}(\bar{X})$. Now, Lemma 3.7 yields that

$$(w_{I^*}(\bar{X}), t_{I^*}(\bar{X})) = (w_{I^* \cup \{j\}}(\bar{X}), t_{I^* \cup \{j\}}(\bar{X})).$$

Therefore, we may have replaced the index set $I^* \in \Phi$ realizing the maximum in (23) from the very beginning by the larger index set $I^* \cup \{j\} \in \Phi$ (with $\#(I^* \cup \{j\}) = \#I^* + 1 \leq d$) realizing the same value in (23). Calling this larger index set I^* again, we may proceed by adding further indices from I_0 to I^* until $I^* = I_0$, in which case we are back to the situation we already dealt with above, or until $\#I^* = d$, which will be the setting we follow next.

Case 2.1: $(w^*, t^*) = (w_{I^*}(\bar{X}), t_{I^*}(\bar{X}))$. The definition of I_0 and Proposition 3.4 yield that

$$\langle w^*, \bar{x}^{(i)} \rangle = t^* = t_{I^*}(\bar{X}) > 0 \quad \forall i \in I_0.$$

Consequently, given an arbitrary $\varepsilon > 0$, we may choose $\delta > 0$ such that for all $X \in \mathbb{B}_\delta(\bar{X}) \subseteq \mathcal{O}$,

$$\langle w^*, x^{(i)} \rangle \geq t^*/2 > 0 \quad \forall i \in I_0; \mu^{\text{Cap}}(t_I(X)) < \mu^{\text{Cap}}(t_I(\bar{X})) + \varepsilon \quad \forall I \in \Phi, \quad (43)$$

$$\left. \begin{aligned} \langle w_I(\bar{X}), x^{(i)} \rangle > t_I(\bar{X}) &\Rightarrow \langle w_I(X), x^{(i)} \rangle > t_I(X) \\ \langle w_I(\bar{X}), x^{(i)} \rangle < t_I(\bar{X}) &\Rightarrow \langle w_I(X), x^{(i)} \rangle < t_I(X) \end{aligned} \right\} \quad \forall i \in \{1, \dots, N\} \quad \forall I \in \Phi, \quad (44)$$

where the continuity of μ^{Cap} and of the w_I, t_I has been exploited. In order to verify (42), we fix an arbitrary $X \in \mathbb{B}_\delta(\bar{X})$ and find J, w, t as required there. To this aim, denote by P the convex hull of the point set $\{x^{(i)}\}_{i \in I_0}$.

Case 2.1. a): $\text{int } P \neq \emptyset$. Clearly, $0 \notin P$ due to the first relation of (43). It is wellknown from the theory of polyhedra (see, e.g., [[15], Theorem 2.15 (7)]), that there exists a representation

$$P = \{x \in \mathbb{R}^d \mid \langle v_k, x \rangle \geq \tau_k \ (k = 1, \dots, m)\} \quad (v_k \in \mathbb{S}^{d-1}, \tau_k \in \mathbb{R}), \quad (45)$$

such that for $k' = 1, \dots, m$ each set

$$F_{k'} := \{x \in \mathbb{R}^d \mid \langle v_{k'}, x \rangle = \tau_{k'}, \langle v_k, x \rangle \geq \tau_k \ (k = 1, \dots, m, k \neq k')\}$$

is a facet of P . There exists some $k_0 \in \{1, \dots, m\}$ such that $\tau_{k_0} > 0$ because otherwise the contradiction $0 \in P$ would result. As a facet of a bounded polyhedron $P \subseteq \mathbb{R}^d$ with $\text{int } P \neq \emptyset$, F_{k_0} must contain at least d vertices of P . Since the vertices of P are contained in the set $\{x^{(i)}\}_{i \in I_0}$, there exists a subset $J \subseteq I_0$ with $\#J = d$ and $\{x^{(i)}\}_{i \in J} \subseteq F_{k_0}$. Hence, by definition of F_{k_0} , $X_J^T v_{k_0} = \tau_{k_0} \mathbf{1}$. By Proposition 3.4, $X_J^T w_J(X) = t_J(X) \mathbf{1}$, with X_J^T being a regular (d, d) -matrix by genericity of X . Then,

$$v_{k_0} = \tau_{k_0} (X_J^T)^{-1} \mathbf{1}, \quad w_J(X) = t_J(X) (X_J^T)^{-1} \mathbf{1}.$$

Since $\|v_{k_0}\| = \|w_J(X)\| = 1$ and $\tau_{k_0}, t_J(X) > 0$, (see Proposition 3.4), it follows that $(w_J(X), t_J(X)) = (v_{k_0}, \tau_{k_0})$. Now, (45) yields that

$$\langle w_J(X), x^{(i)} \rangle = \langle v_{k_0}, x^{(i)} \rangle \geq \tau_{k_0} = t_J(X) \quad \forall i \in I_0. \quad (46)$$

With $\bar{X}_J^T w_J(\bar{X}) = t_J(\bar{X})$ (by Proposition 3.4) and $\bar{X}_J^T w^* = t^*$ (by $J \subseteq I_0$), the same reasoning as before yields that $(w_J(X), t_J(X)) = (w^*, t^*)$. After having fixed $J \in \Phi$, we also fix $(w, t) := (w_J(X), t_J(X))$ as required in (42). Then, by (46), (44) and by the definitions of I_0, I_1 , one gets that

$$\begin{aligned} \mu^{Emp}(X, w, t) &= \mu^{Emp}(X, w_J(X), t_J(X)) \\ &= N^{-1} \# \{i \in \{1, \dots, N\} \mid \langle w_J(X), x^{(i)} \rangle \geq t_J(X)\} \\ &= N^{-1} (\# I_0 + \# I_1) = \mu^{Emp}(X, w^*, t^*), \end{aligned} \quad (47)$$

which is the first desired relation in (42). The second one follows immediately from the second relation in (43).

Case 2.1.b): $\text{int } P = \emptyset$. Then, P as a polytope must be contained in some hyperplane H :

$$P \subseteq H := \{x \in \mathbb{R}^d \mid \langle \hat{w}, x \rangle = \hat{t}\} (\hat{w} \in \mathbb{S}^{(d-1)}, \hat{t} \in \mathbb{R}).$$

We may assume that $\hat{t} \geq 0$. In particular, $\langle \hat{w}, x^{(i)} \rangle = \hat{t}$ for all $i \in I_0$, or, $X_{I_0}^T \hat{w} = \hat{t} \mathbf{1}$, for short. Since also $X_{I^*}^T w_{I^*}(X) = t_{I^*}(X) \mathbf{1}$ (by Proposition 3.4), and recalling that $\#I^* = d$ the same reasoning as above (46) yields that $(w_{I^*}(X), t_{I^*}(X)) = (\hat{w}, \hat{t})$. This implies that the choice $J := I^*$ satisfies (46) (actually as an equation) so that in view of $(w^*, t^*) = (w_{I^*}(X), t_{I^*}(X))$ we may repeat the reasoning after (46) and (47) in order to derive the two relations in (42) in that alternative case too.

Case 2.2): $(w^*, t^*) = (-w_{I^*}(\bar{X}), -t_{I^*}(\bar{X}))$. We observe that, in case of $\tau_k \geq 0$ for all $k = 1, \dots, m$, P must be unbounded according to (45) because from $x^{(1)} \in P$ it would follow that

$$\langle v_k, \lambda x^{(1)} \rangle = \lambda \langle v_k, x^{(1)} \rangle \geq \lambda \tau_k \geq \tau_k \forall \lambda \geq 1 \forall k = 1, \dots, m.$$

This entails that $\lambda x^{(1)} \in P$ for all $\lambda \geq 1$, whence P would be unbounded because $x^{(1)} \neq 0$ thanks to the genericity of X . However, P is bounded as a convex combination of finitely many points. Therefore, there exists some $k_1 \in 1, \dots, m$ with $\tau_{k_1} < 0$. Now, we can repeat exactly the argumentation from the first case above (referring to $\tau_{k_0} > 0$), just with reversed signs.

FUNDING

The authors gratefully acknowledge the support by the German Research Foundation (DFG) within the project B04 of CRC/Transregio 154 and by the FMJH Program Gaspard Monge in optimization and operations research including support to this program by EDF

REFERENCES

- [1] Aistleitner, C., Brauchart, J. S., and Dick, J. Point sets on the sphere S^2 with small spherical cap discrepancy. *Discrete & Computational Geometry* (2012).
- [2] Bakhshizadeh, M., Kamalinejad, A., and Latifi, M. A practical algorithm to calculate cap discrepancy. *arXiv:2010.10454* (2020).
- [3] Brauchart, J., Saff, E., Sloan, I., and Womersley, R. QMC designs: Optimal order Quasi Monte Carlo integration schemes on the sphere. *Mathematical Computation* 83 (2014), 2821–2851.
- [4] Clarke, F. *Optimization and Nonsmooth Analysis*. Wiley New York, 1983.
- [5] Etayo, U. Spherical cap discrepancy of the diamond ensemble. *Discrete & Computational Geometry* 66 (2021), 1218–1238.
- [6] Grabner, P. J., and Tichy, R. F. Spherical designs, discrepancy and numerical integration. *Mathematics of Computation* 60 (1993), 327–336.
- [7] Heitsch, H., and Henrion, R. An enumerative formula for the spherical cap discrepancy. *Journal of Computational and Applied Mathematics* 390 (2021), 113409.
- [8] Heitsch, H., and Henrion, R. On the Lipschitz continuity of the spherical cap discrepancy around generic point sets. *WIAS Preprint No. 3192* (2025).
- [9] Li, S. Concise formulas for the area and volume of a hyperspherical cap. *Asian Journal of Mathematics & Statistics* 4 (2011), 66–70.
- [10] Magnus, J. R., and Neudecker, H. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. John Wiley, 1999.
- [11] Mordukhovich, B. S. *Variational Analysis and Generalized Differentiation I*. Springer Berlin Heidelberg, 2006.
- [12] Nguyen, K., Barileto, N., and Ho, N. Quasi-Monte Carlo for 3d sliced wasserstein. *arXiv:2309.11713* (2023).
- [13] Scholtes, S. *Introduction to Piecewise Differentiable Equations*. Springer New York, 2012.
- [14] van Ackooij, W., and Henrion, R. Gradient formulae for nonlinear probabilistic constraints with Gaussian and Gaussian-like distributions. *SIAM Journal on Optimization* 24 (2014), 1864–1889.

- [15] Ziegler, G. M. Lectures on polytopes, vol. 152 of Graduate texts in mathematics. Springer-Verlag, 1995.

Received April 8, 2025

Accepted May 30, 2025

Holger Heitsch

Weierstrass Institute for Applied Analysis and Stochastics, Mohrenstraße 39, 10117 Berlin, GERMANY

E-mail: heitsch@wias-berlin.de

René Henrion

Weierstrass Institute for Applied Analysis and Stochastics, Mohrenstraße 39, 10117 Berlin, GERMANY

E-mail: henrion@wias-berlin.de