

# Clustering in Hilbert Geometry Two Case Studies: The Probability Simplex and the Correlation Elliptope

Frank Nielsen and Ke Sun

**Abstract** Clustering categorical distributions in the probability simplex is a fundamental primitive often met in applications dealing with histograms or mixtures of multinomials. Traditionally, the differential-geometric structure of the probability simplex has been used either by (i) setting the Riemannian metric tensor to the Fisher information matrix of the categorical distributions, or (ii) defining the information-geometric structure induced by a smooth dissimilarity measure, called a divergence. In this work, we introduce a novel computationally-friendly non-differential framework for modeling the probability simplex: Hilbert simplex geometry. We discuss the pros and cons of those three statistical modelings, and compare them experimentally for clustering tasks.

**Keywords:** Fisher-Riemannian geometry, information geometry, Hilbert simplex geometry, Finsler geometry, center-based clustering.

## 1 Introduction

The multinomial distribution is an important representation in machine learning that is often met in applications [34, 19] as normalized histograms (with non-empty bins). A multinomial distribution (or categorical distribution)  $p \in \Delta^d$  can be thought as a point lying in the probability simplex  $\Delta^d$  (standard simplex) with coordinates  $p = (\lambda_p^0, \dots, \lambda_p^d)$  such that  $\lambda_p^i > 0$  and  $\sum_{i=0}^d \lambda_p^i = 1$ . The open probability simplex  $\Delta^d$  sits in  $\mathbb{R}^{d+1}$  on the hyperplane  $H_{\Delta^d} : \sum_{i=0}^d x^i = 1$ . We consider the task of clustering a set  $\Lambda = \{p_1, \dots, p_n\}$  of  $n$  categorical distributions in  $\Delta^d$  [19] using center-based  $k$ -means++ or  $k$ -center clustering algorithms [6, 25], which rely on a dissim-

---

Frank Nielsen  
Sony Computer Science Laboratories, Tokyo, Japan, e-mail: Frank.Nielsen@acm.org

Ke Sun  
CSIRO Data61, Sydney, Australia, e-mail: Ke.Sun@data61.csiro.au

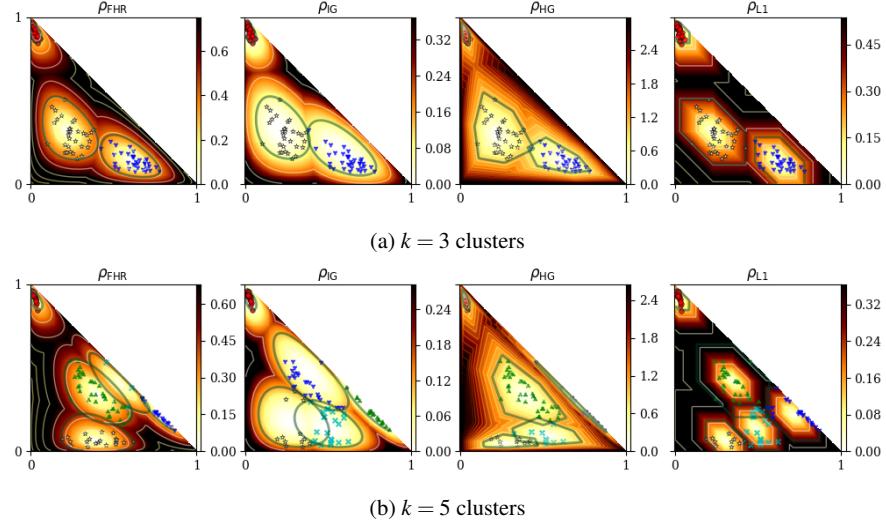


Fig. 1:  $k$ -Center clustering results on a toy dataset in the space of trinomials  $\Delta^2$ . The color density maps indicate the distance from any point to its nearest cluster center.

ilarity measure (loosely called distance or divergence when smooth) between any two categorical distributions. In this work, we mainly consider three distances with their underlying geometries: (1) Fisher-Hotelling-Rao distance  $\rho_{\text{FHR}}$ , (2) Kullback-Leibler divergence  $\rho_{\text{IG}}$ , and (3) Hilbert distance  $\rho_{\text{HG}}$ . The geometric structures are necessary in algorithms, for example, to define midpoint distributions. Figure 1 displays the  $k$ -center clustering results obtained with these three geometries as well as the Euclidean  $L^1$  distance  $\rho_{L1}$  on toy datasets in  $\Delta^2$ . We shall now explain the Hilbert simplex geometry applied to the probability simplex, describe how to perform  $k$ -center clustering in Hilbert geometry, and report experimental results that demonstrate superiority of the Hilbert geometry when clustering multinomials.

The rest of this paper is organized as follows: Section 2 formally introduces the distance measures in  $\Delta^d$ . Section 3 introduces how to efficiently compute the Hilbert distance. Section 4 presents algorithms for Hilbert minimax centers and Hilbert clustering. Section 5 performs an empirical study of clustering multinomial distributions, comparing Riemannian geometry, information geometry and Hilbert geometry. Section 6 presents a second use case of Hilbert geometry in machine learning: clustering correlation matrices. Finally, section 7 concludes this work by summarizing the pros and cons of each geometry. Although some contents require prior knowledge on geometric structures, we will present the detailed algorithms so that general audience can still benefit from this work.

## 2 Three distances with their underlying geometries

### 2.1 Fisher-Hotelling-Rao geometry

The Rao distance between two multinomial distributions is [30, 34]:

$$\rho_{\text{FHR}}(p, q) = 2 \arccos \left( \sum_{i=0}^d \sqrt{\lambda_p^i \lambda_q^i} \right). \quad (1)$$

It is a Riemannian metric length distance (satisfying the symmetric and triangular inequality axioms) obtained by setting the metric tensor  $g$  to the *Fisher information matrix* (FIM)  $\mathcal{I}(p) = (g_{ij}(p))_{d \times d}$  wrt the coordinate system  $(\lambda_p^1, \dots, \lambda_p^d)$ , where

$$g_{ij}(p) = \frac{\delta_{ij}}{\lambda_p^i} + \frac{1}{\lambda_p^0}.$$

We term this geometry the Fisher-Hotelling-Rao (FHR) geometry [28, 61, 53, 54]. The metric tensor  $g$  allows to define an inner product on each tangent plane  $T_p$  of the probability simplex manifold:  $\langle u, v \rangle_p = u^\top g(p)v$ . When  $g$  is everywhere the identity matrix, we recover the Euclidean (Riemannian) geometry with the inner product being the scalar product:  $\langle u, v \rangle = u^\top v$ . The geodesics  $\gamma(p, q; \alpha)$  are defined by the Levi-Civita metric connection [2, 15]. The FHR manifold can be embedded in the positive orthant of an Euclidean  $d$ -sphere in  $\mathbb{R}^{d+1}$  by using the *square root representation*  $\lambda \mapsto \sqrt{\lambda}$  [30]. Therefore the FHR manifold modeling of  $\Delta^d$  has constant *positive* curvature: It is a spherical geometry restricted to the positive orthant with the metric distance measuring the arc length on a great circle.

### 2.2 Information geometry

A divergence  $D$  is a smooth  $C^3$  differentiable dissimilarity measure [3] that allows to define a dual structure in Information Geometry (IG; [60, 15, 2]). A  $f$ -divergence is defined for a strictly convex function  $f$  with  $f(1) = 0$  by:

$$I_f(p : q) = \sum_{i=0}^d \lambda_p^i f \left( \frac{\lambda_q^i}{\lambda_p^i} \right).$$

It is a *separable* divergence since the  $d$ -variate divergence can be written as a sum of  $d$  univariate divergences:  $I_f(p : q) = \sum_{i=0}^d I_f(\lambda_p^i : \lambda_q^i)$ . The class of  $f$ -divergences plays an essential role in information theory since they are provably the *only* separable divergences that satisfy the *information monotonicity* property [2, 37]. That is, by coarse-graining the histograms we obtain lower-dimensional multinomials, say  $p'$  and  $q'$ , such that  $0 \leq I_f(p' : q') \leq I_f(p : q)$  [2]. The Kullback-Leibler (KL)

divergence  $\rho_{\text{IG}}$  is a  $f$ -divergence obtained for  $f(u) = -\log u$ :

$$\rho_{\text{IG}}(p, q) = \sum_{i=0}^d \lambda_p^i \log \frac{\lambda_p^i}{\lambda_q^i}. \quad (2)$$

It is an asymmetric non-metric distance:  $\rho_{\text{IG}}(p, q) \neq \rho_{\text{IG}}(q, p)$ . In differential geometry, the structure of a manifold is defined by two components:

1. A *metric tensor*  $g$  that allows to define an inner product  $\langle \cdot, \cdot \rangle_p$  at each tangent space (for measuring vector lengths and angles between vectors);
2. A *connection*  $\nabla$  that defines *parallel transport*  $\Pi_{p,q}^\nabla$ , i.e., a way to move a tangent vector from one tangent plane  $T_p$  to any other one  $T_q$ .

In FHR geometry, the implicitly-used connection is called the Levi-Civita connection that is induced by the metric  $g$ :  $\nabla^{LC} = \nabla(g)$ . It is a metric connection since it ensures that  $\langle u, v \rangle_p = \langle \Pi_{p,q}^{\nabla^{LC}} u, \Pi_{p,q}^{\nabla^{LC}} v \rangle_q$ . The underlying information-geometric structure of KL is characterized by a pair of *dual* connections [2]  $\nabla = \nabla^{(-1)}$  (mixture connection) and  $\nabla^* = \nabla^{(1)}$  (exponential connection) that induces a corresponding pair of dual geodesics (technically,  $\pm 1$ -autoparallel curves, [15]). Those connections are said *flat* as they define two dual affine coordinate systems  $\theta$  and  $\eta$  on which the  $\theta$ - and  $\eta$ -geodesics are straight line segments, respectively. For multinomials, the *expectation parameters* are:  $\eta = (\lambda^1, \dots, \lambda^d)$  and they one-to-one correspond to the *natural parameters*:  $\theta = \left( \log \frac{\lambda^1}{\lambda^0}, \dots, \log \frac{\lambda^d}{\lambda^0} \right)$ . Thus in IG, we have two kinds of midpoint multinomials of  $p$  and  $q$ , depending on whether we perform the (linear) interpolation on the  $\theta$ - or the  $\eta$ -geodesics. Informally speaking, the dual connections  $\nabla^{(\pm 1)}$  are said coupled to the FIM since we have  $\frac{\nabla + \nabla^*}{2} = \nabla(g) = \nabla^{LC}$ . Those dual connections are not metric connections but enjoy the following property:  $\langle u, v \rangle_p = \langle \Pi_{p,q} u, \Pi_{p,q}^* v \rangle_q$ , where  $\Pi = \Pi^\nabla$  and  $\Pi^* = \Pi^{\nabla^*}$  are the corresponding induced dual parallel transports. The geometry of  $f$ -divergences [3] is the  $\alpha$ -geometry (for  $\alpha = 3 + 2f'''(1)$ ) with the dual  $\pm \alpha$ -connections, where  $\nabla^{(\alpha)} = \frac{1+\alpha}{2} \nabla^* + \frac{1-\alpha}{2} \nabla$ . The Levi-Civita metric connection is  $\nabla^{LC} = \nabla^{(0)}$ . More generally, it was shown how to build a dual information-geometric structure for *any* divergence [3]. For example, we can build a dual structure from the symmetric Cauchy-Schwarz divergence [29]:

$$\rho_{\text{CS}}(p, q) = -\log \frac{\langle \lambda_p, \lambda_q \rangle}{\sqrt{\langle \lambda_p, \lambda_p \rangle \langle \lambda_q, \lambda_q \rangle}}. \quad (3)$$

### 2.3 Hilbert simplex geometry

In Hilbert geometry (HG; [27]), we are given a bounded convex domain  $\mathcal{C}$  (here,  $\mathcal{C} = \Delta^d$ ), and the distance between any two points  $M, M'$  of  $\mathcal{C}$  is defined as follows: Consider the two intersection points  $AA'$  of the line  $(MM')$  with  $\mathcal{C}$ , and order them on the line so that we have  $A, M, M', A'$ . Then the Hilbert metric distance [14] is

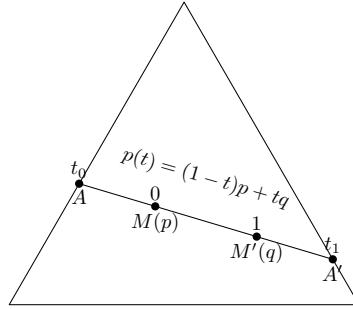


Fig. 2: Computing the Hilbert distance for trinomials on the 2D probability simplex.

defined by:

$$\rho_{\text{HG}}(M, M') = \begin{cases} \left| \log \frac{|A'M||AM'|}{|A'M'||AM|} \right|, & M \neq M', \\ 0 & M = M'. \end{cases} \quad (4)$$

It is also called the Hilbert cross-ratio metric distance [20, 35]. Notice that we take the absolute value of the logarithm since the Hilbert distance is a *signed distance* [56]. When  $\mathcal{C}$  is the unit ball, HG lets us recover the Klein hyperbolic geometry [35]. When  $\mathcal{C}$  is a quadric bounded convex domain, we obtain the Cayley-Klein hyperbolic geometry [12] which can be studied with the Riemannian structure and the corresponding metric distance called the curved Mahalanobis distances [42, 41]. Cayley-Klein hyperbolic geometries have negative curvature.

In Hilbert geometry, the geodesics are *straight* Euclidean lines making them convenient for computation. Furthermore, the domain boundary  $\partial\mathcal{C}$  needs not to be smooth: One may also consider bounded polytopes [11]. This is particularly interesting for modeling  $\Delta^d$ , the  $d$ -dimensional open standard simplex. We call this geometry the *Hilbert simplex geometry*. In Figure 2, we show that the Hilbert distance between two multinomial distributions  $p(M)$  and  $q(M')$  can be computed by finding the two intersection points of the line  $(1-t)p + tq$  with  $\partial\Delta^d$ , denoted as  $t_0 \leq 0$  and  $t_1 \geq 1$ . Then

$$\rho_{\text{HG}}(p, q) = \left| \log \frac{(1-t_0)t_1}{(-t_0)(t_1-1)} \right| = \log \left( 1 - \frac{1}{t_0} \right) - \log \left( 1 - \frac{1}{t_1} \right).$$

The shape of balls in polytope-domain HG is Euclidean polytopes [35], as depicted in Figure 3. Furthermore, the Euclidean shape of the balls do not change with the radius. Hilbert balls have hexagons shapes in 2D [48], rhombic dodecahedra shapes in 3D, and are polytopes [35] with  $d(d+1)$  facets in dimension  $d$ . When the polytope domain is not a simplex, the combinatorial complexity of balls depends on the center location [48], see Figure 4. The HG of the probability simplex yields a non-Riemannian geometry, because at infinitesimal radius, the balls are polytopes and not ellipsoids (corresponding to squared Mahalanobis distance balls used to

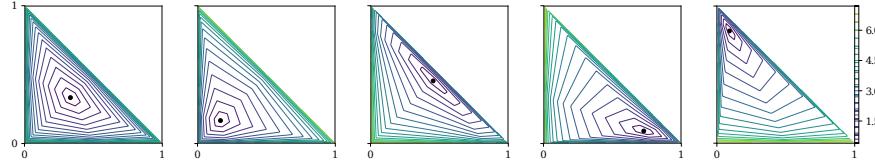


Fig. 3: Balls in the Hilbert simplex geometry  $\Delta^2$  have polygonal Euclidean shapes of constant combinatorial complexity. At infinitesimal scale, the balls have hexagonal shapes, showing that the Hilbert geometry is not Riemannian.

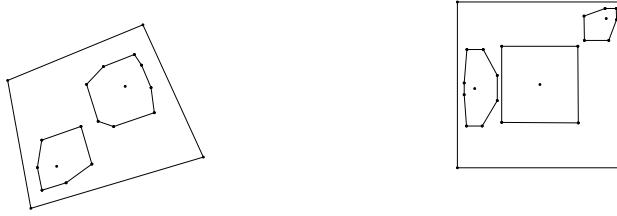


Fig. 4: Hilbert Balls in quadrangle domains have combinatorial complexity depending on the center location.

Table 1: Comparing the three geometric modelings of the probability simplex  $\Delta^d$ .

Riemannian Geometry	Information Rie. Geo.	Non-Rie. Hilbert Geo.
Structure $(\Delta^d, g, \nabla^{LC} = \nabla(g))$ Levi-Civita $\nabla^{LC} = \nabla^{(0)}$	$(\Delta^d, g, \nabla^{(\alpha)}, \nabla^{(-\alpha)})$ dual connections $\nabla^{(\pm\alpha)}$ so that $\frac{\nabla^{(\alpha)} + \nabla^{(-\alpha)}}{2} = \nabla^{(0)}$	$(\Delta^d, \rho)$ connection of $\mathbb{R}^d$
Distance Rao distance (metric) Property invariant to reparameterization	$\alpha$ -divergence (non-metric) KL or reverse KL for $\alpha = \pm 1$ information monotonicity	Hilbert distance (metric) isometric to a normed space
Calculation closed-form Geodesic shortest path	closed-form straight either in $\theta/\eta$	easy (Alg. 1) straight
Smoothness manifold Curvature positive	manifold dually flat	non-manifold negative

visualize metric tensors [33]). The isometries in Hilbert polyhedral geometries are studied in [36]. In Appendix 9, we recall that any Hilbert geometry induces a Finslerian structure that becomes Riemannian iff the boundary is an ellipsoid (yielding the hyperbolic Cayley-Klein geometries [56]). Notice that in Hilbert simplex/polytope geometry, the geodesics are not unique (see Figure 2 of [20]).

Table 1 summarizes the characteristics of the three introduced geometries: FHR, IG, and HG. Let us conclude this introduction by mentioning the Cramér-Rao lower bound and its relationship with information geometry [39]: Consider an unbiased estimator  $\hat{\theta} = T(X)$  of a parameter  $\theta$  estimated from measurements distributed ac-

cording to a smooth density  $p(x; \theta)$  (i.e.,  $X \sim p(x; \theta)$ ). The Cramér-Rao Lower Bound (CRLB) states that the variance of  $T(X)$  is greater or equal to the inverse of the FIM  $\mathcal{I}(\theta)$ :  $V_\theta[T(X)] \succ \mathcal{I}^{-1}(\theta)$ . For regular parametric families  $\{p(x; \theta)\}_\theta$ , the FIM is a positive-definite matrix and defines a metric tensor, called the Fisher metric in Riemannian geometry. The FIM is the cornerstone of information geometry [2] but requires the differentiability of the probability density function (pdf).

A better lower bound that does not require the pdf differentiability is the Hammersley-Chapman-Robbins Lower Bound [26, 18] (HCRLB):

$$V_\theta[T(X)] \geq \sup_{\Delta} \frac{\Delta^2}{E_\theta \left[ \left( \frac{p(x; \theta + \Delta) - p(x; \theta)}{p(x; \theta)} \right)^2 \right]}. \quad (5)$$

By introducing the  $\chi^2$ -divergence,  $\chi^2(P : Q) = \int \left( \frac{dP - dQ}{dQ} \right)^2 dQ$ , we rewrite the HCRLB using the  $\chi^2$ -divergence in the denominator as follows:

$$V_\theta[T(X)] \geq \sup_{\Delta} \frac{\Delta^2}{\chi^2(P(x; \theta + \Delta) : P(x; \theta))}. \quad (6)$$

Note that the FIM is not defined for non-differentiable pdfs, and therefore the Cramér-Rao lower bound does not exist in that case.

### 3 Computing Hilbert distance in $\Delta^d$

Let us start by the simplest case: The 1D probability simplex  $\Delta^1$ , the space of Bernoulli distributions. Any Bernoulli distribution can be represented by the activation probability of the random bit  $x$ :  $\lambda = p(x = 1) \in \Delta^1$ , corresponding to a point in the interval  $\Delta^1 = (0, 1)$ . We write the Bernoulli manifold as an exponential family as

$$p(x) = \exp(x\theta - F(\theta)), \quad x \in \{0, 1\},$$

where  $F(\theta) = \log(1 + \exp(\theta))$ . Therefore  $\lambda = \frac{\exp(\theta)}{1 + \exp(\theta)}$  and  $\theta = \log \frac{\lambda}{1 - \lambda}$ .

#### 3.1 1D probability simplex of Bernoulli distributions

By definition, the Hilbert distance has the closed form:

$$\rho_{\text{HG}}(p, q) = \left| \log \frac{\lambda_q(1 - \lambda_p)}{\lambda_p(1 - \lambda_q)} \right| = \left| \log \frac{\lambda_p}{1 - \lambda_p} - \log \frac{\lambda_q}{1 - \lambda_q} \right|.$$

Note that  $\theta_p = \log \frac{\lambda_p}{1 - \lambda_p}$  is the canonical parameter of the Bernoulli distribution.

The FIM of the Bernoulli manifold in the  $\lambda$ -coordinates is given by:  $g = \frac{1}{\lambda} + \frac{1}{1-\lambda} = \frac{1}{\lambda(1-\lambda)}$ . The FHR distance is obtained by integration as:

$$\rho_{\text{FHR}}(p, q) = 2 \arccos \left( \sqrt{\lambda_p \lambda_q} + \sqrt{(1 - \lambda_p)(1 - \lambda_q)} \right).$$

Notice that  $\rho_{\text{FHR}}(p, q)$  has finite values on  $\partial\Delta^1$ .

The KL divergence of the  $\pm 1$ -geometry is:

$$\rho_{\text{IG}}(p, q) = \lambda_p \log \frac{\lambda_p}{\lambda_q} + (1 - \lambda_p) \log \frac{1 - \lambda_p}{1 - \lambda_q}.$$

The KL divergence belongs to the family of  $\alpha$ -divergences [2].

### 3.2 Arbitrary dimension case

Given  $p, q \in \Delta^d$ , we first need to compute the intersection of line  $(pq)$  with the border of the  $d$ -dimensional probability simplex to get the two intersection points  $p'$  and  $q'$  so that  $p', p, q, q'$  are ordered on  $(pq)$ . Once this is done, we simply apply the formula in Eq. 4 to get the Hilbert distance.

A  $d$ -dimensional simplex consists of  $d+1$  vertices with their corresponding  $(d-1)$ -dimensional facets. For the probability simplex  $\Delta^d$ , let  $e_i = (\underbrace{0, \dots, 0}_i, 1, 0, \dots, 0)$

denote the  $d+1$  vertices of the standard simplex embedded in the hyperplane  $H_\Delta : \sum_{i=0}^d \lambda^i = 1$  in  $\mathbb{R}^{d+1}$ . Let  $f_{\setminus j}$  denote the simplex facets that is the convex hull of all vertices except  $e_j$ :  $f_{\setminus j} = \text{hull}(e_0, \dots, e_{j-1}, e_{j+1}, \dots, e_d)$ . Let  $H_{\setminus j}$  denote the hyperplane supporting this facet, which is the affine hull  $f_{\setminus j} = \text{affine}(e_0, \dots, e_{j-1}, e_{j+1}, \dots, e_d)$ .

To compute the two intersection points of  $(pq)$  with  $\Delta^d$ , a naive algorithm consists in computing the unique intersection point  $r_j$  of the line  $(pq)$  with each hyperplane  $H_{\setminus j}$  ( $j = 0, \dots, d$ ) and checking whether  $r_j$  belongs to  $f_{\setminus j}$ .

A much more efficient implementation given by Alg. (1) calculates the intersection point of the line  $x(t) = (1-t)p + tq$  with each  $H_{\setminus j}$  ( $j = 0, \dots, d$ ). These intersection points are represented using the coordinate  $t$ . For example,  $x(0) = p$  and  $x(1) = q$ . Due to convexity, any intersection point with  $H_{\setminus j}$  must satisfy either  $t \leq 0$  or  $t \geq 1$ . Then, the two intersection points with  $\partial\Delta^d$  are obtained by  $t_0 = \max\{t : \exists j, x(t) \in H_{\setminus j} \text{ and } t \leq 0\}$  and  $t_1 = \min\{t : \exists j, x(t) \in H_{\setminus j} \text{ and } t \geq 1\}$ . This algorithm only requires  $O(d)$  time and  $O(1)$  memoery.

**Lemma 1.** *The Hilbert distance in the probability simplex can be computed in optimal  $\Theta(d)$  time.*

**Algorithm 1:** Computing the Hilbert distance

---

**Data:** Two points  $p = (\lambda_p^0, \dots, \lambda_p^d), q = (\lambda_q^0, \dots, \lambda_q^d)$  in the  $d$ -dimensional simplex  $\Delta^d$

**Result:** Their Hilbert distance  $\rho_{\text{HG}}(p, q)$

```

1 begin
2    $t_0 \leftarrow -\infty; t_1 \leftarrow +\infty;$ 
3   for  $i = 0 \dots d$  do
4     if  $\lambda_p^i \neq \lambda_q^i$  then
5        $t \leftarrow \lambda_p^i / (\lambda_p^i - \lambda_q^i);$ 
6       if  $t_0 < t \leq 0$  then
7          $t_0 \leftarrow t;$ 
8       else if  $1 \leq t < t_1$  then
9          $t_1 \leftarrow t;$ 
10      if  $t_0 = -\infty$  or  $t_1 = +\infty$  then
11        Output  $\rho_{\text{HG}}(p, q) = 0;$ 
12      else if  $t_0 = 0$  or  $t_1 = 1$  then
13        Output  $\rho_{\text{HG}}(p, q) = \infty;$ 
14      else
15        Output  $\rho_{\text{HG}}(p, q) = \left| \log(1 - \frac{1}{t_0}) - \log(1 - \frac{1}{t_1}) \right|;$ 

```

---

Once an arbitrary distance  $\rho$  is chosen, we can define a ball centered at  $c$  and of radius  $r$  as  $B_\rho(c, r) = \{x : \rho(c, x) \leq r\}$ . Figure 3 displays the hexagonal shapes of the Hilbert balls for various center locations in  $\Delta^2$ .

**Theorem 1 (Balls in a simplicial Hilbert geometry [35]).** *A ball in the Hilbert simplex geometry has a Euclidean polytope shape with  $d(d+1)$  facets.*

Note that when the domain is not simplicial, the Hilbert balls can have varying combinatorial complexity depending on the center location. In 2D, the Hilbert ball can have  $s \sim 2s$  edges inclusively, where  $s$  is the number of edges of the boundary of the Hilbert domain  $\partial\mathcal{C}$ .

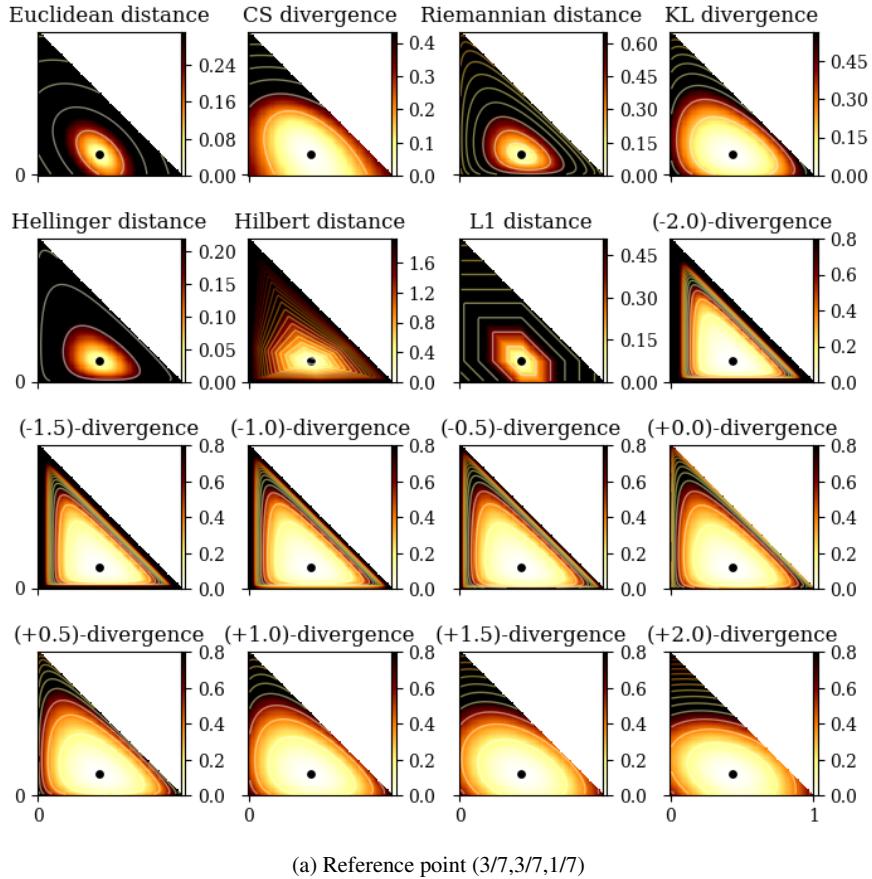
Since a Riemannian geometry is locally defined by a metric tensor, at infinitesimal scales, Riemannian balls have Mahalanobis smooth ellipsoidal shapes:  $B_\rho(c, r) = \{x : (x - c)^\top g(c)(x - c) \leq r^2\}$ . This property allows one to visualize Riemannian metric tensors [33]. Thus we conclude that:

**Lemma 2 ([35]).** *Hilbert simplex geometry is a non-manifold metric length space.*

As a remark, let us notice that slicing a simplex with a hyperplane does not always produce a lower-dimensional simplex. For example, slicing a tetrahedron by a plane yields either a triangle or a quadrilateral. Thus the restriction of a  $d$ -dimensional ball  $B$  in a Hilbert simplex geometry  $\Delta^d$  to a hyperplane  $H$  is a  $(d-1)$ -dimensional ball  $B' = B \cap H$  of varying combinatorial complexity, corresponding to a ball in the induced Hilbert sub-geometry in the convex sub-domain  $H \cap \Delta^d$ .

### 3.3 Visualizing distance profiles

Figure 5 displays the distance profile from any point in the probability simplex to a fixed reference point (trinomial) based on the following common distance measures [15]: Euclidean (metric) distance, Cauchy-Schwarz (CS) divergence, Hellinger (metric) distance, Fisher-Rao (metric) distance, KL divergence and Hilbert simplicial (metric) distance. The Euclidean and Cauchy-Schwarz divergence are clipped to  $\Delta^2$ . The Cauchy-Schwarz distance is projective so that  $\rho_{\text{CS}}(\lambda p, \lambda' q) = \rho_{\text{CS}}(p, q)$  for any  $\lambda, \lambda' > 0$  [49].



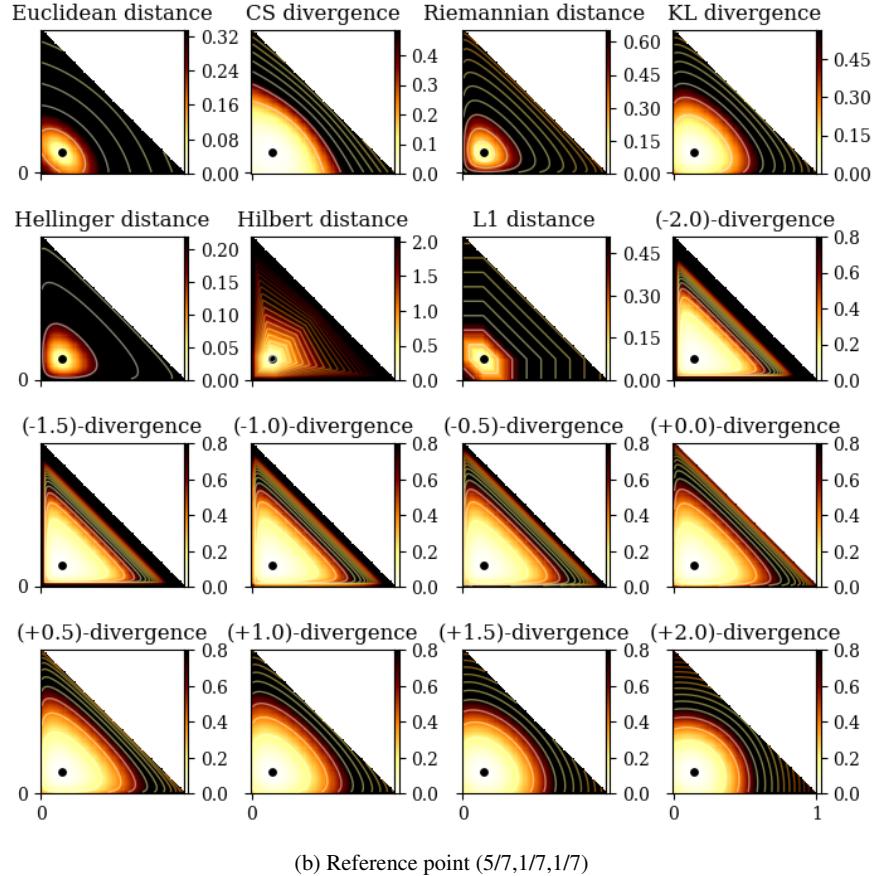


Fig. 5: A comparison of different distance measures on  $\Delta^2$ . The distance is measured from  $\forall p \in \Delta^2$  to a fixed reference point (the black dot). Lighter color means shorter distance. Darker color means longer distance. The contours show equal distance curves with a precision step of 0.2.

#### 4 Center-based clustering

We concentrate on comparing the efficiency of Hilbert simplex geometry for clustering multinomials. We shall compare the experimental results of  $k$ -means++ and  $k$ -center multinomial clustering for the three distances: Rao and Hilbert metric distances, and KL divergence. We describe how to adapt those clustering algorithms to the Hilbert distance.

#### 4.1 *k-means++ clustering*

The celebrated  $k$ -means clustering minimizes the sum of within-cluster variances, where each cluster has a center representative element. When dealing with  $k = 1$  cluster, the center (also called centroid or cluster prototype) is the center of mass defined as the minimizer of

$$E_D(\Lambda, c) = \frac{1}{n} \sum_{i=1}^n D(p_i : c),$$

where  $D(\cdot : \cdot)$  is a dissimilarity measure. For an arbitrary  $D$ , the centroid  $c$  may not be available in closed form. Nevertheless, using a generalization of the  $k$ -means++ initialization [6] (picking randomly seeds), one can bypass the centroid computation, and yet guarantee probabilistically a good clustering.

Let  $C = \{c_1, \dots, c_k\}$  denote the set of  $k$  cluster centers. Then the generalized  $k$ -means energy to be minimized is defined by:

$$E_D(\Lambda, C) = \frac{1}{n} \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} D(p_i : c_j).$$

By defining the distance  $D(p, C) = \min_{j \in \{1, \dots, k\}} D(p : c_j)$  of a point to a set, we can rewrite the objective function as  $E_D(\Lambda, C) = \frac{1}{n} \sum_{i=1}^n D(p_i, C)$ . Let  $E_D^*(\Lambda, k) = \min_{C : |C|=k} E_D(\Lambda, C)$  denote the global minimum of  $E_D(\Lambda, C)$  wrt some given  $\Lambda$  and  $k$ .

The  $k$ -means++ seeding proceeds for an arbitrary divergence  $D$  as follows: Pick uniformly at random at first seed  $c_1$ , and then iteratively choose the  $(k - 1)$  remaining seeds according to the following probability distribution:

$$\Pr(c_j = p_i) = \frac{D(p_i, \{c_1, \dots, c_{j-1}\})}{\sum_{i=1}^n D(p_i, \{c_1, \dots, c_{j-1}\})} \quad (2 \leq j \leq k).$$

Since its inception (2007), this  $k$ -means++ seeding has been extensively studied [7]. We state the general theorem established by [46]:

**Theorem 2 (Generalized  $k$ -means++ performance, [46]).** *Let  $\kappa_1$  and  $\kappa_2$  be two constants such that  $\kappa_1$  defines the quasi-triangular inequality property:*

$$D(x : z) \leq \kappa_1 (D(x : y) + D(y : z)), \quad \forall x, y, z \in \Delta^d,$$

*and  $\kappa_2$  handles the symmetry inequality:*

$$D(x : y) \leq \kappa_2 D(y : x), \quad \forall x, y \in \Delta^d.$$

*Then the generalized  $k$ -means++ seeding guarantees with high probability a configuration  $C$  of cluster centers such that:*

$$E_D(\Lambda, C) \leq 2\kappa_1^2(1 + \kappa_2)(2 + \log k)E_D^*(\Lambda, k). \tag{7}$$

The ratio  $\frac{E_D(A,C)}{E_D^*(A,k)}$  is called the *competitive factor*. The seminal result of ordinary  $k$ -means++ was shown [6] to be  $8(2 + \log k)$ -competitive. When evaluating  $\kappa_1$ , one has to note that squared metric distances are not metric because they do not satisfy the triangular inequality. For example, the squared Euclidean distance is not a metric but it satisfies the 2-quasi triangular inequality with  $\kappa_1 = 2$ .

We state the following general performance theorem:

**Theorem 3 ( $k$ -means++ performance in a metric space).** *In any metric space  $(\mathcal{X}, d)$ , the  $k$ -means++ wrt the squared metric distance  $d^2$  is  $16(2 + \log k)$ -competitive.*

*Proof.* Since a metric distance is symmetric, it follows that  $\kappa_2 = 1$ . Consider the quasi-triangular inequality property for the squared non-metric dissimilarity  $d^2$ :

$$\begin{aligned} d(p, q) &\leq d(p, q) + d(q, r), \\ d^2(p, q) &\leq (d(p, q) + d(q, r))^2, \\ d^2(p, q) &\leq d^2(p, q) + d^2(q, r) + 2d(p, q)d(q, r). \end{aligned}$$

Let us apply the inequality of arithmetic and geometric means<sup>1</sup>:

$$\sqrt{d^2(p, q)d^2(q, r)} \leq \frac{d^2(p, q) + d^2(q, r)}{2}.$$

Thus we have

$$d^2(p, q) \leq d^2(p, q) + d^2(q, r) + 2d(p, q)d(q, r) \leq 2(d^2(p, q) + d^2(q, r)).$$

That is, the squared metric distance satisfies the 2-approximate triangle inequality, and  $\kappa_1 = 2$ . The result is straightforward from Theorem 2.

**Theorem 4 ( $k$ -means++ performance in a normed space).** *In any normed space  $(\mathcal{X}, \|\cdot\|)$ , the  $k$ -means++ with  $D(x : y) = \|x - y\|^2$  is  $16(2 + \log k)$ -competitive.*

*Proof.* In any normed space  $(\mathcal{X}, \|\cdot\|)$ , we have both  $\|x - y\| = \|y - x\|$  and the triangle inequality:

$$\|x - z\| \leq \|x - y\| + \|y - z\|.$$

The proof is very similar to the proof of Theorem 3 and is omitted.

Since any inner product space  $(\mathcal{X}, \langle \cdot, \cdot \rangle)$  has an induced norm  $\|x\| = \sqrt{\langle x, x \rangle}$ , we have the following corollary.

**Corollary 1.** *In any inner product space  $(\mathcal{X}, \langle \cdot, \cdot \rangle)$ , the  $k$ -means++ with  $D(x : y) = \langle x - y, x - y \rangle$  is  $16(2 + \log k)$ -competitive.*

We need to report a bound for the squared Hilbert symmetric distance ( $\kappa_2 = 1$ ). In [35] (Theorem 3.3), it was shown that Hilbert geometry of a bounded convex domain  $\mathcal{C}$  is isometric to a normed vector space iff  $\mathcal{C}$  is an open simplex:

---

<sup>1</sup> For positive values  $a$  and  $b$ , the arithmetic-geometric mean inequality states that  $\sqrt{ab} \leq \frac{a+b}{2}$ .

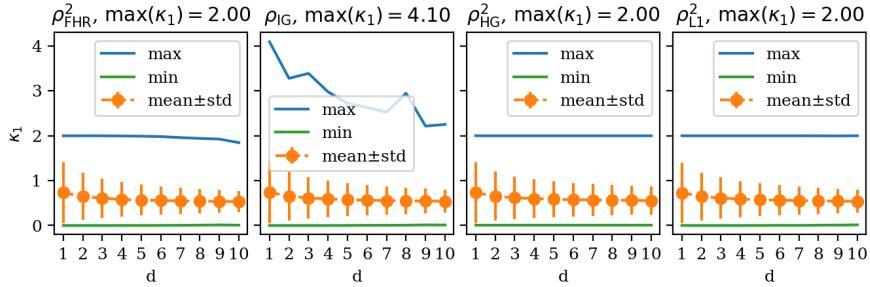


Fig. 6: The maximum, mean, standard deviation, and minimum of  $\kappa_1$  on  $10^6$  randomly generated tuples  $(x, y, z)$  in  $\Delta^d$  for  $d = 1, \dots, 10$ .

$(\Delta^d, \rho_{\text{HG}}) \simeq (V^d, \|\cdot\|_{\text{NH}})$ , where  $\|\cdot\|_{\text{NH}}$  is the corresponding norm. Therefore  $\kappa_1 = 2$ . We write “NH” for short for this equivalent normed Hilbert geometry. Appendix 8 recalls the construction due to [20], and shows the squared Hilbert distance fails the triangle inequality and it is not a distance induced by an inner product.

As an empirical study, we randomly generate  $n = 10^6$  tuples  $(x, y, z)$  based on the uniform distribution in  $\Delta^d$ . For each tuple  $(x, y, z)$ , we evaluate the ratio

$$\kappa_1 = \frac{D(x : z)}{D(x : y) + D(y : z)}.$$

Figure 6 shows the statistics for four different choices of  $D$ : (1)  $D(x : y) = \rho_{\text{FHR}}^2(x, y)$ ; (2)  $D(x : y) = \frac{1}{2}\text{KL}(x : y) + \frac{1}{2}\text{KL}(y : x)$ ; (3)  $D(x : y) = \rho_{\text{HG}}^2(x, y)$ ; (4)  $D(x : y) = \rho_{L1}^2(x, y)$ . We find experimentally that  $\kappa_1$  is upper bounded by 2 for  $\rho_{\text{FHR}}^2$ ,  $\rho_{\text{HG}}^2$  and  $\rho_{L1}^2$ , while the average  $\kappa_1$  value is smaller than 0.5. For all the compared distances,  $\kappa_2 = 1$ . Therefore  $\rho_{\text{FHR}}$  and  $\rho_{\text{HG}}$  have better  $k$ -means++ performance guarantee as compared to  $\rho_{\text{IG}}$ .

We get by applying Theorem 4:

**Corollary 2 ( $k$ -means++ in Hilbert simplex geometry).** *The  $k$ -means++ seeding in a Hilbert simplex geometry in fixed dimension is  $16(2 + \log k)$ -competitive.*

Figure 7 displays the clustering results of  $k$ -means++ in Hilbert simplex geometry as compared to the other geometries for  $k \in \{3, 5\}$ .

The KL divergence can be interpreted as a separable Bregman divergence [1]. The Bregman  $k$ -means++ performance has been studied in [1, 38], and a competitive factor of  $O(\frac{1}{\mu})$  is reported using the notion of Bregman  $\mu$ -similarity (that is suited for data-sets on a compact domain).

In [23], spherical  $k$ -means++ is studied wrt the distance  $d_S(x, y) = 1 - \langle x, y \rangle$  for any pair of points  $x, y$  on the unit sphere. Since  $\langle x, y \rangle = \|x\|_2\|y\|_2\cos(\theta_{x,y}) = \cos(\theta_{x,y})$ , we have  $d_S(x, y) = 1 - \cos(\theta_{x,y})$ , where  $\theta_{x,y}$  denotes the angle between a pair of unit vectors  $x$  and  $y$ . This distance is called the cosine distance since it amounts to one minus the cosine similarity. Notice that the cosine distance is related

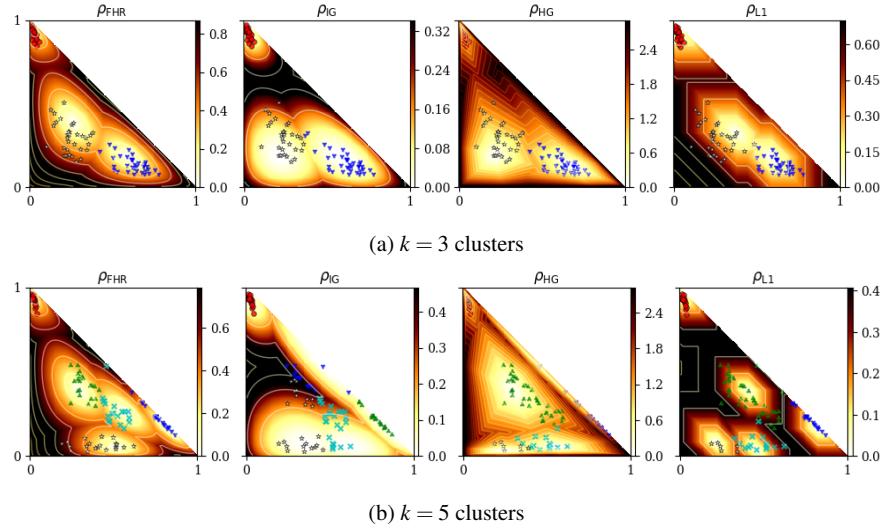


Fig. 7:  $k$ -Means++ clustering results on a toy dataset in the space of trinomials  $\Delta^2$ . The color density maps indicate the distance from any point to its nearest cluster center.

to the squared Euclidean distance via the identity:  $d_S(x, y) = \frac{1}{2} \|x - y\|^2$ . The cosine distance is different from the spherical distance that relies on the arccos function.

Since divergences may be asymmetric, one can further consider mixed divergence  $M(p : q : r) = \lambda D(p : q) + (1 - \lambda)D(q : r)$  for  $\lambda \in [0, 1]$ , and extend the  $k$ -means++ seeding procedure and analysis [47].

For a given data set, we can compute  $\kappa_1$  or  $\kappa_2$  by inspecting triples and pairs of points, and get data-dependent competitive factor improving the bounds mentioned above.

#### 4.2 $k$ -Center clustering

Let  $\Lambda$  be a finite point set. The cost function for a  $k$ -center clustering with centers  $C$  ( $|C| = k$ ) is:

$$f_D(\Lambda, C) = \max_{p_i \in \Lambda} \min_{c_j \in C} D(p_i : c_j).$$

The farthest first traversal heuristic [25] has a guaranteed approximation factor of 2 for any metric distance (see Algorithm 3).

In order to use the  $k$ -center clustering algorithm described in Algorithm 2, we need to be able to compute the 1-center (or minimax center) for the Hilbert sim-

**Algorithm 2:**  $k$ -Center clustering

---

**Data:** A set of points  $p_1, \dots, p_n \in \Delta^d$ . A distance measure  $\rho$  on  $\Delta^d$ . The maximum number  $k$  of clusters. The maximum number  $T$  of iterations.

**Result:** A clustering scheme assigning each  $p_i$  a label  $l_i \in \{1, \dots, k\}$

```

1 begin
2   Randomly pick  $k$  cluster centers  $c_1, \dots, c_k$  using the kmeans++ heuristic;
3   for  $t = 1, \dots, T$  do
4     for  $i = 1, \dots, n$  do
5        $l_i \leftarrow \arg \min_{l=1}^k \rho(p_i, c_l);$ 
6     for  $l = 1, \dots, k$  do
7        $c_l \leftarrow \arg \min_c \max_{i: l_i=l} \rho(p_i, c);$ 
8   Output  $\{l_i\}_{i=1}^n;$ 

```

---

**Algorithm 3:** A 2-approximation of the  $k$ -center clustering for any metric distance  $\rho$ .

---

**Data:** A set  $\Lambda$ ; a number  $k$  of clusters; a metric distance  $\rho$ .

**Result:** A 2-approximation of the  $k$ -center clustering

```

1 begin
2    $c_1 \leftarrow \text{ARandomPointOf}(\Lambda);$ 
3    $C \leftarrow \{c_1\};$ 
4   for  $i = 2, \dots, k$  do
5      $c_i \leftarrow \arg \max_{p \in \Lambda} \rho(p, C);$ 
6      $C \leftarrow C \cup \{c_i\};$ 
7   Output  $C;$ 

```

---

plex geometry, that is the Minimum Enclosing Ball (MEB, also called the Smallest Enclosing Ball, SEB).

We may consider the SEB equivalently either in  $\Delta^d$  or in the normed space  $V^d$ . In both spaces, the shapes of the balls are convex. Let  $\Lambda = \{p_1, \dots, p_n\}$  denote the point set in  $\Delta^d$ , and  $\mathcal{V} = \{v_1, \dots, v_n\}$  the equivalent point set in the normed vector space (following the mapping explained in Appendix 8). Then the SEBs  $B_{\text{HG}}(\Lambda)$  in  $\Delta^d$  and  $B_{\text{NH}}(\mathcal{V})$  in  $V^d$  have respectively radii  $r_{\text{HG}}^*$  and  $r_{\text{NH}}^*$  defined by:

$$\begin{aligned} r_{\text{HG}}^* &= \min_{c \in \Delta^d} \max_{i \in \{1, \dots, n\}} \rho_{\text{HG}}(p_i, c), \\ r_{\text{NH}}^* &= \min_{v \in V^d} \max_{i \in \{1, \dots, n\}} \|v_i - v\|_{\text{NH}}. \end{aligned}$$

The SEB in the normed vector space  $(V^d, \|\cdot\|_{\text{NH}})$  amounts to find the minimum covering norm polytope of a finite point set. This problem has been well-studied in computational geometry [57, 13, 51]. By considering the equivalent Hilbert norm polytope with  $d(d+1)$  facets, we state the result of [57]:

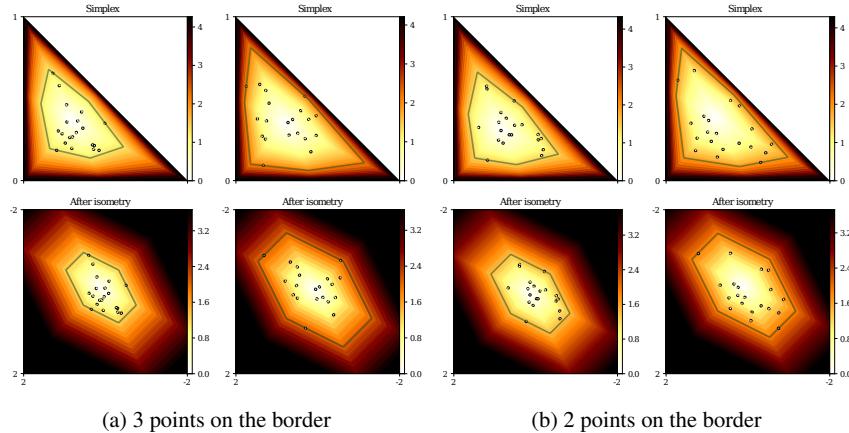


Fig. 8: Computing the SEB in Hilbert simplex geometry amounts to compute the SEB in the corresponding normed vector space.

**Theorem 5 (SEB in Hilbert polytope normed space, [57]).** A  $(1 + \varepsilon)$ -approximation of the SEB in  $V^d$  can be computed in  $O(d^3 \frac{n}{\varepsilon})$  time.

We shall now report two algorithms for computing the SEBs: One exact algorithm in  $V^d$  that does not scale well in high dimensions, and one approximation in  $\Delta^d$  that works well for large dimensions.

#### 4.2.1 Exact smallest enclosing ball in a Hilbert simplex geometry

Given a finite point set  $\{p_1, \dots, p_n\} \in \Delta^d$ , the SEB in Hilbert simplex geometry is centered at

$$c^* = \arg \min_{c \in \Delta^d} \max_{i \in \{1, \dots, n\}} \rho_{\text{HG}}(c, x_i),$$

with radius

$$r^* = \min_{c \in \Delta^d} \max_{i \in \{1, \dots, n\}} \rho_{\text{HG}}(c, x_i).$$

An equivalent problem is to find the SEB in the isometric normed vector space  $V^d$  via the mapping reported in Appendix 8. Each simplex point  $p_i$  corresponds to a point  $v_i$  in the  $V^d$ .

Figure 8 displays some examples of the exact smallest enclosing balls in the Hilbert simplex geometry and in the corresponding normed vector space.

To compute the SEB, one may also consider the generic LP-type randomized algorithm [44]. We notice that an enclosing ball for a point set in general has a number  $k$  of points on the border of the ball, with  $2 \leq k \leq \frac{d(d+1)}{2}$ . Let  $D = \frac{d(d+1)}{2}$  denote the varying size of the combinatorial basis, then we can apply the LP-type frame-

work (we check the axioms of locality and monotonicity, [58]) to solve efficiently the SEBs.

**Theorem 6 (Smallest Enclosing Hilbert Ball is LP-type, [63, 58]).** *The smallest enclosing Hilbert ball amounts to find the smallest enclosing ball in a vector space wrt a polytope norm that can be solved using a LP-type randomized algorithm.*

The Enclosing Ball Decision Problem (EBDP, [45]) asks for a given value  $r$ , whether  $r \geq r^*$  or not. The decision problem amounts to find whether a set  $\{rB_V + v_i\}$  of translates can be stabbed by a point [45]: That is, whether  $\cap_{i=1}^n (rB_V + v_i)$  is empty or not. Since these translates are polytopes with  $d(d+1)$  facets, this can be solved in linear time using *Linear Programming*.

**Theorem 7 (Enclosing Hilbert Ball Decision Problem).** *The decision problem to test whether  $r \geq r^*$  or not can be solved by Linear Programming.*

This yields a simple scheme to approximate the optimal value  $r^*$ : Let  $r_0 = \max_{i \in \{2, \dots, n\}} \|v_i - v_1\|_{NH}$ . Then  $r^* \in [\frac{r_0}{2}, r_0] = [a_0, b_0]$ . At stage  $i$ , perform a dichotomic search on  $[a_i, b_i]$  by answering the decision problem for  $r_{i+1} = \frac{a_i+b_i}{2}$ , and update the radius range accordingly [45].

However, the LP-type randomized algorithm or the decision problem-based algorithm do not scale well in high dimensions. Next, we introduce a simple approximation algorithm that relies on the fact that the line segment  $[pq]$  is a geodesic in Hilbert simplex geometry. (Geodesics are not unique. See Figure 2 of [20].)

#### 4.2.2 Geodesic bisection approximation heuristic

In Riemannian geometry, the 1-center can be arbitrarily finely approximated by a simple geodesic bisection algorithm [9, 5]. This algorithm can be extended to HG straightforwardly as detailed in Algorithm 4.

---

**Algorithm 4:** Geodesic walk for approximating the Hilbert minimax center, generalizing [9]

---

**Data:** A set of points  $p_1, \dots, p_n \in \Delta^d$ . The maximum number  $T$  of iterations.

**Result:**  $c \approx \arg \min_c \max_i \rho_{HG}(p_i, c)$

```

1 begin
2    $c_0 \leftarrow \text{ARandomPointOf}(\{p_1, \dots, p_n\});$ 
3   for  $t = 1, \dots, T$  do
4      $p \leftarrow \arg \max_{p_i} \rho_{HG}(p_i, c_{t-1});$ 
5      $c_t \leftarrow c_{t-1} \#_{1/(t+1)}^\rho p;$ 
6   Output  $c_T$ ;

```

---

The algorithm first picks up a point  $c_0$  at random from  $\Lambda$  as the initial center, then computes the farthest point  $p$  (with respect to the distance  $\rho$ ), and then walk

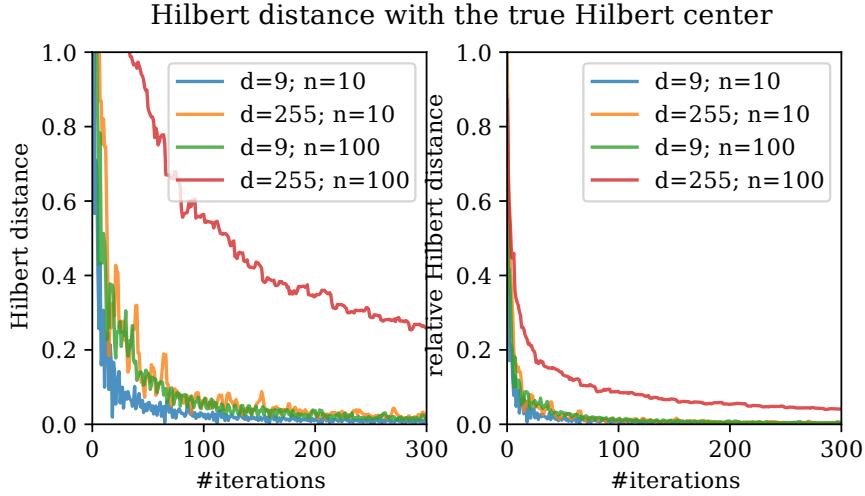


Fig. 9: Convergence rate of Alg. (4) measured by the Hilbert distance between the current minimax center and the true center (left) or their Hilbert distance divided by the Hilbert radius of the dataset (right). The plot is based on 100 random points in  $\Delta^9/\Delta^{255}$ .

on the geodesic from  $c_0$  to  $p$  by a certain amount to define  $c_1$ , etc. For an arbitrary distance  $\rho$ , we define the operator  $\#_\alpha^\rho$  as follows:

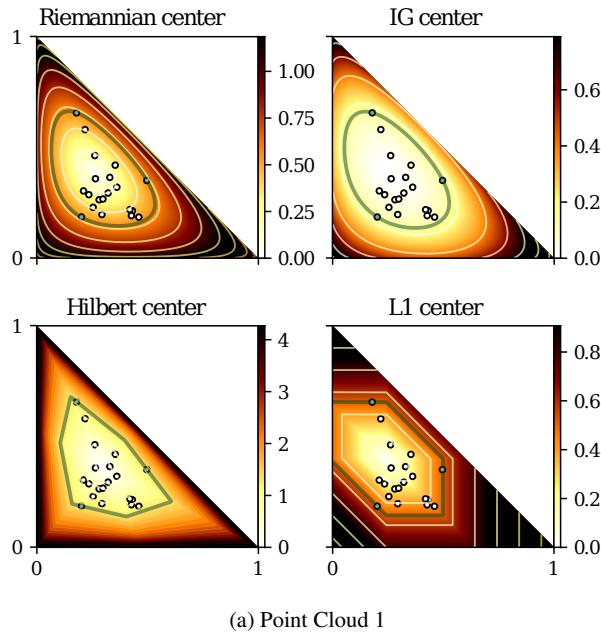
$$p \#_\alpha^\rho q = v = \gamma(p, q, \alpha), \quad \rho(p : v) = \alpha \rho(p : q),$$

where  $\gamma(p, q, \alpha)$  is the geodesic passing through  $p$  and  $q$ , and parameterized by  $\alpha$  ( $0 \leq \alpha \leq 1$ ). When the equations of the geodesics are explicitly known, we can either get a closed form solution for  $\#_\alpha^\rho$  or perform a bisection search to find  $v'$  such that  $\rho(p : v') \approx \alpha \rho(p : q)$ . See [40] for an extension and analysis in hyperbolic geometry. See Fig. (9) to get an intuitive idea on the *experimental* convergence rate of Algorithm 4. See Fig. (10) for visualizations of centers wrt different geometries.

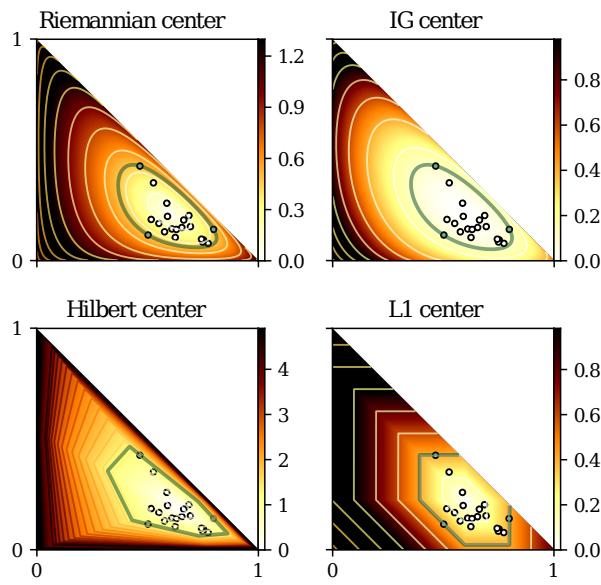
Furthermore, this iterative algorithm implies a core-set [10] (namely, the set of farthest points visited during the geodesic walks) that is useful for clustering large data-sets [8]. See [13] for core-set results on containment problems wrt a convex homothetic object (the equivalent Hilbert polytope norm in our case).

A simple algorithm dubbed MINCON [51] can find an approximation of the Minimum Enclosing Polytope. The algorithm induces a core-set of size  $O(\frac{1}{\epsilon^2})$  although the theorem is challenged in [13].

Thus by combining the  $k$ -center seeding [25] with the Lloyd-like batched iterations, we get an efficient  $k$ -center clustering algorithm for the FHR and Hilbert metric geometries. When dealing with the Kullback-Leibler divergence, we use the fact that KL is a Bregman divergence, and use the 1-center algorithm ([50, 43] for



(a) Point Cloud 1



(b) Point Cloud 2

approximation in any dimension, or [44] which is exact but limited to small dimensions).

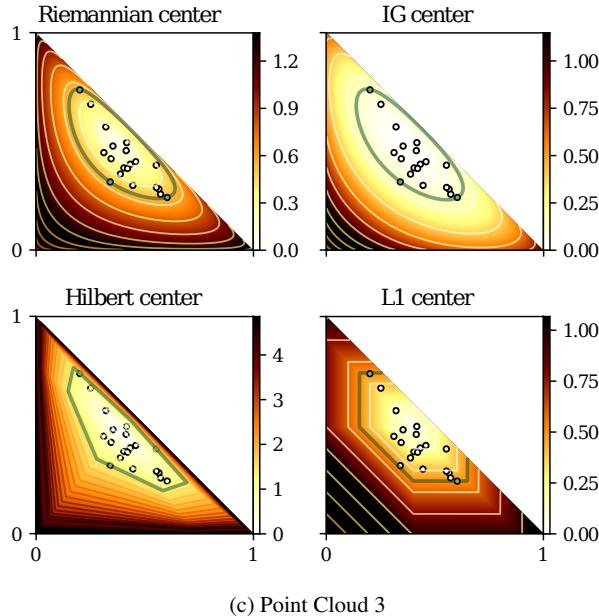


Fig. 10: The Riemannian/IG/Hilbert/ $L_1$  minimax centers of three point clouds in  $\Delta^2$  based on Alg. (4). The color maps show the distance from  $\forall p \in \Delta^2$  to the corresponding center.

Since Hilbert simplex geometry is isomorphic to a normed vector space [35] with a polytope norm with  $d(d + 1)$  facets, the Voronoi diagram in Hilbert geometry of  $\Delta^d$  amounts to compute a Voronoi diagram wrt a polytope norm [32, 55, 21].

## 5 Experiments

We generate a dataset consisting of a set of clusters in a high dimensional statistical simplex  $\Delta^d$ . Each cluster is generated independently as follows. We first pick a random center  $c = (\lambda_c^0, \dots, \lambda_c^d)$  based on the uniform distribution on  $\Delta^d$ . Then any random sample  $p = (\lambda^0, \dots, \lambda^d)$  associated with  $c$  is independently generated by

$$\lambda^i = \frac{\exp(\log \lambda_c^i + \sigma \varepsilon^i)}{\sum_{i=0}^d \exp(\log \lambda_c^i + \sigma \varepsilon^i)},$$

where  $\sigma > 0$  is a noise level parameter, and each  $\varepsilon^i$  follows independently a standard Gaussian distribution (generator 1) or the Student's  $t$ -distribution with five degrees of freedom (generator 2). Let  $\sigma = 0$ , we get  $\lambda^i = \lambda_c^i$ . Therefore  $p$  is randomly

distributed around  $c$ . We repeat generating random samples for each cluster center, and make sure that different clusters have almost the same number of samples. Then we perform clustering based on the configurations  $n \in \{50, 100\}$ ,  $d \in \{9, 255\}$ ,  $\sigma \in \{0.5, 0.9\}$ ,  $\rho \in \{\rho_{\text{FHR}}, \rho_{\text{IG}}, \rho_{\text{HG}}, \rho_{\text{EUC}}, \rho_{\text{L1}}\}$ . For simplicity, the number of clusters  $k$  is set to the ground truth. For each configuration, we repeat the clustering experiment based on 300 different random datasets. The performance is measured by the normalized mutual information (NMI), which is a scalar indicator in the range  $[0, 1]$  (the larger the better).

The results of  $k$ -means++ and  $k$ -centers are shown in Table 2 and Table 3, respectively. The large variance of NMI is because that each experiment is performed on random datasets wrt different random seeds. Generally, the performance deteriorates as we increase the number of clusters, increase the noise level or decrease the dimensionality, which have the same effect to reduce the inter-cluster gap.

The key comparison is the three columns  $\rho_{\text{FHR}}$ ,  $\rho_{\text{HG}}$  and  $\rho_{\text{IG}}$ , as they are based on exactly the same algorithm with the only difference being the underlying geometry. We see clearly that in general their clustering performance presents the order  $\text{HG} > \text{FHR} > \text{IG}$ . The performance of HG is superior to the other two geometries, especially when the noise level is large. Intuitively, the Hilbert balls are more compact in size and therefore can better capture the clustering structure (see Fig. (1)).

The column  $\rho_{\text{EUC}}$  is based on the Euclidean enclosing ball. It shows the worst scores because the intrinsic geometry of the probability simplex is far from the Euclidean geometry.

## 6 Hilbert geometry of the space of correlation matrices

In this section, we present the Hilbert geometry to the space of correlation matrices

$$\mathcal{C}^d = \{C_{d \times d} : C \succ 0; C_{ii} = 1, \forall i\}.$$

If  $C_1, C_2 \in \mathcal{C}$ , then  $(1 - \lambda)C_1 + \lambda C_2 \in \mathcal{C}$  for  $0 < \lambda < 1$ . Therefore  $\mathcal{C}$  is a convex set, known as an *elliptope* embedded in the p.s.d. cone. See Fig. (11) for an intuitive view of  $\mathcal{C}_3$ , where the coordinate system  $(x, y, z)$  is the off-diagonal entries of  $C \in \mathcal{C}_3$ .

In order to compute the Hilbert distance  $\rho_{\text{HG}}(C_1, C_2)$ , we need to compute the intersection of the line  $(C_1, C_2)$  with  $\partial\mathcal{C}$ , denoted as  $C'_1$  and  $C'_2$ , then we have

$$\rho_{\text{HG}}(C_1, C_2) = \left| \log \frac{\|C_1 - C'_2\| \|C'_1 - C_2\|}{\|C_1 - C'_1\| \|C_2 - C'_2\|} \right|.$$

Unfortunately there is no closed form solution of  $C'_1$  and  $C'_2$ . Instead, we apply a binary searching algorithm. Note a necessary condition for  $C \in \mathcal{C}$  is that  $C$  has a positive spectrum. If  $C$  has at least one non-positive eigenvalue, then  $C \notin \mathcal{C}$ . To determine whether a given  $C$  is inside the elliptope requires a spectral decomposition of  $C$ . Therefore the computation of  $C'_1$  and  $C'_2$  is in general expensive.

Table 2:  $k$ -means++ clustering accuracy in NMI on randomly generated datasets based on different geometries. The table shows the mean and standard deviation after 300 independent runs for each configuration.  $\rho$  is the distance measure.  $n$  is the sample size.  $d$  is the dimensionality of  $\Delta^d$ .  $\sigma$  is noise level.

$k$	$n$	$d$	$\sigma$	$\rho_{\text{FHR}}$	$\rho_{\text{IG}}$	$\rho_{\text{HG}}$	$\rho_{\text{EUC}}$	$\rho_{L1}$
3	50	9	0.5	0.76 ± 0.22	0.76 ± 0.24	<b>0.81 ± 0.22</b>	0.64 ± 0.23	0.70 ± 0.22
			0.9	0.44 ± 0.20	0.44 ± 0.20	<b>0.57 ± 0.22</b>	0.31 ± 0.17	0.38 ± 0.18
	100	9	0.5	0.80 ± 0.24	0.81 ± 0.24	<b>0.88 ± 0.21</b>	0.74 ± 0.25	0.79 ± 0.24
			0.9	0.65 ± 0.27	0.66 ± 0.28	<b>0.72 ± 0.27</b>	0.46 ± 0.24	0.63 ± 0.27
	255	9	0.5	0.76 ± 0.22	0.76 ± 0.21	<b>0.82 ± 0.22</b>	0.60 ± 0.21	0.69 ± 0.23
			0.9	0.42 ± 0.19	0.41 ± 0.18	<b>0.54 ± 0.22</b>	0.27 ± 0.14	0.34 ± 0.16
5	50	9	0.5	0.82 ± 0.23	0.82 ± 0.24	<b>0.89 ± 0.20</b>	0.74 ± 0.24	0.80 ± 0.25
			0.9	0.66 ± 0.26	0.66 ± 0.28	<b>0.72 ± 0.26</b>	0.45 ± 0.25	0.64 ± 0.27
	100	9	0.5	0.75 ± 0.14	0.74 ± 0.15	<b>0.81 ± 0.13</b>	0.61 ± 0.13	0.68 ± 0.13
			0.9	0.44 ± 0.13	0.42 ± 0.13	<b>0.55 ± 0.15</b>	0.31 ± 0.11	0.36 ± 0.12
	255	9	0.5	0.83 ± 0.15	0.83 ± 0.15	<b>0.88 ± 0.14</b>	0.77 ± 0.16	0.82 ± 0.15
			0.9	0.71 ± 0.17	0.70 ± 0.19	<b>0.75 ± 0.17</b>	0.50 ± 0.17	0.68 ± 0.18
10	50	9	0.5	0.74 ± 0.13	0.74 ± 0.14	<b>0.80 ± 0.14</b>	0.60 ± 0.13	0.67 ± 0.13
			0.9	0.42 ± 0.11	0.40 ± 0.12	<b>0.55 ± 0.15</b>	0.29 ± 0.09	0.35 ± 0.11
	100	9	0.5	0.83 ± 0.14	0.83 ± 0.15	<b>0.88 ± 0.13</b>	0.77 ± 0.15	0.81 ± 0.15
			0.9	0.69 ± 0.18	0.69 ± 0.18	<b>0.73 ± 0.17</b>	0.48 ± 0.17	0.67 ± 0.18
	255	9	0.5	0.75 ± 0.14	0.74 ± 0.15	<b>0.81 ± 0.13</b>	0.61 ± 0.13	0.68 ± 0.13
			0.9	0.44 ± 0.13	0.42 ± 0.13	<b>0.55 ± 0.15</b>	0.31 ± 0.11	0.36 ± 0.12

(a) generator 1

$k$	$n$	$d$	$\sigma$	$\rho_{\text{FHR}}$	$\rho_{\text{IG}}$	$\rho_{\text{HG}}$	$\rho_{\text{EUC}}$	$\rho_{L1}$
3	50	9	0.5	0.62 ± 0.22	0.60 ± 0.22	<b>0.71 ± 0.23</b>	0.45 ± 0.20	0.54 ± 0.22
			0.9	0.29 ± 0.17	0.27 ± 0.16	<b>0.39 ± 0.19</b>	0.17 ± 0.13	0.25 ± 0.15
	100	9	0.5	0.70 ± 0.25	0.69 ± 0.26	<b>0.74 ± 0.25</b>	0.37 ± 0.29	0.70 ± 0.26
			0.9	<b>0.42 ± 0.25</b>	0.35 ± 0.20	0.40 ± 0.19	0.03 ± 0.08	<b>0.44 ± 0.26</b>
	255	9	0.5	0.63 ± 0.22	0.61 ± 0.22	<b>0.71 ± 0.22</b>	0.46 ± 0.19	0.56 ± 0.20
			0.9	0.29 ± 0.15	0.26 ± 0.14	<b>0.38 ± 0.20</b>	0.18 ± 0.12	0.24 ± 0.14
5	50	9	0.5	0.71 ± 0.26	0.69 ± 0.27	<b>0.75 ± 0.25</b>	0.31 ± 0.28	0.70 ± 0.27
			0.9	0.41 ± 0.26	0.33 ± 0.20	0.38 ± 0.18	0.02 ± 0.06	<b>0.43 ± 0.26</b>
	100	9	0.5	0.64 ± 0.15	0.61 ± 0.14	<b>0.70 ± 0.14</b>	0.48 ± 0.14	0.57 ± 0.15
			0.9	0.31 ± 0.12	0.29 ± 0.12	<b>0.41 ± 0.15</b>	0.20 ± 0.09	0.26 ± 0.10
	255	9	0.5	0.74 ± 0.17	0.72 ± 0.17	<b>0.77 ± 0.16</b>	0.41 ± 0.20	0.74 ± 0.17
			0.9	0.44 ± 0.17	0.37 ± 0.16	0.44 ± 0.15	0.04 ± 0.06	<b>0.47 ± 0.17</b>
10	50	9	0.5	0.62 ± 0.14	0.61 ± 0.14	<b>0.71 ± 0.14</b>	0.46 ± 0.13	0.54 ± 0.14
			0.9	0.30 ± 0.10	0.27 ± 0.11	<b>0.40 ± 0.13</b>	0.19 ± 0.08	0.25 ± 0.09
	100	9	0.5	0.73 ± 0.18	0.70 ± 0.18	<b>0.75 ± 0.16</b>	0.37 ± 0.20	0.73 ± 0.17
			0.9	0.43 ± 0.16	0.35 ± 0.14	0.41 ± 0.12	0.03 ± 0.06	<b>0.46 ± 0.18</b>
	255	9	0.5	0.75 ± 0.14	0.74 ± 0.15	<b>0.81 ± 0.13</b>	0.61 ± 0.13	0.68 ± 0.13
			0.9	0.44 ± 0.13	0.42 ± 0.13	<b>0.55 ± 0.15</b>	0.31 ± 0.11	0.36 ± 0.12

(b) generator 2

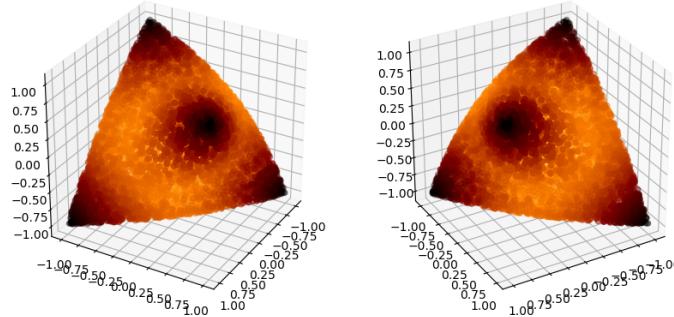
Table 3:  $k$ -center clustering accuracy in NMI on randomly generated datasets based on different geometries. The table shows the mean and standard deviation after 300 independent runs for each configuration.  $\rho$  is the distance measure.  $n$  is the sample size.  $d$  is the dimensionality of the statistical simplex.  $\sigma$  is noise level.

$k$	$n$	$d$	$\sigma$	$\rho_{\text{FHR}}$	$\rho_{\text{IG}}$	$\rho_{\text{HG}}$	$\rho_{\text{EUC}}$	$\rho_{L1}$
3	50	9	0.5	0.87 ± 0.19	0.85 ± 0.19	<b>0.92 ± 0.16</b>	0.72 ± 0.22	0.80 ± 0.20
			0.9	0.54 ± 0.21	0.51 ± 0.21	<b>0.70 ± 0.23</b>	0.36 ± 0.17	0.44 ± 0.19
		255	0.5	0.93 ± 0.16	0.92 ± 0.18	<b>0.95 ± 0.14</b>	0.89 ± 0.18	0.90 ± 0.19
	100	9	0.5	0.76 ± 0.24	0.72 ± 0.26	<b>0.82 ± 0.24</b>	0.50 ± 0.28	0.76 ± 0.25
			0.9	0.88 ± 0.17	0.86 ± 0.18	<b>0.93 ± 0.14</b>	0.70 ± 0.20	0.80 ± 0.20
		255	0.5	0.53 ± 0.20	0.49 ± 0.19	<b>0.70 ± 0.22</b>	0.33 ± 0.14	0.41 ± 0.18
5	50	9	0.5	0.93 ± 0.16	0.92 ± 0.17	<b>0.95 ± 0.13</b>	0.88 ± 0.19	0.93 ± 0.16
			0.9	0.81 ± 0.22	0.75 ± 0.24	<b>0.83 ± 0.22</b>	0.47 ± 0.28	0.79 ± 0.22
		255	0.5	0.50 ± 0.13	0.47 ± 0.13	<b>0.66 ± 0.15</b>	0.34 ± 0.11	0.40 ± 0.12
	100	9	0.5	<b>0.92 ± 0.11</b>	<b>0.91 ± 0.12</b>	<b>0.93 ± 0.11</b>	0.87 ± 0.13	<b>0.92 ± 0.12</b>
			0.9	0.77 ± 0.15	0.71 ± 0.17	<b>0.85 ± 0.17</b>	0.54 ± 0.19	0.74 ± 0.16
		255	0.5	0.48 ± 0.12	0.46 ± 0.12	<b>0.66 ± 0.15</b>	0.33 ± 0.09	0.39 ± 0.10
	100	9	0.5	<b>0.93 ± 0.10</b>	<b>0.92 ± 0.11</b>	<b>0.94 ± 0.09</b>	0.89 ± 0.11	0.92 ± 0.11
		255	0.9	0.81 ± 0.14	0.74 ± 0.15	<b>0.84 ± 0.16</b>	0.52 ± 0.19	0.79 ± 0.14

(a) generator 1

$k$	$n$	$d$	$\sigma$	$\rho_{\text{FHR}}$	$\rho_{\text{IG}}$	$\rho_{\text{HG}}$	$\rho_{\text{EUC}}$	$\rho_{L1}$
3	50	9	0.5	0.68 ± 0.22	0.67 ± 0.22	<b>0.80 ± 0.20</b>	0.48 ± 0.22	0.60 ± 0.22
			0.9	0.32 ± 0.18	0.29 ± 0.17	<b>0.45 ± 0.21</b>	0.20 ± 0.14	0.26 ± 0.15
		255	0.5	0.79 ± 0.24	0.75 ± 0.24	<b>0.82 ± 0.22</b>	0.13 ± 0.23	<b>0.81 ± 0.24</b>
	100	9	0.5	0.35 ± 0.27	0.35 ± 0.21	<b>0.42 ± 0.19</b>	0.00 ± 0.02	0.32 ± 0.30
			0.9	0.30 ± 0.16	0.28 ± 0.14	<b>0.42 ± 0.19</b>	0.20 ± 0.12	0.26 ± 0.14
		255	0.5	0.78 ± 0.25	0.76 ± 0.24	<b>0.82 ± 0.21</b>	0.05 ± 0.14	0.77 ± 0.27
5	50	9	0.5	0.29 ± 0.28	0.29 ± 0.20	<b>0.39 ± 0.20</b>	0.00 ± 0.02	0.22 ± 0.25
			0.9	0.69 ± 0.14	0.66 ± 0.14	<b>0.77 ± 0.13</b>	0.50 ± 0.13	0.61 ± 0.14
		255	0.5	0.34 ± 0.12	0.30 ± 0.12	<b>0.46 ± 0.15</b>	0.22 ± 0.09	0.28 ± 0.10
	100	9	0.5	0.42 ± 0.21	0.38 ± 0.16	<b>0.46 ± 0.15</b>	0.00 ± 0.02	0.39 ± 0.22
			0.9	0.66 ± 0.13	0.64 ± 0.14	<b>0.77 ± 0.14</b>	0.47 ± 0.13	0.57 ± 0.13
		255	0.5	0.31 ± 0.11	0.28 ± 0.10	<b>0.44 ± 0.13</b>	0.21 ± 0.08	0.25 ± 0.09
	100	9	0.5	0.32 ± 0.19	0.30 ± 0.15	<b>0.41 ± 0.13</b>	0.00 ± 0.01	0.26 ± 0.18

(b) generator 2

Fig. 11: The ellipope  $\mathcal{C}_3$  (two different perspectives).Table 4: NMI (mean $\pm$ std) of  $k$ -means++ clustering based on different distance measures in the ellipope (500 independent runs)

$v_1$	$v_2$	$\rho_{HG}$	$\rho_{EUC}$	$\rho_{L1}$	$\rho_{LD}$
4	10	<b><math>0.62 \pm 0.22</math></b>	$0.57 \pm 0.21$	$0.56 \pm 0.22$	$0.58 \pm 0.22$
4	30	<b><math>0.85 \pm 0.18</math></b>	$0.80 \pm 0.20$	$0.81 \pm 0.19$	$0.82 \pm 0.20$
4	50	<b><math>0.89 \pm 0.17</math></b>	$0.87 \pm 0.17$	$0.86 \pm 0.18$	$0.88 \pm 0.18$
5	10	<b><math>0.50 \pm 0.21</math></b>	$0.49 \pm 0.21$	$0.48 \pm 0.20$	$0.47 \pm 0.21$
5	30	<b><math>0.77 \pm 0.20</math></b>	$0.75 \pm 0.21$	$0.75 \pm 0.21$	$0.75 \pm 0.21$
5	50	<b><math>0.84 \pm 0.19</math></b>	$0.82 \pm 0.19$	$0.82 \pm 0.20$	<b><math>0.84 \pm 0.18</math></b>

We compare the Hilbert ellipope geometry with commonly used distance measures including the  $L_2$  distance  $\rho_{EUC}$ ,  $L_1$  distance  $\rho_{L1}$ , and the square root of the log-det divergence

$$\rho_{LD}(C_1, C_2) = \text{tr}(C_1 C_2^{-1}) - \log |C_1 C_2^{-1}| - d.$$

Due to the high computational complexity, we only investigate  $k$ -means++ clustering. The investigated dataset consists of 100 matrices forming 3 clusters in  $\mathcal{C}_3$  with almost identical size. Each cluster is independently generated according to

$$\begin{aligned} P &\sim \mathcal{W}^{-1}(I_{3 \times 3}, v_1), \\ C_i &\sim \mathcal{W}^{-1}(P, v_2), \end{aligned}$$

where  $\mathcal{W}^{-1}(A, v)$  denotes the inverse Wishart distribution with scale matrix  $A$  and  $v$  degrees of freedom, and  $C_i$  is a point in the cluster associated with  $P$ . Table 4 shows the  $k$ -means++ clustering performance in terms of NMI. Again Hilbert geometry is favorable as compared to alternatives, showing that the good performance of Hilbert clustering is generalizable.

## 7 Conclusion

We introduced the Hilbert metric distance and its underlying non-Riemannian geometry for modeling the space of multinomials or the open probability simplex. We compared experimentally in simulated clustering tasks this geometry with the traditional differential geometric modelings (either FHR metric connection or dually coupled non-metric affine connections of information geometry [2]).

The main feature of HG is that it is a metric non-manifold geometry, where geodesics are straight (Euclidean) line segments. In simplex domains, the Hilbert balls have fixed combinatorial (Euclidean) polytope structures, and HG is known to be isometric to a normed space [20, 24]. This latter isometry allows one to generalize easily the standard proofs of clustering (*e.g.*,  $k$ -means or  $k$ -center). We demonstrated it for the  $k$ -means++ competitive performance analysis, and for the convergence of the  $k$ -center heuristic [9] (smallest enclosing Hilbert ball allows one to implement efficiently the  $k$ -center clustering). Our experimental  $k$ -means++ or  $k$ -center comparisons of HG algorithms with the manifold modeling approach yield superior performance. This may be intuitively explained by the sharpness of Hilbert balls as compared to the FHR/IG ball profiles.

Chentsov [17] defined statistical invariance on a probability manifold under Markov morphisms, and proved that the Fisher Information Metric is the unique Riemannian metric (up to rescaling) for multinomials. However, this does not rule out that other distances (with underlying geometric structures) may be used to model statistical manifolds (*e.g.*, Finsler statistical manifolds [16, 59], or the total variation distance — the only metric  $f$ -divergence [31]). Defining statistical invariance related to geometry is the cornerstone problem of information geometry that can be tackled from many directions (see [22] and references therein for a short review).

In this paper, we introduced Hilbert geometries in machine learning by considering clustering tasks in the probability simplex and in the ellipope. Hilbert geometries proved computationally handy since geodesics are straight lines. One future direction is to consider the Hilbert metric for regularization and sparsity in machine learning (due to its equivalence with a polytope normed distance).

Our Python codes are freely available online for reproducible research:

<https://www.lix.polytechnique.fr/~nielsen/HSG/>

## References

1. Ackermann, M.R., Blömer, J.: Bregman clustering for separable instances. In: Scandinavian Workshop on Algorithm Theory, pp. 212–223. Springer (2010)
2. Amari, S.i.: Information Geometry and Its Applications, *Applied Mathematical Sciences*, vol. 194. Springer Japan (2016)
3. Amari, S.i., Cichocki, A.: Information geometry of divergence functions. Bulletin of the Polish Academy of Sciences: Technical Sciences **58**(1), 183–195 (2010)

4. Arnaudon, M., Nielsen, F.: Medians and means in Finsler geometry. *LMS Journal of Computation and Mathematics* **15**, 23–37 (2012)
5. Arnaudon, M., Nielsen, F.: On approximating the Riemannian 1-center. *Computational Geometry* **46**(1), 93–104 (2013)
6. Arthur, D., Vassilvitskii, S.:  $k$ -means++: The advantages of careful seeding. In: ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 1027–1035 (2007)
7. Bachem, O., Lucic, M., Hassani, S.H., Krause, A.: Approximate  $k$ -means++ in sublinear time. In: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, pp. 1459–1467 (2016)
8. Bachem, O., Lucic, M., Krause, A.: Scalable and distributed clustering via lightweight coresets (2017). arXiv:1702.08248 [stat.ML]
9. Bădoiu, M., Clarkson, K.L.: Smaller core-sets for balls. In: ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 801–802 (2003)
10. Bădoiu, M., Clarkson, K.L.: Optimal core-sets for balls. *Computational Geometry* **40**(1), 14–22 (2008)
11. Bernig, A.: Hilbert geometry of polytopes. *Archiv der Mathematik* **92**(4), 314–324 (2009)
12. Bi, Y., Fan, B., Wu, F.: Beyond Mahalanobis metric: Cayley-Klein metric learning. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2339–2347 (2015)
13. Brandenberg, R., König, S.: No dimension-independent core-sets for containment under homothetics. *Discrete & Computational Geometry* **49**(1), 3–21 (2013)
14. Busemann, H.: The Geometry of Geodesics, *Pure and Applied Mathematics*, vol. 6. Elsevier Science (1955)
15. Calin, O., Udriste, C.: Geometric Modeling in Probability and Statistics. Mathematics and Statistics. Springer International Publishing (2014)
16. Cena, A.: Geometric structures on the non-parametric statistical manifold. Ph.D. thesis, University of Milano (2002)
17. Cencov, N.N.: Statistical Decision Rules and Optimal Inference, *Translations of mathematical monographs*, vol. 53. American Mathematical Society (2000)
18. Chapman, D.G., Robbins, H.: Minimum variance estimation without regularity assumptions. *The Annals of Mathematical Statistics* pp. 581–586 (1951)
19. Chaudhuri, K., McGregor, A.: Finding metric structure in information theoretic clustering. In: Conference on Learning Theory (COLT), pp. 391–402 (2008)
20. de la Harpe, P.: On Hilbert's metric for simplices. In: Geometric Group Theory, vol. 1, pp. 97–118. Cambridge Univ. Press (1991)
21. Deza, M., Sikirić, M.D.: Voronoi polytopes for polyhedral norms on lattices. *Discrete Applied Mathematics* **197**, 42–52 (2015)
22. Dowty, J.G.: Chentsov's theorem for exponential families (2017). arXiv:1701.08895 [math.ST]
23. Endo, Y., Miyamoto, S.: Spherical  $k$ -means++ clustering. In: Modeling Decisions for Artificial Intelligence, pp. 103–114. Springer (2015)
24. Foertsch, T., Karlsson, A.: Hilbert metrics and Minkowski norms. *Journal of Geometry* **83**(1-2), 22–31 (2005)
25. Gonzalez, T.F.: Clustering to minimize the maximum intercluster distance. *Theoretical Computer Science* **38**, 293–306 (1985)
26. Hammersley, J.: On estimating restricted parameters. *Journal of the Royal Statistical Society. Series B (Methodological)* **12**(2), 192–240 (1950)
27. Hilbert, D.: Über die gerade linie als kürzeste verbindung zweier punkte. *Mathematische Annalen* **46**(1), 91–96 (1895)
28. Hotelling, H.: Spaces of statistical parameters. In: Bulletin AMS, vol. 36, p. 191 (1930)
29. Jenssen, R., Principe, J.C., Erdogmus, D., Eltoft, T.: The Cauchy-Schwarz divergence and Parzen windowing: Connections to graph theory and mercer kernels. *Journal of the Franklin Institute* **343**(6), 614–629 (2006)
30. Kass, R.E., Vos, P.W.: Geometrical Foundations of Asymptotic Inference. Wiley Series in Probability and Statistics. Wiley-Interscience (1997)

31. Khosravifard, M., Fooladivanda, D., Gulliver, T.A.: Conflicton of the convexity and metric properties in  $f$ -divergences. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences* **90**(9), 1848–1853 (2007)
32. Körner, M.C.: Minisum hyperspheres, *Springer Optimization and Its Applications*, vol. 51. Springer New York (2011)
33. Laidlaw, D.H., Weickert, J.: Visualization and Processing of Tensor Fields: Advances and Perspectives. Mathematics and Visualization. Springer-Verlag Berlin Heidelberg (2009)
34. Lebanon, G.: Learning Riemannian metrics. In: Conference on Uncertainty in Artificial Intelligence (UAI), pp. 362–369 (2002)
35. Lemmens, B., Nussbaum, R.: Birkhoff's version of Hilberts metric and its applications in analysis. *Handbook of Hilbert Geometry* pp. 275–303 (2014)
36. Lemmens, B., Walsh, C.: Isometries of polyhedral Hilbert geometries. *Journal of Topology and Analysis* **3**(02), 213–241 (2011)
37. Liang, X.: A note on divergences. *Neural Computation* **28**(10), 2045–2062 (2016)
38. Manthey, B., Röglin, H.: Worst-case and smoothed analysis of  $k$ -means clustering with Bregman divergences. *Journal of Computational Geometry* **4**(1), 94–132 (2013)
39. Nielsen, F.: Cramér-Rao lower bound and information geometry. In: *Connected at Infinity II*, pp. 18–37. Springer (2013)
40. Nielsen, F., Hadjeres, G.: Approximating covering and minimum enclosing balls in hyperbolic geometry. In: International Conference on Networked Geometric Science of Information, pp. 586–594. Springer (2015)
41. Nielsen, F., Muzellec, B., Nock, R.: Classification with mixtures of curved Mahalanobis metrics. In: IEEE International Conference on Image Processing (ICIP), pp. 241–245 (2016)
42. Nielsen, F., Muzellec, B., Nock, R.: Large margin nearest neighbor classification using curved Mahalanobis distances (2016). arXiv:1609.07082 [cs.LG]
43. Nielsen, F., Nock, R.: On approximating the smallest enclosing Bregman balls. In: Proceedings of the twenty-second annual symposium on Computational geometry, pp. 485–486. ACM (2006)
44. Nielsen, F., Nock, R.: On the smallest enclosing information disk. *Information Processing Letters* **105**(3), 93–97 (2008)
45. Nielsen, F., Nock, R.: Approximating smallest enclosing balls with applications to machine learning. *International Journal of Computational Geometry & Applications* **19**(05), 389–414 (2009)
46. Nielsen, F., Nock, R.: Total Jensen divergences: Definition, properties and  $k$ -means++ clustering (2013). arXiv:1309.7109 [cs.IT]
47. Nielsen, F., Nock, R., Amari, S.i.: On clustering histograms with  $k$ -means by using mixed  $\alpha$ -divergences. *Entropy* **16**(6), 3273–3301 (2014)
48. Nielsen, F., Shao, L.: On balls in a polygonal Hilbert geometry. In: 33st International Symposium on Computational Geometry (SoCG 2017). Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik (2017)
49. Nielsen, F., Sun, K., Marchand-Maillet, S.: On Hölder projective divergences. *Entropy* **19**(3) (2017)
50. Nock, R., Nielsen, F.: Fitting the smallest enclosing Bregman ball. In: ECML, pp. 649–656. Springer (2005)
51. Panigrahy, R.: Minimum enclosing polytope in high dimensions (2004). arXiv:cs/0407020 [cs.CG]
52. Papadopoulos, A., Troyanov, M.: From Funk to Hilbert geometry (2014). arXiv:1406.6983 [math.MG]
53. Rao, C.R.: Information and accuracy attainable in the estimation of statistical parameters. *Bull. Cal. Math. Soc.* **37**(3), 81–91 (1945)
54. Rao, C.R.: Information and the accuracy attainable in the estimation of statistical parameters. In: Breakthroughs in statistics, pp. 235–247. Springer (1992)
55. Reem, D.: The geometric stability of Voronoi diagrams in normed spaces which are not uniformly convex (2012). arXiv:1212.1094 [cs.CG]

56. Richter-Gebert, J.: Perspectives on projective geometry: A guided tour through real and complex geometry. Springer-Verlag Berlin Heidelberg (2011)
57. Saha, A., Vishwanathan, S., Zhang, X.: New approximation algorithms for minimum enclosing convex shapes. In: ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 1146–1160 (2011)
58. Sharir, M., Welzl, E.: A combinatorial bound for linear programming and related problems. STACS 92 pp. 567–579 (1992)
59. Shen, Z.: Riemann-Finsler geometry with applications to information geometry. Chinese Annals of Mathematics-Series B **27**(1), 73–94 (2006)
60. Shima, H.: The Geometry of Hessian Structures. World Scientific (2007)
61. Stigler, S.M., et al.: The epic story of maximum likelihood. Statistical Science **22**(4), 598–620 (2007)
62. Vernicos, C.: Introduction aux géométries de Hilbert. Séminaire de théorie spectrale et géométrie **23**, 145–168 (2004)
63. Welzl, E.: Smallest enclosing disks (balls and ellipsoids). New results and new trends in computer science pp. 359–370 (1991)

## 8 Isometry of Hilbert simplex geometry to a normed vector space

Consider the Hilbert simplex metric space  $(\Delta^d, \rho_{\text{HG}})$  where  $\Delta^d$  denotes the  $d$ -dimensional open probability simplex and  $\rho_{\text{HG}}$  the Hilbert cross-ratio metric. Let us recall the isometry ([20], 1991) of the open standard simplex to a normed vector space  $(V^d, \|\cdot\|_{\text{NH}})$ . Let  $V^d = \{v \in \mathbb{R}^{d+1} : \sum_i v^i = 0\}$  denote the  $d$ -dimensional vector space sitting in  $\mathbb{R}^{d+1}$ . Map a point  $p = (\lambda^0, \dots, \lambda^d) \in \Delta^d$  to a point  $v(x) = (v^0, \dots, v^d) \in V^d$  as follows:

$$v^i = \frac{1}{d+1} \left( d \log \lambda^i - \sum_{j \neq i} \log \lambda^j \right) = \log \lambda^i - \frac{1}{d+1} \sum_j \log \lambda^j.$$

We define the corresponding norm  $\|\cdot\|_{\text{NH}}$  in  $V^d$  by considering the shape of its unit ball  $B_V = \{v \in V^d : |v^i - v^j| \leq 1, \forall i \neq j\}$ . The unit ball  $B_V$  is a symmetric convex set containing the origin in its interior, and thus yields a *polytope norm*  $\|\cdot\|_{\text{NH}}$  (Hilbert norm) with  $2\binom{d+1}{2} = d(d+1)$  facets. Reciprocally, let us notice that a norm induces a unit ball centered at the origin that is convex and symmetric around the origin.

The distance in the normed vector space between  $v \in V^d$  and  $v' \in V^d$  is defined by:

$$\rho_V(v, v') = \|v - v'\|_{\text{NH}} = \inf \{ \tau : v' \in \tau(B_V \oplus \{v\}) \},$$

where  $A \oplus B = \{a + b : a \in A, b \in B\}$  is the Minkowski sum.

The reverse map from the normed space  $V^d$  to the probability simplex  $\Delta^d$  is given by:

$$\lambda^i = \frac{\exp(v^i)}{\sum_j \exp(v^j)}.$$

Thus we have  $(\Delta^d, \rho_{\text{HG}}) \cong (V^d, \|\cdot\|_{\text{NH}})$ . In 1D,  $(V^1, \|\cdot\|_{\text{NH}})$  is isometric to the Euclidean line.

Note that computing the distance in the normed vector space requires naively  $O(d^2)$  time.

Unfortunately, the norm  $\|\cdot\|_{\text{NH}}$  does not satisfy the parallelogram law.<sup>2</sup> Notice that a norm satisfying the parallelogram law can be associated with an inner product via the polarization identity. Thus the isometry of the Hilbert geometry to a normed vector space is not equipped with an inner product. However, all norms in a finite dimensional space are equivalent. This implies that in finite dimension,  $(\Delta^d, \rho_{\text{HG}})$  is *quasi-isometric* to the Euclidean space  $\mathbb{R}^d$ . An example of Hilbert geometry in infinite dimension is reported in [20]. Hilbert spaces are not CAT spaces except when  $\mathcal{C}$  is an ellipsoid [62].

---

<sup>2</sup> Consider  $A = (1/3, 1/3, 1/3)$ ,  $B = (1/6, 1/2, 1/3)$ ,  $C = (1/6, 2/3, 1/6)$  and  $D = (1/3, 1/2, 1/6)$ . Then  $2AB^2 + 2BC^2 = 4.34$  but  $AC^2 + BD^2 = 3.84362411135$ .

## 9 Hilbert geometry with Finslerian/Riemannian structures

In a Riemannian geometry, each tangent plane  $T_p M$  of a  $d$ -dimensional manifold  $M$  is equivalent to  $\mathbb{R}^d$ :  $T_p M \simeq \mathbb{R}^d$ . The inner product at each tangent plane  $T_p M$  can be visualized by an ellipsoid shape, a convex symmetric object centered at point  $p$ . In a *Finslerian geometry*, a norm  $\|\cdot\|_p$  is defined in each tangent plane  $T_p M$ , and this norm is visualized as a symmetric convex object with non-empty interior. Finslerian geometry thus generalizes Riemannian geometry by taking into account generic symmetric convex objects instead of ellipsoids for inducing norms at each tangent plane. Any Hilbert geometry induced by a compact convex domain  $\mathcal{C}$  can be expressed by an equivalent Finslerian geometry by defining the norm in  $T_p$  at  $p$  as follows [62]:

$$\|v\|_p = F_{\mathcal{C}}(p, v) = \frac{\|v\|}{2} \left( \frac{1}{pp^+} + \frac{1}{pp^-} \right),$$

where  $F_{\mathcal{C}}$  is the *Finsler metric*,  $\|\cdot\|$  is an *arbitrary norm* on  $\mathbb{R}^d$ , and  $p^+$  and  $p^-$  are the intersection points of the line passing through  $p$  with direction  $v$ :

$$p^+ = p + t^+ v, \quad p^- = p + t^- v.$$

A geodesic  $\gamma$  in a Finslerian geometry satisfies:

$$d_{\mathcal{C}}(\gamma(t_1), \gamma(t_2)) = \int_{t_1}^{t_2} F_{\mathcal{C}}(\gamma(t), \dot{\gamma}(t)) dt.$$

In  $T_p M$ , a ball of center  $c$  and radius  $r$  is defined by:

$$B(c, r) = \{v : F_{\mathcal{C}}(c, v) \leq r\}.$$

Thus any Hilbert geometry induces an equivalent Finslerian geometry, and since Finslerian geometries include Riemannian geometries, one may wonder which Hilbert geometries induce Riemannian structures? The only Riemannian geometries induced by Hilbert geometries are the *hyperbolic Cayley-Klein geometries* [56, 42, 41] with the domain  $\mathcal{C}$  being an ellipsoid. The Finslerian modeling of information geometry has been studied in [16, 59].

There is not a canonical way of defining measures in a Hilbert geometry since Hilbert geometries are Finslerian but not necessary Riemannian geometries [62]. The Busemann measure is defined according to the Lebesgue measure  $\lambda$  of  $\mathbb{R}^d$ : Let  $B_p$  denote the unit ball wrt. to the Finsler norm at point  $p \in \mathcal{C}$ , and  $B_e$  the Euclidean unit ball. Then the Busemann measure for a Borel set  $\mathcal{B}$  is defined by [62]:

$$\mu_{\mathcal{C}}(\mathcal{B}) = \int_{\mathcal{B}} \frac{\lambda(B_e)}{\lambda(B_p)} d\lambda(p).$$

The existence and uniqueness of center points of a probability measure in Finsler geometry have been investigated in [4].

## 10 Bounding Hilbert norm with other norms

Let us show that  $\|v\|_{\text{NH}} \leq \beta_{d,c} \|v\|_c$ , where  $\|\cdot\|_c$  is any norm. Let  $v = \sum_{i=0}^d e_i x_i$ , where  $\{e_i\}$  is a basis of  $\mathbb{R}^{d+1}$ . We have:

$$\|v\|_c \leq \sum_{i=0}^d |x_i| \|e_i\|_c \leq \|x\|_2 \underbrace{\sqrt{\sum_{i=0}^d \|e_i\|_c^2}}_{\beta_d},$$

where the first inequality comes from the triangle inequality, and the second inequality is from the Cauchy-Schwarz inequality. Thus we have:

$$\|v\|_{\text{NH}} \leq \beta_d \|x\|_2,$$

with  $\beta_d = \sqrt{d+1}$  since  $\|e_i\|_{\text{NH}} \leq 1$ .

Let  $\alpha_{d,c} = \min_{\{v : \|v\|_c=1\}} \|v\|_{\text{NH}}$ . Consider  $u = \frac{v}{\|v\|_c}$ . Then  $\|u\|_c = 1$  so that  $\|v\|_{\text{NH}} \geq \alpha_{d,c} \|v\|_c$ . To find  $\alpha_d$ , we consider the unit  $\ell_2$  ball in  $V^d$ , and find the smallest  $\lambda > 0$  so that  $\lambda B_V$  fully contains the Euclidean ball.

Therefore, we have overall:

$$\alpha_d \|x\|_2 \leq \|v\|_{\text{NH}} \leq \sqrt{d+1} \|x\|_2$$

In general, note that we may consider two arbitrary norms  $\|\cdot\|_l$  and  $\|\cdot\|_u$  so that:

$$\alpha_{d,l} \|x\|_l \leq \|v\|_{\text{NH}} \leq \beta_{d,u} \|x\|_u.$$

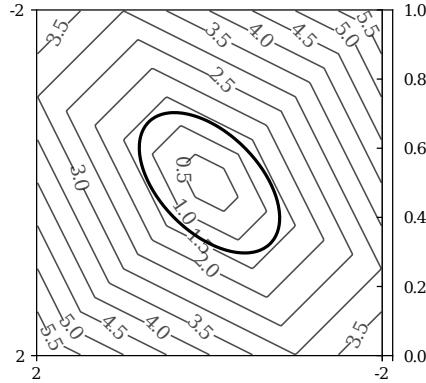


Fig. 12: Polytope balls  $B_V$  and the Euclidean unit ball  $B_E$ . From the figure the smallest polytope ball has radius  $\approx 1.5$ .

## 11 Funk directed metrics and Funk balls

The Funk metric [52] wrt a convex domain  $\mathcal{C}$  is defined by

$$F_{\mathcal{C}}(x, y) = \log \left( \frac{\|x - a\|}{\|y - a\|} \right),$$

where  $a$  is the intersection of the domain boundary and the affine ray  $R(x, y)$  starting from  $x$  and passing through  $y$ . Correspondingly, the reverse Funk metric is

$$F_{\mathcal{C}}(y, x) = \log \left( \frac{\|y - b\|}{\|x - b\|} \right),$$

where  $b$  is the intersection of  $R(y, x)$  with the boundary. The Funk metric is *not* a metric distance.

The Hilbert metric is simply the arithmetic symmetrization:

$$H_{\mathcal{C}}(x, y) = \frac{F_{\mathcal{C}}(x, y) + F_{\mathcal{C}}(y, x)}{2}.$$

It is interesting to explore clustering based on the Funk geometry, which we leave as a future work.