

SPARSE HYPERBOLIC REPRESENTATION LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Reducing the space complexity of representations while minimizing the loss of information makes data science procedures computationally efficient. For the entities with the tree structure, hyperbolic-space-based representation learning (HSBRL) has successfully reduced the space complexity of representations by using low-dimensional space. Nevertheless, it has not minimized the space complexity of each representation since it has used the same dimension for all representations and has not selected the best dimension for each representation. Hence, this paper aims to minimize representations' space complexity for HSBRL. For minimizing each representation's space complexity, sparse learning has been effective in the context of linear-space-based machine learning; however, no sparse learning has been proposed for HSBRL. It is non-trivial to propose sparse learning for HSBRL because (i) sparse learning methods designed for linear space cause non-uniform sparseness in hyperbolic space, and (ii) existing Riemannian gradient descent methods fail to obtain sparse representations owing to an oscillation problem. This paper, for the first time, establishes a sparse learning scheme for hyperbolic space, overcoming the above issues with our novel sparse regularization term and optimization algorithm. Our regularization term achieves uniform sparseness since it is defined based on geometric distance from subspaces inducing sparsity. Our optimization algorithm successfully obtains sparse representations, avoiding the oscillation problem by realizing the shrinkage-thresholding idea in a general Riemannian manifold. Numerical experiments demonstrate that our scheme can obtain sparse representations with smaller information loss than traditional methods, successfully avoiding the oscillation problem.

1 INTRODUCTION

Data science applies a composition of mathematical operations, called an algorithm, to solve real-life issues. For mathematical operations to handle real-world entities, we need their mathematical representations. Representation learning (RL) aims to obtain those entities' mathematical representations that reflect their semantic meaning. As an interface between the real world and data science, RL has been applied to various areas, e.g., machine translation and sentiment analysis for natural language (Mikolov et al., 2013; Pennington et al., 2014; Bojanowski et al., 2017; Ganea et al., 2018a), and community detection and link prediction for social network data (Hoff et al., 2002; Perozzi et al., 2014; Tang et al., 2015b;a; Grover & Leskovec, 2016), pathway prediction of biochemical network (Dale et al., 2010; MA Basher & Hallam, 2021), and link prediction and triplet classification for knowledge base (Nickel et al., 2011; Bordes et al., 2013; Riedel et al., 2013; Nickel et al., 2016; Trouillon et al., 2016; Ebisu & Ichise, 2018). The fact that data science handles real entities throughout mathematical representations implies that reducing the space complexity of representations will benefit the whole of data science. Recent studies have shown that hyperbolic-space-based RL (HSBRL) can obtain effective low-dimensional representations, typical low space complexity representations, if the entities have the tree structure (Nickel & Kiela, 2017; Ganea et al., 2018b; Nickel & Kiela, 2018; Tay et al., 2018; Sala et al., 2018; Gu et al., 2019; Suzuki et al., 2019; Tifrea et al., 2019; Yu & De Sa, 2019; Law et al., 2019; Mathieu et al., 2019; Balazevic et al., 2019; Tabaghi & Dokmanic, 2020; Chami et al., 2020; Sonthalia & Gilbert, 2020; Kyriakis et al., 2021; Suzuki et al., 2021a;b). Still, existing HSBRL methods have not minimized the space complexity of each representation since they have used the same dimension for all representations and have not selected the best space complexity for each entity. Specifically, it would be more efficient to allocate high space complexity only for representations in non-tree-like parts of the data, if the data structure has tree-like parts and non-tree-like

parts. For minimizing each representation’s space complexity, sparse learning has been effective for linear-space-based machine learning. Sparse learning aims to obtain representations filled with the most zero elements to the extent that it does not lose the essential information. Typically, we have implemented sparse learning by the 1-norm regularization, which optimizes the sum of the original loss function of the problem and the 1-norm regularization term. The 1-norm regularization has mainly been developed in many application areas, such as time-frequency data analysis (Donoho & Johnstone, 1995; Chen et al., 2001), natural image processing (Olshausen & Field, 1996), and visual tracking (Liu et al., 2010). In the context of the linear model, the 1-norm regularization’s statistical property has also been well studied (Tibshirani, 1996; Chen et al., 2001; Friedman et al., 2010) as well as geometric property (Donoho & Tanner, 2009a;b; 2010). Also, in statistics and machine learning communities, it has been intensively studied from theoretical perspectives (Ng, 2004; Shalev-Shwartz & Tewari, 2009; Tomioka et al., 2011; Zhang et al., 2016b) and applications for precision matrix learning (Friedman et al., 2008), neural-network-based learning (Liu et al., 2017b; Alizadeh et al., 2020), Ising model learning (Kuang et al., 2017), and transfer learning (Takada & Fujisawa, 2020). However, the above work is all for linear space; sparse learning has not been proposed for HSBRL. This paper aims to establish a sparse learning scheme for HSBRL. We have the following three challenges in establishing the scheme.

Firstly, defining sparsity in hyperbolic space is non-trivial. In linear space, we can define a point’s sparsity as the number of zero elements in the coordinate representing the point. This definition’s advantage is that it directly indicates representations’ space complexity since a computer uses a coordinate system to represent points. Hence, we want to define the sparsity of a point in hyperbolic space so that it inherits this property. The issue is that the straightforward application of the definition in hyperbolic space is ill-defined since we cannot determine a unique coordinate system for a general manifold. To avoid this issue, we must have a geometric definition, i.e., one independent of the coordinate system, and we expect it to be interpretable as computational complexity at the same time.

Secondly, we need to design a continuous regularization term to induce the sparsity of representations. If the objective function includes the 0-norm, the number of the non-zero elements, it involves combinatorial optimization since the 0-norm is not a continuous function. Hence, direct use of the 0-norm does not scale to big data. For this reason, we relax the 0-norm to the 1-norm in optimizing a function in linear space. However, defining a hyperbolic counterpart of the 1-norm is, again, non-trivial. For example, applying the 1-norm in a coordinate system of hyperbolic space naïvely has two drawbacks below. First, we will get different results depending on which coordinate system we use. Second, the regularization’s strength cannot be uniform because a coordinate does not always reflect the distance structure in hyperbolic space.

Thirdly, we need an effective optimization algorithm to obtain sparse representations. We can apply the Riemannian gradient descent (RGD) algorithm in theory to the objective function with our 1-norm regularization term since it is a continuous function on the Riemannian manifold. However, the RGD fails to get sparse representations since it oscillates around the true sparse solution.

This paper proposes a novel sparse learning scheme for HSBRL to solve these issues. It is regularization-based, so applicable to any HSBRL using a continuous loss function. The **core contributions** of our scheme are the following three, each of which addresses one of the above issues:

1. The ***Hyperbolic sparsity***, a novel geometric definition of a point’s sparsity. We can also interpret the value as the space complexity of a point in hyperbolic space.
2. The ***Hyperbolic 1-norm***, a novel sparse regularization term for continuous optimization. It achieves a uniform sparseness strength since it measures the geometric distance from the point to the subspaces with higher hyperbolic sparsity.
3. The ***Hyperbolic iterative shrinkage-thresholding algorithm (HISTA)***, a novel optimization algorithm free from the oscillating issue, realizing the shrinkage-thresholding idea in a Riemannian manifold (Bruck Jr, 1977; Chambolle et al., 1998; Figueiredo & Nowak, 2003; Daubechies et al., 2004; Vonesch & Unser, 2007; Hale et al., 2007; Wright et al., 2009; Beck & Teboulle, 2009). Despite its oscillating-free property, HISTA’s time/space computational complexity with respect to the space dimension is linear, the same as the RGD.

The above definitions are all geometric, i.e., free from the coordinate system selection problem. Also, we give closed form formulae for them in hyperbolic space. We can calculate them in linear time complexity with respect to the space dimension.

Our experiments demonstrate that our scheme can obtain sparse representations without the oscillation problem. Also, our comparisons show that our scheme can obtain representations with smaller loss of information than existing low-dimensional representations with the same space complexity. Our experiments also show that our optimization algorithm HISTA works better than the RGD if the true solution is sparse.

The rest of the paper is organized as follows. Section 2 geometrically defines the sparsity on a Riemannian manifold, including hyperbolic space. Section 3 proposes our sparse learning scheme on a Riemannian manifold. Section 4 specializes the scheme for hyperbolic space and gives specific formulae. Section 5 evaluates our scheme and Section 6 concludes this paper.

1.1 RELATED WORK

This paper proposes multiple operations in hyperbolic space. There has been much work on defining operations for hyperbolic space in the mathematical context e.g., (Gromov, 1987; Ungar, 1994; 1996; Sabinin et al., 1998), and machine learning context e.g., the work cited above and general Riemannian optimization work (Absil et al., 2009; Qi et al., 2010; Zhang & Sra, 2016; Zhang et al., 2016a; Liu et al., 2017a; Becigneul & Ganeva, 2019; Kasai et al., 2019; Zhou et al., 2019) and neural hyperbolic network papers (Ganeva et al., 2018c; Chami et al., 2019; Liu et al., 2019; Tan et al., 2021; Shimizu et al., 2021; Takeuchi et al., 2022). Still, none of them have proposed the operations to get sparse representations. The work closest to our motivation is the Riemannian proximal gradient (RPG) descent (Huang & Wei, 2022) since the iterative shrinkage-thresholding algorithm (ISTA) can be regarded as a special case of the proximal gradient descent method. Despite its solid theoretical background, however, the closed form solution of the RPG operation for the sparse regularization has not been known except for the linear space case, where the RPG is reduced to the ISTA. Hence, the RPG cannot obtain sparse representations. Conversely, our algorithm can give sparse representations and computationally cheap since the closed form algorithm is given for hyperbolic space. Still, both of the RPG and our algorithm can be regarded as generalizations of the ISTA since they are reduced to the ISTA for optimizing the 1-norm regularization problem in linear space.

1.2 NOTATION

Linear space We denote the set of real values, nonnegative real values, positive real values, and nonnegative integers by \mathbb{R} , $\mathbb{R}_{\geq 0}$, $\mathbb{R}_{> 0}$, and $\mathbb{Z}_{\geq 0}$, respectively. For $D \in \mathbb{Z}_{\geq 0}$, we denote the D -dimensional real coordinate space (RCS) by \mathbb{R}^D . Also, for $M, N \in \mathbb{Z}_{\geq 0}$ we denote the set of $M \times N$ -real matrices by $\mathbb{R}^{M \times N}$. In this paper, an element in RCS is denoted by a bold lower letter, e.g., \mathbf{p}, \mathbf{q} , and a real matrix is denoted by a bold upper letter, e.g., \mathbf{G} . We denote the transpose of $\mathbf{p} \in \mathbb{R}^D$ by \mathbf{p}^\top and define $|\mathbf{p}| := \sqrt{\mathbf{p}^\top \mathbf{p}}$. We denote the D -dimensional zero vector and one vector by $\mathbf{0}_D$ and $\mathbf{1}_D$, respectively, and $D \times D$ identity matrix by \mathbf{I}_D . In this paper, a mathematical constant is denoted by an upright letter, e.g., $\mathbf{0}_D$ and \mathbf{I}_D , whereas a variable is denoted by an italic letter, e.g., $D, \mathbf{p}, \mathbf{q}$. We call the inner product space $(\mathbb{R}^D, \langle \cdot, \cdot \rangle)$ the D -dimensional Euclidean vector space (EVS), where $\langle \cdot, \cdot \rangle : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ by $\langle \mathbf{p}, \mathbf{q} \rangle = \mathbf{p}^\top \mathbf{q}$. For a scalar-valued function f of a scalar, typically a hyperbolic function in this paper, we extend it to a vector-valued function of a vector, defined by $f\left(\begin{bmatrix} p_1 & p_2 & \cdots & p_D \end{bmatrix}^\top\right) := [f(p_1) \ f(p_2) \ \cdots \ f(p_D)]^\top$.

Riemannian manifold Let \mathcal{M} and \mathcal{M}' be C^∞ -Riemannian manifolds. We denote by $C^\infty(\mathcal{M})$ the set of real C^∞ -functions on \mathcal{M} , by $T_p\mathcal{M}$ the tangent space at point p , and by $\langle \cdot, \cdot \rangle_p^\mathcal{M} : T_p\mathcal{M} \times T_p\mathcal{M} \rightarrow \mathbb{R}$ the metric tensor at $T_p\mathcal{M}$. For C^∞ -map $\varphi : \mathcal{M} \rightarrow \mathcal{M}'$, the differential $d\varphi_p : T_p\mathcal{M} \rightarrow T_{\varphi(p)}\mathcal{M}'$ of φ at p is defined by $(d\varphi_p(v))(f) = v(f \circ \varphi)$. We call a diffeomorphism $\varphi : \mathcal{M} \rightarrow \mathcal{M}'$ an *isometry* if $\langle u, v \rangle_p^\mathcal{M} = \langle d\varphi_p(u), d\varphi_p(v) \rangle_{\varphi(p)}^{\mathcal{M}'}$ for any point p and tangent vectors $u, v \in T_p\mathcal{M}$. We denote by $\Delta_\mathcal{M}(p, q)$ the Riemannian distance between $p \in \mathcal{M}$ and $q \in \mathcal{M}$, which is defined as the length of the shortest piecewise smooth curve connecting p and q . We denote by $\exp_p : T_p\mathcal{M} \rightarrow \mathcal{M}$ the *exponential map* at point $p \in \mathcal{M}$. Here, we define the exponential map based on the Levi-Civita connection. Intuitively, $\exp_p(v)$ indicates the point reached in one unit time by the motion with a zero acceleration (called a *geodesic*) starting at the point p with the initial velocity v . If the exponential map \exp_p is bijective, then we can define its inverse map, called *logarithmic map* $\log_p : \mathcal{M} \rightarrow T_p\mathcal{M}$. For definitions of these terms of the Riemannian geometry, we refer readers to (Lee, 2018).

Coordinate system This paper mainly discusses manifolds diffeomorphic to \mathbb{R}^D for some $D \in \mathbb{Z}_{\geq 0}$. We can determine such a C^∞ -Riemannian manifold by a pair $(\mathbb{U}, \mathbf{G}_\cdot)$ of an open set $\mathbb{U} \subset \mathbb{R}^D$ and the coordinate representation matrix $\mathbf{G}_p \in \mathbb{R}^{D \times D}$ of the Riemannian metric at point p for all $p \in \mathbb{U}$. We denote by $(\mathbb{U}, \mathbf{G}_\cdot)$ the Riemannian manifold determined this way. For the open set $\mathbb{U} \subset \mathbb{R}^D$, the coordinate system (x^1, x^2, \dots, x^D) for \mathbb{U} , and $d = 1, 2, \dots, D$, we define the partial derivative operator $\partial_d|_p : C^\infty(\mathbb{U}) \rightarrow \mathbb{R}$ by $\partial_d|_p f := \frac{\partial}{\partial x^d} f(\mathbf{p})$, where $C^\infty(\mathbb{U})$ is the set of real C^∞ -functions on \mathbb{U} . For $\mathbf{v} := [v^1 \ v^2 \ \dots \ v^D]^\top \in \mathbb{R}^D$, we denote $v^1 \partial_1|_p + v^2 \partial_2|_p + \dots + v^D \partial_D|_p$ by $\mathbf{v}^\top \boldsymbol{\partial}|_p$. Here, we have that $T_p \mathcal{M} = \{\mathbf{v}^\top \boldsymbol{\partial}|_p \mid \mathbf{v} \in \mathbb{R}^D\}$. Using the coordinate representation \mathbf{G}_p at $p \in \mathbb{U}$ of the Riemannian metric $\langle \cdot, \cdot \rangle_{(\mathbb{U}, \mathbf{G}_\cdot)}$, we can calculate the inner product $\langle \mathbf{u}^\top \boldsymbol{\partial}|_p, \mathbf{v}^\top \boldsymbol{\partial}|_p \rangle_{(\mathbb{U}, \mathbf{G}_\cdot)}$ of two tangent vectors $\mathbf{u}^\top \boldsymbol{\partial}|_p, \mathbf{v}^\top \boldsymbol{\partial}|_p \in T_p(\mathbb{U}, \mathbf{G}_\cdot)$ by $\langle \mathbf{u}^\top \boldsymbol{\partial}|_p, \mathbf{v}^\top \boldsymbol{\partial}|_p \rangle_{(\mathbb{U}, \mathbf{G}_\cdot)} = \mathbf{u}^\top \mathbf{G}_p \mathbf{v}$.

The following examples are the Riemannian manifolds we mainly discuss in the paper.

Example 1. (a) The Riemannian manifold $(\mathbb{R}^D, \mathbf{I}_D)_\cdot$ is the D -dimensional Euclidean space, where $\mathbf{I}_{D,p} = \mathbf{I}_D$ for all $p \in \mathbb{R}^D$.

(b) Let \mathbb{D}^D be the D -dimensional open ball $\mathbb{D}^D := \{p \in \mathbb{R}^D \mid p^\top p < 1\}$. For $p \in \mathbb{D}^D$, we define $\mathbf{G}_p^p, \mathbf{G}_p^K \in \mathbb{R}^{D \times D}$ by $\mathbf{G}_p^p = \frac{1-p^\top p}{4} \mathbf{I}_D$ and $\mathbf{G}_p^K = (1 - p^\top p)(\mathbf{I}_D - pp^\top)$. The Riemannian manifolds $(\mathbb{D}^D, \mathbf{G}_\cdot^p)$ and $(\mathbb{D}^D, \mathbf{G}_\cdot^K)$ are isometric to each other. Specifically, there is an isometry $\varphi_{p,K} : \mathbb{D}^D \rightarrow \mathbb{D}^D$ defined by $\varphi_{p,K}(p) := \tanh(\operatorname{atanh}(p)) \frac{p}{|p|}$. The Riemannian manifold determined by these two is called the *D -dimensional hyperbolic space*. When we distinguish these two coordinate systems of the D -dimensional hyperbolic space, we call $(\mathbb{D}^D, \mathbf{G}_\cdot^p)$ and $(\mathbb{D}^D, \mathbf{G}_\cdot^K)$ the *Poincaré model* and *Klein model*, respectively.

2 THE RIEMANNIAN SPARSITY

In this section, we define the sparsity of a point in Riemannian manifold geometrically. Our sparsity also indicates the space computational complexity to represent a point in hyperbolic space. In the definition of the sparsity of a point in RCS, the origin and canonical bases play an essential role as a “reference point” and “reference direction.” Hence, we need to begin with the definition of the origin and the bases to discuss the sparsity of a point in a manifold.

Definition 1 (Riemannian manifold with an origin and orthonormal bases (RMOO)). Let \mathcal{M} be a D -dimensional Riemannian manifold. For a point $o \in \mathcal{M}$ and orthonormal bases (ONBs) $e_1, e_2, \dots, e_D \in T_o \mathcal{M}$, the triple $(\mathcal{M}, o, (e_1, e_2, \dots, e_D))$ is called a *D -dimensional Riemannian manifold with an origin and orthonormal bases (RMOO)*. Here, the point o and set (e_1, e_2, \dots, e_D) are called the origin and ONBs of the RMOO, respectively.

Definition 2 (Isometric RMOOs). Let $(\mathcal{M}, o, (e_1, e_2, \dots, e_D))$ and $(\mathcal{M}', o', (e'_1, e'_2, \dots, e'_D))$ be RMOOs and $\varphi : \mathcal{M} \rightarrow \mathcal{M}'$ is an isometry. If $\varphi(o) = o'$ and $d\varphi_o(e_d) = e'_d$ for $d = 1, 2, \dots, D$, then φ is called an *isometry* from $(\mathcal{M}, o, (e_1, e_2, \dots, e_D))$ and $(\mathcal{M}', o', (e'_1, e'_2, \dots, e'_D))$. If there exists an isometry between two RMOOs, we say that the RMOOs are *isometric* to each other.

Example 2. (a) The triple $((\mathbb{R}^D, \mathbf{I}_D)_\cdot, \mathbf{0}, (\partial_1|_0, \partial_2|_0, \dots, \partial_D|_0))$ is an RMOO. We call it the *D -dimensional Euclidean vector RMOO (EVRMOO)*. Here, we can identify $\mathbf{v}^\top \boldsymbol{\partial}|_0 \in T_0 \mathbb{R}^D$ with $\mathbf{v} \in \mathbb{R}^D$ since $\exp_0(\mathbf{v}^\top \boldsymbol{\partial}|_0) = \mathbf{v}$ holds. Under these identifications, we can say that the canonical bases e_1, e_2, \dots, e_D of the D -dimensional RCS are the ONBs of the EVRMOO. Although the D -dimensional EVRMOO is no more than Riemannian reformulation of the D -dimensional RCS, it helps us to see the correspondence between the D -dimensional RCS and hyperbolic space.

(b) $((\mathbb{D}^D, \mathbf{G}_\cdot^p), \mathbf{0}, (\frac{1}{2}\partial_1|_0, \frac{1}{2}\partial_2|_0, \dots, \frac{1}{2}\partial_D|_0))$ and $((\mathbb{D}^D, \mathbf{G}_\cdot^K), \mathbf{0}, (\partial_1|_0, \partial_2|_0, \dots, \partial_D|_0))$ are RMOOs isometric to each other. We call it the *D -dimensional hyperbolic RMOO*. Here, the map $\varphi_{p,K}$ is an isometry.

Now, we define the sparsity of a point with respect to a basis before considering all the bases. In EVS, a vector has sparsity with respect to a basis if the point is orthogonal to the basis. We can generalize this definition to RMOOs as follows.

Definition 3 (Sparse hyperplane (SHP)). Let $(\mathcal{M}, o, (e_1, e_2, \dots, e_D))$ be a D -dimensional RMOO, and assume that $\exp_o : T_o \mathcal{M} \rightarrow \mathcal{M}$ is bijective. The d -th *sparse hyperplane* (SHP), denoted by $\Pi_d^{\mathcal{M}}$, is defined by $\Pi_d^{\mathcal{M}} := \{p \in \mathcal{M} \mid \langle e_d, \log_o(p) \rangle_o = 0\}$.

Once we have defined the sparseness of a point with respect to a basis, we can define that with respect to the all bases as follows.

Definition 4 (Riemannian sparseness and 0 norm). Let $(\mathcal{M}, \mathbf{o}, (e_1, e_2, \dots, e_D))$ be a D -dimensional RMOO, and assume that $\exp_{\mathbf{o}} : T_{\mathbf{o}}\mathcal{M} \rightarrow \mathcal{M}$ is bijective. We define the **Riemannian sparseness** $\text{sp}^{\mathcal{M}}(p)$ of a point $p \in \mathcal{M}$ as the number of SHPs that includes p , i.e., $\text{sp}^{\mathcal{M}}(p) := |\{d = 1, 2, \dots, D \mid p \in \Pi_d\}|$. Also, we define the **Riemannian 0-norm** $\|p\|_0^{\mathcal{M}}$ of a point $p \in \mathcal{M}$ by $\|p\|_0^{\mathcal{M}} = D - \text{sp}^{\mathcal{M}}(p)$.

Remark 1. (i) The Riemannian 0-norm does not depend on the coordinate system. Specifically, for two isometric RMOOs $(\mathcal{M}, \mathbf{o}, (e_1, e_2, \dots, e_D))$ and $(\mathcal{M}', \mathbf{o}', (e'_1, e'_2, \dots, e'_D))$ and an isometry $\varphi : \mathcal{M} \rightarrow \mathcal{M}'$, and we have that $\|p\|_0^{\mathcal{M}} = \|\varphi(p)\|_0^{\mathcal{M}'}$ for any point $p \in \mathcal{M}$.

(ii) Since the SHPs, Riemannian sparseness, and Riemannian 0-norm depend on the origin \mathbf{o} and ONBs (e_1, e_2, \dots, e_D) , we could clarify these by, e.g., $\Pi_d^{((\mathcal{M}, \langle \cdot, \cdot \rangle, \cdot), \mathbf{o}, (e_1, e_2, \dots, e_D))}$. Nevertheless, we omit these from notation for simplicity since they are clear for each Riemannian manifold and coordinate system in this paper.

Example 3. For the D -dimensional EVRMOO $((\mathbb{R}^D, \mathbf{I}_D, \cdot), \mathbf{0}, (\partial_1|_{\mathbf{0}}, \partial_2|_{\mathbf{0}}, \dots, \partial_D|_{\mathbf{0}}))$, the d -th SHP is given by $\Pi_d^{\mathcal{M}} = \{p \in \mathbb{R}^D \mid p^\top e_D = 0\} = \{p \in \mathbb{R}^D \mid [p]_d = 0\}$, where $[p]_d$ indicates the d -th element of p . Hence, we have $\|p\|_0^{(\mathbb{R}^D, \mathbf{I}_D, \cdot)} = \|p\|_0$. The above results directly follow the fact that $\langle \log_{\mathbf{0}}(p), \log_{\mathbf{0}}(q) \rangle_{\mathbf{0}}^{(\mathbb{R}^D, \mathbf{I}_D, \cdot)} = p^\top q$.

Example 4. In D -dimensional hyperbolic space, the specific form of the 0-norm is given by $\|p\|_0^{(\mathbb{D}^D, G^p)} = \|p\|_0$, and $\|p\|_0^{(\mathbb{D}^D, G^k)} = \|p\|_0$.

We call $\|p\|_0^{(\mathbb{D}^D, G^p)}$ and $\|p\|_0^{(\mathbb{D}^D, G^k)}$ the **hyperbolic 0-norm** of $p \in \mathbb{D}^D$ in the Poincaré model and in the Klein model, respectively.

Remark 2. (i) Example 4 implies that the Riemannian 0-norm, which is defined geometrically in Definition 4, also indicates the number of non-zero elements to represent the point in each coordinate system on a computer. This fact justifies our Riemannian 0-norm definition for HSBRL in both geometric and computational senses.

(ii) Results similar to Example 4 hold for the proper velocity model, hyperboloid model, and hemisphere model. Note that it does NOT hold for the upper half space model. In the first place, the upper plane model does not use the coordinate whose 0-norm is larger than two; hence the model is not suitable for discussing the sparsity of a point.

3 SPARSE LEARNING SCHEME ON A RIEMANNIAN MANIFOLD

We have defined the Riemannian sparseness and 0-norm. We are ready to discuss the scheme to obtain sparse representations in the sense of the Riemannian sparseness. This section discusses the general RMOO case. We will discuss the hyperbolic RMOO case in the next section as a special case of this section's discussion.

3.1 THE RIEMANNIAN 1-NORM

The 0-norm in RCS is not continuous as a function of a point. Hence, we cannot optimize the 0-norm regularized function by a gradient method. To solve this issue, we have often relaxed the 0-norm to the 1-norm, which is a continuous function. Likewise, we define the Riemannian 1-norm of a point on an RMOO as a relaxation of the Riemannian 0-norm so that we can use it as a regularization term for gradient methods. To define an appropriate counterpart of the 1-norm on RMOO, we discuss D -dimensional RCS again. The 1-norm $\|p\|_1$ of a point $p = [p_1 \ p_2 \ \dots \ p_D]^\top \in \mathbb{R}^D$ is given by the sum of the absolute element values $|p_1|, |p_2|, \dots, |p_D|$. The absolute value $|p_d|$ of each element of a vector is natural as a relaxed regularization term to induce sparsity with respect to the d -th axis since the element indicates the distance $\Delta_{\mathcal{M}}(p, \Pi_d) := \inf_{p' \in \Pi_d} \Delta_{\mathcal{M}}(p, p')$ between the point and the d -th SHP. The above observation inspires our definition below of the 1-norm in an RMOO. We begin with defining the (signed) distance between a point and an SHP, since it plays an essential role to define the 1-norm as we have seen in RCS.

Definition 5 (Signed SHP distance map (SSDM)). Let $(\mathcal{M}, o, (e_1, e_2, \dots, e_D))$ be a D -dimensional RMOO, and assume that $\exp_o : T_o\mathcal{M} \rightarrow \mathcal{M}$ is bijective. We define the **signed distance** $\delta_d^{\mathcal{M}}(p) \in \mathbb{R}$ of a point $p \in \mathcal{M}$ from the d -th SHP by $\delta_d^{\mathcal{M}}(p) := \text{sgn}(\langle e_d, \log_o(p) \rangle_o) \Delta_{\mathcal{M}}(p, \Pi_{\hat{d}})$, where $\Delta_{\mathcal{M}}(p, \Pi_{\hat{d}}) := \inf_{p' \in \Pi_{\hat{d}}} \Delta_{\mathcal{M}}(p, p')$. We also define the **signed SHP distance map (SSDM)** $\delta^{\mathcal{M}} : \mathcal{M} \rightarrow \mathbb{R}^D$ by $\delta^{\mathcal{M}}(p) := [\delta_1^{\mathcal{M}}(p) \ \delta_2^{\mathcal{M}}(p) \ \dots \ \delta_D^{\mathcal{M}}(p)]^\top$.

Remark 3. We consider the signed version in Definition 5 since we often have a bijective SSDM as a result of considering the signed version. We can specify a point by the signed distances if the SSDM is bijective. Hence, a bijective SSDM plays an essential role to develop our optimization algorithm controlling the signed distances, as we see later.

Appendix A has the visualization of the SSDM. Based on the SSDM, we define the Riemannian 1-norm, which is this subsection's objective.

Definition 6 (Riemannian 1-norm). Let $(\mathcal{M}, o, (e_1, e_2, \dots, e_D))$ be a D -dimensional RMOO, and assume that $\exp_o : T_o\mathcal{M} \rightarrow \mathcal{M}$ is bijective. We define the **Riemannian 1-norm** $\|p\|_1^{\mathcal{M}}$ of a point $p \in \mathcal{M}$ by $\|p\|_1^{\mathcal{M}} := \|\delta^{\mathcal{M}}(p)\|_1 = \sum_{d=1}^D |\delta_d^{\mathcal{M}}(p)|$. We call it the **Riemannian 1-norm regularization** to add the Riemannian 1-norms of the points multiplied by some factor in optimizing a loss function of points in an RMOO.

Example 5. Consider the D -dimensional EVRMOO $((\mathbb{R}^D, \mathbf{G}^{\mathbb{R}^D}), \mathbf{0}, (\partial_1|_0, \partial_2|_0, \dots, \partial_D|_0))$. Then, the SSDM is given by $\delta^{(\mathbb{R}^D, \mathbf{G}^{\mathbb{R}^D})}(p) = p$. The 1-norm is given by $\|p\|_1 = \sum_{d=1}^D |p_d|$, which is equivalent to the 1-norm as a real numeric vector.

3.2 THE RIEMANNIAN ITERATIVE SHRINKAGE-THRESHOLDING ALGORITHM

To optimize a function differentiable almost everywhere, one might use (sub)gradient methods. However, the direct application of gradient methods in the 1-norm regularization causes residuals to the sparse solution. See the following simplest example. Details are also in Appendix B

Example 6. Suppose that we optimize $f(p) = |p|$ by the gradient descent method with learning rate $\alpha > 0$ and initial point $p^{(0)} \neq 0$. Then the algorithm ends up oscillating between $p^{(0)} - \alpha n$ and $p^{(0)} - \alpha(n+1)$ and cannot achieve the true sparse solution $p^{(0)} = 0$, unless $p^{(0)}$ is an integral multiple of α . Here, $n = \left\lfloor \frac{p^{(0)}}{\alpha} \right\rfloor$ is the maximum integer that is no greater than $\frac{p^{(0)}}{\alpha}$.

Example 6 is also an example of an oscillation in an RMOO optimization problem since it is equivalent to optimization of the 1-norm function of 1-dimensional Euclidean space or hyperbolic space. Note that Euclidean space and hyperbolic space are isometric to each other for 1-dimensional case.

To address the above issue to obtain the true sparse solution, the **iterative shrinkage-thresholding algorithm (ISTA)** (e.g., Bruck Jr, 1977) has been used. The ISTA is designed to optimize a function in RCS in the following form:

$$J(p) := L(p) + \lambda \|p\|_1, \quad (1)$$

where $L : \mathbb{R}^D \rightarrow \mathbb{R}$ is an almost everywhere differentiable function and $\lambda \in \mathbb{R}_{\geq 0}$ is the regularization weight. The ISTA iterates the following two steps: (i) $q \leftarrow p - \alpha \frac{\partial}{\partial p} L(p)$, (ii) $p \leftarrow \mathcal{T}_{\alpha\lambda}(q)$. Here, $\mathcal{T}_\beta : \mathbb{R}^D \rightarrow \mathbb{R}^D$ is the **soft-thresholding operator (STO)** defined by

$$\mathcal{T}_\beta \left([p_1 \ p_2 \ \dots \ p_D]^\top \right) := [\tau_\beta(p_1) \ \tau_\beta(p_2) \ \dots \ \tau_\beta(p_D)]^\top, \quad (2)$$

where $\tau_\beta : \mathbb{R}^D \rightarrow \mathbb{R}^D$ is defined by $\tau_\beta(p) := \text{sgn}(p) \max\{|p| - \beta, 0\}$.

This section aims to generalize the ISTA for functions in an RMOO. We begin with discussing why in \mathbb{R}^D the ISTA, and in particular, the STO, works well to obtain the sparse solution in the 1-norm regularization. The STO decreases the absolute value of each coordinate element by β , leaving the sign unchanged. If the absolute value is no larger than β , then the element will be shrunk to zero. This is why the ISTA causes sparsity of the solution. This is in contrast to the gradient descent method. Since each SHP's volume is zero, we cannot expect gradient descent method to induce sparsity unless the function has a special property. Let us evaluate the strength of the shrinkage effect geometrically. Recall that each element of a vector indicates the signed distance from an SHP. Hence, the following Proposition holds.

Proposition 1. For $p \in \mathbb{R}^D$, $\|\mathcal{T}_\beta(p)\|_0 < \|p\|_0$ if and only if $\min_{d=1,2,\dots,D} \Delta_{\mathbb{R}^D}(p, \Pi_{\hat{d}}) \leq \beta$.

Remark 4. Proposition 1 clarifies when the STO increases the sparsity of a point. It only depends on the distance of a point from the SHPs and β , and does not directly depend on the position of the point. This means that the strength of the sparsity induction effect by the STO is “uniform” everywhere in \mathbb{R}^D and only β determines the strength. Since $\beta = \alpha\lambda$ in the ISTA, the strength of the STO is proportional to λ . This is compatible to the aim of the objective function J in (1), where we expect λ to control the regularization strength.

As explained in Remark 4, Proposition 1 is the core property of the STO for the ISTA. This motivates us to keep the property in designing the STO for an RMOO and the Riemannian ISTA based on that. Our idea is to shrink the signed distance from each SHP directly the same way as the STO in RCS does. We can formulate this idea using the SSDM and its inverse if the SSDM is bijective, as follows:

Definition 7 (Riemannian STO (RSTO)). Let $(\mathcal{M}, \circ, (e_1, e_2, \dots, e_D))$ be a D -dimensional RMOO, and assume that $\exp_o : T_o\mathcal{M} \rightarrow \mathcal{M}$ and the SSDM $\delta : \mathcal{M} \rightarrow \mathbb{R}^D$ are bijective. We define the **Riemannian STO (RSTO)** $\mathcal{T}_\beta^\mathcal{M} : \mathcal{M} \rightarrow \mathcal{M}$ on the RMOO by $\mathcal{T}_\beta^\mathcal{M}(p) := \delta^{-1}(\mathcal{T}_\beta(\delta(p)))$.

The RSTO applies the STO on the signed distances from SHPs. By this definition, the RSTO has the uniform strength property corresponding to Proposition 1.

Proposition 2. Let $(\mathcal{M}, \circ, (e_1, e_2, \dots, e_D))$ be a D -dimensional RMOO, and assume that $\exp_o : T_o\mathcal{M} \rightarrow \mathcal{M}$ is bijective. Also, we assume that the SSDM $\delta^\mathcal{M} : \mathcal{M} \rightarrow \mathbb{R}^D$ is bijective. For all $p \in \mathcal{M}$, $\|\mathcal{T}_\beta^\mathcal{M}(p)\|_0 < \|p\|_0$ if and only if $\min_{d=1,2,\dots,D} \Delta_{\mathbb{R}^D}(p, \Pi_{\hat{d}}) \leq \beta$.

Remark 5. Similar to Proposition 1, Proposition 2 states that the RSTO has a “uniform” sparsity inducing strength determined only by β everywhere in \mathcal{M} . See also Remark 4.

We can define the **Riemannian ISTA (RISTA)**, the RMOO version of the ISTA, based on Definition 7. Recall that the Riemannian gradient descent methods, e.g., (Zhang & Sra, 2016), update the point by the exponential map: $p^{(t+1)} \leftarrow \exp_{p^{(t)}}(-\alpha \text{grad}_{p^{(t)}} L)$, where $J : \mathcal{M} \rightarrow \mathbb{R}$ is the objective function, $\text{grad}_{p^{(t)}} J$ is its Riemannian gradient at $p^{(t)}$, and $\alpha \in \mathbb{R}_{>0}$ is the learning rate. This operation is the Riemannian counterpart of the “negative gradient addition” in the gradient descent method in linear space. Replacing the “negative gradient addition” and the STO by the update by the exponential map and the RSTO, we obtain the RISTA as in Algorithm 1. As in the ISTA for linear space, the STO’s parameter is proportional to $\alpha\lambda$. We have established the sparse learning scheme for a general RMOO. In the next section, we will derive the specific formulae of the 1-norm and STO for hyperbolic space, which is this paper’s primary interest.

4 SPARSE LEARNING SCHEME FOR HYPERBOLIC SPACE

This section derives specific formulae of the core operations of our sparse learning scheme, the Riemannian 1-norm and RSTO, for hyperbolic space. As we can see in the previous sections, once we have the formula of the SSDM, we can immediately calculate the Riemannian 1-norm and RSTO. We provide the SSDM formula for hyperbolic space below (The proof is in Appendix C).

Theorem 1. The SSDM $\delta^{(\mathbb{D}^D, \mathcal{G}^p)} : \mathbb{D}^D \rightarrow \mathbb{R}^D$ is bijective. We can calculate the SSDM $\delta^{(\mathbb{D}^D, \mathcal{G}^p)}$ and its inverse maps $(\delta^{(\mathbb{D}^D, \mathcal{G}^p)})^{-1} : \mathbb{R}^D \rightarrow \mathbb{D}^D$ by

$$\delta^{(\mathbb{D}^D, \mathcal{G}^p)}(p) = \text{asinh}\left(\frac{2p}{1 - p^\top p}\right), \quad (\delta^{(\mathbb{D}^D, \mathcal{G}^p)})^{-1}(\sigma) = \frac{\sinh(\sigma)}{\sqrt{1 + (\sinh \sigma)^\top (\sinh \sigma) + 1}}, \quad (3)$$

whose time and space computational complexities are $O(D)$.

Remark 6. The SSDM formula in Theorem 1 is not surprising since point-hyperplane distance formulae have been given (e.g., Cho et al., 2019; Chien et al., 2021). Theorem 1’s significance is that

Algorithm 1 Riemannian ISTA (RISTA)

Input: $p_{\text{init}} \in \mathcal{M}$: initial point,
 $\alpha \in \mathbb{R}_{>0}$: learning rate,
 $T \in \mathbb{Z}_{>0}$: # iterations.
Output: $p_{\text{output}} \in \mathcal{M}$
 $p^{(0)} \leftarrow p_{\text{init}}$
for $t \leftarrow 1, 2, \dots, T$ **do**
 $q^{(t)} \leftarrow \exp_{p^{(t-1)}}(-\alpha \text{grad}_{p^{(t-1)}} L)$
 $p^{(t)} \leftarrow \mathcal{T}_{\alpha\lambda}(q^{(t)})$
end for
 $p_{\text{output}} \leftarrow p^{(T)}$

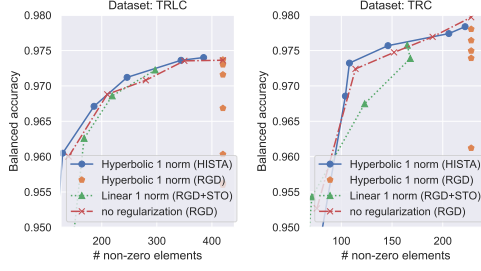


Figure 1: The trade-off between the representation quality (balanced accuracy) and the space complexity (the 0-norm). Left: **TRLC**, right: **TRC**. The closer to the left upper corner are the graphs, the better.

it calculates the distance to D SHPs at once in $O(D)$, while the straightforward application of the existing formula requires time complexity $O(D^2)$. The inverse SSDM formula is also novel.

Corollary 1. The Riemannian 1-norm $\|p\|_1^{(\mathbb{D}^D, \mathcal{G}^p)}$ of any $p \in \mathbb{D}^D$, which we call the **hyperbolic 1-norm** of p , and the RSTO $\mathcal{T}_\beta^{(\mathbb{D}^D, \mathcal{G}^p)}$, which we call the **hyperbolic STO (HSTO)**, are given by

$$\|p\|_1^{(\mathbb{D}^D, \mathcal{G}^p)} = \left\| \operatorname{asinh} \left(\frac{2p}{1 - p^\top p} \right) \right\|_1, \quad \mathcal{T}_\beta^{(\mathbb{D}^D, \mathcal{G}^p)}(p) = \frac{\sinh(\sigma)}{\sqrt{1 + (\sinh \sigma)^\top (\sinh \sigma) + 1}}, \quad (4)$$

where $\sigma \in \mathbb{R}^D$ is given by $\sigma := \operatorname{asinh} \left(\frac{2p}{1 - p^\top p} \right) - \beta \mathbf{1}_D$.

Now, by substituting the RSTO in Algorithm 1 by HSTO, we obtain the RISTA for hyperbolic space, which we call the **hyperbolic ISTA (HISTA)**. We also an explicit pseudocode of HISTA in Appendix D for the implementation convenience.

Lastly, we evaluate the non-uniform strength of the STO applied to the Poincaré model. Let $\epsilon = 10^3$ and $p, q = [0, \epsilon]^\top, [1 - \epsilon, \epsilon]^\top$. The distances from these two points to the 2nd SHP are significantly different: $\Delta_{(\mathbb{D}^D, \mathcal{G}^p)}(p, \Pi_2) \approx 0.002$, $\Delta_{(\mathbb{D}^D, \mathcal{G}^p)}(q, \Pi_2) \approx 0.882$. Nevertheless, both points are on the border of the area where the sparsity is increased by \mathcal{T}_ϵ . This example shows the non-uniformness of the STO. In contrast, Proposition 2 guarantees that the RSTO does not cause this issue.

5 EXPERIMENTS

Our numerical experiments aim to confirm when our sparse learning scheme works effectively. The objectives are (i) to confirm that our sparse learning scheme actually obtains sparse representations, (ii) to analyze when the sparse representations obtained by our scheme is effective, and (iii) to analyze when the HISTA effectively works as an optimization algorithm. We achieve the (i) and (ii) in our graph embedding experiments, and (iii) in our square distance minimization experiments.

5.1 SPARSE REPRESENTATION QUALITY EVALUATION IN GRAPH EMBEDDING

Our objective here is NOT to maximize the quality of representations, but to compare our sparse learning scheme with possible alternatives. Hence, we consider a simple edge labeling problem for a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} is the vertex set and \mathcal{E} is the edge set. We explain the outline of the experimental settings and leave the details in Appendix E, with additional results supporting our discussion. For an edge $\{u, v\}$, where $u, v \in \mathcal{V}$, we define the label $y_{u,v}$ as +1 if $\{u, v\} \in \mathcal{V}$ and -1 otherwise. We predict the label from the representations $(z_v)_{v \in \mathcal{V}}$ learned by optimizing a hinge-loss-based loss function with a regularization term. Specifically, the predicted label $\hat{y}_{u,v}$ is given by +1 if $\Delta_{(\mathbb{D}^D, \mathcal{G}^p)}(z_u, z_v) < \theta$ and -1 otherwise. Here, $\theta \in \mathbb{R}_{>0}$ is the threshold hyperparameter.

We get representations by optimizing a simple hinge-loss-based loss function with a regularization term. We compare the following three regularization functions: **hyperbolic 1-norm** $\|z\|_{1, (\mathbb{D}^D, \mathcal{G}^p)}$, **linear 1-norm** (in the Poincaré model) $\|z\|_1$, and **no regularization**. In this section, the 1-norm $\|\cdot\|_1$ is called the linear 1-norm to distinguish it with the hyperbolic 1-norm explicitly. We evaluate

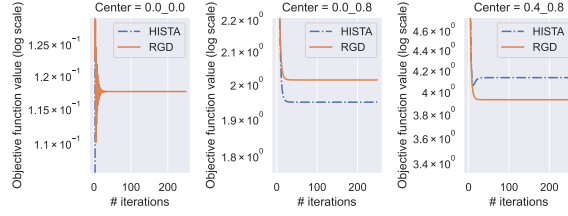


Figure 2: The optimization performance comparison. $z' = [0.0 \ 0.0]^\top, [0.0 \ 0.8]^\top, [0.4 \ 0.8]^\top$ from left to right. The HISTA's plot hits the bottom for $z' = [0.0 \ 0.0]^\top$ since it achieves the zero value, i.e., the true solution.

the trade-off between the quality and the space complexity of the representations. We measure the quality by the balanced accuracy and the space complexity by the 0-norm, i.e., the number of zero elements. Specifically, we vary the regularization weight in $\lambda = \{1.0, 2.0, 5.0\} \times 10^{\{-3, -2, -1\}}$ for regularization methods and the space dimension $D = 2, 3, 4, 5, 6$ for the no regularization method, to observe how each method can achieve a high accuracy with a low 0-norm. Here, the dimension $D = 6$ is fixed for the two regularization methods.

We consider the following two graphs with partial tree-like structures, to see the difference between the hyperbolic 1-norm regularization and linear 1-norm regularization. **Tree-RootLeafCubes (TRLIC)** consists of a m -dimensional cube graph connecting the roots of two n -ary trees with height h , and four m -dimensional cube graphs connected to leaves of the trees. **Tree-RootCubes (TRC)** is similar to the above but cubes connected to leaves are omitted. Uniform regularization is needed for **TRLIC** since it has cubes both at the root and around the leaves. Conversely, strong regularization around the boundary and weak regularization could work well for **TRC** since it has a cube only at the root. Hence, our natural expectation is that hyperbolic 1-norm regularization works better for **TRLIC** than linear 1-norm regularization, but the tendency is not clear for **TRC**.

Figure 1 shows how the accuracy and the 0-norm changes by varying λ or D . The closer to the left upper corner are graphs, the better. Here, we show the range where the sum of the 0-norms no smaller than $2|\mathcal{V}|$; otherwise the mean 0-norm would be lower than two, which would be meaningless as representations. We have also plotted the results of the hyperbolic 1-norm regularization optimized by RGD, which shows that RGD fails to get sparse representations, while shrinkage-thresholding operators succeeded. As we have expected, the hyperbolic 1-norm regularization outperforms the others in **TRLIC**. It shows that our Riemannian 1-norm regularization can select the dimension of each representations efficiently. In **TRC**, the linear 1-norm regularization outperforms others around where the sum of the 0-norm is 75, as we have expected. Interestingly, our hyperbolic 1-norm regularization outperforms the linear 1-norm regularizations where the sum of the 0-norm is larger. One possible reason is that the linear 1-norm regularizations tend to be unstable since the STO changes the representations dramatically around the ball boundary. Although comparing the optimization process is not trivial since they optimizes different functions, clarifying the reason of the low performance of the linear regularization would be interesting future work.

5.2 OPTIMIZATION PERFORMANCE EVALUATION ON THE SQUARE DISTANCE MINIMIZATION

This subsection evaluates the optimization algorithms' performance. Although the RGD tends to oscillate and fail to get a sparse result, it is possible that the RGD is better if our motivation is in the optimization of the function. Since the ISTA can optimize the objective function in linear space more efficiently than the gradient descent if the regularization term is not smooth around the solution (e.g., Vonesch & Unser, 2007), our interest is whether we have similar results in hyperbolic space. Hence, we compare the HISTA and RGD for sparse solution cases and non-sparse solution cases.

We consider minimizing the square distance with the hyperbolic 1-norm regularization: $L(z) = \left[\Delta_{(\mathbb{D}^2, \mathcal{G}^p)}(z, z') \right]^2 + \lambda \|z\|_{1, (\mathbb{D}^2, \mathcal{G}^p)}$. Here, we set $z' = [0.0 \ 0.0]^\top, [0.0 \ 0.8]^\top, [0.4 \ 0.8]^\top$. We expect that the true solution is sparse for the first two cases and non-sparse for the last case, though we do not know the analytic solution for the latter two. We set $\lambda = 1.0$ and $\alpha = 0.1$ for all cases.

Figure 2 shows that the HISTA outperforms the RGD for $z' = [0.0 \ 0.0]^\top, [0.0 \ 0.4]^\top$ in terms of the objective function's value as well as obtaining a sparse solution. For $z' = [0.4 \ 0.8]^\top$, the RGD can outperform the HISTA. We also observe a "bounce back" effect by the HISTA, which could be a drawback. Still, the HISTA is stable for all the cases, while the oscillation of the RGD is significant for $z' = [0.0 \ 0.0]^\top$. Our results confirm that the superiority of the HISTA over the RGD in terms of the function value for sparse solution cases. Nevertheless, detecting the cause of the bounce back effect by the HISTA would be interesting future work.

6 CONCLUSION

This paper has established a novel sparse learning scheme for HSBRL for the first time, based on geometric definitions of the sparsity, regularization term, and the optimization algorithm, with effective closed-form formulae. The HISTA is the fundamental optimization algorithm for sparse learning, and its extension to accelerated methods could be interesting and realistic future work.

REFERENCES

- P-A Absil, Robert Mahony, and Rodolphe Sepulchre. Optimization algorithms on matrix manifolds. In *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2009.
- Milad Alizadeh, Arash Behboodi, Mart van Baalen, Christos Louizos, Tijmen Blankevoort, and Max Welling. Gradient ℓ_1 regularization for quantization robustness. In *International Conference on Learning Representations*, 2020.
- Ivana Balazevic, Carl Allen, and Timothy Hospedales. Multi-relational poincaré graph embeddings. *Advances in Neural Information Processing Systems*, 32, 2019.
- Gary Becigneul and Octavian-Eugen Ganea. Riemannian adaptive optimization methods. In *International Conference on Learning Representations*, 2019.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomáš Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. *Advances in Neural Information Processing Systems*, 26, 2013.
- Ronald E Bruck Jr. On the weak convergence of an ergodic iteration for the solution of variational inequalities for monotone operators in hilbert space. *Journal of Mathematical Analysis and Applications*, 61(1):159–164, 1977.
- Antonin Chambolle, Ronald A De Vore, Nam-Yong Lee, and Bradley J Lucier. Nonlinear wavelet image processing: variational problems, compression, and noise removal through wavelet shrinkage. *IEEE Transactions on Image Processing*, 7(3):319–335, 1998.
- Ines Chami, Zhitao Ying, Christopher Ré, and Jure Leskovec. Hyperbolic graph convolutional neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- Ines Chami, Albert Gu, Vaggos Chatziafratis, and Christopher Ré. From trees to continuous embeddings and back: Hyperbolic hierarchical clustering. *Advances in Neural Information Processing Systems*, 33:15065–15076, 2020.
- Scott Shaobing Chen, David L Donoho, and Michael A Saunders. Atomic decomposition by basis pursuit. *SIAM review*, 43(1):129–159, 2001.
- Eli Chien, Chao Pan, Puoya Tabaghi, and Olgica Milenkovic. Highly scalable and provably accurate classification in poincaré balls. In *2021 IEEE International Conference on Data Mining (ICDM)*, pp. 61–70. IEEE, 2021.
- Hyunghoon Cho, Benjamin DeMeo, Jian Peng, and Bonnie Berger. Large-margin classification in hyperbolic space. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1832–1840. PMLR, 2019.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- Joseph M Dale, Liviu Popescu, and Peter D Karp. Machine learning methods for metabolic pathway prediction. *BMC bioinformatics*, 11(1):1–14, 2010.
- Ingrid Daubechies, Michel Defrise, and Christine De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 57(11):1413–1457, 2004.
- David Donoho and Jared Tanner. Counting faces of randomly projected polytopes when the projection radically lowers dimension. *Journal of the American Mathematical Society*, 22(1):1–53, 2009a.

- David Donoho and Jared Tanner. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 367(1906): 4273–4293, 2009b.
- David L Donoho and Iain M Johnstone. Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association*, 90(432):1200–1224, 1995.
- David L Donoho and Jared Tanner. Counting the faces of randomly-projected hypercubes and orthants, with applications. *Discrete & Computational Geometry*, 43(3):522–541, 2010.
- Takuma Ebisu and Ryutaro Ichise. TorusE: Knowledge graph embedding on a Lie group. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, pp. 1819–1826. AAAI Press, 2018.
- Mário AT Figueiredo and Robert D Nowak. An em algorithm for wavelet-based image restoration. *IEEE Transactions on Image Processing*, 12(8):906–916, 2003.
- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.
- Jerome Friedman, Trevor Hastie, and Rob Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1):1, 2010.
- Octavian Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic entailment cones for learning hierarchical embeddings. In *International Conference on Machine Learning*, pp. 1646–1655. PMLR, 2018a.
- Octavian-Eugen Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic entailment cones for learning hierarchical embeddings. In *the 35th International Conference on Machine Learning*, volume 80 of *Machine Learning Research*, pp. 1632–1641. PMLR, 2018b.
- Octavian-Eugen Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic neural networks. In *the 32nd Conference on Neural Information Processing Systems*, pp. 5350–5360, 2018c.
- Mikhael Gromov. Hyperbolic groups. In *Essays in Group Theory*, pp. 75–263. Springer, 1987.
- Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 855–864. ACM, 2016.
- Albert Gu, Frederic Sala, Beliz Gunel, and Christopher Ré. Learning mixed-curvature representations in product spaces. In *International Conference on Learning Representations*, 2019.
- Elaine T Hale, Wotao Yin, and Yin Zhang. A fixed-point continuation method for l_1 -regularized minimization with applications to compressed sensing. *CAAM TR07-07, Rice University*, 43:44, 2007.
- Peter D Hoff, Adrian E Raftery, and Mark S Handcock. Latent space approaches to social network analysis. *Journal of the American Statistical Association*, 97(460):1090–1098, 2002.
- Wen Huang and Ke Wei. Riemannian proximal gradient methods. *Mathematical Programming*, 194(1):371–413, 2022.
- HiroYuki Kasai, Pratik Jawanpuria, and Bamdev Mishra. Riemannian adaptive stochastic gradient algorithms on matrix manifolds. In *International Conference on Machine Learning*, pp. 3262–3271. PMLR, 2019.
- Zhaobin Kuang, Sinong Geng, and David Page. A screening rule for l_1 -regularized ising model estimation. *Advances in Neural Information Processing Systems*, 30, 2017.
- Panagiotis Kyriakis, Iordanis Fostropoulos, and Paul Bogdan. Learning hyperbolic representations of topological features. In *International Conference on Learning Representations*, 2021.

- Marc Law, Renjie Liao, Jake Snell, and Richard Zemel. Lorentzian distance learning for hyperbolic representations. In *International Conference on Machine Learning*, pp. 3672–3681. PMLR, 2019.
- John M Lee. *Introduction to Riemannian manifolds*, volume 176. Springer, 2018.
- Baiyang Liu, Lin Yang, Junzhou Huang, Peter Meer, Leiguang Gong, and Casimir Kulikowski. Robust and fast collaborative tracking with two stage sparse optimization. In *European Conference on Computer Vision*, pp. 624–637. Springer, 2010.
- Qi Liu, Maximilian Nickel, and Douwe Kiela. Hyperbolic graph neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yuanyuan Liu, Fanhua Shang, James Cheng, Hong Cheng, and Licheng Jiao. Accelerated first-order methods for geodesically convex optimization on riemannian manifolds. *Advances in Neural Information Processing Systems*, 30, 2017a.
- Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2736–2744, 2017b.
- Abdur Rahman MA Basher and Steven J Hallam. Leveraging heterogeneous network embedding for metabolic pathway prediction. *Bioinformatics*, 37(6):822–829, 2021.
- Emile Mathieu, Charline Le Lan, Chris J Maddison, Ryota Tomioka, and Yee Whye Teh. Continuous hierarchical representations with poincaré variational auto-encoders. *Advances in neural Information Processing Systems*, 32, 2019.
- Tomás Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. In *the 27th Conference on Neural Information Processing Systems*, pp. 3111–3119, 2013.
- Andrew Y Ng. Feature selection, l_1 vs. l_2 regularization, and rotational invariance. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 78, 2004.
- Maximilian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. In *the 31st Conference on Neural Information Processing Systems*, pp. 6338–6347, 2017.
- Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In *the 28th International Conference on Machine Learning*, pp. 809–816. Omnipress, 2011.
- Maximilian Nickel, Lorenzo Rosasco, and Tomaso A. Poggio. Holographic embeddings of knowledge graphs. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*, pp. 1955–1961. AAAI Press, 2016.
- Maximilian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz model of hyperbolic geometry. In *International Conference on Machine Learning*, pp. 3779–3788. PMLR, 2018.
- Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *the 2014 Conference on Empirical Methods in Natural Language Processing*, pp. 1532–1543. ACL, 2014.
- Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. DeepWalk: online learning of social representations. In *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New York, NY, USA - August 24 - 27, 2014*, pp. 701–710. ACM, 2014.
- Chunhong Qi, Kyle A Gallivan, and P-A Absil. Riemannian bfgs algorithm with applications. In *Recent advances in optimization and its applications in engineering*, pp. 183–192. Springer, 2010.

- Sebastian Riedel, Limin Yao, Andrew McCallum, and Benjamin M. Marlin. Relation extraction with matrix factorization and universal schemas. In *Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics, Proceedings, June 9-14, 2013, Westin Peachtree Plaza Hotel, Atlanta, Georgia, USA*, pp. 74–84. The Association for Computational Linguistics, 2013.
- Lev V Sabinin, Ludmila L Sabinina, and Larissa V Sbitneva. On the notion of gyrogroup. *aequationes mathematicae*, 56(1):11–17, 1998.
- Frederic Sala, Christopher De Sa, Albert Gu, and Christopher Ré. Representation tradeoffs for hyperbolic embeddings. In *the 35th International Conference on Machine Learning*, volume 80 of *Machine Learning Research*, pp. 4457–4466. PMLR, 2018.
- Shai Shalev-Shwartz and Ambuj Tewari. Stochastic methods for l_1 regularized loss minimization. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pp. 929–936, 2009.
- Ryohei Shimizu, YUSUKE Mukuta, and Tatsuya Harada. Hyperbolic neural networks++. In *International Conference on Learning Representations*, 2021.
- Rishi Sonthalia and Anna Gilbert. Tree! i am no tree! i am a low dimensional hyperbolic embedding. *Advances in Neural Information Processing Systems*, 33:845–856, 2020.
- Atsushi Suzuki, Jing Wang, Feng Tian, Atsushi Nitanda, and Kenji Yamanishi. Hyperbolic ordinal embedding. In *the 11th Asian Conference on Machine Learning*, volume 101 of *Machine Learning Research*, pp. 1065–1080. PMLR, 2019.
- Atsushi Suzuki, Atsushi Nitanda, Jing Wang, Linchuan Xu, Kenji Yamanishi, and Marc Cavazza. Generalization error bound for hyperbolic ordinal embedding. *arXiv preprint arXiv:2105.10475*, 2021a.
- Atsushi Suzuki, Atsushi Nitanda, Linchuan Xu, Kenji Yamanishi, Marc Cavazza, et al. Generalization bounds for graph embedding using negative sampling: Linear vs hyperbolic. *Advances in Neural Information Processing Systems*, 34:1243–1255, 2021b.
- Puoya Tabaghi and Ivan Dokmanic. Hyperbolic distance matrices. In *the 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1728–1738. ACM, 2020.
- Masaaki Takada and Hironori Fujisawa. Transfer learning via ℓ_1 regularization. *Advances in Neural Information Processing Systems*, 33:14266–14277, 2020.
- Jun Takeuchi, Noriki Nishida, and Hideki Nakayama. Neural networks in a product of hyperbolic spaces. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop*, pp. 211–221, 2022.
- Xingwei Tan, Gabriele Pergola, and Yulan He. Extracting event temporal relations via hyperbolic geometry. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 8065–8077, 2021.
- Jian Tang, Meng Qu, and Qiaozhu Mei. PTE: predictive text embedding through large-scale heterogeneous text networks. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, pp. 1165–1174. ACM, 2015a.
- Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. LINE: large-scale information network embedding. In *Proceedings of the 24th International Conference on World Wide Web, WWW 2015, Florence, Italy, May 18-22, 2015*, pp. 1067–1077. ACM, 2015b.
- Yi Tay, Luu Anh Tuan, and Siu Cheung Hui. Hyperbolic representation learning for fast and efficient neural question answering. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, pp. 583–591, 2018.

- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Alexandru Tifrea, Gary Becigneul, and Octavian-Eugen Ganea. Poincare glove: Hyperbolic word embeddings. In *International Conference on Learning Representations*, 2019.
- Ryota Tomioka, Taiji Suzuki, and Masashi Sugiyama. Super-linear convergence of dual augmented lagrangian algorithm for sparsity regularized estimation. *Journal of Machine Learning Research*, 12(5), 2011.
- Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pp. 2071–2080. JMLR.org, 2016.
- Abraham A Ungar. The holomorphic automorphism group of the complex disk. *aequationes mathematicae*, 47(2):240–254, 1994.
- Abraham A Ungar. Extension of the unit disk gyrogroup into the unit ball of any real inner product space. *Journal of Mathematical Analysis and Applications*, 202(3):1040–1057, 1996.
- Cédric Vonesch and Michael Unser. A fast iterative thresholding algorithm for wavelet-regularized deconvolution. In *Wavelets XII*, volume 6701, pp. 135–139. SPIE, 2007.
- Stephen J Wright, Robert D Nowak, and Mário AT Figueiredo. Sparse reconstruction by separable approximation. *IEEE Transactions on Signal Processing*, 57(7):2479–2493, 2009.
- Tao Yu and Christopher M De Sa. Numerically accurate hyperbolic embeddings using tiling-based models. *Advances in Neural Information Processing Systems*, 32, 2019.
- Hongyi Zhang and Suvrit Sra. First-order methods for geodesically convex optimization. In *Conference on Learning Theory*, pp. 1617–1638. PMLR, 2016.
- Hongyi Zhang, Sashank J Reddi, and Suvrit Sra. Riemannian svrg: Fast stochastic optimization on riemannian manifolds. *Advances in Neural Information Processing Systems*, 29, 2016a.
- Yuchen Zhang, Jason D Lee, and Michael I Jordan. l1-regularized neural networks are improperly learnable in polynomial time. In *International Conference on Machine Learning*, pp. 993–1001. PMLR, 2016b.
- Pan Zhou, Xiao-Tong Yuan, and Jiashi Feng. Faster first-order methods for stochastic non-convex optimization on riemannian manifolds. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 138–147. PMLR, 2019.

A VISUALIZATION OF THE SSDM

. See Figure 3.

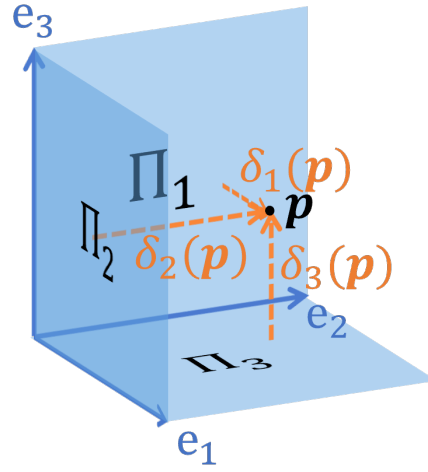


Figure 3: The SSDM’s visualization. It measures the distance from each SHP.

B DETAILED EXPLANATION OF EXAMPLE 6

The function $f(p) = |p|$ is differentiable at $p \neq 0$ and the derivative is given by $\frac{d}{dp}f(p) = \text{sgn}(p)$. Suppose that the learning rate is $\alpha > 0$ and the initial point is $p^{(0)} \neq 0$. By the symmetry about the origin, we can assume that $p^{(0)} > 0$ without loss of generality. Then the gradient descent generates the series $p^{(0)}, p^{(1)}, \dots$ of points according to the following recursion:

$$p^{(t+1)} \leftarrow p^{(t)} - \alpha \frac{d}{dp}f(p^{(t)}) = \begin{cases} p^{(t)} - \alpha & \text{if } p^{(t)} > 0, \\ p^{(t)} + \alpha & \text{if } p^{(t)} < 0. \end{cases} \quad (5)$$

Here, we usually set $p^{(t+1)} \leftarrow p^{(t)}$ if $p^{(t)} = 0$, which we can justify as a subgradient method. We can see from (5) that the algorithm ends up oscillating between $p^{(0)} - \alpha n$ and $p^{(0)} - \alpha(n+1)$ unless $p^{(0)}$ is an integral multiple of α , where $n = \left\lceil \frac{p^{(0)}}{\alpha} \right\rceil$ is the maximum integer that is no greater than $\frac{p^{(0)}}{\alpha}$.

C PROOF OF THEOREM 1

Proof. We prove for the SSDM $\delta^{(\mathbb{D}^D, G^p)}$. It suffices to prove that the absolute values are correct since the logarithmic map at the origin of the Poincaré model does not change the sign of each element. Let $\mathbf{h} \in \mathbb{D}^2$ be the foot of the geodesic pass through \mathbf{p} on Π_d . Note that \mathbf{h} is unique according to Gauss-Bonnet theorem. We have that $\mathbf{h} = \text{argmin}_{\mathbf{q}} \Delta_{(\mathbb{D}^D, G^p)}(\mathbf{p}, \mathbf{q})$ from hyperbolic Pythagorean theorem. In the following, we regard the ball of the Poincaré model as a unit ball in Euclidean space and discuss using elementary geometry. A geodesic in hyperbolic space is now an arc orthogonal to the unit ball and Π_d and passing through \mathbf{p} . Define $\mathbf{p}' = \frac{\mathbf{p}}{p^\top \mathbf{p}}$ and $\mathbf{h}' = \frac{\mathbf{h}}{h^\top \mathbf{h}}$. Also denote by \mathbf{m} the midpoint of \mathbf{p} and \mathbf{p}' and by \mathbf{j} the midpoint of \mathbf{h} and \mathbf{h}' . Note that \mathbf{j} is the center of the arc drawn by the geodesic. Since the arc is orthogonal to the unit ball, it also passes through \mathbf{p}' and \mathbf{h}' according to the power of a point theorem. The subplane including the arc also contains \mathbf{p} , \mathbf{h} , and \mathbf{p}' . Hence, the following discussion is on the subplane. We regard the axis in the subplane on the intersection of the subplane and Π_d as x -axis, and the other axis orthogonal to Π_d as y -axis. We indicate the coordinate of the \mathbf{p} in the subplane by $[x \ y]^\top$ and that of \mathbf{h} by $[h \ 0]^\top$. The coordinates of \mathbf{p}' and \mathbf{h}' are $\frac{1}{x^2+y^2}[x \ y]^\top$ and $[1/h \ 0]^\top$, respectively. See also Figure 4. We have that $|\mathbf{m}| = \frac{1}{2} \left(\sqrt{x^2 + y^2} + \frac{1}{\sqrt{x^2 + y^2}} \right)$ and $|\mathbf{j}| = \frac{1}{2} \left(h + \frac{1}{h} \right)$. By similarity of two right triangles, we have that $\frac{\sqrt{x^2+y^2}}{x} = \frac{|\mathbf{j}|}{|\mathbf{m}|}$. Hence, $|\mathbf{j}| = \frac{x^2+y^2+1}{2x}$. Noting that $h < 1 < \frac{1}{h}$, we have that $h = \frac{x^2+y^2 - \sqrt{(x^2+y^2)^2 - 4x^2}}{2x}$. We get the expected result by substituting this to the distance formula of the Poincaré model: $\Delta_{(\mathbb{D}^2, G^p)}(\mathbf{p}, \mathbf{q}) = \text{acosh} \left(1 + \frac{2|\mathbf{p}-\mathbf{q}|^2}{(1-|\mathbf{p}|^2)(1-|\mathbf{q}|^2)} \right)$. Specifically,

$$\begin{aligned} \Delta_{(\mathbb{D}^2, G^p)}([x \ y], [h \ 0]) &= \text{acosh} \left(1 + \frac{2((x-h)^2 + y^2)}{(1-(x^2+y^2))(1-h^2)} \right) \\ &= \text{acosh} \left(\sqrt{1 + \frac{4y^2}{(1-(x^2+y^2))^2}} \right) \\ &= \text{asinh} \left(\frac{2y}{(1-(x^2+y^2))^2} \right). \end{aligned} \quad (6)$$

We complete the proof by recalling that $y = p_d$ and $(x^2 + y^2) = \mathbf{p}^\top \mathbf{p}$. \square

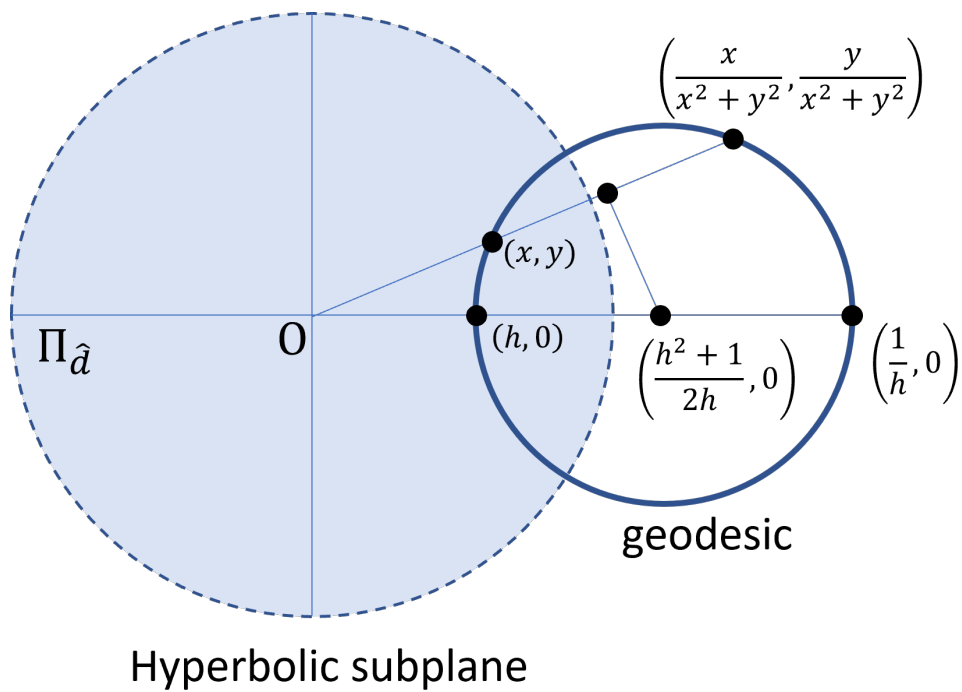


Figure 4: Hyperbolic subdisk.

D EXPLICIT PSEUDOCODE OF THE HISTA

Input:

$\mathbf{p}_{\text{init}} \in \mathbb{D}^D$: initial point,
 $\alpha \in \mathbb{R}_{>0}$: learning rate,
 $T \in \mathbb{Z}_{\geq 0}$: # iterations.

Output: $\mathbf{p}_{\text{output}} \in \mathbb{D}^D$

```

 $\mathbf{p}^{(0)} \leftarrow \mathbf{p}_{\text{init}}$ 
for  $t \leftarrow 1, 2, \dots, T$  do
   $\boldsymbol{\gamma}^{(t)} \leftarrow \partial|_{\mathbf{p}^{(t-1)}} J$ 
   $\rho^{(t)} \leftarrow \frac{4}{(1 - |\mathbf{p}^{(t-1)}|^2)^2}$ 
   $\mathbf{g}^{(t)} \leftarrow (\rho^{(t)})^2 \boldsymbol{\gamma}^{(t)}$ 
   $\mathbf{q}^{(t-1)} \leftarrow \rho^{(t)} \cdot \frac{[\cosh(|-\alpha \mathbf{g}^{(t)}|) - \rho^{(t)} \alpha (\mathbf{g}^{(t)})^\top (\mathbf{p}^{(t-1)})] \mathbf{p}^{(t-1)} + \sinh(|\alpha \mathbf{g}^{(t)}|) \mathbf{g}^{(t)}}{1 + (\rho^{(t)} - 1) \cosh(|-\alpha \mathbf{g}^{(t)}|) - (\rho^{(t)})^2 \alpha (\mathbf{g}^{(t)})^\top (\mathbf{p}^{(t-1)}) \sinh(|\alpha \mathbf{g}^{(t)}|)}$ 
   $\boldsymbol{\sigma}^{(t)} \leftarrow \text{asinh}\left(\frac{2\mathbf{p}^{(t-1)}}{1 - \mathbf{p}^{(t-1)\top} \mathbf{p}^{(t-1)}}\right) - \alpha \lambda \mathbf{1}_D$ 
   $\mathbf{p}^{(t)} \leftarrow \frac{\sinh(\boldsymbol{\sigma}^{(t)})}{\sqrt{1 + (\sinh \boldsymbol{\sigma}^{(t)})^\top (\sinh \boldsymbol{\sigma}^{(t)})} + 1}$ 
end for
 $\mathbf{p}_{\text{output}} \leftarrow \mathbf{p}^{(T)}$ 

```

E EXPERIMENTAL SETTING DETAILS OF SECTION 5.1

We denote a graph by $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} is the vertex set and \mathcal{E} is the edge set. Since edges are the most fundamental form of information about entity relations, graph embedding has wide applications. Hence, we choose graph embedding as the problem on which we evaluate the performance of our sparse learning scheme. Our objective here is NOT to maximize the quality of representations, but to compare our sparse learning scheme with possible alternatives. Hence, we use the following simplest graph embedding setting. Let $\mathcal{C}(\mathcal{V}, 2)$ be the set of subsets of \mathcal{V} , whose size is two. That is, $\mathcal{C}(\mathcal{V}, 2)$ is the set of unordered vertex pairs. Define the label of a vertex pair $y_{u,v} \in \{-1, +1\}$ and the sample weight $w_{u,v} \in \mathbb{R}_{>0}$ for $u, v \in \mathcal{V}$ such that $u \notin v$ by

$$y_{u,v} := \begin{cases} +1 & \text{if } \{u, v\} \in \mathcal{E}, \\ -1 & \text{if } \{u, v\} \notin \mathcal{E}, \end{cases} \quad w_{u,v} := 2 \cdot \frac{|\{\{u', v'\} \in \mathcal{C}(\mathcal{V}, 2) \mid y_{u',v'} = y_{u,v}\}|}{|\mathcal{C}(\mathcal{V}, 2)|}. \quad (7)$$

Also, define the 0-1 loss function $l_{0-1} : \mathbb{R} \times \{-1, +1\} \rightarrow \{0, +1\}$ by

$$l_{0-1}(\hat{y}, y) := \begin{cases} 0 & \text{if } \text{sgn}(\hat{y}) = y, \\ +1 & \text{otherwise.} \end{cases} \quad (8)$$

The objective of our experimental setting is to minimize the following balanced 0-1 loss:

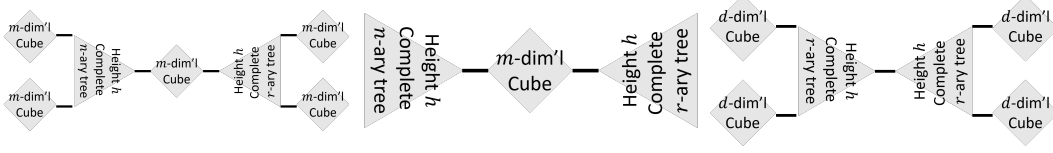
$$L_{0-1}((\mathbf{z}_v)_{v \in \mathcal{V}}; \mathcal{G}) := \sum_{\{u,v\} \in \mathcal{C}(\mathcal{V}, 2)} w_{u,v} l_{0-1} \left(\left[\Delta_{(\mathbb{D}^D, \mathbf{G}^p)}(\mathbf{z}_u, \mathbf{z}_v) \right]^2 - \theta, y_{u,v} \right), \quad (9)$$

where the hyperparameter $\theta \in \mathbb{R}_{>0}$ determines the threshold in labeling the pair to be positive or negative.

The above function L_{0-1} is not easy to optimize since it is not continuous. Hence, in the optimization step, we replace l_{0-1} by the hinge loss function $l_{\text{hinge}} : \mathbb{R} \times \{-1, +1\} \rightarrow \mathbb{R}_{\leq 0}$ defined by $l_{\text{hinge}}(\hat{y}, y) = \max -\hat{y}y + 1, 0$, widely used in machine learning area, e.g., support vector machines (Cortes & Vapnik, 1995). That is, the loss function in the optimization step is

$$L((\mathbf{z}_v)_{v \in \mathcal{V}}; \mathcal{G}) := \sum_{\{u,v\} \in \mathcal{C}(\mathcal{V}, 2)} w_{u,v} l \left(\left[\Delta_{(\mathbb{D}^D, \mathbf{G}^p)}(\mathbf{z}_u, \mathbf{z}_v) \right]^2 - \theta, y_{u,v} \right). \quad (10)$$

Also, we add the regularization term $\lambda \sum_{v \in \mathcal{V}} r(\mathbf{z}_v)$ to the objective function, where the regularization function $r : \mathbb{D}^D \rightarrow \mathbb{R}_{\geq 0}$ is the object of the comparison in the experiments and varies for each

Figure 5: The datasets' structure. Left: **TRLC**, center: **TRC**, right: **TLC**.

method. Note that $\lambda \mathbb{R}_{\geq 0}$ determines the regularization strength. To wrap up, we optimize the function $J((z_v)_{v \in \mathcal{V}}; \mathcal{G}) := L((z_v)_{v \in \mathcal{V}}; \mathcal{G}) + \lambda \sum_{v \in \mathcal{V}} r(z_v)$.

As a regularization function r , we compare the following three:

$$r(z) = \begin{cases} \|z\|_{1,(\mathbb{D}^D, \mathcal{G}^p)} & \text{Riemannian 1-norm,} \\ \|z\|_1 & \text{linear 1-norm,} \\ 0 & \text{no regularization.} \end{cases} \quad (11)$$

We use the HISTA for the hyperbolic 1-norm. For the linear-norm, we apply Riemannian gradient descent and traditional shrinkage-thresholding operator. We evaluate the balanced accuracy $1 - L_{0.1}((z_v)_{v \in \mathcal{V}}; \mathcal{G})$ and the sum $\sum_{v \in \mathcal{V}} \|z_v\|_{0,(\mathbb{D}^D, \mathcal{G}^p)} = \sum_{v \in \mathcal{V}} \|z_v\|_0$ of the 0-norms of the representations. The higher accuracy and lower 0-norm, the better, but there is a trade-off between the accuracy and the 0-norm. Specifically, the stronger the regularization is, the lower accuracy and lower 0-norm it gets, and vice versa. Hence, we vary the regularization weight λ and observe how the accuracy and the 0-norm changes. For the no regularization method, we vary D instead of λ . In this case, the lower D is, the lower accuracy and lower 0-norm it gets, and vice versa.

As a graph, we consider tree-like structures that are not completely tree, which are our main focus. To see the difference between the Riemannian 1-norm regularization and Linear 1-norm regularization, we consider the following two structures.

Tree-RootLeafCubes (TRLC) consisting of two complete n -ary trees with height h and five m -dimensional cubes. One cube is in between the roots of the two trees, where each vertex of a hyperbody diagonal pair (a most distant pair) in the cube has an edge to the root of a tree. The other four cubes are connected to a leaf of a tree. Two cubes are connected to one tree and the other two cubes are connected to the other tree. Here, each of the former two cubes have one edge to a leaf of the tree, where the two leaves connected to a cube are most distant to each other. The same holds true for the other tree and the latter two cubes.

Tree-RootCubes (TRC) is similar to the above but without the cube connected to the leaves of the trees.

To support the discussion, we also show the results on the following graph.

Tree-LeafCubes (TLC) is similar to the above but without the cube connecting the roots of the trees and the roots are directly connected.

Uniform regularization is needed for **TRLC** since it has cubes both at the root and around the leaves. Conversely, strong regularization around the boundary and weak regularization could work well for **TRC** since it has a cube only at the root. For **TLC**, strong regularization around the center is needed. Hence, we expect the linear 1-norm regularization does not work well. Still, the non-tree-like area is smaller than **TRLC**, we expect that the advantage of the hyperbolic 1-norm regularization is smaller. Hence, our natural expectation is that Riemannian 1-norm regularization works better for **TRLC** than Linear 1-norm regularization, but the tendency is not clear for **TRC**. Figure 5 visualizes these graph structures.

Other experimental settings are as follows. We set $n = 2$, $h = 3$, and $m = 3$. The dimension $D = 6$ is fixed for the two regularization methods, while $D = 2, 3, 4, 5, 6$ for the no regularization method. The regularization strength varies among $\lambda = \{1.0, 2.0, 5.0\} \times 10^{\{-3, -2, -1\}}$ for the two regularization methods. The learning rate that achieved the best accuracy is selected from $\lambda = \{1.0, 2.0, 5.0\} \times 10^{\{1, 2, 3\}}$. The threshold hyperparameter θ is set to 1.0. The number of iterations is set to $T = 10000$.

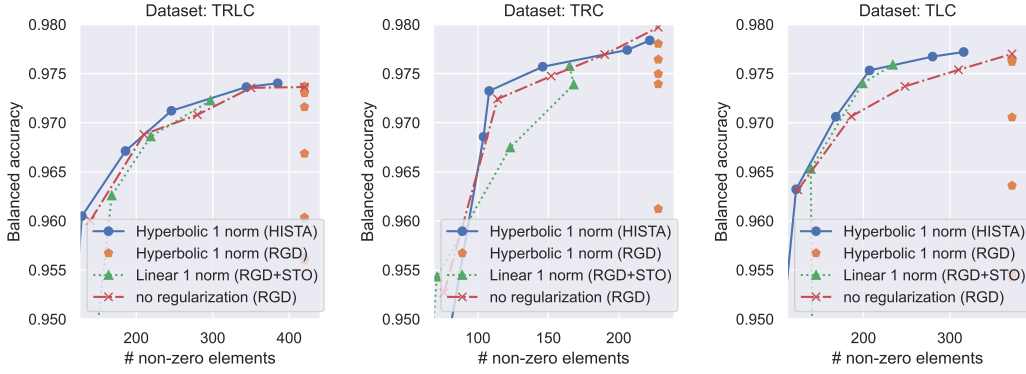


Figure 6: The trade-off between the representation quality (balanced accuracy) and the space complexity (the 0-norm). Left: **TRLC**, center: **TRC**, right: **TLC**. The closer to the left upper corner are the graphs, the better.

Figure 1 shows how the accuracy and the 0-norm changes by varying λ or D . The closer to the left upper corner are graphs, the better. Here, we show the range where the sum of the 0-norms no smaller than $2|\mathcal{V}|$; otherwise the mean 0-norm would be lower than two, which would be meaningless as representations. We have also plotted the results of the hyperbolic 1-norm regularization optimized by RGD, which shows that RGD fails to get sparse representations, while shrinkage-thresholding operators succeeded. As we have expected, the hyperbolic 1-norm regularization outperforms the others in **TRLC**. It shows that our Riemannian 1-norm regularization can select the dimension of each representations efficiently. In **TRC**, the linear 1-norm regularization outperforms others around where the sum of the 0-norm is 75, as we have expected. Interestingly, our hyperbolic 1-norm regularization outperforms the linear 1-norm regularizations where the sum of the 0-norm is larger. One possible reason is that the linear 1-norm regularizations tend to be unstable since the STO changes the representations dramatically around the ball boundary. Although comparing the optimization process is not trivial since they optimizes different functions, clarifying the reason of the low performance of the linear regularization would be interesting future work. In **TLC**, the hyperbolic 1-norm works the best as expected.