

# Supplementary materials

## A VISUALIZATION OF THE SSDM

This section gives a visualization of the SSDM. See Figure 1. Recall that the  $d$ -th *sparse hyperplane* (SHP), denoted by  $\Pi_d^{\mathcal{M}}$ , is defined by  $\Pi_d^{\mathcal{M}} := \{p \in \mathcal{M} \mid \langle e_d, \log_o(p) \rangle_o = 0\}$ . This means that  $\Pi_d^{\mathcal{M}} = \exp_o(\text{span}\{e_1, e_2, \dots, e_{d-1}, e_{d+1}, e_{d+2}, \dots, e_D\})$ . For example,  $\Pi_1^{\mathcal{M}} = \Pi_1$  in Figure 1 is the image under the exponential map  $\exp_o$  of a linear subspace spanned by  $e_2$  and  $e_3$ . The  $d$ -th element  $\delta_d^{\mathcal{M}}(p) = \delta_d(p)$  of SSDM measures the signed distance from the  $d$ -th SHP to the point  $p$ . For example,  $\delta_1(p)$  is the signed distance from  $\Pi_1$  to  $p$  in Figure 1.

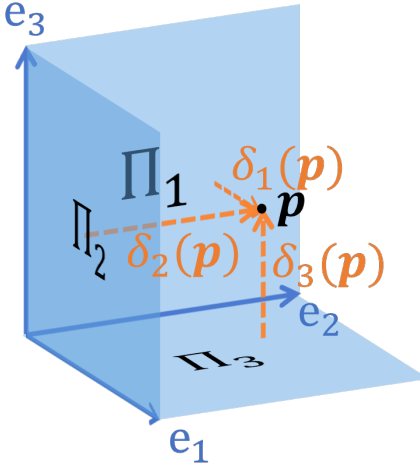


Figure 1: The SSDM’s visualization for a 3-dimensional CHMOO case. The hyperplane  $\Pi_d^{\mathcal{M}} = \Pi_d$  is the  $d$ -th sparse hyperplane (SHP). The  $d$ -th element  $\delta_d^{\mathcal{M}}(p) = \delta_d(p)$  of SSDM measures the signed distance from the  $d$ -th SHP to the point  $p$ .

## B DETAILED EXPLANATION OF EXAMPLE 7

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:

$$p^{(t+1)} \leftarrow p^{(t)} - \alpha \frac{d}{dp}f(p^{(t)}) = \begin{cases} p^{(t)} - \alpha & \text{if } p^{(t)} > 0, \\ p^{(t)} + \alpha & \text{if } p^{(t)} < 0. \end{cases} \quad (6)$$

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 (6) 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  $\alpha$ , where  $n = \left\lfloor \frac{p^{(0)}}{\alpha} \right\rfloor$  is the maximum integer that is no greater than  $\frac{p^{(0)}}{\alpha}$ .

## C PROOF OF THEOREM 2

*Proof.* We prove for the SSDM  $\delta^{(\mathbb{D}^D, \mathcal{G}^p)}$ . 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  $\Pi_d$ . Note that  $\mathbf{h}$  is unique according to Gauss-Bonnet theorem. We have that  $\mathbf{h} = \text{argmin}_{\mathbf{q}} \Delta_{(\mathbb{D}^D, \mathcal{G}^p)}(\mathbf{p}, \mathbf{q})$  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  $\Pi_d$  and passing through  $\mathbf{p}$ . Define  $\mathbf{p}' = \frac{\mathbf{p}}{p^T \mathbf{p}}$  and  $\mathbf{h}' = \frac{\mathbf{h}}{h^T \mathbf{h}}$ . 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  $\Pi_d$  as  $x$ -axis, and the other axis orthogonal to  $\Pi_d$  to  $y$ -axis. We indicate the coordinate of the  $\mathbf{p}$  in the subplane by  $[x \ y]^T$  and that of  $\mathbf{h}$  by  $[h \ 0]^T$ . The coordinates of  $\mathbf{p}'$  and  $\mathbf{h}'$  are  $\frac{1}{x^2+y^2}[x \ y]^T$  and  $[1/h \ 0]^T$ , respectively. See also Figure 2. 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

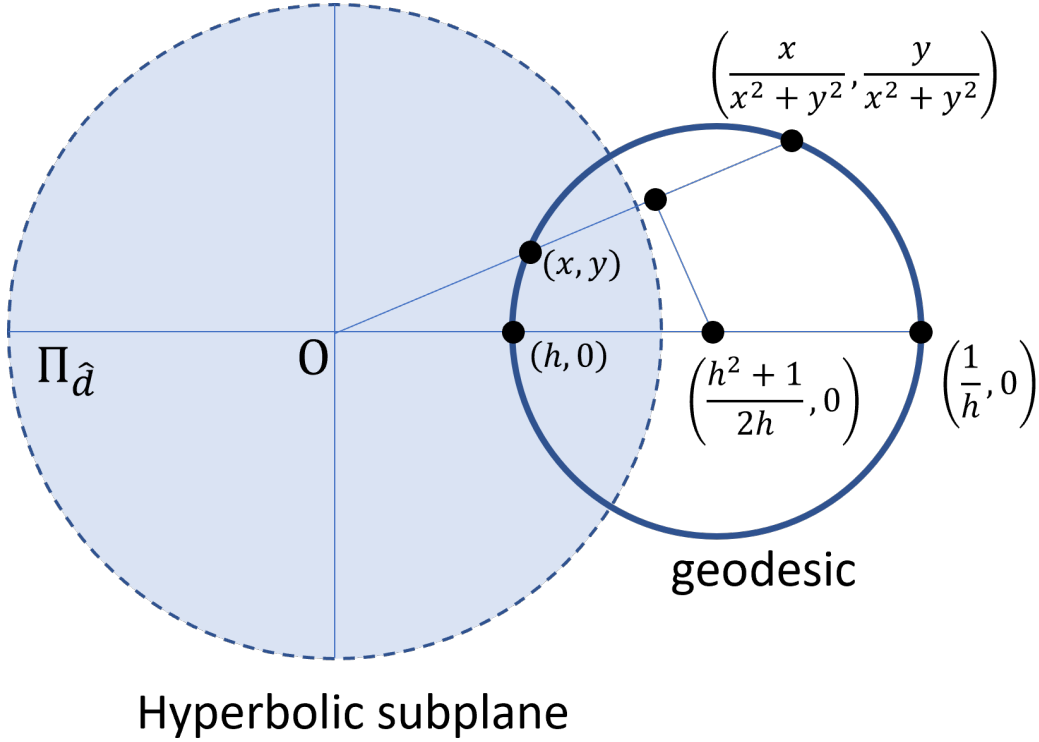


Figure 2: Hyperbolic subdisk.

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 Poincarè model:  $\Delta_{(\mathbb{D}^2, \mathcal{G}^p)}(\mathbf{p}, \mathbf{q}) = \operatorname{acosh} \left( 1 + \frac{2|\mathbf{p}-\mathbf{q}|^2}{(1-|\mathbf{p}|^2)(1-|\mathbf{q}|^2)} \right)$ . Specifically,

$$\begin{aligned}
 & \Delta_{(\mathbb{D}^2, \mathcal{G}^p)}([x \ y], [h \ 0]) \\
 &= \operatorname{acosh} \left( 1 + \frac{2((x-h)^2 + y^2)}{(1-(x^2+y^2))(1-h^2)} \right) \\
 &= \operatorname{acosh} \left( \sqrt{1 + \frac{4y^2}{(1-(x^2+y^2))^2}} \right) \\
 &= \operatorname{asinh} \left( \frac{2y}{(1-(x^2+y^2))^2} \right).
 \end{aligned} \tag{7}$$

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

## D EXPLICIT PSEUDOCODE OF THE HISTA

Algorithm 2 shows the explicit form of HISTA on the Poincaré model. Here,  $\text{sinhc}$  is defined by

$$\text{sinhc}(x) := \begin{cases} 1 & \text{if } x = 0, \\ \frac{\sinh x}{x} & \text{if } x \neq 0. \end{cases} \quad (8)$$

---

### Algorithm 2 HISTA (Explicit form)

---

**Require:**  $\mathbf{p}_{\text{init}} \in \mathbb{D}^D$ : initial point,

$\alpha \in \mathbb{R}_{>0}$ : learning rate,

$T \in \mathbb{Z}_{\geq 0}$ : # iterations.

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

$\mathbf{p}^{(0)} \leftarrow \mathbf{p}_{\text{init}}$

**for**  $t \leftarrow 1, 2, \dots, T$  **do**

$\boldsymbol{\gamma}^{(t)} \leftarrow \partial|_{\mathbf{p}^{(t-1)}} J$

$\rho^{(t)} \leftarrow \frac{4}{(1 - |\mathbf{p}^{(t-1)}|^2)^2}$

$\mathbf{g}^{(t)} \leftarrow (\rho^{(t)})^2 \boldsymbol{\gamma}^{(t)}$

$\mathbf{q}^{(t-1)} \leftarrow \rho^{(t)} \cdot \frac{[\cosh(|-\alpha \mathbf{g}^{(t)}|) - \rho^{(t)} \alpha (\mathbf{g}^{(t)})^\top (\mathbf{p}^{(t-1)})] \mathbf{p}^{(t-1)} + \text{sinhc}(|\alpha \mathbf{g}^{(t)}|) \mathbf{g}^{(t)}}{1 + (\rho^{(t)} - 1) \cosh(|-\alpha \mathbf{g}^{(t)}|) - (\rho^{(t)})^2 \alpha (\mathbf{g}^{(t)})^\top (\mathbf{p}^{(t-1)}) \text{sinhc}(|\alpha \mathbf{g}^{(t)}|)}$

$\boldsymbol{\sigma}^{(t)} \leftarrow \text{asinh}\left(\frac{2\mathbf{p}}{1 - \mathbf{p}^\top \mathbf{p}}\right) - \alpha \lambda \mathbf{1}_D$

$\mathbf{p}^{(t)} \leftarrow \frac{\sinh(\boldsymbol{\sigma}^{(t)})}{\sqrt{1 + (\sinh \boldsymbol{\sigma}^{(t)})^\top (\sinh \boldsymbol{\sigma}^{(t)}) + 1}}$

**end for**

$\mathbf{p}_{\text{output}} \leftarrow \mathbf{p}^{(T)}$

---

## E FORMULAE FOR A PRODUCT MANIFOLD.

We discuss the formulae of our sparse representation learning scheme for product manifolds. We do not give detailed proofs, but they can easily be proved using a specific coordinate space. Let  $\mathcal{M}^{[1]}, \mathcal{M}^{[2]}, \dots, \mathcal{M}^{[M]}$  are Riemannian manifolds. The Riemannian product manifold  $\mathcal{M} = \mathcal{M}^{[1]} \times \mathcal{M}^{[2]} \times \dots \times \mathcal{M}^{[M]}$  is given as the topological product manifold of  $\mathcal{M}^{[1]}, \mathcal{M}^{[2]}, \dots, \mathcal{M}^{[M]}$  equipped with the metric tensor defined as follows. For  $p = (p^{[1]}, p^{[2]}, \dots, p^{[M]}) \in \mathcal{M}$ , the metric tensor and a direct some decomposition of the tangent space  $T_p \mathcal{M}$  is given as follows. For tangent vectors  $v, v' \in T_p \mathcal{M}$ , let  $c : I \rightarrow \mathcal{M}$  and  $c' : I' \rightarrow \mathcal{M}$  be  $C^\infty$  curves on  $\mathcal{M}$  tangent to  $v$  and  $v'$  respectively, where  $I, I' \subset \mathbb{R}$  is open intervals such that  $0 \in I \cap I'$ . Here,  $c$  being tangent to  $u$  means that  $c(0) = p$  and  $v = \dot{c}|_0$ , where  $\dot{c}|_t : C^\infty(\mathcal{M}) \rightarrow \mathbb{R}$  for  $t \in I$  is defined by  $\dot{c}|_{t'} f := \frac{d}{dt}(c(t))|_{t'}$ . There exist  $C^\infty$  curves  $c^{[m]} : I \rightarrow \mathcal{M}$  and  $c'^{[m]} : I' \rightarrow \mathcal{M}$  for  $m = 1, 2, \dots, M$  such that  $c^{[m]}(t) = c'^{[m]}(t) = p^{[m]}$  for  $m = 1, 2, \dots, M$  and  $c(t) = (c^{[1]}(t), c^{[2]}(t), \dots, c^{[M]}(t))$  and  $c'(t) = (c'^{[1]}(t), c'^{[2]}(t), \dots, c'^{[M]}(t))$ . Then we can define  $\tilde{v}^{[m]} := c|_0 \in T_p \mathcal{M}$  and  $\tilde{v}'^{[m]} := c'|_0 \in T_p \mathcal{M}$  for  $m = 1, 2, \dots, M$ . We can prove that it does not depend on the choice of  $c$ . In the following, we denote the linear operation to obtain  $\tilde{v}^{[m]} \in T_{p^{[m]}} \mathcal{M}^{[m]}$  from  $v \in T_p \mathcal{M}$  by  $\pi_p^{[m]} : T_p \mathcal{M} \rightarrow T_{p^{[m]}} \mathcal{M}^{[m]}$ . We define the metric tensor on  $T_p \mathcal{M}$  by

$$\langle v, v' \rangle_p^{\mathcal{M}} := \sum_{m=1}^M \langle \tilde{v}^{[m]}, \tilde{v}'^{[m]} \rangle_p^{\mathcal{M}}. \quad (9)$$

We can prove that the above metric is symmetric and positive-definite.

A tangent vector  $\tilde{v}^{[m]} \in T_{p^{[m]}}\mathcal{M}^{[m]}$  can be identified with  $v^{[m]} \in T_p\mathcal{M}$  that satisfies

$$\begin{cases} \pi_p^{[m']} (v^{[m]}) = \tilde{v}^{[m]} & \text{if } m' = m, \\ \pi_p^{[m']} (v^{[m]}) = 0 & \text{if } m' \neq m. \end{cases} \quad (10)$$

We can prove that the above  $v^{[m]}$  is uniquely defined. By the above identification, we can identify each  $T_{p^{[m]}}\mathcal{M}^{[m]}$  with  $T_p^{[m]}\mathcal{M} := \left(\pi_p^{[m]}\right)^{-1} (T_{p^{[m]}}\mathcal{M}^{[m]}) \subset T_p\mathcal{M}$ , which is a linear subspace of  $T_p\mathcal{M}$  with the same dimension as  $T_{p^{[m]}}\mathcal{M}^{[m]}$ . Also, we can prove that  $T_{p^{[m]}}\mathcal{M}^{[m]} = \bigoplus_{m=1}^M T_p^{[m]}\mathcal{M}$  is the orthogonal direct sum decomposition. Since the restriction  $\pi_p^{[m]}|_{T_{p^{[m]}}\mathcal{M}^{[m]}}$  is one-to-one, we can

define its inverse map  $\left(\pi_p^{[m]}|_{T_{p^{[m]}}\mathcal{M}^{[m]}}\right)^{-1} : T_{p^{[m]}}\mathcal{M}^{[m]} \rightarrow T_p^{[m]}\mathcal{M}$ .

Based on the identification, if we have multiple CHMOOs  $\left(\mathcal{M}^{[m]}, \mathfrak{o}^{[m]}, \left(\tilde{e}_1^{[m]}, \tilde{e}_2^{[m]}, \dots, \tilde{e}_{D^{[m]}}^{[m]}\right)\right)_{m=1}^M$ , we can see that  $\left(e_1^{[1]}, e_2^{[1]}, \dots, e_{D^{[1]}}^{[1]}, e_1^{[2]}, e_2^{[2]}, \dots, e_{D^{[2]}}^{[2]}, \dots, e_1^{[M]}, e_2^{[M]}, \dots, e_{D^{[M]}}^{[M]}\right)$  is an ONB in  $T_p\mathcal{M}$ , where  $e_d^{[m]} = \left(\pi_o^{[m]}|_{T_o^{[m]}\mathcal{M}^{[m]}}\right)^{-1} \left(\tilde{e}_d^{[m]}\right)$  for  $m = 1, 2, \dots, M$  and  $d = 1, 2, \dots, D^{[m]}$ ,  $\mathcal{M} = \mathcal{M}^{[1]} \times \mathcal{M}^{[2]} \times \dots \times \mathcal{M}^{[M]}$ , and  $\mathfrak{o} = (\mathfrak{o}^{[1]}, \mathfrak{o}^{[2]}, \dots, \mathfrak{o}^{[M]}) \in \mathcal{M}$ . Hence, we can define the product CHMOO  $\left(\mathcal{M}, \mathfrak{o}, \left(e_1^{[1]}, e_2^{[1]}, \dots, e_{D^{[1]}}^{[1]}, e_1^{[2]}, e_2^{[2]}, \dots, e_{D^{[2]}}^{[2]}, \dots, e_1^{[M]}, e_2^{[M]}, \dots, e_{D^{[M]}}^{[M]}\right)\right)$ .

We are interested in the SHP, Riemannian 0-norm, SSDM, its inverse, CH 1-norm, and CHSTO of the product CHMOO. In the following, we derive the formula for the SSDM, which enable us to calculate the others. To achieve that, we review the basic property of the product Riemannian manifold.

Suppose that  $p = (p^{[1]}, p^{[2]}, \dots, p^{[M]}) \in \mathcal{M}$  and  $v \in T_p\mathcal{M}$  and we consider the unique decomposition  $v = \sum_m^M v^{[m]}$ , where  $v^{[m]} \in T_p^{[m]}\mathcal{M}$  for  $m = 1, 2, \dots, M$ . Define  $\tilde{v}^{[m]} = \pi_p^{[m]}|_{T_p^{[m]}\mathcal{M}^{[m]}} (v^{[m]})$  for  $m = 1, 2, \dots, M$ . We can see that if we can define  $\exp_{p^{[m]}}(\tilde{v}^{[m]}) \in \mathcal{M}^{[m]}$  for  $m = 1, 2, \dots, M$ , it follows that  $\exp_p(v) = (\exp_{p^{[1]}}(\tilde{v}^{[1]}), \exp_{p^{[2]}}(\tilde{v}^{[2]}), \dots, \exp_{p^{[M]}}(\tilde{v}^{[M]}))$ . We can calculate the logarithmic map similarly.

The distance between two points  $p = (p^{[1]}, p^{[2]}, \dots, p^{[M]})$  and  $q = (q^{[1]}, q^{[2]}, \dots, q^{[M]})$  is given by

$$\Delta_{\mathcal{M}}(p, q) = \sqrt{\sum_{m=1}^M [\Delta_{\mathcal{M}^{[m]}}(p^{[m]}, q^{[m]})]^2}. \quad (11)$$

From the above property, we can confirm that the SHP  $\Pi_{[m],d}^{\mathcal{M}}$  of the product CHMOO corresponding to  $e_d^{[m]} \in T_o^{[m]}\mathcal{M}^{[m]} \subset T_o\mathcal{M}$  is given by

$$\Pi_{[m],d}^{\mathcal{M}} = \left\{ \left( p^{[1]}, p^{[2]}, \dots, p^{[M]} \right) \left| p^{[m']} \in \Pi_d^{\mathcal{M}^{[m]}} \text{ if } m' = m, \quad p^{[m']} \in \mathcal{M}^{[m]} \text{ if } m' \neq m, \right. \right\}. \quad (12)$$

Hence, we immediately get the SSDM formula.

$$\delta^{\mathcal{M}}(p) = \begin{bmatrix} \delta^{\mathcal{M}^{[1]}}(p^{[1]}) \\ \delta^{\mathcal{M}^{[2]}}(p^{[2]}) \\ \vdots \\ \delta^{\mathcal{M}^{[M]}}(p^{[M]}) \end{bmatrix}. \quad (13)$$

The above formula enables us to calculate the SHP, Riemannian 0-norm, SSDM, its inverse, CH 1-norm, and CHSTO for the product CHMOO if we can calculate them for each component CHMOO.

For example, we can calculate them for the product of the EVCHMOO and hyperbolic CHMOOs. Note that the product of EVCHMOOs is again a higher dimensional EVCHMOO, while the product of hyperbolic CHMOOs is not a hyperbolic CHMOO since the product of hyperbolic space is not a hyperbolic space.

## F NUMERICAL EXPERIMENTS: HISTA AVOIDS THE OSCILLATION

Remark 9 states that our motivation in this paper has almost been achieved. The last thing we need to do is to confirm by numerical experiments that our HISTA avoids the oscillation issue. Hence, we compare the HISTA and RGD for sparse solution and non-sparse solution cases.

We consider minimizing the square distance with the hyperbolic 1-norm regularization:  $L(z) = [\Delta_{(\mathbb{D}^2, \mathcal{G}^P)}(z, z')]^2 + \lambda \|z\|_{1, (\mathbb{D}^2, \mathcal{G}^P)}$ . 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 3 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  $z' = [0.0 \ 0.0]^\top$ . See also Figure 4. Our results confirm that the superiority of the HISTA over the RGD in the function value for sparse solution cases. Still, detecting the cause of the bounce back effect by the HISTA would be interesting future work.

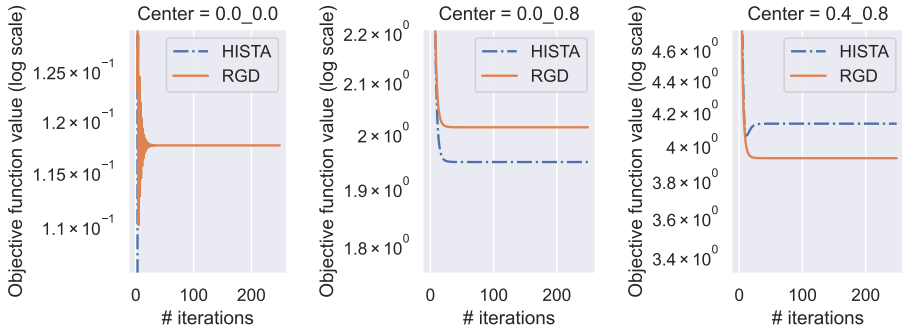


Figure 3: The optimization performances of HISTA and RGD in minimizing the square distance from a fixed point  $z'$  with the hyperbolic 1-norm, where  $z' = [0.0, 0.0]^\top$  (Left),  $z' = [0.0, 0.8]^\top$  (Center), and  $z' = [0.4, 0.8]^\top$  (Right). Note that in (Left), the blue dashed line indicating HISTA goes infinitely downward because the objective function value reaches 0. See Figure 4 for details.

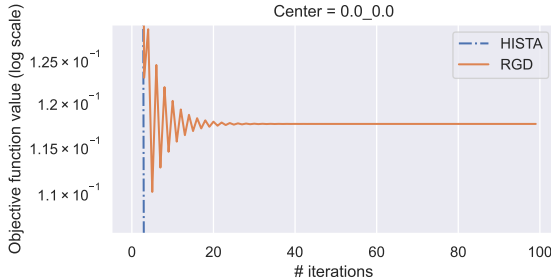


Figure 4: The optimization performances of HISTA and RGD in minimizing the square distance from a fixed point  $z'$  with the hyperbolic 1-norm, where  $z' = [0.0, 0.0]^\top$ . Here we focus on the first 100 iterations to see the oscillation issue of the RGD.

## G NUMERICAL EXPERIMENTS ON GRAPH EMBEDDING SETTING

This section gives numerical-experimental results. Note that the objective is not to achieve state-of-the-art representations, but to show how our sparse learning scheme influences HSBRL.

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  $C(\mathcal{V}, 2)$  be the set of subsets of  $\mathcal{V}$ , whose size is two. That is,  $C(\mathcal{V}, 2)$  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

$$\begin{aligned} y_{u,v} &:= \begin{cases} +1 & \text{if } \{u, v\} \in \mathcal{E}, \\ -1 & \text{if } \{u, v\} \notin \mathcal{E}, \end{cases} \\ w_{u,v} &:= 2 \cdot \frac{|\{\{u', v'\} \in C(\mathcal{V}, 2) \mid y_{u', v'} = y_{u,v}\}|}{|C(\mathcal{V}, 2)|}. \end{aligned} \quad (14)$$

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

$$l_{0-1}(\hat{y}, y) := \begin{cases} 0 & \text{if } \text{sgn}(\hat{y}) = y, \\ +1 & \text{otherwise.} \end{cases} \quad (15)$$

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

$$L_{0-1}((z_v)_{v \in \mathcal{V}}; \mathcal{G}) := \sum_{\{u,v\} \in C(\mathcal{V}, 2)} w_{u,v} l_{0-1} \left( [\Delta_{(\mathbb{D}^D, \mathcal{G}^p)}(z_u, z_v)]^2 - \theta, y_{u,v} \right), \quad (16)$$

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

$$L((z_v)_{v \in \mathcal{V}}; \mathcal{G}) := \sum_{\{u,v\} \in C(\mathcal{V}, 2)} w_{u,v} l \left( [\Delta_{(\mathbb{D}^D, \mathcal{G}^p)}(z_u, z_v)]^2 - \theta, y_{u,v} \right). \quad (17)$$

Also, we add the regularization term  $\lambda \sum_{v \in \mathcal{V}} r(z_v)$  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 method. Note that  $\lambda \mathbb{R}_{>0}$  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:

$$r(z) = \begin{cases} \|z\|_{1, (\mathbb{D}^D, \mathcal{G}^p)} & \text{hyperbolic 1-norm (H 1-norm),} \\ \|z\|_1 & \text{linear 1-norm,} \\ 0 & \text{no regularization.} \end{cases} \quad (18)$$

We use the HISTA for the hyperbolic 1-norm. For the linear-norm, we apply Riemannian gradient descent and traditional shrinkage-thresholding operator. Note that this is also what we propose for a baseline. For the no regularization case, these two are the same. Strictly speaking, we need to regard the problem as the optimization of a function of the product of  $|\mathcal{V}|$  hyperbolic spaces since we consider  $|\mathcal{V}|$  points in hyperbolic space. The rigorous discussion for the product manifold is given in Appendix E. Still, it shows that what we need to do is to calculate a partial derivative for each point, convert it into a Riemannian gradient, and apply the HISTA or RGD for each point, like existing papers do.

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



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

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  $\lambda$  and observe how the accuracy and the 0-norm changes. For the no regularization method, we vary  $D$  instead of  $\lambda$ . 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 CH 1-norm regularization and Linear 1-norm regularization, we experiment on synthetic datasets defined below.

- **TREE-ROOTLEAF CUBES (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-ROOT CUBES (TRC)** is similar to TRLC but without the cube connected to the leaves.

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. Hence, our natural expectation is that CH 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.

To show our sparse learning scheme's behavior in real applications, we conduct the same experiment on the following real datasets.

- **ENRON-EMAIL** is an email network reflecting the hierarchical tree structure of the company. At the same time, it also contains edges corresponding to cross-departmental communications, which might be an omen of Enron's bankruptcy in 2001. Since it is a mixture of a tree-like structure and a non-tree-like structure, we expect that our sparse learning scheme works better than the **no regularization** method in the ENRON-EMAIL.
- **CORA** is a citation network, which shows a highly tree-like structure. Since it is highly tree-like, we expect it to be more effective to apply a (non-sparse) hyperbolic embedding method in low-dimensional space than to apply our sparse learning scheme in high-dimensional space. Hence, we do NOT expect that our sparse learning scheme works so effectively in CORA as in ENRON-EMAIL.

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  $\theta$  is set to 1.0. The number of iterations is set to  $T = 10000$ .

Figure 6 shows how the accuracy and the 0-norm changes by varying  $\lambda$  or  $D$ . The **closer to the left upper corner**, the better. Here, we show the range where the sum of the 0-norms is 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 **H 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 **H 1-norm** regularization outperforms the others in TRLC. It shows that our **H 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

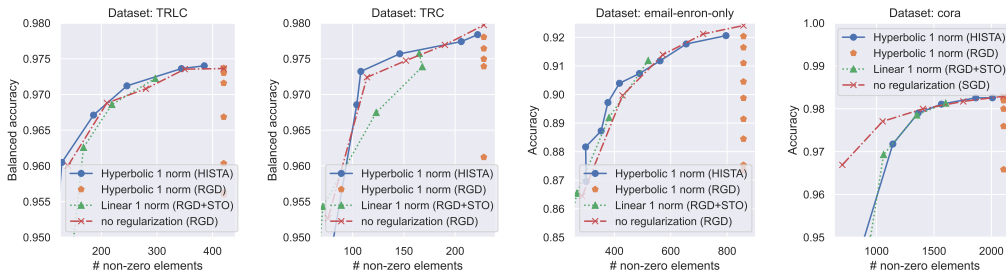


Figure 6: The trade-off between the representation quality (balanced accuracy) and the space complexity (the 0-norm). From left to right: TRLC, TRC, EMAIL-ENRON, **Cora**. The closer to the left upper corner are the graphs, the better.

75, as we have expected. Interestingly, our **H 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 optimize different functions, clarifying the reason for the low performance of the linear regularization would be interesting future work. In ENRON-EMAIL, our **H 1-norm** outperformed other methods both in the balanced accuracy and the sum of 0-norm, when the sum of 0-norm is around 400. While our method’s advantage is clear where the sum of 0-norm is small, no method outperforms the remaining in both balanced accuracy and the sum of 0-norm where the sum of 0-norm is large. The investigation of this phenomenon is interesting future work. In CORA, **no regularization** method outperformed both the regularization methods. The reason is that CORA is highly tree-like, so it is most suitable for low-dimensional hyperbolic space. Hence, the result is consistent with our expectations.

## H ACRONYM TABLE

To increase readability, we provide the table of acronyms used in this paper in Table 1

RL	representation learning
HSBRL	hyperbolic-space-based representation learning
RGD	Riemannian gradient descent
HISTA	hyperbolic iterative shrinkage-thresholding algorithm
RCS	real coordinate space
CHM	Cartan-Hadamard manifold
CHMOO	CHM with an origin and orthonormal bases
ONB	orthonormal basis
EVCHMOO	Euclidean vector CHMOO
SHP	sparse hyperplane
SSDM	signed SHP distance map
CH	Cartan-Hadamard
H	hyperbolic
ISTA	iterative shrinkage-thresholding algorithm
STO	soft-thresholding operator
CHSTO	Cartan-Hadamard STO
CHISTA	Cartan-Hadamard ISTA

Table 1: Acronyms