003 004

010 011

012

013

014

015

016

017

018

019

021

## MIXED-CURVATURE DECISION TREES AND RANDOM FORESTS

Anonymous authors

Paper under double-blind review

### ABSTRACT

Decision trees (DTs) and their random forest (RF) extensions are workhorses of classification and regression in Euclidean spaces. However, algorithms for learning in non-Euclidean spaces are still limited. We extend DT and RF algorithms to product manifolds: Cartesian products of several hyperbolic, hyperspherical, or Euclidean components. Such manifolds handle heterogeneous curvature while still factorizing neatly into simpler components, making them compelling embedding spaces for complex datasets. Our novel angular reformulation of DTs respects the geometry of the product manifold, yielding splits that are geodesically convex, maximum-margin, and composable. In the special cases of singlecomponent manifolds, our method simplifies to its Euclidean or hyperbolic counterparts, or introduces hyperspherical DT algorithms, depending on the curvature. We benchmark our method on various classification, regression, and link prediction tasks on synthetic data, graph embeddings, mixed-curvature variational autoencoder latent spaces, and empirical data. Compared to six other classifiers, product DTs and RFs ranked first on 21 of 22 single-manifold benchmarks and 18 of 35 product manifold benchmarks, and placed in the top 2 on 53 of 57 benchmarks overall. This highlights the value of product DTs and RFs as straightforward yet powerful new tools for data analysis in product manifolds.

027 028 029

030

025

026

### 1 INTRODUCTION

While much of machine learning focuses on Euclidean spaces, these can fail to capture the true
structure of complex datasets. For example, hierarchical structures, which are common in taxonomy
(e.g., phylogenetic trees) are better represented in hyperbolic space due to its exponential volume
growth, which naturally mirrors tree-like data (Sonthalia & Gilbert, 2020). Similarly, cyclical structures, often encountered in time-series data with periodic patterns (e.g., seasonal trends, neuronal
spiking dynamics), can benefit from spherical representations (Ding & Regev, 2021).

However, many real-world datasets don't conform to a single geometric structure. Any constantcurvature manifold—whether hyperbolic, spherical, or Euclidean—would struggle to represent all the nuances of such data simultaneously. Product manifolds, as proposed by Gu et al. (2018), offer
a solution. By combining multiple constant-curvature component manifolds (spherical, Euclidean, and hyperbolic spaces) into a single product manifold, they can better capture the complexity of such mixed-structure data. This flexibility reduces distortion when modeling pairwise distances and enables a more accurate representation of the underlying data structure.

Despite their advantages, product manifolds have seen limited adoption in machine learning, particularly for inference tasks like classification and regression. Existing work has primarily focused on applications in biology (McNeela et al., 2024) and knowledge graphs (Wang et al., 2021). However, tools for leveraging product manifold representations in downstream tasks remain scarce.

In this paper, we introduce mixed-curvature decision trees (DTs) and random forests (RFs), expanding the toolkit for analyzing product manifold data. By enabling inference directly on product manifold coordinates, our approach is well-suited for datasets that combine hierarchical, cyclical, and other complex geometric patterns. This framework provides a principled way to learn from such structures, achieving more accurate results than competing models. These contributions offer new possibilities for applying product manifold representations in fields ranging from biological modeling to temporal-spatial analysis.



Figure 1: An illustration of the product manifold DT in action. We consider a sample of labeled points  $(\mathbf{X}, \mathbf{y})$  from one of the simplest possible product manifolds: the torus  $\mathcal{P} = \mathbb{S}^1 \times \mathbb{S}^1$ . Since we know the signature for  $\mathcal{P}$ , we can factorize  $\mathbf{X}$  into coordinates on two circles. Our DT splits these factorized coordinates to a maximum depth of 2, partitioning  $\mathcal{P}$  into a total of  $2^2 = 4$  disjoint decision areas (colored **positive** or **negative** to reflect the classes).

### Our contributions:

072

073

074

075

076 077 078

079

081

082

084

085

087

090 091

092

- 1. We generalize DTs and RFs to *all* constant-curvature manifolds. Unlike existing methods, we represent data and splits as angles in two-dimensional subspaces. This guarantees splits are geodesically convex, maximum-margin, and composable. In the single-manifold case, this extends existing Euclidean and hyperbolic models or introduces *hyperspherical* DTs and RFs.
  - 2. We introduce novel DT and RF algorithms for product manifolds.
  - 3. We extend techniques for sampling distributions in non-Euclidean manifolds to describe mixtures of Gaussians in product manifolds.
  - 4. We show how problems like link prediction in graphs and signal analysis can be recast as inference problems on product manifolds.
  - 5. We demonstrate the effectiveness of our component- and product-manifold algorithms over competing algorithms on a suite of 57 diverse non-Euclidean benchmarks.

### 1.1 Related work

Non-Euclidean representation learning. Important background on manifolds in machine learning is given in Cayton (2005) and Bengio et al. (2014). Much of the work on product manifolds is indebted to early works on hyperbolic spaces, including Nickel & Kiela (2017); Chamberlain et al. (2017), and Ganea et al. (2018).

Machine learning in product manifolds. Tabaghi et al. (2021) describe linear classifiers, including perceptron and support vector machines; Tabaghi et al. (2024) adapt principal component analysis; and Cho et al. (2023) generalize Transformer architectures to product manifolds.

Computationally tractable manifolds. Besides product manifolds, other methods for representing data with heterogeneous curvature also exist: Borde & Kratsios (2023) is based on fractals, while Cruceru et al. (2020) is based on matrix manifolds.

Product manifold-derived features. Tsagkrasoulis & Montana (2017) train RF classifiers on distance matrices from arbitrary manifolds, e.g. product manifolds. Sun et al. (2021) and Borde et al. (2023b) use product manifolds to compute rich similarity measures as features for classification. Giovanni et al. (2022) introduce a heterogeneous variant of product manifolds; Borde et al. (2024) combine quasi-metrics and partial orders in a product manifold for graph representations.



Figure 2: Decision boundaries in any constant-curvature manifold are found by 2-D projections into 2-dimensional subspaces. In this perspective, both data and splits are parameterized by an angle  $\theta$ . Each **split** divides the manifold into **positive-class regions** and **negative-class regions**.

Hyperbolic random forests. Our method is inspired by recent work by Doorenbos et al. (2023)
and Chlenski et al. (2024) extending RFs to hyperbolic space. In particular, our angular-split perspective synthesizes the ideas in Chlenski et al. (2024) and Tabaghi et al. (2021).

Applications of product manifolds. Product manifolds are popular for embedding knowledge graphs Wang et al. (2021); Li et al. (2024); Nguyen-Van et al. (2023). In biology, they have been used to represent pathway graphs (McNeela et al., 2024), cryo-EM images (Zhang et al., 2021), and single-cell transcriptomic profiles (Tabaghi et al., 2021). Skopek et al. (2020) also embed image datasets into product manifolds.

130 131

132 133

134

135

136 137

### 2 PRELIMINARIES

We review relevant details of different Riemannian manifolds (Euclidean spaces, hyperspheres, hyperboloids, and product manifolds), along with key properties of the Euclidean and hyperbolic variants of DTs and RFs.

## 137 2.1 RIEMANNIAN MANIFOLDS138

We will begin by reviewing key details of hyperspheres, hyperboloids, and Euclidean spaces. Formore details, readers can consult Do Carmo (1992).

Each space described is a Riemannian manifold, meaning that it is locally isomorphic to Euclidean
 space and equipped with a distance metric. The shortest paths between two points u and v on a
 manifold are called geodesics. As all three spaces we consider have constant Gaussian curvature,
 we define simple closed forms for geodesic distances in each of the following subsections in lieu of
 a more general discussion of geodesic distances in arbitrary Riemannian manifolds.

Any constant-curvature manifold  $\mathcal{M}$  is parameterized by a dimensionality D and a curvature K. They can also all be considered embedded in an ambient space  $\mathbb{R}^{D+1}$ . Finally, for each point  $\mathbf{x} \in \mathcal{M}$ , the tangent plane at  $\mathbf{x}, T_{\mathbf{x}}\mathcal{M}$ , is the space of all tangent vectors at  $\mathbf{x}$ :

$$T_{\mathbf{x}}\mathcal{M} = \{ \mathbf{x}' \in \mathcal{M} : \langle x', x \rangle_{\mathcal{M}} = 0 \}.$$
<sup>(1)</sup>

150 151

154

159

### 2.1.1 EUCLIDEAN SPACE

Euclidean spaces are naturally understood as  $\mathbb{R}^D$ , but we will use the notation  $\mathbb{E}^D = \mathbb{R}^D$  when treating Euclidean spaces as manifolds. In contrast, we will continue to use  $\mathbb{R}^D$  to refer to ambient spaces. Euclidean spaces use the familiar inner product (dot product), norm ( $\ell_2$  norm), and distance function (Euclidean distance):

$$\langle \mathbf{u}, \mathbf{v} \rangle = u_0 v_0 + u_1 v_1 + \ldots + u_2 v_2, \tag{2}$$

$$\|\mathbf{u}\| = \sqrt{\langle \mathbf{u}, \mathbf{u} \rangle},\tag{3}$$

$$\delta_{\mathbb{R}}(\mathbf{u}, \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\|. \tag{4}$$

## 162 2.1.2 HYPERSPHERICAL SPACE

Hyperspheres can be viewed as surfaces *embedded* in a higher-dimensional, Euclidean ambient space. Hyperspherical space uses the same inner products as Euclidean space. The hypersphere is the set of points in the ambient space having a Euclidean norm equal to some radius inversely proportional to the curvature K > 0:

$$\mathbb{S}^{D,K} = \{ \mathbf{x} \in \mathbb{R}^{D+1} : \|x\| = 1/K \}.$$
(5)

Because shortest paths between two points  $\mathbf{u}$  and  $\mathbf{v}$  in  $\mathbb{S}^{D,K}$  through the ambient space leave the surface of the manifold, we must define the hyperspherical distance function for the shortest path entirely in  $\mathbb{S}^{D,K}$  between  $\mathbf{u}$  and  $\mathbf{v}$ :

183

188 189

194

197

203

213 214

168

$$\delta_{\mathbb{S}}(\mathbf{u}, \mathbf{v}) = \cos^{-1}(K^2 \langle \mathbf{u}, \mathbf{v} \rangle) / K.$$
(6)

# 175 2.1.3 HYPERBOLIC SPACE

177 Hyperbolic space is characterized by constant negative metric curvature. This has several conse-178 quences: for instance, the angles in any triangle sum to less than  $\pi$ , many lines through a point can 179 be parallel to any given line, and neighborhoods grow exponentially with radius.

There are several equivalent models of hyperbolic space. For our purposes, we will describe the
hyperbolic space from the perspective of the hyperboloid model. First, we must define the ambient
Minkowski space. This is a vector space equipped with the Minkowski inner product:

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}} = -u_0 v_0 + u_1 v_1 + \ldots + u_n v_n. \tag{7}$$

Similar to the Euclidean case, we let  $\|\mathbf{u}\|_{\mathcal{L}} = \langle \mathbf{u}, \mathbf{u} \rangle_{\mathcal{L}}$  (we do not wish to take the square root of a negative number). The hyperboloid of dimension D and curvature K < 0, written  $\mathbb{H}^{D,K}$ , is a set of points with constant Minkowski norm:

$$\mathbb{H}^{D,K} = \{ \mathbf{x} \in \mathbb{R}^D : \|\mathbf{x}\|_{\mathcal{L}} = -1/K^2, \ x_0 > 0 \},$$
(8)

Finally, the hyperbolic distance function for geodesic distances between  $\mathbf{u}, \mathbf{v} \in \mathbb{H}^{D,K}$  is given by

$$\delta_{\mathbb{H}}(\mathbf{u}, \mathbf{v}) = -\cosh^{-1}(K^2 \langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}}) / K.$$
(9)

#### 2.1.4 MIXED-CURVATURE PRODUCT MANIFOLDS

<sup>195</sup> We reiterate the definition of product manifolds from Gu et al. (2018). A product manifold  $\mathcal{P}$  is the <sup>196</sup> Cartesian product of one or more spherical, Euclidean, and hyperbolic manifolds:

$$\mathcal{P} = \mathbb{S}^{s_1, K_1} \times \mathbb{S}^{s_2, K_2} \times \dots \times \mathbb{S}^{s_n, K_n} \times \mathbb{H}^{h_1, K_1} \times \dots \times \mathbb{H}^{h_m, K_m} \times \mathbb{R}^d$$
(10)

The total number of dimensions is  $\sum_{i}^{n} s_{i} + \sum_{j}^{m} h_{j} + d$ . Each individual manifold is called a component manifold, and the decomposition of the product manifold into component manifolds is called the signature. Informally, the signature can be considered a list of dimensionalities and curvatures for each component manifold.

Distances in  $\mathcal{P}$  decompose as the  $\ell_2$  norm of the distances in each of the component manifolds:

$$\delta_{\mathcal{P}}(\mathbf{u}, \mathbf{v}) = \sqrt{\sum_{\mathcal{M} \in \mathcal{P}} \delta_{\mathcal{M}}(\mathbf{u}_{\mathcal{M}}, \mathbf{v}_{\mathcal{M}})^2},\tag{11}$$

where  $\mathbf{u}_{\mathcal{M}}$  and  $\mathbf{v}_{\mathcal{M}}$  denotes the restriction of  $\mathbf{u}$  and  $\mathbf{v}$  to their components in  $\mathcal{M}$  and  $\delta_{\mathcal{M}}$  refers the distance function appropriate to  $\mathcal{M}$ .

For  $\mathbf{x} \in \mathcal{P}$ , the tangent plane at  $\mathbf{x}$ ,  $T_{\mathbf{x}}\mathcal{P}$ , is the concatenation (denoted by the direct sum  $\bigoplus$ ) of all component tangent planes:

$$T_{\mathbf{x}}\mathcal{P} = \bigoplus_{\mathcal{M} \in \mathcal{P}} T_{\mathbf{x}_{\mathcal{M}}}\mathcal{M}.$$
 (12)

We additionally define the origin of  $\mathcal{P}$ ,  $\mu_0$ , as the concatenation of the origins of each respective manifold. The origin is (1/|K|, 0, ...) for  $\mathbb{H}^{D,K}$  and  $\mathbb{S}^{D,K}$ , and (0, 0, ...) for  $\mathbb{E}^D$ .

# 216 2.2 DECISION TREES AND RANDOM FORESTS

The Classification and Regression Trees (CART) (Breiman, 2017) algorithm fits a DT  $\mathcal{T}$  to a set of labeled data (X, y). Specifically, it greedily selects a split at each set to partition the dataset in such a way as to maximize the information gain,

221

222 223

229 230

$$\operatorname{IG}(\mathbf{y}) = C(\mathbf{y}) - \frac{|\mathbf{y}^+|}{|\mathbf{y}|}C(\mathbf{y}^+) - \frac{|\mathbf{y}^-|}{|\mathbf{y}|}C(\mathbf{y}^-).$$
(13)

In this case,  $C(\cdot)$  is some sort of impurity function (we use Gini impurity for classification and variance for regression). Some splitting function  $S(\cdot)$  is used to partition the *labels* y into two classes,  $y^+$  and  $y^-$ ; however,  $S(\cdot)$  also partitions the *input space* (corresponding to some X that does not appear in Eq. 13) into decision regions. Classically,  $S(\cdot)$  is a thresholding function and thus breaks the input space into high-dimensional boxes given some dimension d and threshold  $\theta$ :

$$S(\mathbf{x}) = \mathbb{I}\{x_d > \theta\}.$$
(14)

This algorithm is applied recursively to each decision region until a stopping condition is met (e.g., maximum number of splits is reached). The result is a fitted DT,  $\mathcal{T}$ , which can be used for inference. During inference, an unseen point x is passed through  $\mathcal{T}$  until it reaches a leaf node corresponding to some decision region. For classification, the point is then assigned the majority label inside that region; for regression, it is assigned the mean value inside that region.

Finally, a RF is an ensemble of DTs, typically trained on a bootstrapped subsample of the points and features in X (Breiman, 2001).

238 239 240

### 2.2.1 HYPERBOLIC DECISION TREE ALGORITHMS

The hyperplane perspective on DTs is helpful background for understanding our method: mathematically, thresholding x on a dimension is equivalent to taking its dot product with the normal vector of a separating hyperplane  $\mathbb{P}$ , even in hyperbolic space. Although this is easy to compute for classical thresholding boundaries, which are zero in all dimensions but *d*, this perspective principally admits *any hyperplane*  $\mathbb{P}$  as a valid decision boundary.

Naturally, considerations around choosing an appropriate (and computationally efficient)  $\mathbb{P}$  abound. 246 To this end, Chlenski et al. (2024) impose homogeneity and sparsity constraints on the hyperplanes 247 they consider for hyperbolic DTs. In hyperbolic space, homogenous hyperplanes— hyperplanes 248 that contain the origin of the ambient space—intersect  $\mathbb{H}^{D,K}$  at geodesic submanifolds: that is, 249  $\mathbb{P} \cap \mathbb{H}^{D,K}$  is closed under shortest paths *according to*  $\delta_{\mathbb{H}}$ . The sparsity constraint enforces that the 250 normal vectors of  $\mathbb{P}$  must be nonzero only in two positions: the timelike coordinate  $x_0$  and some 251 other  $x_d$ , which ensures that only  $\mathcal{O}(nd)$  candidate hyperplanes are considered per split, and each 252 decision can be computed in  $\mathcal{O}(1)$  time using sparse dot products. 253

254 255

256

257

258

259

260

261

262

263 264 265

266 267 268

269

## 3 MIXED-CURVATURE DECISION TREES

For any DT, we must transform the input **X** into a set of candidate hyperplanes. To this end, we reframe and generalize the hyperplane approach of hyperbolic DTs. First, we observe that homogenous hyperplanes are geodesically convex in *any constant-curvature manifold*; therefore, we can extend the hyperbolic DT approach to  $\mathbb{E}$  and  $\mathbb{S}$ . Second, we observe that fitting sparse, homogenous DTs is equivalent to thresholding on angles under 2-dimensional projections.

We consider the set of all projections onto the basis  $\{x_0, x_d\}$ , which can be computed in  $\mathcal{O}(1)$  time per projection by coordinate selection. First, we compute the angles in each projection:<sup>1</sup>

$$\theta(\mathbf{x},d) = \tan^{-1}(x_0/x_d). \tag{15}$$

Next, we use a modified splitting criterion to account for the geometry of angular splits:

$$S(\mathbf{x}, d, \theta) = \mathbb{I}\{\theta(\mathbf{x}, d) \in [\theta, \theta + \pi)\}.$$
(16)

<sup>&</sup>lt;sup>1</sup>Note that, in our implementation, we use the PyTorch arctan2 function to ensure that we can recover the full range of angles in  $[0, 2\pi)$ . This is essential for properly specifying decision boundaries in S.

270 Once the best angle is selected, we must compute angular midpoints to select  $\mathbb{P}$  that intersects  $\mathcal{M}$  at 271 a point *geodesically equidistant* the two points to either side of it (Euclidean DTs do this by sampling 272 averaging the threshold values). Angular midpoints for each component manifold are described in 273 the following sections and summarized in Table 4 in the Appendix.

With the angular features and manifold-informed midpoint modifications in place, the rest of the algorithm follows Section 2.2 unmodified.

#### 277 278 3.1 EUCLIDEAN DECISION TREES

While the intersections of homogenous hyperplanes in  $\mathbb{R}^D$  with  $\mathbb{E}^D$  are (trivially) convex, these lack the expressiveness of an ambient-space formulation. Thus, we embed  $\mathbb{E}^D$  in  $\mathbb{R}^{D+1}$  by a trivial lift:

$$\phi : \mathbb{E}^D \to \mathbb{R}^{D+1}, \ \phi(\mathbf{u}) = (1, \mathbf{u}). \tag{17}$$

For two points  $\mathbf{u}, \mathbf{v} \in \mathbb{E}^D$ , the midpoint angles in  $\mathbb{E}^D$  can be described in terms of the coordinates of  $\mathbf{u}$  and  $\mathbf{v}$  or their respective projection angles  $(\theta_u, \theta_v)$  as

$$m_{\mathbb{E}}(\mathbf{u}, \mathbf{v}) = \tan^{-1}(2/(u_d + v_d))$$
(18)

$$= \tan^{-1} \left( \frac{\tan^{-1}(\theta_{\mathbf{u}}) \tan^{-1}(\theta_{\mathbf{v}})}{\tan^{-1}(\theta_{\mathbf{u}}) + \tan^{-1}(\theta_{\mathbf{v}})} \right).$$
(19)

While this presentation of Euclidean DTs is unconventional, it is completely equivalent to thresholding in the basis dimensions. See Appendix C for the proof.

#### 3.2 HYPERBOLIC DECISION TREES

For two points  $\mathbf{u}, \mathbf{v} \in \mathbb{H}^{D,K}$ , we compute  $\theta_{\mathbf{u}}$  and  $\theta_{\mathbf{v}}$  according to Eq 15 and follow Chlenski et al. (2024) in computing the hyperbolic midpoint angle in  $\mathbb{H}^{D,K}$  as:

$$V := \frac{\sin(2\theta_{\mathbf{u}} - 2\theta_{\mathbf{v}})}{\sin(\theta_{\mathbf{u}} + \theta_{\mathbf{v}})\sin(\theta_{\mathbf{v}} - \theta_{\mathbf{u}})},\tag{20}$$

$$m_{\mathbb{H}}(\mathbf{u}, \mathbf{v}) = \begin{cases} \cot^{-1}(V - \sqrt{V^2 - 1}) & \text{if } \theta_{\mathbf{u}} + \theta_{\mathbf{v}} < \pi\\ \cot^{-1}(V + \sqrt{V^2 - 1}) & \text{otherwise.} \end{cases}$$
(21)

300 301 302

303

308

309

310 311 312

313

279

280

281 282

283 284

290

291 292

293 294

#### 3.3 HYPERSPHERICAL DECISION TREES

The hyperspherical case is quite simple, except that unlike hyperbolic space and the "lifted" Euclidean space after applying Eq 17, we lack a natural choice of  $x_0$ . We adopt the convention of fixing the first dimension of the embedding space as  $x_0$ , which intuitively corresponds to fixing a "north pole" at the origin  $\mu_0 = (1/|K|, 0, ...)$ .

Angular midpoints are particularly well-behaved in hyperspherical manifolds: given  $\mathbf{u}, \mathbf{v} \in \mathbb{S}^{D,K}$ , the hyperspherical midpoint angle by finding  $\theta_{\mathbf{u}}$  and  $\theta_{\mathbf{v}}$  using Eq 15 and taking their mean:

$$m_{\mathbb{S}}(\mathbf{u}, \mathbf{v}) = (\theta_{\mathbf{u}} + \theta_{\mathbf{v}})/2.$$
(22)

#### 3.4 PRODUCT DECISION TREE ALGORITHM

Intuitively, the transition from DTs in a single component manifold to a product manifold is that we
 now iterate over all preprocessed angles together, using the angular midpoint formula appropriate to
 each component. The complete pseudocode for this algorithm is in Appendix B.

Allowing for a single DT to span all components—as opposed to, e.g., an ensemble of DTs, each operating in a single component,—allows the model to independently allocate its splits across components according to their relevance to the task at hand. Recasting DT learning in terms of angular comparisons has three major advantages over finding planar decision boundaries directly:

1. We can consider angles under *arbitrary* linear projections (not just projections onto basis dimensions) while maintaining  $\mathcal{O}(1)$  decision complexity. For instance, we can easily search over all  $\binom{D}{2}$  2-dimensional projections if we wish.

- 2. As there is no longer any need to enforce the constraints in Equations 5 and 8 at inference time, it becomes possible to subsample the features (precomputed angles) in RFs.
  - 3. Product manifolds can always represent additional features in a new Euclidean manifold. For instance, this can be useful for incorporating metadata into DT training.

### 4 BENCHMARKS

We carried out benchmarks to evaluate which model, given a labeled set of mixed-curvature embeddings, achieves the lowest validation error. While we produced embeddings using a range of datasets and embedding techniques, our results focus only on performance on downstream tasks. We describe our data generation/embedding methods in more detail in the Appendix.

We summarize our benchmark results, with references to specific figures and tables, in Table 1. Our full results can be found in Table 5 in the Appendix.

Table 1: Benchmarks summary. "#Top-k" columns count how often product DTs or RFs were among the top k predictors for a given set of benchmarks.

| Manifold type    | Task           | Reference | #Top-1    | #Top-2    | Total |
|------------------|----------------|-----------|-----------|-----------|-------|
| Single-curvature | Classification | Figure 3  | 10 (91%)  | 11 (100%) | 11    |
| Single-curvature | Regression     | Figure 4  | 11 (100%) | 11 (100%) | 11    |
| Product manifold | Classification | Table 2   | 11 (46%)  | 22 (92%)  | 24    |
| Product manifold | Regression     | Table 3   | 7 (64%)   | 9 (82%)   | 11    |
| Total            |                |           | 39 (68%)  | 53 (93%)  | 57    |

#### 348 349 350

351

361 362

363

324

325

326

327

328

330 331

332

333

334

335

336

337 338

339

### 4.1 EXPERIMENT DETAILS

**Problem setup.** Given a dataset  $\mathbf{X}$ , a set of labels  $\mathbf{y}$ , and a product manifold  $\mathcal{P}$ , we evaluate a variety of classifiers on their ability to predict  $\mathbf{y}$  from  $\mathbf{X}$ . We apply an identical 80:20 train-test split to all of our data, train our models on the training set, and evaluate performance on the test set.

**Results reporting.** We report 95% confidence intervals for micro-averaged  $F_1$  scores for classification and root mean squared error (RMSE) for regression benchmarks. Pairwise statistical significance is determined by the Wilcoxon signed-rank test comparing all same-type classifiers (i.e. trees to trees and forests to forests). We also apply a Bonferroni correction: starting with a critical value of .05, we divide by the total number of comparisons carried out *for a given signature*: since we compare 5 different models, our critical value becomes .05/10 = .005.

### 4.2 DATASETS

**Synthetic data.** We develop a novel method to sample mixtures of Gaussians in  $\mathcal{P}$  to generate classification and regression datasets. For classification, we generate 8 classes using 32 clusters. For regression, we generate a single scalar response variable using 32 clusters with randomly-generated intercepts. Our full method is described in Appendix Section A

368 Graph embeddings. For classification and regression on graph datasets, we generate embeddings that approximate shortest-path distances in the graph using the method de-369 scribed in Gu et al. (2018). We select the optimal signature from the candidate set 370  $\{(\mathbb{H}^2)^2, \mathbb{H}^2\mathbb{E}^2, \mathbb{H}^2\mathbb{S}^2, \mathbb{S}^2\mathbb{E}^2, (\mathbb{S}^2)^2, \mathbb{H}^4, \mathbb{E}^4, \mathbb{S}^4\}$  by generating embeddings in each signature and se-371 lecting the signature with the lowest metric distortion. For link prediction, we embed all datasets in 372  $\mathcal{P} = (\mathbb{S}^2 \mathbb{E}^2 \mathbb{H}^2)$ , then create a binary classification dataset by associating each pair of nodes with a 373 point in  $\mathcal{P}^2\mathbb{E}^1$ , where each pair of points is included and the last Euclidean dimension is the man-374 ifold distance  $\delta_{\mathcal{P}}(\mathbf{x_i}, \mathbf{x_i})$ ; labels are simply whether there is an edge between nodes i and j. Full 375 details on graph embeddings are described in Appendix Section E.2. 376

377 **Mixed-curvature VAE latent space.** We follow Skopek et al. (2020) in training variational autoencoders (VAEs) whose latent space is  $\mathcal{P}$ . Once the VAE is trained, we use its encoder to generate



Figure 3: Classification benchmark comparison of DTs (top) and RFs (bottom). We report microaveraged  $F_1$  scores on a synthetic data classification task involving mixtures of 8 Gaussians in manifolds of varying constant curvatures K. We compare DTs and RFs in the **product manifold**, **the ambient space**, and **the tangent plane**, along with *k*-nearest neighbors on distances in  $\mathcal{P}$ . Statistical significance (Bonferroni-corrected p < 0.05) is marked with an asterisk (\*). We omit product space perceptrons, which never achieved competitive results.



Figure 4: Regression benchmarks (RMSE) for single-curvature manifolds. We follow the conventions of Figure 3 and mark Bonferroni-corrected significance with an asterisk (\*).

embeddings for our dataset and classify these embeddings. Full details on VAE training and down-stream inference are described in Appendix Section E.3.

**Empirical datasets.** Some datasets can be represented in a non-Euclidean geometry without generating embeddings: for instance, geospatial data lives in  $\mathbb{S}^2$ , while cyclic time series embed in  $\mathbb{S}^1$ . We describe our approach to generating embeddings for these empirical datasets in Appendix Section E.4.

| 432 | Table 2: $F_1$ scores for all classification and link prediction (LP) benchmarks. Best predic-            |
|-----|---|
| 433 | tors are shown in <b>bold</b> , while second-best predictors are <u>underlined</u> . For brevity, we omit |
| 434 | columns for low-performing methods and merge DT and RF columns (e.g. "Ambient" means                      |
| 435 | max(mean(Ambient DT), mean(Ambient RF)).)   |

| - |      | Dataset          | Signature  | k-NN           | Ambient        | Product          |
|---|------|------------------|--|----------------|----------------|------------------|
| - |      | Gaussian         | $(S^2)^2$  | 35.7±.5        | 32.3±.5        | 33.1±.5          |
|   | -K   |                  | $\mathbb{E}^4$   | 34.9±.5        | 30.0±.4        | 31.3±.4          |
|   | ilti |                  | $\mathbb{H}^2\mathbb{E}^2$                             | 40.3±.4        | 34.7±.5        | 36.1±.5          |
|   | Ĩ    |                  | $\mathbb{H}^2 \mathbb{S}^2$                            | 38.4±.5        | 33.3±.5        | <u>35.2±.5</u>   |
|   | ic ( |                  | $\mathbb{H}^4$   | 47.5±.5        | 35.9±.4        | <u>40.0±.4</u>   |
|   | het  |                  | $\mathbb{S}^2 \mathbb{E}^2$                            | 35.8±.5        | 32.7±.5        | $33.5 \pm .4$    |
|   | ynt  |                  | $\mathbb{S}^4$   | 33.1±.4        | 27.6±.4        | <u>28.0±.5</u>   |
|   | S.   |                  | $(\mathbb{H}^2)^2$                                     | 41.5±.5        | 34.5±.5        | <u>37.0±.5</u>   |
|   |      | CiteSeer         | $\mathbb{H}^2\mathbb{S}^2$                             | <u>25.9±.5</u> | 26.1±.7        | 25.8±.6          |
|   | So   | Cora             | $\mathbb{H}^4$   | $20.7 \pm .4$  | $28.9 \pm .5$  | $28.9 \pm .4$    |
|   | lin  | PolBlogs         | $(\mathbb{S}^2)^2$                                     | 93.5±.4        | <u>93.2±.5</u> | 92.9±.4          |
|   | edc  | AdjNoun (LP)     | $(\mathbb{S}^2\mathbb{E}^2\mathbb{H}^2)^2\mathbb{E}^1$ | 93.3±1.1       | 93.7±1.1       | 93.7±1.1         |
|   | nb   | Dolphins (LP)    | $(\mathbb{S}^2\mathbb{E}^2\mathbb{H}^2)^2\mathbb{E}^1$ | 96.6±.3        | 92.3±.9        | 96.6±.3          |
|   | 1ei  | Football (LP)    | $(\mathbb{S}^2\mathbb{E}^2\mathbb{H}^2)^2\mathbb{E}^1$ | 79.8±3.3       | 85.7±3.6       | 85.7±3.6         |
|   | apł  | Karate Club (LP) | $(\mathbb{S}^2\mathbb{E}^2\mathbb{H}^2)^2\mathbb{E}^1$ | 95.1±1.5       | 88.6±2.2       | 95.1±1.5         |
|   | G    | Les Mis (LP)     | $(\mathbb{S}^2\mathbb{E}^2\mathbb{H}^2)^2\mathbb{E}^1$ | 95.7±.7        | 92.7±.9        | <u>95.6±.8</u>   |
|   |      | PolBooks (LP)    | $(\mathbb{S}^2\mathbb{E}^2\mathbb{H}^2)^2\mathbb{E}^1$ | 95.8±.4        | 92.9±.6        | 95.8±.4          |
| - |      | Blood            | $\mathbb{S}^2\mathbb{E}^2(\mathbb{H}^2)^3$             | 17.4±.5        | 19.3±.5        | 20.1±.5          |
|   | ш    | CIFAR-100        | $(\mathbb{S}^2)^4$                                     | 8.6±.4         | <u>11.5±.5</u> | $12.0 \pm .3$    |
|   | /AJ  | Lymphoma         | $(\mathbb{S}^2)^2$                                     | 77.8±1.4       | 81.7±1.2       | 83.7±1.2         |
|   | -    | MNIST            | $\mathbb{S}^2 \mathbb{E}^2 \mathbb{H}^2$               | 41.9±3.7       | 35.7±2.8       | <u>39.4±2.3</u>  |
| - | r    | Landmasses       | $\mathbb{S}^2$   | 91.4±.2        | 83.5±.3        | <u>84.2±.3</u>   |
|   | the  | Neuron 33        | $(S^{1})^{5}$  | $50.5 \pm .5$  | $76.2 \pm .4$  | 77 <b>.0±.</b> 4 |
| _ | 0    | Neuron 46        | $(S^1)^5$  | 50.2±.2        | <u>61.1±.3</u> | 61.2±.3          |

Table 3: Regression results (RMSE) for all benchmarks. We follow the conventions of Table 2. CS PhDs is a graph embedding dataset, whereas Temperature and Traffic are empirical.

| Dataset                      | Signature                      | k-Neighbors      | Ambient          | Product          |
|------------------------------|--------------------------------|------------------|------------------|------------------|
| Synthetic (multi- <i>K</i> ) | $(S^2)^2$                      | .196±.002        | .191±.002        | .191±.002        |
| -                            | $\mathbb{E}^4$                 | .194±.003        | .191±.002        | .190±.002        |
|                              | $\mathbb{H}^2\mathbb{E}^2$     | .196±.003        | .194±.002        | .193±.002        |
|                              | $\mathbb{H}^2 \mathbb{S}^2$    | .197±.003        | <u>.194±.002</u> | .193±.002        |
|                              | $\mathbb{H}^4$                 | .175±.003        | .184±.003        | <u>.178±.003</u> |
|                              | $\mathbb{S}^2 \mathbb{E}^2$    | .199±.003        | .194±.002        | .194±.002        |
|                              | $\mathbb{S}^4$                 | .194±.002        | .188±.002        | .189±.002        |
|                              | $(\mathbb{H}^2)^2$             | .193±.003        | $.193 \pm .002$  | .191±.002        |
| CS PhDs                      | $\mathbb{H}^4$                 | .053±.005        | .052±.005        | .041±.004        |
| Temperature                  | $\mathbb{S}^2\mathbb{S}^1$     | 7.198±.212       | 4.531±.187       | $7.130 \pm .123$ |
| Traffic                      | $\mathbb{E}^1(\mathbb{S}^1)^4$ | <u>.510±.003</u> | $.505 \pm .003$  | $.534 \pm .003$  |

### 4.3 BASELINES

We use Scikit-Learn (Pedregosa et al., 2011) DTs and RFs in both the ambient space  $\mathbb{R}^{D+1}$  and the tangent plane  $T_{\mu_0}\mathcal{P}$  as baselines. Ambient space models operate directly on ambient space coordinates. Tangent plane models project points from  $\mathcal{P}$  to  $\mathcal{T}_{\mu_0}\mathcal{P}$  by applying the logarithmic map at  $\mu_0$ as a preprocessing step. We use Scikit-Learn k-nearest neighbor (k-NN) classifiers and regressors with precomputed pairwise distance matrices according to  $\delta_{\mathcal{P}}$  (Eq. 11). Finally, we implemented the product space perceptron algorithm described in Tabaghi et al. (2021). For our own models, we set hyperparameters identically to Scikit-Learn DTs and RFs, except we consider all  $\binom{D}{2}$  projections— for a total of 3 features per 2-dimensional component manifold, just like ambient space methods use. Full details for each model can be found in Appendix E.1.

#### 486 Decision boundaries visualized: land vs water 487 488 **Product Space RF Euclidean RF** 1.0 489 490 491 492 493 P(Land) 494 495 **Tangent RF** k-Nearest Neighbors 496 497 498 499 500 0.0501

Figure 5: We color a world map with each model's predicted  $\mathbb{P}(\text{Land})$  for the "Landmasses" dataset, a land vs. water classification benchmark in  $\mathbb{S}^2$ . Each RF consists of 12 DTs with a max depth of 3. Note the artifacts learned by Euclidean, tangent RFs, and *k*-NN models.

## 4.4 **RESULTS**

For single-curvature synthetic datasets, our method was the best classifier in 10 out of 11 signatures (Figure 3) and the best regressor (Figure 4 for all signatures. In Tables 2 and 3, we demonstrate consistently good performance across a diverse range of benchmarks.

Further experiments can be found in the Appendix: we provide ablations in F, detailed tables and latent space visualizations in G, comparisons to MLP and GNN models in I, runtime and computational complexity analysis in J, and interpretability experiments in K.

502

503

504

505 506 507

508 509

510

511

512

513

## 5 CONCLUSION

We present strong preliminary evidence favoring mixed-curvature DTs and RFs. In particular, we motivate and describe our entire algorithm and demonstrate its effectiveness across a highly diverse set of 57 benchmarks covering a variety of tasks and geometries.

Product manifold DTs and RFs offer a valuable balance of expressiveness and simplicity, positioned
between extremely legible but underpowered linear classifiers and powerful but uninterpretable neural networks operating in product manifolds. We believe that these qualities, combined with their
demonstrated performance across our benchmark datasets, are compelling evidence of our method's
usefulness in a non-Euclidean data analysis toolkit.

527 **Limitations.** While we view our work as downstream of signature selection and embedding generation, its value heavily depends on the availability of good product manifold embeddings. There are 528 challenges in selecting appropriate signatures (Borde et al., 2023a), and product manifolds are not 529 able to represent all patterns in data (Borde & Kratsios, 2023). Furthermore, it is computationally 530 intensive to generate of embeddings. There are also tradeoffs between DTs and RFs and other high-531 performing methods, especially graph neural networks when topologies are known. The complexity 532 of working with non-Euclidean data could pose a potential barrier to adoption. Finally, the lack of 533 a privileged basis (Elhage et al., 2023) in non-Euclidean embeddings makes the inductive bias of 534 decision trees less well-motivated. 535

**Future work.** It may be possible to exploit non-privileged basis dimensions using approaches such as rotation forests (Bagnall et al., 2020), random 2-D subspace angles, or oblique decision trees. A continuous unification, such as the  $\kappa$ -stereographic model described in (Skopek et al., 2020), may be more robust and elegant. Extensions to simplex geometry (Mishra et al., 2020) are also worth considering.

## 540 REFERENCES

570

581

582

583

586

588

589

- 10x Genomics. Hodgkin's Lymphoma, Dissociated Tumor: Targeted, Immunology Panel, 2020a.
   URL https://www.10xgenomics.com/datasets/hodgkins-lymphoma-disso
   ciated-tumor-targeted-immunology-panel-3-1-standard-4-0-0.
- 10x Genomics. PBMCs from a Healthy Donor: Targeted-Compare, Immunology Panel, 2020b.
   URL https://www.10xgenomics.com/datasets/pbm-cs-from-a-healthy-d
   onor-targeted-compare-immunology-panel-3-1-standard-4-0-0.
- Lada A. Adamic and Natalie Glance. The political blogosphere and the 2004 U.S. election: divided they blog. In *Proceedings of the 3rd international workshop on Link discovery*, LinkKDD '05, pp. 36–43, New York, NY, USA, August 2005. Association for Computing Machinery. ISBN 978-1-59593-215-0. doi: 10.1145/1134277. URL https://doi.org/10.1145/1134271.1134277.
- Gregor Bachmann, Gary Bécigneul, and Octavian-Eugen Ganea. Constant Curvature Graph Convolutional Networks, May 2020. URL http://arxiv.org/abs/1911.05076.
   arXiv:1911.05076 [cs].
- A. Bagnall, M. Flynn, J. Large, J. Line, A. Bostrom, and G. Cawley. Is rotation forest the best classifier for problems with continuous features?, April 2020. URL http://arxiv.org/ab s/1809.06705. arXiv:1809.06705 [cs].
- Gary Becigneul and Octavian-Eugen Ganea. Riemannian Adaptive Optimization Methods. September 2018. URL https://openreview.net/forum?id=rleiqi09K7.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation Learning: A Review and New
   Perspectives, April 2014. URL http://arxiv.org/abs/1206.5538. arXiv:1206.5538.
- Haitz Saez de Ocariz Borde, Alvaro Arroyo, Ismael Morales, Ingmar Posner, and Xiaowen Dong.
   Neural Latent Geometry Search: Product Manifold Inference via Gromov-Hausdorff-Informed
   Bayesian Optimization, October 2023a. URL http://arxiv.org/abs/2309.04810.
- Haitz Sáez de Ocáriz Borde and Anastasis Kratsios. Neural Snowflakes: Universal Latent Graph
   Inference via Trainable Latent Geometries, October 2023. URL http://arxiv.org/abs/
   2310.15003. arXiv:2310.15003.
- Haitz Sáez de Ocáriz Borde, Anees Kazi, Federico Barbero, and Pietro Liò. Latent Graph Inference using Product Manifolds, June 2023b. URL http://arxiv.org/abs/2211.16199. arXiv:2211.16199 [cs].
- Haitz Sáez de Ocáriz Borde, Anastasis Kratsios, Marc T. Law, Xiaowen Dong, and Michael Bronstein. Neural Spacetimes for DAG Representation Learning, August 2024. URL http: //arxiv.org/abs/2408.13885. arXiv:2408.13885.
  - Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, October 2001. ISSN 1573-0565. doi: 10.1023/A:1010933404324. URL https://doi.org/10.1023/A:1010933404324.
- Leo Breiman. *Classification and Regression Trees*. Routledge, New York, October 2017. ISBN 978-1-315-13947-0. doi: 10.1201/9781315139470.
- Lawrence Cayton. Algorithms for manifold learning. 2005.
  - Benjamin Paul Chamberlain, James Clough, and Marc Peter Deisenroth. Neural Embeddings of Graphs in Hyperbolic Space, May 2017. URL http://arxiv.org/abs/1705.10359. arXiv:1705.10359.
- Philippe Chlenski, Ethan Turok, Antonio Moretti, and Itsik Pe'er. Fast hyperboloid decision tree algorithms, March 2024. URL http://arxiv.org/abs/2310.13841. arXiv:2310.13841 [cs].

| 594<br>595<br>596               | Sungjun Cho, Seunghyuk Cho, Sungwoo Park, Hankook Lee, Honglak Lee, and Moontae Lee.<br>Curve Your Attention: Mixed-Curvature Transformers for Graph Representation Learning,<br>September 2023. URL http://arxiv.org/abs/2309.04082. arXiv:2309.04082 [cs].   |
|---------------------------------|--|
| 597<br>598<br>599<br>600        | Calin Cruceru, Gary Bécigneul, and Octavian-Eugen Ganea. Computationally Tractable Riemannian Manifolds for Graph Embeddings, June 2020. URL http://arxiv.org/abs/2002.08665. 65. arXiv:2002.08665.  |
| 601<br>602<br>603<br>604        | Jiarui Ding and Aviv Regev. Deep generative model embedding of single-cell RNA-Seq profiles<br>on hyperspheres and hyperbolic spaces. <i>Nature Communications</i> , 12(1):2554, May 2021. ISSN<br>2041-1723. doi: 10.1038/s41467-021-22851-4. URL https://www.nature.com/artic<br>les/s41467-021-22851-4. Publisher: Nature Publishing Group.                             |
| 605<br>606<br>607               | Manfredo Do Carmo. <i>Riemannian Geometry</i> . Springer US, 1992. URL https://link.springer.com/book/9780817634902.   |
| 608<br>609<br>610               | Lars Doorenbos, Pablo Márquez-Neila, Raphael Sznitman, and Pascal Mettes. Hyperbolic Random Forests, August 2023. URL http://arxiv.org/abs/2308.13279. arXiv:2308.13279 [cs].  |
| 612<br>613<br>614               | Nelson Elhage, Robert Lasenby, and Christopher Olah. Privileged Bases in the Transformer Residual Stream, 2023. URL https://transformer-circuits.pub/2023/privileged-basis/index.html.   |
| 615<br>616<br>617               | Fedesoriano. Traffic Prediction Dataset, 2020. URL https://www.kaggle.com/dataset s/fedesoriano/traffic-prediction-dataset.  |
| 618<br>619<br>620               | Octavian-Eugen Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic Entailment Cones for<br>Learning Hierarchical Embeddings, June 2018. URL http://arxiv.org/abs/1804.018<br>82. arXiv:1804.01882.   |
| 621<br>622<br>623<br>624<br>625 | C. Lee Giles, Kurt D. Bollacker, and Steve Lawrence. CiteSeer: an automatic citation indexing system. In <i>Proceedings of the third ACM conference on Digital libraries</i> , DL '98, pp. 89–98, New York, NY, USA, May 1998. Association for Computing Machinery. ISBN 978-0-89791-965-4. doi: 10.1145/276675.276685. URL https://dl.acm.org/doi/10.1145/276 675.276685. |
| 626<br>627<br>628<br>629        | Francesco Di Giovanni, Giulia Luise, and Michael Bronstein. Heterogeneous manifolds for curvature-aware graph embedding, February 2022. URL http://arxiv.org/abs/2202.01185. arXiv:2202.01185.   |
| 630<br>631<br>632<br>633        | M. Girvan and M. E. J. Newman. Community structure in social and biological networks. <i>Proceedings of the National Academy of Sciences of the United States of America</i> , 99(12):7821–7826, June 2002. ISSN 0027-8424. doi: 10.1073/pnas.122653799. URL https://www.ncbi.nlm.nih.gov/pmc/articles/PMC122977/.   |
| 634<br>635<br>636<br>637        | Albert Gu, Frederic Sala, Beliz Gunel, and Christopher Ré. Learning Mixed-Curvature Representa-<br>tions in Product Spaces. September 2018. URL https://openreview.net/forum?id=<br>HJxeWnCcF7.  |
| 638<br>639<br>640               | Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring Network Structure, Dynamics, and Function using NetworkX. pp. 11–15, Pasadena, California, June 2008. doi: 10.25080/TCW V9851. URL https://doi.curvenote.com/10.25080/TCWV9851.  |
| 641<br>642<br>643<br>644<br>645 | John D. Hunter. Matplotlib: A 2D Graphics Environment. <i>Computing in Science &amp; Engineering</i> , 9(3):90–95, May 2007. ISSN 1558-366X. doi: 10.1109/MCSE.2007.55. URL https://ieeexplore.ieee.org/document/4160265. Conference Name: Computing in Science & Engineering.   |
| 646<br>647                      | David S. Johnson. The genealogy of theoretical computer science: a preliminary report. <i>SIGACT</i><br><i>News</i> , 16(2):36–49, July 1984. ISSN 0163-5700. doi: 10.1145/1008959.1008960. URL<br>https://dl.acm.org/doi/10.1145/1008959.1008960.   |

| 648<br>649<br>650<br>651        | Allan R. Jones, Caroline C. Overly, and Susan M. Sunkin. The Allen Brain Atlas: 5 years and beyond. <i>Nature Reviews Neuroscience</i> , 10(11):821–828, November 2009. ISSN 1471-0048. doi: 10.1038/nrn2722. URL https://www.nature.com/articles/nrn2722. Publisher: Nature Publishing Group.  |
|---------------------------------|---|
| 652<br>653<br>654               | Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization, January 2017.<br>URL http://arxiv.org/abs/1412.6980. arXiv:1412.6980 [cs].   |
| 655<br>656<br>657               | Donald Ervin Knuth. <i>The Stanford GraphBase : a platform for combinatorial computing</i> . New York, N.Y. : ACM Press ; Reading, Mass. : Addison-Wesley, 1993. ISBN 978-0-201-54275-2. URL http://archive.org/details/stanfordgraphbas00knut.   |
| 658<br>659<br>660               | Max Kochurov, Rasul Karimov, and Serge Kozlukov. Geoopt: Riemannian Optimization in PyTorch, July 2020. URL http://arxiv.org/abs/2005.02819. arXiv:2005.02819 [cs].   |
| 661                             | Valdis Krebs. Books about US politics, 2004. URL http://www.orgnet.com.   |
| 663                             | Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. 2009.   |
| 664<br>665<br>666<br>667        | Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. <i>Proceedings of the IEEE</i> , 86(11):2278–2324, November 1998. ISSN 1558-2256. doi: 10.1109/5.726791. URL https://ieeexplore.ieee.org/document/726791. Conference Name: Proceedings of the IEEE.  |
| 669<br>670<br>671<br>672        | Kaibei Li, Yihao Zhang, Junlin Zhu, Xiaokang Li, and Xibin Wang. Multi-space interaction learn-<br>ing for disentangled knowledge-aware recommendation. <i>Expert Systems with Applications</i> , 254:<br>124458, November 2024. ISSN 0957-4174. doi: 10.1016/j.eswa.2024.124458. URL https:<br>//www.sciencedirect.com/science/article/pii/S0957417424013241.  |
| 673<br>674<br>675<br>676<br>677 | David Lusseau, Karsten Schneider, Oliver J. Boisseau, Patti Haase, Elisabeth Slooten, and Steve M. Dawson. The bottlenose dolphin community of Doubtful Sound features a large proportion of long-lasting associations. <i>Behavioral Ecology and Sociobiology</i> , 54(4):396–405, September 2003. ISSN 1432-0762. doi: 10.1007/s00265-003-0651-y. URL https://doi.org/10.1 007/s00265-003-0651-y.   |
| 678<br>679<br>680               | Daniel McNeela, Frederic Sala, and Anthony Gitter. Product Manifold Representations for Learning on Biological Pathways, January 2024. URL http://arxiv.org/abs/2401.15478. arXiv:2401.15478 [cs, q-bio].   |
| 682<br>683<br>684               | Bamdev Mishra, Hiroyuki Kasai, and Pratik Jawanpuria. Riemannian optimization on the simplex of positive definite matrices, November 2020. URL http://arxiv.org/abs/1906.10436. arXiv:1906.10436.   |
| 685<br>686<br>687               | Yoshihiro Nagano, Shoichiro Yamaguchi, Yasuhiro Fujita, and Masanori Koyama. A Wrapped<br>Normal Distribution on Hyperbolic Space for Gradient-Based Learning, May 2019. URL http:<br>//arxiv.org/abs/1902.02992. arXiv:1902.02992 [cs, stat].  |
| 688<br>689<br>690               | M. E. J. Newman. Finding community structure in networks using the eigenvectors of matrices,<br>May 2006. URL https://arxiv.org/abs/physics/0605087v3.  |
| 691<br>692<br>693               | Tuc Nguyen-Van, Dung D. Le, and The-Anh Ta. Improving Heterogeneous Graph Learning with Weighted Mixed-Curvature Product Manifold, July 2023. URL http://arxiv.org/abs/2307.04514. arXiv:2307.04514 [cs].   |
| 695<br>696                      | Maximilian Nickel and Douwe Kiela. Poincaré Embeddings for Learning Hierarchical Representa-<br>tions, May 2017. URL http://arxiv.org/abs/1705.08039. arXiv:1705.08039.   |
| 697<br>698<br>699<br>700<br>701 | Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Edouard Duchesnay. Scikit-learn: Machine Learning in Python. <i>Journal of Machine Learning Research</i> , 12 (85):2825–2830, 2011. ISSN 1533-7928. URL http://jmlr.org/papers/v12/pedr egosalla.html. |

| 702<br>703<br>704<br>705                                    | Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad.<br>Collective Classification in Network Data. <i>AI Mag.</i> , 29(3):93–106, September 2008. ISSN 0738-4602. doi: 10.1609/aimag.v29i3.2157. URL https://doi.org/10.1609/aimag.v29i3.2157.   |
|---|---|
| 706<br>707<br>708<br>709                                    | Ondrej Skopek, Octavian-Eugen Ganea, and Gary Bécigneul. Mixed-curvature Variational Autoen-<br>coders, February 2020. URL http://arxiv.org/abs/1911.08411. arXiv:1911.08411<br>[cs, stat].   |
| 710<br>711<br>712   | Rishi Sonthalia and Anna C. Gilbert. Tree! I am no Tree! I am a Low Dimensional Hyperbolic Em-<br>bedding, October 2020. URL http://arxiv.org/abs/2005.03847. arXiv:2005.03847<br>[cs, math, stat].   |
| 713<br>714<br>715   | Li Sun, Zhongbao Zhang, Junda Ye, Hao Peng, Jiawei Zhang, Sen Su, and Philip S. Yu. A Self-<br>supervised Mixed-curvature Graph Neural Network, December 2021. URL http://arxiv.<br>org/abs/2112.05393. arXiv:2112.05393 [cs].  |
| 716<br>717<br>718   | Puoya Tabaghi, Chao Pan, Eli Chien, Jianhao Peng, and Olgica Milenkovic. Linear Classifiers in<br>Product Space Forms, February 2021. URL http://arxiv.org/abs/2102.10204.<br>arXiv:2102.10204 [cs, stat] version: 1.   |
| 719<br>720<br>721<br>722                                    | Puoya Tabaghi, Michael Khanzadeh, Yusu Wang, and Sivash Mirarab. Principal Component<br>Analysis in Space Forms, July 2024. URL http://arxiv.org/abs/2301.02750.<br>arXiv:2301.02750 [cs, eess, math, stat].  |
| 723<br>724<br>725   | Dimosthenis Tsagkrasoulis and Giovanni Montana. Random Forest regression for manifold-<br>valued responses, February 2017. URL http://arxiv.org/abs/1701.08381.<br>arXiv:1701.08381 [stat].   |
| 726<br>727<br>728   | Marco Virgolin. Time complexity for different machine learning algorithms, February 2021. URL https://marcovirgolin.github.io/extras/details_time_complexity_m achine_learning_algorithms/.   |
| 729<br>730<br>731<br>732<br>733<br>734                      | Shen Wang, Xiaokai Wei, Cicero Nogueira Nogueira dos Santos, Zhiguo Wang, Ramesh Nallapati,<br>Andrew Arnold, Bing Xiang, Philip S. Yu, and Isabel F. Cruz. Mixed-Curvature Multi-Relational<br>Graph Neural Network for Knowledge Graph Completion. In <i>Proceedings of the Web Conference</i><br>2021, WWW '21, pp. 1761–1771, New York, NY, USA, June 2021. Association for Computing<br>Machinery. ISBN 978-1-4503-8312-7. doi: 10.1145/3442381.3450118. URL https://doi.<br>org/10.1145/3442381.3450118.  |
| 735<br>736<br>737<br>738                                    | Wikipedia. List of cities by average temperature, August 2024. URL https://en.wikipedia<br>.org/w/index.php?title=List_of_cities_by_average_temperature&old<br>id=1241784795#cite_note-1. Page Version ID: 1241784795.  |
| 739<br>740<br>741<br>742                                    | Wayne W. Zachary. An Information Flow Model for Conflict and Fission in Small Groups. <i>Journal of Anthropological Research</i> , 33(4):452–473, 1977. ISSN 0091-7710. URL https://www.jstor.org/stable/3629752. Publisher: [University of New Mexico, University of Chicago Press].   |
| 743<br>744<br>745<br>746                                    | Sharon Zhang, Amit Moscovich, and Amit Singer. Product Manifold Learning. In <i>Proceedings</i> of The 24th International Conference on Artificial Intelligence and Statistics, pp. 3241–3249. PMLR, March 2021. URL https://proceedings.mlr.press/v130/zhang21j.html. ISSN: 2640-3498.   |
| 747<br>748<br>749<br>750<br>751<br>752<br>753<br>754<br>755 | <ul> <li>Grace X. Y. Zheng, Jessica M. Terry, Phillip Belgrader, Paul Ryvkin, Zachary W. Bent, Ryan Wilson, Solongo B. Ziraldo, Tobias D. Wheeler, Geoff P. McDermott, Junjie Zhu, Mark T. Gregory, Joe Shuga, Luz Montesclaros, Jason G. Underwood, Donald A. Masquelier, Stefanie Y. Nishimura, Michael Schnall-Levin, Paul W. Wyatt, Christopher M. Hindson, Rajiv Bharadwaj, Alexander Wong, Kevin D. Ness, Lan W. Beppu, H. Joachim Deeg, Christopher McFarland, Keith R. Loeb, William J. Valente, Nolan G. Ericson, Emily A. Stevens, Jerald P. Radich, Tarjei S. Mikkelsen, Benjamin J. Hindson, and Jason H. Bielas. Massively parallel digital transcriptional profiling of single cells. <i>Nature Communications</i>, 8(1):14049, January 2017. ISSN 2041-1723. doi: 10.1038/ncomms14049. URL https://www.nature.com/articles/ncomms14049. Publisher: Nature Publishing Group.</li> </ul> |

# A GAUSSIAN MIXTURE DETAILS

#### 758 A.1 OVERALL STRUCTURE

The structure of our sampling algorithm is as follows. Note that, rather than letting  $\mathcal{M}$  be a manifold of arbitrary curvature, we force its curvature to be one of  $\{-1, 0, 1\}$  for implementation reasons. This necessitates rescaling steps, which take place in Equations 29, 33, and 39. The result is equivalent to performing the equivalent steps, without rescaling, on a manifold of the proper curvature. 1. Generate c, a vector that divides m samples into n clusters:

 $\mathbf{p_{raw}} = \langle p_0, p_1, \dots, p_{n-1} \rangle \tag{23}$ 

$$p_i \sim \text{Uniform}(0, 1)$$
 (24)

$$\mathbf{p_{norm}} = \frac{\mathbf{p_{raw}}}{\sum_{i=0}^{n-1} p_i} \tag{25}$$

 $\sum_{i=0}^{n} p_i$   $\mathbf{c} = \langle c_0, c_1, \dots c_{m-1} \rangle$ (26)

$$e_i \sim \text{Categorical}(n, \mathbf{p_{norm}})$$
 (27)

2. Sample  $\mathbf{M}_{euc}$ , an  $n \times D$  matrix of n class means:

$$\mathbf{M}_{\mathbf{euc}} = \langle \mathbf{m}_{\mathbf{0}}, \mathbf{m}_{\mathbf{1}}, \dots, \mathbf{m}_{\mathbf{n-1}} \rangle^{T}$$
(28)

$$\mathbf{m_i} \sim \mathcal{N}(0, \sqrt{K\mathbf{I}}). \tag{29}$$

3. Move  $\mathbf{M}_{euc}$  into  $T_0\mathcal{M}$ , the tangent plane at the origin of  $\mathcal{M}$ , by applying  $\psi : \mathbf{x} \to (0, \mathbf{x})$  per-row to  $\mathbf{M}_{euc}$ :

$$\mathbf{M}_{\mathbf{tan}} = \langle \psi(\mathbf{m_0}), \psi(\mathbf{m_1}), \dots \psi(\mathbf{m_{n-1}}) \rangle^T,$$
(30)

$$\psi: \mathbb{R}^D \to \mathbb{R}^{D+1}, \ \mathbf{x} \to \langle 0, \mathbf{x} \rangle.$$
(31)

4. Project  $\mathbf{M}_{tan}$  onto  $\mathcal{M}$  using the exponential map from  $T_0 \mathcal{M}$  to  $\mathbf{M}_{tan}$ :

$$\mathbf{M} = \exp_0(\mathbf{M_{tan}}). \tag{32}$$

5. For  $0 \le i < n$ , sample a corresponding covariance matrix. Here,  $\sigma$  is a variance scale parameter that can be set:

$$\Sigma_{i} \sim \text{Wishart}(\sigma \sqrt{K} \mathbf{I}, D)$$
 (33)

6. For  $0 \le j < m$ , sample  $\mathbf{X}_{euc}$ , a matrix of m points according to their clusters' covariance matrices:

$$\mathbf{X}_{\mathbf{euc}} = \langle \mathbf{x}_0, \mathbf{x}_1, \dots \mathbf{x}_{\mathbf{m}-1} \rangle^T$$
(34)

$$x_j \sim \mathcal{N}(0, \Sigma_{\mathbf{c}_j}).$$
 (35)

7. Apply  $\psi(\cdot)$  from Eq 31 to each  $\mathbf{x}_i$  to move it into  $T_0\mathcal{M}$ :

$$\mathbf{X}_{\mathbf{tan}} = \langle \psi(\mathbf{x_0}), \psi(\mathbf{x_1}), \dots \psi(\mathbf{x_{m-1}}) \rangle^T.$$
(36)

8. For each row in  $\mathbf{X}_{tan}$ , apply parallel transport from  $T_0 \mathcal{M}$  to its class mean:

$$\mathbf{X}_{\mathbf{PT}} = \langle \mathbf{x}_{\mathbf{0},\mu}, \mathbf{x}_{\mathbf{1},\mu}, \dots, \mathbf{x}_{\mathbf{m-1},\mu} \rangle \tag{37}$$

$$\mathbf{x}_{\mathbf{j},\mu} = PT_{0\to\mathbf{m}_{\mathbf{c}_{\mathbf{i}}}}(\mathbf{x}_{\mathbf{j}}) \tag{38}$$

9. Use the exponential map at  $T_{\mu}\mathcal{M}$  to move the points onto the manifold:

$$\mathbf{X}_{\mathcal{M}} = \left\langle \mathbf{x}_{\mathbf{0},\mathcal{M}}, \mathbf{x}_{\mathbf{1},\mathcal{M}}, \dots, \mathbf{x}_{\mathbf{m}-\mathbf{1},\mathcal{M}} \right\rangle$$
(39)

$$\mathbf{x}_{\mathbf{j},\mathcal{M}} = \frac{\exp_{\mathbf{m}_{\mathbf{c}_{\mathbf{j}}}}(\mathbf{x}_{\mathbf{j},\mu})}{\sqrt{K}}$$
(40)

10. Repeat steps 2–9 for as many manifolds as desired; produce a final embedding by concatenating all component embeddings column-wise:

$$\mathbf{X} = \langle \mathbf{X}_{\mathcal{M}_0}, \mathbf{X}_{\mathcal{M}_1}, \dots \mathbf{X}_{\mathcal{M}_p} \rangle \tag{41}$$

#### A.2 EQUATIONS FOR MANIFOLD OPERATIONS

First, we provide the forms of the parallel transport operation in hyperbolic, hyperspherical, and Euclidean spaces:

$$PT^{\mathbb{H}}_{\nu \to \mu}(\mathbf{v}) = \mathbf{v} + \frac{\langle \mu - \alpha \nu, \nu \rangle_{\mathcal{L}}}{\alpha + 1}(\nu + \mu)$$
(42)

$$\alpha = -\langle \nu, \mu \rangle_{\mathcal{L}} \tag{43}$$

$$PT^{\mathbb{S}}_{\nu \to \mu}(\mathbf{v}) = \mathbf{v}\cos(d) + \frac{\sin(d)}{d}(\mu - \cos(d)\nu)$$

$$d = \cos^{-1}(\nu \cdot \mu)$$
(44)
(45)

$$\mathrm{PT}_{\nu \to \mu}^{\mathbb{E}}(\mathbf{v}) = \mathbf{v} + \mu - \nu. \tag{46}$$

(45)

The exponential map is defined as follows in each of the three spaces:

$$\exp_{\mu}^{\mathbb{H}}(\mathbf{u}) = \cosh(\|\mathbf{u}\|_{\mathcal{L}})\mu + \sinh(\|\mathbf{u}\|_{\mathcal{L}})\frac{\mathbf{u}}{\|\mathbf{u}\|_{\mathcal{L}}}$$
(47)

$$\exp^{\mathbb{S}}_{\mu}(\mathbf{u}) = \cos(\|\mathbf{u}\|)\mu + \sin(\|\mathbf{u}\|)\frac{\mathbf{u}}{\|\mathbf{u}\|}$$
(48)

$$\mathrm{sp}^{\mathbb{E}}_{\mu}(\mathbf{u}) = \mathbf{u}.\tag{49}$$

#### A.3 GENERATING CLASSIFICATION TARGETS

ez

To generate classification targets covering  $p \le n$  classes, all we need to do is map clusters to classes. To ensure that each class has at least one associated cluster, we arbitrarily assign the first p clusters to the first p classes. In the p = n case, this is equal to the p-dimensional identity matrix, and we conclude. In the p < n case, we assign the remaining n - p by drawing assignments from a uniform categorical distribution over the p classes. 

#### A.4 GENERATING REGRESSION TARGETS

To generate regression targets, we draw per-cluster slopes and intercepts: 

$$\beta_{i,k} \sim Uniform(-1,1) \tag{50}$$

$$\alpha_i \sim Uniform(-10, 10 \tag{51})$$

We then multiply each  $x_i \in \mathbf{X}_{\mathbf{e}}\mathbf{uc}$  (i.e. the pre-transport samples from the normal distribution) by  $\beta$  and add  $\alpha$ :

$$y_j = \mathbf{x}_j \beta + \alpha + \varepsilon \tag{52}$$

To make the regression task more constrained and, therefore, to make the RMSEs across samples more comparable, we further normalize the labels to the range [0,1] by subtracting the minimum y value and dividing by the range.

#### A.5 RELATIONSHIP TO OTHER WORK

Nagano et al. (2019) developed the overall technique used for a single cluster and a single manifold, i.e. steps 6-9. Chlenski et al. (2024) modified this method to work for mixtures of Gaussians in  $\mathbb{H}^{d,1}$ , and deployed it for  $d \in \{2, 4, 8, 16\}$ . This corresponds to steps 1–5 of our procedure (although note that our covariance matrices are sampled differently in step 5). Thus, our contribution is simply to add step 10, modify step 5 to use the Wishart distribution, to add curvature-related scaling factors in Equations 29, 33, and 39, and to generate classification and regression targets as described in the preceding sections. 

We apply this to *hyperspherical* manifolds, for which the von Mises-Fisher (VMF) distribution is typically preferred. This is an unconventional choice, but has been employed previously by Skopek et al. (2020) in their mixed-curvature VAE formulation. We do not argue for the superiority of our approach over the VMF distribution in general; however, we prefer to use ours for these benchmarks, as it allows us to draw simpler parallels between manifolds of different curvatures.

## **B PRODUCT SPACE DECISION TREE PSEUDOCODE**

|                         | Procedure FIT:   |
|-------------------------|--|
| 2:                      | $\mathcal{P}$ (signature of) product manifold  |
| 3:                      | X data points  |
| 4:                      | y target labels  |
| 5:                      | Initialize:  |
| 6:                      | $\mathcal{T}$ an empty tree  |
| 7:                      | return $FITTREE(\mathbf{X}, \mathbf{y}, 0)$  |
| 8:                      |  |
| 9:                      | Procedure FITTREE:   |
| 10:                     | X data points  |
| 11:                     | y target labels  |
| 12:                     | t current depth of the tree.   |
| 13:                     | Initialize:  |
| 14:                     | $d_{\text{best}}$ dimension of best split,   |
| 15:                     | $\theta_{\text{best}}$ angle of best split,  |
| 10:                     | $I_{Gbest}$ information gain of best split.  |
| 1/:                     | <b>10</b> reach $a \in \mathbf{D}^{*}$ <b>do</b>   |
| 18:                     | $\mathcal{M} \leftarrow \text{component manifold for dimension } a$  |
| 19:                     | $\Theta \leftarrow \text{GETCANDIDATES}(\mathcal{M}, \mathbf{A}, a)$   |
| 20:                     | Destition <b>X</b> winto $\mathbf{X}^+$ $\mathbf{X}^ \mathbf{x}^+$ $\mathbf{x}^-$ via Eq. 16   |
| 21:                     | Faltition $\mathbf{A}$ , y find $\mathbf{A}^+$ , $\mathbf{A}^-$ , $\mathbf{y}^+$ , $\mathbf{y}^-$ via Eq. 10.                                    |
| 22.<br>73.              | if $IC = \sum IC$ , then   |
| 23.<br>24.              | $d = \theta = IC = d \theta IC$  |
| 2 <del>4</del> .<br>25. | end if   |
| 25.<br>26·              | end for  |
| 20.<br>27.              | end for  |
| 28:                     | if no valid split was found then   |
| 29:                     | <b>return</b> $\mathcal{N}$ , a new <b>leaf node</b> with <b>v</b> probabilities.  |
| 30:                     | else   |
| 31:                     | Create $\mathcal{N}$ , a decision node with $d_{\text{best}}$ and $\theta_{\text{best}}$   |
| 32:                     | $\mathcal{N}_L \leftarrow \text{FitTree}(\mathbf{X}^-, \mathbf{y}^-, t+1)$   |
| 33:                     | $\mathcal{N}_R \leftarrow FitTree(\mathbf{X}^+, \mathbf{y}^+, t+1)$  |
| 34:                     | <b>return</b> $\mathcal{N}$ with left child $\mathcal{N}_L$ and right child $\mathcal{N}_R$  |
| 35:                     | end if   |
| 36:                     |  |
| 37:                     | Procedure GETCANDIDATES:   |
| 38:                     | $\mathcal{M}$ A component manifold   |
| 39:                     | $\mathbf{X}$ A dataset of points in $\mathcal{M}$  |
| 40:                     | d A dimension index  |
| 41:                     | if $d$ is the special dimension then   |
| 42:                     | return empty array []  |
| 43:                     | end if   |
| 44:                     | $\Theta \leftarrow$ Angles of X via Eq. 15   |
| 45:                     | $\Theta \leftarrow$ sort and deduplicate $\Theta$  |
| 46:                     | <b>return</b> $[\theta_m \text{ for } \theta_i, \theta_{i+1} \in \Theta \text{ via Eq. 18, 20, or 22]}$ (depending on curvature of $\mathcal{M}$ |

# 918 C PROOF OF EQUIVALENCE FOR EUCLIDEAN CASE

A classical CART tree splits data points according to whether their value in a given dimension is greater than or less than some threshold value t. Midpoints are simple arithmetic means. This can be written as:

$$S'(\mathbf{x}, d, t) = \begin{cases} 1 \text{ if } x_d > t, \\ 0 \text{ otherwise.} \end{cases}$$
(53)

$$m_{DT}(\mathbf{u}, \mathbf{v}) = \frac{u_d + v_d}{2}.$$
(54)

In our transformed DT, we lift the data points by applying  $\phi : \mathbf{x} \to (1, \mathbf{x})$  and then check which side of an axis-inclined hyperplane they fall on. The splitting function is based on the angle  $\theta$  of inclination with respect to the (0, d) plane, i.e.,  $\langle 1, x_d \rangle$ . Our midpoints are computed to ensure equidistance in the original manifold:

$$S(\mathbf{x}, d, \theta) = \operatorname{sign}(\sin(\theta)x_d - \cos(\theta)x_0)$$
(55)

$$m_{\mathbb{E}}(\mathbf{u}, \mathbf{v}) = \tan^{-1}\left(\frac{2}{u_d + v_d}\right)$$
(56)

To demonstrate the equivalence of the classical DT formulation to our transformed algorithm in  $\mathbb{E}$ , we will show that Equation 53 is equivalent to Equation 55 and Equation 54 is equivalent to Equation 56 under

$$\theta = \cot^{-1}(t). \tag{57}$$

#### C.1 EQUIVALENCE OF SPLITS

First, we show that Equations 53 and 55 are equivalent, assuming  $t \neq 0$ :

$$S(\mathbf{x}, d, \theta) = \operatorname{sign}(\sin(\theta)x_d - \cos(\theta)x_0) = 1$$
(58)

$$\iff \sin(\theta) x_d - \cos(\theta) > 0 \tag{59}$$

$$\iff \frac{\sin(\theta)}{\cos(\theta)} x_d = \tan(\theta) x_d > 1 \tag{60}$$

$$\iff x_d/t > 1$$
 (61)

$$\iff x_d > t \tag{62}$$

$$\iff S'(\mathbf{x}, d, t) = 1 \tag{63}$$

#### C.2 EQUIVALENCE OF MIDPOINTS

Now, we show that Equations 54 and 56 are equivalent:

$$\cot^{-1}(m_{DT}(\mathbf{u}, \mathbf{v})) = \cot^{-1}\left(\frac{u_d + v_d}{2}\right)$$
(64)

$$=\tan^{-1}\left(\frac{2}{u_d - v_d}\right) \tag{65}$$

$$=m_{\mathbb{E}}(\mathbf{u},\mathbf{v})$$
 (66)

### D SUMMARY OF ANGULAR MIDPOINT FORMULAS

Table 4: Distance functions and midpoint angle formulas for each component manifold type.

| Manifold $\mathcal{M}$ | Distance $\delta_{\mathcal{M}}(\mathbf{u}, \mathbf{v})$                          | Midpoint angle $\theta_{\mathcal{M}}(\mathbf{u}, \mathbf{v})$  |
|------------------------|--|--|
| $\mathbb{S}^{D,K}$     | $\cos^{-1}\left(\frac{K^2\langle \mathbf{u}, \mathbf{v} \rangle}{K}\right)$      | $\frac{\theta_{\mathbf{u}} + \theta_{\mathbf{v}}}{2}$  |
| $\mathbb{E}^{D,0}$     | $\sqrt{\langle {f u},{f v} angle}$   | $\tan^{-1}\left(\frac{2}{u_d + v_d}\right)$  |
| $\mathbb{H}^{D,K}$     | $\frac{-\cosh^{-1}(K^2\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}})}{K}$ | $\cot^{-1}(V - \sqrt{V^2 - 1}) \text{ if } \theta_{\mathbf{u}} + \theta_{\mathbf{v}} < \pi,$                       |
|                        |  | $V := \frac{\sin(2\theta_{\mathbf{u}} - 2\theta_{\mathbf{v}})}{\sin(2\theta_{\mathbf{u}} - 2\theta_{\mathbf{v}})}$ |
|                        |  | $2\sin(\theta_{\mathbf{u}}+\theta_{\mathbf{v}})\sin(\theta_{\mathbf{v}}-\theta_{\mathbf{u}})$                      |

#### E FULL EXPERIMENTAL HYPERPARAMETERS

#### E.1 SCIKIT-LEARN HYPERPARAMETERS

**Random forests and decision trees.** For fairness, we set all DT and RF hyperparameters identically. Specifically, we set the following hyperparameters for both DTs and RFs:

```
• max_depth = 5
```

```
• min_samples_split = 2
```

```
• min_samples_leaf = 1
```

```
• min_impurity_decrease = 0.0
```

For RFs, we also set the following hyperparameters:

```
• n_estimators = 12
```

```
• max_features = "sqrt"
```

```
• bootstrap = True (subsamples the training data)
```

• max\_samples = None (draws *n* samples from a set of *n* points)

Because the scikit-learn implementation differs substantially from ours, subsamples vary even whenthe random seed is set. Nevertheless, we also employ the same random seed for all RF models.

k-nearest neighbor models. For k-nearest neighbors, we use default hyperparameters.

Product space perceptrons and SVMs. Product space perceptrons only have one hyperparameter,
 which is the relative weight assigned to each component manifold. We elect to give each component manifold equal weight.

Neither the SVM code provided by Tabaghi et al. (2021) nor our own reimplementation would run on our datasets. In particular, we had issues satisfying the convexity constraints described in their paper, causing the solve to crash. Correcting this mistake and augmenting our benchmarks with SVM evaluations is a direction for future research.

```
Product space decision trees and random forests. For our models, we set the n_features
= "n_choose_2" parameter. This means that we consider all \binom{n}{2} linear projections. We do this
because we restrict ourselves to 2-dimensional component manifolds, and therefore we only observe
\binom{3}{2} = 3 total angles, equal to the number of features used by ambient space Euclidean methods.
```

- - B E.2 GRAPH EMBEDDINGS
- **Learning embeddings.** We reimplement the method in Gu et al. (2018) to learn graph embeddings. In particular, we use the NetworkX package (Hagberg et al., 2008) to load the graph, extract

the largest connected component, and compute pairwise distances between nodes using the Floyd-Warshall algorithm. For embedding purposes, we treat all graphs as undirected. Pairwise distances were normalized into the range [0, 1] by dividing by the maximum distance.

Embedding hyperparameters. Embeddings were learned using Riemannian Adam (Becigneul & Ganea, 2018) implemented in Geoopt (Kochurov et al., 2020). For each signature, we train 10 randomly-initialized embeddings for 3,000 epochs each. We treat the first 300 epochs as a burn-in period, during which the learning rate is .01 and the curvature of each manifold is fixed. For the remaining epochs, we train embedding coordinates with a learning rate of 0.1 and scale factors with a learning rate of 0.01. These hyperparameters were chosen based on their stability and convergence in exploratory experiments.

Train-test split. Because embeddings must be learned per-node, it was not possible to perform a train-test split prior to the embedding step; however, we performed the train-test split at the node level for all tasks including link prediction. This means that we discarded all edges between test and training nodes from our dataset. While we acknowledge the embeddings step could be a source of leakage, we have no reason to believe this would bias evaluations in favor of any particular model. Future work should focus on developing methods to learn node embeddings in phases or to use masked gradients to minimize leakage at the embedding step.

**Evaluations.** Since it was not clear *a priori* which signature would embed each graph the best, we learned 10 embeddings for each candidate signature and took the one with the best  $D_{avg}$  to be the benchmark signature. Our reasoning is that the lowest-distortion embedding of the graph is the most appropriate benchmark for evaluating the geometrical appropriateness of a classifier. Thus, scores for the lowest-distortion signature appear in Tables 2 and 3, whereas scores for all signatures can be found in Table 5.

Link prediction. To generate link prediction datasets, we trained 100 randomly initialized sets of node embeddings in  $\mathbb{S}^2 \times \mathbb{E}^2 \times \mathbb{H}^2$ . If we let **X** be our original node embeddings and  $\mathcal{E}$  be the ground-truth edges of the graph, we then generated the following dataset:

$$\mathbf{X}_{LP} = \{ \langle \mathbf{x}_i, \mathbf{x}_j \rangle \text{ for } (\mathbf{x}_i, \mathbf{x}_j, \delta_{\mathcal{P}}(\mathbf{x}_i, \mathbf{x}_j) \in \mathbf{X} \}$$
(67)

$$\mathbf{y}_{LP} = \{ \mathbb{I}\{ (\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{E} \} \text{ for } (\mathbf{x}_i, \mathbf{x}_j) \in \mathbf{X} \}$$
(68)

1055 The corresponding signature is  $(\mathcal{P})^2 \times \mathbb{E}^1$ ; in the case of our embeddings, that is  $(\mathbb{S}^2 \times \mathbb{E}^2 \times \mathbb{H}^2)^2 \times \mathbb{E}^1$ .

1057 E.3 VAE TRAINING

1053 1054

1056

1058

1061

1062

1072

1073

1074

**Encoder/decoder architectures.** Following Tabaghi et al. (2021), we use the following encoder/decoder architectures:

- Lymphoma dataset: Two 200-dimensional hidden layers, 500 epochs
- Blood cell dataset: Three 400-dimensional hidden layers, 200 epochs
- Omniglot and MNIST: 400-dimensional latent
- CIFAR-100: 4 × 4 convolutional kernels with stride 2 and padding 1. Encoder: 3 CNN layers of 64, 128, and 512 channels. Decoder: 2048-dimensional dense layer, followed by 2 CNN layers of 256, 64, and 3 channels.

**Training hyperparameters.** Our VAEs were trained using the Adam optimizer (Kingma & Ba, 2017) with default parameters (learning rate .001,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ . In all models, each layer except the last is followed by a ReLU activation function. Curvatures were trained identically, except using a learning rate of .0001, after 100 burn-in epochs. Because some training details were omitted from the original papers, we additionally chose the following hyperparameters:

- Batch size: 4,096
  - Number of samples per point: 64
- $\beta$  (weight for KL-divergence in VAE loss): 1

Train-test split. To minimize the risk of data leakage, we trained our VAEs on only the training data, then used the trained VAEs to generate embeddings for the training and test data. Embeddings were generated by running points through the VAE encoder and taking the returned mean parameter.

**Evaluations.** To conserve memory, we randomly subsampled 1,000 points from the training and test sets for each evaluation. We ran 10 trials per dataset in total.

# 1080 E.4 EMPIRICAL DATASETS

**Landmasses.** We generated a geospatial classification dataset for land versus water prediction by sampling 1,000 points from an evenly sampled grid of 10,000 longitudes and latitudes, transforming them to 3-dimensional coordinates, and assigning a "land" or "water" label to each point using the Basemap library in Matplotlib (Hunter, 2007). For classification, we associate the 3-dimensional coordinates with the signature  $\mathbb{S}^2$ .

**Neural spiking prediction.** We use patch-clamp electrophysiology datasets downloaded from the Allen Mouse Brain Atlas (Jones et al., 2009). We arbitrarily pick Neurons 33 and 46 for their nontrivial spiking dynamics. To represent signals in product spaces, we apply a Fast Fourier Transform and take the top 5 Fourier coefficients by magnitude. We then take their corresponding frequencies  $f_i$  and represent each time point in  $\mathbb{S}^1$  via the following transformation:

1092

1093 1094

1102

1103 1104

1106

 $\phi: \mathbb{R}^1 \to (\mathbb{S}^1)^5, \ \phi(t) = \left( \cos\left(2\pi \frac{t}{f_i}\right), \ \sin\left(2\pi \frac{t}{f_i}\right) \right) \Big|_{i=1}^5$ (69)

This yields a product space representation in  $(\mathbb{S}^1)^5$ . We plot both signals, along with their reconstruction using their top 5 Fourier components, in Figure 6.

**Global temperature by month.** We downloaded a list of global average monthly temperatures for the 400 largest cities in the world from Wikipedia (Wikipedia, 2024). We transform longitude and latitude into 3-D coordinates to represent our data in  $S^2$ . To convert months to  $S^1$  valued coordinates, we transform ordinal representations of months  $t \in [0, 11]$  via the following transformation:

$$\phi : \mathbb{R}^1 \to \mathbb{S}^1, \ \phi(t) = \left(\cos\left(2\pi\frac{t}{12}\right), \ \sin\left(2\pi\frac{t}{12}\right)\right)$$
(70)

1105 This yields a product space representation of the data in  $\mathbb{S}^2 \times \mathbb{S}^1$ .



Figure 6: The "Neuron 33" and "Neuron 46" datasets, along with their reconstruction using the top 5 Fourier coefficients shown in red.

**Traffic prediction.** We download an automobile traffic prediction dataset from Kaggle (Fedesoriano, 2020). This dataset aggregates readings across four sensors with date and time annotations. We process the date and time annotation into day of year (d), day of week (w), hour (h), and minute (m)labels and transform to  $(\mathbb{S}^1)^4$  analogously to the month timestamps in the global temperature data. Letting l be the (numeric) label of the sensor, we apply the following transformation to our data:

$$\phi: \mathbb{R}^5 \to (\mathbb{S}^1)^5 \times \mathbb{E}^1$$
(71)

 $\phi(d, w, h, m, l) = \left(\cos\left(2\pi \frac{d}{365}\right), \sin\left(2\pi \frac{d}{365}\right), \right.$ 

 $\cos\left(2\pi\frac{w}{7}\right), \sin\left(2\pi\frac{w}{7}\right),$ 

 $\cos\left(2\pi\frac{h}{24}\right),\ \sin\left(2\pi\frac{h}{24}\right),$ 

 $\cos\left(2\pi\frac{m}{60}\right), \sin\left(2\pi\frac{m}{60}\right), l\right)$ 

1141

1145

1148 1149

1150

1151

1165

1166

## F ABLATIONS AND EFFECTS OF HYPERPARAMETERS

For all experiments, we sampled 100 mixtures of 32 Gaussians using the signature  $\mathcal{P} = \mathbb{S}^2 \times \mathbb{E}^2 \times \mathbb{H}^2$ in an 8-class regression setting (analogous to the multi-*K* benchmark in Tables 2 and 3, varying one parameter at a time. Results are plotted in Figure 7.



(a) Changing the number of features seen by each DT/RF from 6 to 9 by including the  $(x_1, x_2)$  angle is massively beneficial.





(72)

(b) Changing feature subsampling approaches in RFs doesn't appear to do much.



(c) Increasing the maximum depth of each DT/RF is massively beneficial, and shows no signs of overfitting even at unrestricted max depth. All within-predictor differences are significant.

(d) Replacing the midpoint-angle computations with arithmetic means has no statistically significant effect on performance for DTs or RFs, surprisingly.

Figure 7: Effects of various hyperparameters on the performance of our algorithms.

1182 1183

1180

- 1184
- 1185
- 1186
- 1187

# <sup>1188</sup> G DETAILED RESULTS



## 1190 G.1 GLOBAL TEMPERATURE PREDICTION PLOTS

Figure 8: Decision boundaries for the temperature prediction task for the months of January, April,July, and October, colored by predicted temperature across four trained predictors.

### G.2 VAE LATENT SPACE VISUALIZATIONS

1216 1217



Figure 9: Visualizations of the latent space for all four of the datasets we embed using a VAE, colored by class. For visualization purposes, we show  $S^2$  components in 2-dimensional polar coordinates, and project  $\mathbb{H}^2$  embeddings to the Poincaré disk.

## 1242 G.3 FULL RESULTS TABLE

1243

1250

Table 5: Full results for all benchmarks. Unlike Tables 2 and 3, this table reports full results for all classifiers. It also includes single-*K* results and all signatures for graph embeddings. Recall the following shorthand: C=Classification ( $F_1$ ), R=Regression (RMSE), LP=Link prediction ( $F_1$ ); **bold** = best predictor, <u>underline</u> = second-best predictor; \* = beating **product space methods**, † = beating **ambient space methods**, ‡ = beating **tangent plane methods**, § = beating *k*-nearest neighbors, ¶ = beating **product space perceptrons**.

|       | Dataset              | Task | Signature   | Perceptron  | k-Neighbors  | Euclidean DT   | Euclidean RF  | Tangent DT  | Tangent RF  | Product DT  | Product RF  |
|-------|----------------------|------|---|---|--|--|---|---|---|---|---|
|       | Gaussian             | С    | $\mathbb{E}^2$  | 19.5 ± .5***  | 28.2 ± .4****  | 28.4 ± .5  | <u>30.5 ± .5</u> <sup>§</sup> 1                       | 28.4 ± .5¶  | <u>30.5 ± .5</u> <sup>§1</sup>                        | 28.5 ± .5   | 30.9 ± .5 <sup>§</sup> 1                                |
|       |                      |      | H <sup>2,0.25</sup><br>H <sup>2,0.5</sup>                                   | $18.4 \pm .6^{19^{+4}}$                             | $27.8 \pm .4^{11}$                                   | $28.5 \pm .5$<br>27.8 ± 5  | 29.7 ± .5 <sup>81</sup><br>28.8 ± 51*                 | $28.7 \pm .5^{1}$<br>27.6 ± 5 <sup>1</sup>            | $\frac{30.1 \pm .5^{\$1}}{28.9 \pm .5^{\$1}}$         | $28.7 \pm .5^{\circ}$<br>28.3 ± .5 <sup>°</sup>     | $30.6 \pm .5^{191}$<br>30.2 ± 4 <sup>1818</sup>         |
|       |                      |      | $\mathbb{H}^{2,1.0}$  | 21.0 ± .5 <sup>†§*‡</sup>                           | $28.5 \pm .5^{1*}$                                   | 27.9 ± .5 <sup>11</sup>  | 28.9 ± .5 <sup>¶*‡</sup>                              | 26.1 ± .5 <sup>*</sup>                                | $\frac{20.9 \pm .5}{27.7 \pm .5}$                     | 28.2 ± .5 <sup>1</sup>                              | 30.8 ± .5 <sup>18</sup>                                 |
|       |                      |      | $\mathbb{H}^{2,2.0}$  | 18.2 ± .5 <sup>†§*‡</sup>                           | $27.0 \pm .5^{1}$                                    | 26.2 ± .5  | 26.6 ± .5   | 24.3 ± .5 <sup>†¶*</sup>                              | 25.6 ± .5 <sup>*¶*</sup>                              | 26.4 ± .5 <sup>1</sup> *                            | 28.1 ± .5 <sup>*¶‡</sup>                                |
|       |                      |      | $\mathbb{H}^{2,4.0}$<br>$\otimes 2,-0.25$                                   | $17.3 \pm .5^{10}$                                  | $\frac{26.1 \pm .5^{1}}{20.4 \pm .5^{12}}$           | $25.3 \pm .5^{13}$   | $25.5 \pm .5^{1+3}$                                   | $22.7 \pm .5$   | $24.3 \pm .4^{11}$                                    | $25.2 \pm .5^{13}$                                  | $27.2 \pm .5^{13}$                                      |
| G     |                      |      | $S^{2,-0.5}$  | $16.3 \pm .5^{10^{+1}}$                             | $30.6 \pm .5^{1*}$                                   | $29.2 \pm .5^{1}$  | $30.4 \pm .5^{1*}$                                    | $29.6 \pm .5^{1}$                                     | $\frac{50.5 \pm .5}{31.4 \pm .5}$                     | $30.2 \pm .5^{1}$                                   | 32.3 ± .5**   |
| gle   |                      |      | $S^{2,-1.0}$  | 16.1 ± .5 <sup>†§*‡</sup>                           | $32.6 \pm .5^{1}$                                    | $30.0 \pm .5^{1}$  | 30.4 ± .5 <sup>¶*</sup>                               | $30.3 \pm .5^{1}$                                     | 32.1 ± .5   | $30.9 \pm .5^{1}$                                   | 33.2 ± .5 <sup>*¶‡</sup>                                |
| (sin  |                      |      | $S^{2,-2.0}$<br>$S^{2,-4.0}$  | $15.6 \pm .6^{19^{+4}}$                             | $\frac{36.1 \pm .5^{1}}{41.6 \pm .5^{1}}$            | $32.3 \pm .51$<br>$32.4 \pm .51$                                 | $32.9 \pm .61^{\circ}$<br>33.0 ± 51 <sup>*</sup>      | $33.8 \pm .6^{1}$<br>$35.2 \pm .5^{1}$                | $35.5 \pm .6^{\circ}$<br>37.6 ± 51°                   | $34.1 \pm .6^{11}$<br>36.3 + 5 <sup>11</sup>        | $36.7 \pm .5^{11}$                                      |
| etic  |                      | R    | $\mathbb{E}^2$  | 10.4 ± .0   | .211 ± .002***                                       | .201 ± .002 <sup>§</sup>   | .199 ± .002 <sup>§*</sup>                             | .201 ± .002 <sup>§</sup>                              | .199 ± .002 <sup>§*</sup>                             | .201 ± .002 <sup>§</sup>                            | .198 ± .002***  |
| dfu - |                      |      | H <sup>2,0.25</sup>   |   | .207 ± .003  | .199 ± .002  | .196 ± .002 <sup>§</sup>                              | .199 ± .002   | .196 ± .002   | .199 ± .002   | .196 ± .002   |
| ŝ     |                      |      | H2,0.0<br>H2,1.0  |   | $.206 \pm .003^{10}$<br>$214 \pm .003^{10}$          | $.198 \pm .002^{\circ}$<br>205 ± .003^{\circ}                    | $\frac{.194 \pm .002^{\circ}}{201 \pm .003^{\circ*}}$ | $.197 \pm .002$<br>206 ± 003                          | $\frac{.194 \pm .002^{\circ}}{203 \pm .003^{\circ}}$  | $.197 \pm .003^{\circ}$<br>205 + 003^{\circ}        | $.193 \pm .002^{18}$<br>201 + 003 <sup>188</sup>        |
|       |                      |      | $\mathbb{H}^{2,2.0}$  |   | .211 ± .003  | $.202 \pm .003^{\$}$   | <u>.199 ± .003</u> <sup>\$‡</sup>                     | .203 ± .003   | .201 ± .003   | $.202 \pm .003$                                     | .198 ± .003 <sup>\$‡</sup>                              |
|       |                      |      | $\mathbb{H}^{2,4.0}$<br>$\mathbb{G}^2 = 0.25$                               |   | .215 ± .003***                                       | .206 ± .003  | $\frac{.202 \pm .003^{\$1}}{.002^{\$1}}$              | .207 ± .003   | .204 ± .003**   | .206 ± .003   | .202 ± .003 <sup>\$‡</sup>                              |
|       |                      |      | S <sup>2</sup> , 0.20<br>S <sup>2</sup> , -0.5                              |   | $.200 \pm .003^{++}$                                 | $.198 \pm .003^{\circ}$<br>$.203 \pm .003^{\circ}$               | $\frac{.195 \pm .003^{\circ}}{.199 \pm .003^{\circ}}$ | $.198 \pm .003^{\circ}$<br>$.201 \pm .003^{\circ}$    | $\frac{.195 \pm .003^{\circ}}{.198 \pm .003^{\circ}}$ | $.198 \pm .003^{\circ}$<br>$.202 \pm .003^{\circ}$  | $.194 \pm .003^{\circ}$<br>.197 ± .003 <sup>\circ</sup> |
|       |                      |      | $\mathbb{S}^{2,-1.0}$   |   | .207 ± .003***                                       | .203 ± .003 <sup>§*</sup>  | .199 ± .003 <sup>§*</sup>                             | .201 ± .003 <sup>§</sup>                              | .197 ± .003 <sup>§</sup>                              | $.201 \pm .003^{\dagger \$}$                        | .197 ± .003*8   |
|       |                      |      | $S^{2,-2,0}$<br>$S^{2}-4,0$   |   | .206 ± .003***                                       | $.205 \pm .003^{*}$  | .201 ± .003   | .201 ± .003   | $\frac{.198 \pm .003^{\$}}{102 \pm .003^{\$}}$        | $.201 \pm .003^{10}$                                | .196 ± .003*  |
|       | Gaussian             | С    | $(S^2)^2$   | 25.4 ± .5   | 35.7 ± .5 <sup>1</sup>                               | 29.6 ± .5  | 32.3 ± .5 <sup>1</sup>                                | 28.9 ± .4 <sup>1*</sup>                               | $\frac{.192 \pm .003}{31.5 \pm .51}$                  | 30.2 ± .5 <sup>13</sup>                             | 33.1 ± .5   |
|       |                      |      | $\mathbb{E}^4$  | 21.2 ± .5 <sup>†§*‡</sup>                           | 34.9 ± .5  | 27.1 ± .5 <sup>¶</sup>   | 30.0 ± .4 <sup>¶*</sup>                               | 27.1 ± .5 <sup>¶</sup>                                | 30.0 ± .4 <sup>¶*</sup>                               | 28.1 ± .5 <sup>¶</sup>                              | <u>31.3 ± .4</u> *¶‡                                    |
|       |                      |      | $\mathbb{H}^2 \mathbb{E}^2$<br>$\mathbb{H}^2 \mathbb{C}^2$                  | $20.2 \pm .5^{10+1}$                                | $40.3 \pm .4^{1}$                                    | $30.8 \pm .5$  | $34.7 \pm .5^{1*}$                                    | $30.9 \pm .5^{10}$                                    | 34.8 ± .5   | $32.0 \pm .5^{11}$                                  | $\frac{36.1 \pm .5}{35.2 \pm .5}$                       |
| _     |                      |      | H <sup>4</sup>  | 15.6 ± .5 <sup>†§*‡</sup>                           | $47.5 \pm .5^{1}$                                    | 32.9 ± .5 <sup>1*</sup>  | 35.9 ± .4 <sup>¶*</sup>                               | $33.0 \pm .4^{1*}$                                    | 37.9 ± .4 <sup>1*</sup>                               | 35.2 ± .5*1   | $\frac{33.2 \pm .3}{40.0 \pm .4}$                       |
| (-K)  |                      |      | $S^2 \mathbb{E}^2$  | 23.7 ± .5***  | 35.8 ± .5  | 29.5 ± .4 <sup>¶</sup>   | 32.7 ± .51  | 29.3 ± .5¶  | 32.4 ± .4¶*   | 30.0 ± .5   | 33.5 ± .4 <sup>¶‡</sup>                                 |
| nulti |                      |      | $(\mathbb{H}^2)^2$  | $18.1 \pm .6^{18^{+3}}$<br>$15.7 \pm .6^{18^{+3}}$  | $33.1 \pm .4^{1}$<br>$41.5 \pm .5^{1}$               | 25.9 ± .5 <sup>13</sup><br>31.7 + .5 <sup>1</sup>                | $27.6 \pm .413$<br>$34.5 \pm .518$                    | $24.0 \pm .4^{11}$<br>31.9 + 5 <sup>11</sup>          | $26.2 \pm .5^{11}$<br>$35.4 \pm .5^{11}$              | $25.1 \pm .5^{1}$<br>$32.3 \pm .5^{1}$              | $\frac{28.0 \pm .5^{11}}{37.0 \pm .5^{11}}$             |
| ic (n |                      | R    | $(S^2)^2$   | 10.7 ± 10   | .196 ± .002***                                       | .196 ± .002  | .191 ± .002 <sup>§</sup>                              | .197 ± .002   | .191 ± .002   | .196 ± .002   | .191 ± .002   |
| uthet |                      |      | E <sup>4</sup><br>112122  |   | .194 ± .003***                                       | .197 ± .003  | $\frac{.191 \pm .002^{\$*}}{.104 \pm .002}$           | .197 ± .003   | $\frac{.191 \pm .002^{\$*}}{.192 \pm .002}$           | .196 ± .003   | .190 ± .002***  |
| Syn   |                      |      | $\mathbb{H}^2 \mathbb{S}^2$   |   | $.190 \pm .003$<br>$.197 \pm .003^{110}$             | $.199 \pm .002$<br>.200 ± .002                                   | $.194 \pm .002$<br>$.194 \pm .002^{\$*}$              | $.199 \pm .002$<br>$.199 \pm .002$                    | .195 ± .002<br>.194 ± .002 <sup>§*</sup>              | .198 ± .002<br>.199 ± .002                          | .195 ± .002<br>.193 ± .002                              |
|       |                      |      | $\mathbb{H}^4$  |   | $.175 \pm .003$                                      | $.189 \pm .003^{\circ}$  | $.184 \pm .003^*$                                     | $.187 \pm .003^*$                                     | .181 ± .003   | .185 ± .003**                                       | .178 ± .003   |
|       |                      |      | $S^2 \mathbb{E}^2$<br>$\mathbb{S}^4$  |   | $.199 \pm .003^{110}$                                | $.200 \pm .003$<br>193 ± 002                                     | $.194 \pm .002^{\$}$<br>188 ± .002 <sup>\$\$</sup>    | $.201 \pm .003$<br>193 $\pm .002$                     | $.195 \pm .003^{\$}$<br>190 ± .002^{18}               | $.199 \pm .003$<br>194 ± .002                       | $.194 \pm .002^{\$}$<br>189 $\pm .002^{\$}$             |
|       |                      |      | $(\mathbb{H}^{2})^{2}$  |   | $.194 \pm .002$                                      | $.193 \pm .002$<br>$.198 \pm .002^*$                             | .193 ± .002*  | $.195 \pm .002$<br>$.198 \pm .002$                    | $.190 \pm .002$                                       | $.194 \pm .002$<br>$.196 \pm .002^{\dagger}$        | $.139 \pm .002^{\circ}$<br>.191 ± .002                  |
|       | CiteSeer             | С    | (S <sup>2</sup> ) <sup>2</sup>  | 13.5 ± .678**                                       | 25.1 ± .61   | 25.3 ± .5  | $\frac{27.0 \pm .7}{27.0 \pm .7}$                     | 25.3 ± .61  | 26.2 ± .71  | 25.8 ± .6   | 27.1 ± .71  |
| ngs   |                      |      | $\mathbb{H}^{4}$<br>$\mathbb{H}^{2}\mathbb{E}^{2}$                          | $13.4 \pm .5^{18^{+}}$<br>$13.4 \pm .4^{18^{+}}$    | $24.1 \pm .4^{1}$<br>$24.7 \pm .7^{1}$               | $24.5 \pm .7^{1}$<br>$25.4 \pm .5^{1}$                           | 25.9 ± .4 <sup>1</sup><br>26.6 ± .6 <sup>1</sup>      | $24.5 \pm .7^{1}$<br>$24.9 \pm .5^{1}$                | 25.0 ± .7 <sup>1</sup><br>25.6 ± .5 <sup>1</sup>      | $23.7 \pm .4^{1}$<br>$25.2 \pm .5^{1}$              | $\frac{25.1 \pm .5^{1}}{26.4 \pm .6^{1}}$               |
| ippa  |                      |      | $\mathbb{H}^2\mathbb{S}^2$  | 13.7 ± .2***  | $25.9 \pm .5^{\parallel}$                            | $25.4 \pm .4^{\circ}$  | $26.1 \pm .7$   | $25.0 \pm .6^{\circ}$                                 | $25.9 \pm .7$   | 25.8 ± .6   | $25.6 \pm .6^{\circ}$                                   |
| šmbé  |                      |      | H <sup>4</sup><br>62172   | 14.1 ± .4 <sup>†§*‡</sup>                           | 26.2 ± .8 <sup>1</sup>                               | $25.5 \pm .8^{1}$  | $26.9 \pm .9^{\circ}$                                 | 24.6 ± .7 <sup>¶</sup>                                | $\frac{26.5 \pm .8}{25.8 \pm .6}$                     | 24.0 ± .4   | 26.0 ± .9 <sup>¶</sup>                                  |
| phe   |                      |      | 3"E"<br>S <sup>4</sup>  | $13.7 \pm .3^{10}$                                  | $25.6 \pm .8^{\circ}$<br>24.7 ± .7 <sup>1</sup>      | $24.9 \pm .6^{\circ}$<br>24.5 ± .9 <sup><math>\circ</math></sup> | $20.3 \pm .0^{\circ}$<br>25.5 ± .7 <sup>1</sup>       | $25.0 \pm .7^{1}$<br>24.4 ± .7 <sup>1</sup>           | $25.8 \pm .6$<br>24.9 ± .7                            | $24.4 \pm .5^{\circ}$<br>$24.2 \pm .7^{\circ}$      | $\frac{25.9 \pm .5}{24.9 \pm .5}$                       |
| . Gra |                      |      | $(\mathbb{H}^{2})^{2}$  | 5.4 ± .2 <sup>†§*‡</sup>                            | 24.9 ± .7¶   | 25.8 ± .6 <sup>1</sup>   | 27.3 ± .51  | 25.3 ± .8¶  | 26.6 ± .5   | 25.4 ± .61  | 26.8 ± .7   |
|       | Cora                 | С    | $(S^2)^2$<br>$\mathbb{R}^4$   | $18.2 \pm 1.2^{+1}$<br>$16.4 \pm 5^{+}$             | 21.4 ± .4 <sup>***</sup><br>21.1 ± .4 <sup>***</sup> | 29.2 ± .5 <sup>81</sup><br>28.4 ± .4 <sup>81</sup>               | 29.9 ± .4 <sup>81</sup><br>29.3 ± .5 <sup>81</sup>    | 29.2 ± .4 <sup>8</sup><br>28.4 + .4 <sup>8</sup>      | $\frac{29.7 \pm .5^{\$1}}{29.7 \pm .6^{\$1}}$         | 27.9 ± 1.2 <sup>81</sup><br>28.9 ± 5 <sup>81</sup>  | 29.6 ± .5 <sup>\$1</sup><br>29.4 ± .6 <sup>\$1</sup>    |
|       |                      |      | $\mathbb{H}^2\mathbb{E}^2$  | 17.3 ± .5 <sup>†*‡</sup>                            | 20.1 ± .6 <sup>†*‡</sup>                             | 28.7 ± .5 <sup>§</sup>   | $29.2 \pm .5^{\$1}$                                   | 28.8 ± .4 <sup>§</sup>                                | $29.4 \pm .5^{\$1}$                                   | 28.7 ± .4 <sup>§</sup>                              | 29.2 ± .5 <sup>§</sup>                                  |
|       |                      |      | $\mathbb{H}^2 \mathbb{S}^2$   | 15.9 ± .2***  | 20.9 ± .5 <sup>***</sup>                             | 28.5 ± .5 <sup>§</sup>   | $\frac{29.8 \pm .6^{\$1}}{29.8 \pm .6^{\$1}}$         | 28.8 ± .4 <sup>§¶</sup>                               | $\frac{29.8 \pm .6^{\$1}}{20.7 \pm .6^{\$1}}$         | 29.0 ± .6 <sup>§</sup>                              | 29.9 ± .5 <sup>§</sup>                                  |
|       |                      |      | $\mathbb{H}^*$<br>$\mathbb{S}^2\mathbb{E}^2$                                | 20.4 ± 1.8°<br>16.7 ± 2 <sup>†§*‡</sup>             | $20.7 \pm .4^{13}$<br>$21.0 \pm 5^{11}$              | $28.5 \pm .6^{8}$<br>$28.0 \pm .6^{81}$                          | 28.9 ± .5 <sup>8</sup><br>28.6 ± .5 <sup>8</sup>      | $28.1 \pm .6^{\circ}$<br>27.8 + 6 <sup>\\$1</sup>     | $28.7 \pm .5^{\circ}$<br>$28.5 \pm .5^{\circ}$        | $27.4 \pm .6^{\circ}$<br>28.5 + $6^{\circ}$         | 28.9 ± .4 <sup>81</sup><br>28.5 ± .5 <sup>81</sup>      |
|       |                      |      | S4  | 16.8 ± .5***  | 20.7 ± .5****  | 27.9 ± .8 <sup>§</sup>   | $28.9 \pm .6^{\$1}$                                   | 27.8 ± .5 <sup>§</sup>                                | 28.9 ± .5 <sup>§</sup>                                | 27.7 ± .7 <sup>§</sup>                              | 28.8 ± .6 <sup>§</sup>                                  |
|       | D-IDI                |      | $(\mathbb{H}^2)^2$  | 5.1 ± .2***   | 20.7 ± .4***   | 29.3 ± .6 <sup>§</sup>   | 30.3 ± .6 <sup>§</sup>                                | 29.6 ± .6 <sup>§</sup>                                | 30.5 ± .6 <sup>§</sup>                                | 29.5 ± .5 <sup>§</sup>                              | 30.3 ± .6 <sup>§</sup>                                  |
|       | POIBIOGS             | C    | (3)))<br>E <sup>4</sup>   | 50.4 ± .8***<br>48.4 ± .8***                        | 93.9 ± .4"   | $\frac{93.2 \pm .5}{92.4 \pm .3}$                                | $93.1 \pm .4^{\circ}$<br>$92.7 \pm .5^{\circ}$        | $90.9 \pm .8^{1}$<br>$92.5 \pm .3^{1}$                | $91.0 \pm .7$<br>$93.1 \pm .3$                        | 69.0 ± 3.3<br>93.3 ± .3                             | 92.9 ± .4"<br>93.6 ± .5 <sup>¶</sup>                    |
|       |                      |      | $\mathbb{H}^2 \mathbb{E}^2$   | 63.6 ± 5.2***                                       | 93.6 ± .4  | 92.9 ± .4¶   | 93.4 ± .3¶  | 92.7 ± .4¶  | 92.9 ± .5   | 93.1 ± .5   | 93.6 ± .2   |
|       |                      |      | $\mathbb{H}^2 \mathbb{S}^2$<br>ur <sup>4</sup>                              | $62.1 \pm 5.6^{\$*1}$                               | $94.0 \pm .5^{1}$                                    | $92.9 \pm .5$<br>$92.1 \pm .7$                                   | $\frac{93.6 \pm .6}{92.5 \pm .5}$                     | $92.2 \pm .5$<br>01.0 + 7                             | $93.2 \pm .6^{1}$                                     | $92.4 \pm .4^{1}$                                   | $93.3 \pm .5^{1}$                                       |
|       |                      |      | $\mathbb{S}^{2}\mathbb{E}^{2}$  | 49.8 ± 2.6***                                       | $93.5 \pm .5^{\circ}$<br>94.3 ± .2 <sup>°</sup>      | $93.4 \pm .6^{1}$  | $92.3 \pm .5^{\circ}$<br>$93.9 \pm .5^{\circ}$        | $91.9 \pm .7$<br>$92.5 \pm 1.1$                       | 93.6 ± .5 <sup>1</sup>                                | $92.7 \pm .5^{\circ}$<br>93.8 ± .6 <sup>1</sup>     | $\frac{75.1 \pm .4}{93.2 \pm .6}$                       |
|       |                      |      | S <sup>4</sup>  | 48.0 ± 1.2***                                       | 93.4 ± .3  | $91.5 \pm .5^{11}$   | $92.9 \pm .5^{1}$                                     | 91.7 ± .8¶  | 92.8 ± .5   | 92.8 ± .3   | $93.1 \pm .4^{1}$                                       |
|       | CS PhDs              | R    | $\frac{(\mathbb{H}^2)^2}{(\mathbb{S}^2)^2}$                                 | 48.0 ± .91878<br>.043 ± .004                        | 93.9 ± .41   | $92.5 \pm .6^{1}$<br>048 + .004                                  | 92.8 ± .7 <sup>1</sup>                                | $92.7 \pm .6^{1}$<br>.046 ± .004                      | 93.3 ± .51<br>.041 ± .003                             | 92.7 ± .4 <sup>1</sup><br>.046 ± .003               | $93.7 \pm .6^{1}$                                       |
|       |                      |      | $\mathbb{E}^4$  | $.035 \pm .003$                                     | $.056 \pm .003$                                      | $.040 \pm .005$  | $.047 \pm .005$                                       | $.040 \pm .005$                                       | $.048 \pm .004$                                       | .038 ± .004   | $.046 \pm .005$   |
|       |                      |      | $\mathbb{H}^2 \mathbb{E}^2$<br>$\mathbb{H}^2 \mathbb{R}^2$                  | $.043 \pm .004$                                     | $.057 \pm .004$<br>048 ± 003                         | $\frac{.040 \pm .004}{.043 \pm .004}$                            | $.051 \pm .003$<br>$.054 \pm .004$                    | .039 ± .005   | $.049 \pm .006$                                       | $.045 \pm .003$                                     | $.044 \pm .004$   |
|       |                      |      | H <sup>4</sup>  | .045 ± .004<br>.044 ± .002                          | .048 ± .005<br>.053 ± .005                           | .045 ± .004<br>.052 ± .002                                       | .054 ± .004<br>.052 ± .005                            | $.045 \pm .005$<br>$.049 \pm .005$                    | .048 ± .008<br>.060 ± .003                            | $.044 \pm .004$<br>.041 ± .004                      | $.055 \pm .005$<br>$.057 \pm .003$                      |
|       |                      |      | $S^2 \mathbb{E}^2$  | .042 ± .005   | .067 ± .006**  | $.040 \pm .002$  | $.050 \pm .004$                                       | .038 ± .004   | .048 ± .007   | <u>.040 ± .004</u>                                  | .053 ± .005   |
|       |                      |      | 3 <sup>*</sup><br>(Ⅲ <sup>2</sup> ) <sup>2</sup>                            | .045 ± .004   | .065 ± .005**  | .040 ± .004 <sup>s</sup><br>.040 ± .003                          | $.050 \pm .004$<br>$.050 \pm .004$                    | $\frac{.041 \pm .003^{\circ}}{.043 \pm .002^{\circ}}$ | $.043 \pm .004$<br>$.052 \pm .006$                    | $.047 \pm .004$<br>$.042 \pm .004$                  | $.048 \pm .003$<br>$.051 \pm .004$                      |
|       | AdjNoun              | LP   | $\mathbb{S}^2 \mathbb{E}^2 \mathbb{H}^2$                                    | 93.7 ± 1.1  | 93.3 ± 1.1   | 93.5 ± .9  | 93.7 ± 1.1  | 93.5 ± .9   | 93.7 ± 1.1  | 93.5 ± .9   | 93.7 ± 1.1  |
|       | Dolphins<br>Football | LP   | $S^2 \mathbb{E}^2 \mathbb{H}^2$<br>$\mathbb{E}^2 \mathbb{E}^2 \mathbb{H}^2$ | 90.7 ± .7 <sup>†*‡</sup>                            | 92.1 ± .6 <sup>†*‡</sup>                             | 96.6 ± .3 <sup>§¶</sup><br>85.7 ± 3.6                            | 92.3 ± .9   | 96.6 ± .3 <sup>§¶</sup><br>85.7 ± 3.6                 | 90.9 ± .8   | 96.6 ± .3 <sup>§¶</sup><br>85.7 ± 3.6               | 90.7 ± .7   |
|       | Karate Club          | LP   | $S^2 \mathbb{E}^2 \mathbb{H}^2$   | 62.0 ± 3.5<br>65.7 ± 10.0                           | 79.8 ± 5.5<br>89.4 ± 1.9                             | 05.7 ± 3.0<br>95.1 ± 1.5   | 63.1 ± 3.5<br>88.6 ± 2.2                              | 05.7 ± 3.0<br>95.1 ± 1.5                              | 65.1 ± 5.5<br>88.8 ± 1.9                              | 05.7 ± 3.0<br>95.1 ± 1.5                            | 82.0 ± 3.3<br>88.8 ± 2.6                                |
|       | Les Mis              | LP   | $S^2 \mathbb{E}^2 \mathbb{H}^2$   | 92.2 ± .9   | 93.8 ± 1.1   | 95.7 ± .7  | 92.7 ± .9   | 95.5 ± .6   | 93.7 ± 1.0  | $\frac{95.6 \pm .8}{25.0 \pm .8}$                   | 92.2 ± .9   |
|       | PolBooks             | LP   | $S^2 \mathbb{E}^2 \mathbb{H}^2$<br>$S^2 \mathbb{E}^2 (\mathbb{H}^2)^3$      | 92.2 ± .6 <sup>7*‡</sup><br>2 8 ± 0 <sup>7§*‡</sup> | 94.8 ± .3  | 95.8 ± .4 <sup>1</sup><br>18.9 ± .4 <sup>8</sup>                 | 92.9 ± .6   | 95.8 ± .4 <sup>1</sup><br>18.7 ± .4 <sup>8</sup>      | 92.1 ± .6   | 95.8 ± .4   | 92.2 ± .6   |
| ш     | CIFAR-100            | č    | $(S^2)^4$   | 5.7 ± .4 <sup>†§*‡</sup>                            | 8.6 ± .4 <sup>*¶*‡</sup>                             | $10.0 \pm .4^{1}$  | <u>11.5 ± .5<sup>§</sup></u>                          | $10.1 \pm .4^{\circ}$                                 | $\frac{11.5 \pm .5^{\$1}}{11.5 \pm .5^{\$1}}$         | 10.8 ± .3 <sup>§</sup>                              | 12.0 ± .3 <sup>§</sup>                                  |
| VAL   | Lymphoma             | С    | $(S^2)^2$   | 78.1 ± .3 <sup>**‡</sup>                            | 77.8 ± 1.4 <sup>***</sup>                            | 81.7 ± 1.2 <sup>§¶*</sup>  | 81.6 ± 1.3 <sup>§</sup>                               | 81.2 ± 1.4 <sup>§*</sup>                              | 81.4 ± 1.3 <sup>§</sup>                               | 83.7 ± 1.2****                                      | $\frac{83.1 \pm 1.2^{\$1}}{20.4 \pm 2.2^{\$1}}$         |
|       | MNIST<br>Landmasses  | C    | 5"E"H"<br>S <sup>2</sup>  | $10.9 \pm 3.0^{18}$                                 | $41.9 \pm 3.7$<br>91.4 ± .2                          | 28.9 ± 1.3<br>81.2 ± .4  | 35.7 ± 2.8"<br>83.5 ± .3 <sup>11</sup>                | $28.6 \pm 1.0^{4}$<br>79.7 ± .3 <sup>11</sup>         | 30.5 ± 2.9*<br>81.8 ± .3***                           | 50.9 ± 1.6 <sup>4</sup><br>83.5 ± .3 <sup>*10</sup> | <u>59.4 ± 2.5</u><br>84.2 ± .3 <sup>13</sup>            |
| Jer   | Neuron 33            | č    | $(S^{1})^{5}$   | 53.9 ± .3***  | 50.5 ± .5 <sup>†*‡</sup>                             | 76.2 ± .4 <sup>§</sup>   | 76.0 ± .5 <sup>§</sup>                                | 76.2 ± .4 <sup>§¶</sup>                               | 75.8 ± .5 <sup>§¶*</sup>                              | 76.0 ± .5 <sup>§</sup>                              | 77.0 ± .4   |
| ē     | Neuron 46            | C    | $(S^1)^5$<br>$C^2 C^1$  | 51.5 ± .1 <sup>†*‡</sup>                            | $50.2 \pm .2^{\dagger * \ddagger}$                   | $60.7 \pm .3^{\$11}$<br>5 200 ± 246 <sup>\$1</sup>               | $\frac{61.1 \pm .3^{\$11}}{4.531 \pm .187^{\$11}}$    | $59.3 \pm .3^{181}$                                   | $59.9 \pm .3^{181}$                                   | 60.8 ± .3 <sup>\$1</sup>                            | 61.2 ± .3 <sup>§</sup>                                  |
|       | Traffic              | R    | $\mathbb{E}^{1}(\mathbb{S}^{1})^{4}$  |   | $1.196 \pm .212$                                     | .521 ± .003 <sup>\$</sup>  | 4.531 ± .187**  | $1.825 \pm .196$<br>.526 ± .003 <sup>†</sup>          | $.515 \pm .003^{\circ}$                               | $1.574 \pm .249$<br>.534 ± .003                     | $.130 \pm .123$<br>.577 ± .005                          |
|       |                      |      | . /   |   |  |  |   |   |   |   |   |

1294

## 1296 H DATASETS AVAILABILITY

| 1304 |   |
|------|---|
| 1005 | Table 6: All of the datasets used in this paper, with download links and citations. CC-BY-SA is short |
| 1305 | for the Creative Commons Attribution-ShareAlike license. Allen TOU is the Allen Institute terms       |
| 1306 | of use found at https://alleninstitute org/terms-of-use/  |
| 1307 | of use, found at neeps . / attentinise reaction of use /.   |

| Dataset        | Link   | License   | Citation                    |
|----------------|--|-----------|-----------------------------|
| CiteSeer       | Network Repository: CiteSeer                     | CC-BY-SA  | Giles et al. (1998)         |
| Cora           | Network Repository: CORA                         | CC-BY-SA  | Sen et al. (2008)           |
| Polblogs       | Network Repository: Polblogs                     | CC-BY-SA  | Adamic & Glance (200        |
| CS PhDs        | Pajek datasets: PhD students in CS               | CC-BY-SA  | Johnson (1984)              |
| Adjnoun        | Network Repository: Adjnoun                      | CC-BY-SA  | Newman (2006)               |
| Dolphins       | Network Repository: Dolphins                     | CC-BY-SA  | Lusseau et al. (2003)       |
| Football       | Network Repository: Football                     | CC-BY-SA  | Girvan & Newman (20         |
| Karate Club    | Network Repository: Karate                       | CC-BY-SA  | Zachary (1977)              |
| Les Mis        | Network Repository: Les Mis                      | CC-BY-SA  | Knuth (1993)                |
| Polbooks       | Network Repository: Polblooks                    | CC-BY-SA  | Krebs (2004)                |
| Blood          | 10x Genomics: CD8+ Cytotoxic T-<br>cells         | CC-BY-SA  | Zheng et al. (2017)         |
| Blood          | CD8+/CD45RA+ Naive Cytotoxic T Cells             | CC-BY-SA  | Zheng et al. (2017)         |
| Blood          | 10x Genomics: CD56+ Natural<br>Killer Cells      | CC-BY-SA  | Zheng et al. (2017)         |
| Blood          | 10x Genomics: CD4+ Helper T<br>Cells             | CC-BY-SA  | Zheng et al. (2017)         |
| Blood          | 10x Genomics: CD4+/CD45RO+<br>Memory T Cells     | CC-BY-SA  | Zheng et al. (2017)         |
| Blood          | 10x Genomics:                                    | CC-BY-SA  | Zheng et al. (2017)         |
| Dioou          | CD4+/CD45RA+/CD25- Naive T                       |           | 2.1011g et ul. (2017)       |
| Dlaad          | CD4+/CD25+ Degulatory T Calls                    | CC DV CA  | Then $\alpha$ at al. (2017) |
| Diood          | 10x Company CD24   Collo                         | CC DV SA  | Zheng et al. $(2017)$       |
| DIOOU<br>D11   | TOX Genomics: CD34+ Cens                         | CC-DI-SA  | Zheng et al. $(2017)$       |
| Blood<br>Dlasd | CD19+ B Cells                                    | CC-BY-SA  | Zheng et al. $(2017)$       |
| BIOOD          | It delain's Lemenhame Disconisted                | CC-BY-SA  | Zheng et al. $(2017)$       |
| Lymphoma       | Tumor: Targeted, Immunology                      | СС-В 1-5А | Tox Genomics (2020a)        |
| Lumphama       | Fallel<br>DDMCa from a Haalthy Donor             | CC DV CA  | $10\pi$ Companying (2020b)  |
| Lymphoma       | Torrected Compare Immunology                     | СС-Б І-ЗА | Tox Genomics (2020b)        |
|                | Targeted-Compare, Ininunology                    |           |                             |
| MAUCT          | Panel  | MIT       | $I_{1}$                     |
| MINIST         | HuggingFace: MINIST                              | MIII      | Lecun et al. (1998)         |
| CIFAR-100      | HuggingFace: CIFAR-100                           | None      | Kriznevsky (2009)           |
| Landmasses     | Basemap 1.4.1: 1s_land                           | None      | Inone                       |
| Neurons        | Allen Brain Atlas                                | Allen TOU | Jones et al. $(2009)$       |
| Iemperature    | wikipedia: List of cities by average temperature | СС-ВҮ-ЅА  | wikipedia (2024)            |
| Traffic        | Kaggle: Traffic Prediction Dataset               | None      | Fedesoriano (2020)          |

# 1350 I COMPARISON TO NEURAL BASELINES

# 1352 I.1 MODELS

Neural networks, especially graph neural networks, are a popular choice for representing and working with mixed-curvature representations (Sun et al., 2021; Cho et al., 2023; Bachmann et al., 2020;
McNeela et al., 2024). Following the typical node classification approach described in the literature,
we compared our method to both deep neural networks (ML) and graph neural networks (GNN).
We find that our methods are generally competitive with neural baselines, especially when less data is available.

1360 We trained each model for 200 epochs using Adam Kingma & Ba (2017) with a learning rate of 1361 .01,  $\beta_1 = 0.9$ , and  $\beta_2 = -0.999$ . Each model used one hidden dimension equal to its ambient 1362 dimension. For classification, we used an output dimension equal to the number of classes and 1363 cross-entropy loss; for regression, we used a single output dimension and mean squared error loss.

We additionally train tangent plane variants of the MLP and GNN models, in which the data is preprocessed using a logarithmic map and subsequently treated as Euclidean. For datasets where the graph topology is not provided (e.g. Gaussian mixtures), we take the pairwise distances in the manifold geometry and transform them into dense weighted edges using the Gaussian kernel:

$$\mathbf{D}_{i,j} = \delta_{\mathcal{P}}(\mathbf{x}_i, \mathbf{x}_j) \tag{73}$$

$$\mathbf{A} = \exp(-\mathbf{D}) \tag{74}$$

For graph datasets with known topology, we instead used the true adjacency matrix and an ablation
in which the adjacency matrix is replaced by the dense Gaussian kernel estimate. Interestingly, for
several of the graph datasets, substituting the true adjacencies with a Gaussian kernel on embedding
distances did not substantially hinder performance; however, it rarely helped.

#### 1376 I.2 DATASETS

1368 1369

1370

1375

1377

1382

1383

1389

1390

1378 Due to time considerations, we ran our benchmarks on a representative sample of the datasets. For 1379 adjacency-free datasets, we chose Gaussian mixtures in single-component signatures (as in Figures 3 1380 and 4) and multiple-component signatures (as in Tables 2 and 3). We also ran each graph dataset on 1381 the lowest  $D_{avg}$  signature (i.e. the signatures reported in Tables 2 and 3).

#### I.3 RESULTS

We tabulate our results as follows: Tables 7 and 9 contain classification and regression benchmarks for single-*K* manifolds; Tables 8 and 10 contain classification and regression benchmarks for product spaces; and Table 11 contains graph dataset benchmarks. For all benchmarks except graph data, we beat both MLPs and GNNs; for datasets with informative graph topologies, GNNs are a sensible alternative.

Table 7: Comparison to neural networks on the constant-curvature classification task.

| K     | Product DT      | Product RF      | $T_{\mu_0} \mathcal{P} \operatorname{MLP}$ | MLP             | $T_{\mu_0} \mathcal{P} \operatorname{GNN}$ | GNN             |
|-------|-----------------|-----------------|--|-----------------|--|-----------------|
| -4    | $0.34 \pm 0.10$ | $0.36 \pm 0.07$ | $0.27 \pm 0.08$                            | $0.17 \pm 0.10$ | $0.23 \pm 0.18$                            | $0.18 \pm 0.15$ |
| -2    | $0.36 \pm 0.10$ | $0.37 \pm 0.11$ | $0.29 \pm 0.12$                            | $0.21 \pm 0.16$ | $0.29 \pm 0.10$                            | $0.21 \pm 0.11$ |
| -1    | $0.32 \pm 0.09$ | $0.34 \pm 0.08$ | $0.26 \pm 0.14$                            | $0.23 \pm 0.09$ | $0.27 \pm 0.11$                            | $0.22 \pm 0.15$ |
| -0.5  | $0.32 \pm 0.09$ | $0.33 \pm 0.10$ | $0.27 \pm 0.10$                            | $0.29 \pm 0.08$ | $0.28 \pm 0.07$                            | $0.23 \pm 0.12$ |
| -0.25 | $0.30 \pm 0.11$ | $0.32 \pm 0.11$ | $0.26 \pm 0.11$                            | $0.25 \pm 0.08$ | $0.26 \pm 0.09$                            | $0.21 \pm 0.08$ |
| 0     | $0.29 \pm 0.11$ | $0.31 \pm 0.08$ | $0.23 \pm 0.11$                            | $0.26 \pm 0.12$ | $0.24 \pm 0.05$                            | $0.24 \pm 0.05$ |
| 0.25  | $0.28 \pm 0.09$ | $0.30 \pm 0.11$ | $0.26 \pm 0.11$                            | $0.28 \pm 0.09$ | $0.27 \pm 0.09$                            | $0.21 \pm 0.14$ |
| 0.5   | $0.25 \pm 0.10$ | $0.29 \pm 0.11$ | $0.26 \pm 0.11$                            | $0.25 \pm 0.13$ | $0.25 \pm 0.12$                            | $0.22 \pm 0.14$ |
| 1     | $0.27 \pm 0.07$ | $0.29 \pm 0.05$ | $0.21 \pm 0.07$                            | $0.22 \pm 0.07$ | $0.22 \pm 0.07$                            | $0.21 \pm 0.07$ |
| 2     | $0.26 \pm 0.07$ | $0.29 \pm 0.07$ | $0.23 \pm 0.09$                            | $0.26 \pm 0.09$ | $0.23 \pm 0.12$                            | $0.23 \pm 0.12$ |
| 4     | $0.25 \pm 0.11$ | $0.26 \pm 0.08$ | $0.22 \pm 0.06$                            | $0.21 \pm 0.07$ | $0.21 \pm 0.09$                            | $0.22 \pm 0.07$ |
|       |                 |                 |  |                 |  |                 |

|                             | I                          |                 |  |                 |  |      |
|-----------------------------|----------------------------|-----------------|--|-----------------|--|------|
| $\mathcal{P}$               | Product DT                 | Product RF      | $T_{\mu_0} \mathcal{P} \operatorname{MLP}$ | MLP             | $T_{\mu_0} \mathcal{P} \operatorname{GNN}$ | GNI  |
| $\mathbb{E}^4$              | $0.31 \pm 0.05$            | $0.35 \pm 0.07$ | $0.30 \pm 0.11$                            | $0.31 \pm 0.11$ | $0.27 \pm 0.12$                            | 0.27 |
| $\mathbb{H}^4$              | $0.39 \pm 0.08$            | $0.43 \pm 0.10$ | $0.38 \pm 0.11$                            | $0.31 \pm 0.09$ | $0.32 \pm 0.13$                            | 0.22 |
| $\mathbb{H}^2\mathbb{E}^2$  | $\overline{0.34 \pm 0.11}$ | $0.38 \pm 0.11$ | $0.35 \pm 0.13$                            | $0.34 \pm 0.14$ | $0.30 \pm 0.13$                            | 0.23 |
| $(\mathbb{H}^2)^2$          | $0.35 \pm 0.09$            | $0.36 \pm 0.08$ | $0.35 \pm 0.06$                            | $0.34 \pm 0.06$ | $0.28 \pm 0.07$                            | 0.23 |
| $\mathbb{H}^2 \mathbb{S}^2$ | $\overline{0.31 \pm 0.09}$ | $0.36 \pm 0.10$ | $0.33 \pm 0.10$                            | $0.35 \pm 0.09$ | $0.27 \pm 0.09$                            | 0.22 |
| $\mathbb{S}^4$              | $0.25 \pm 0.06$            | $0.30 \pm 0.05$ | $0.25 \pm 0.07$                            | $0.23 \pm 0.08$ | $0.20 \pm 0.08$                            | 0.20 |
| $\mathbb{S}^2 \mathbb{E}^2$ | $0.29 \pm 0.09$            | $0.32 \pm 0.08$ | $0.30 \pm 0.09$                            | $0.31 \pm 0.11$ | $0.25 \pm 0.06$                            | 0.24 |
| $(S^2)^2$                   | $0.30 \pm 0.08$            | $0.34 \pm 0.10$ | $0.33 \pm 0.08$                            | $0.34 \pm 0.10$ | $0.23 \pm 0.09$                            | 0.18 |

Table 8: Comparison to neural networks on the mixed-curvature classification task.

Table 9: Comparison to neural networks on the single- and mixed-curvature regression tasks.

| K     | Product DT      | Product RF      | $T_{\mu_0} \mathcal{P} \operatorname{MLP}$ | MLP             | $T_{\mu_0} \mathcal{P} \operatorname{GNN}$ | GNN             |
|-------|-----------------|-----------------|--|-----------------|--|-----------------|
| -4    | $0.20 \pm 0.04$ | $0.19 \pm 0.04$ | $0.55 \pm 0.13$                            | $0.55 \pm 0.13$ | $0.55 \pm 0.13$                            | $0.55 \pm 0.13$ |
| -2    | $0.21 \pm 0.04$ | $0.20\pm0.04$   | $0.55 \pm 0.08$                            | $0.55 \pm 0.08$ | $0.55 \pm 0.08$                            | $0.55 \pm 0.08$ |
| -1    | $0.19 \pm 0.03$ | $0.19 \pm 0.03$ | $0.54 \pm 0.09$                            | $0.54 \pm 0.09$ | $0.54 \pm 0.09$                            | $0.54 \pm 0.09$ |
| -0.5  | $0.19 \pm 0.05$ | $0.18\pm0.05$   | $0.53 \pm 0.12$                            | $0.53 \pm 0.12$ | $0.53 \pm 0.12$                            | $0.53 \pm 0.12$ |
| -0.25 | $0.20 \pm 0.02$ | $0.19 \pm 0.02$ | $0.53 \pm 0.10$                            | $0.53 \pm 0.10$ | $0.53 \pm 0.10$                            | $0.53 \pm 0.10$ |
| 0     | $0.21 \pm 0.03$ | $0.20\pm0.03$   | $0.55 \pm 0.11$                            | $0.55 \pm 0.11$ | $0.55 \pm 0.11$                            | $0.55 \pm 0.11$ |
| 0.25  | $0.21 \pm 0.03$ | $0.20\pm0.03$   | $0.54 \pm 0.08$                            | $0.54 \pm 0.08$ | $0.54 \pm 0.08$                            | $0.54 \pm 0.08$ |
| 0.5   | $0.20 \pm 0.06$ | $0.19 \pm 0.05$ | $0.53 \pm 0.09$                            | $0.53 \pm 0.09$ | $0.53 \pm 0.09$                            | $0.53 \pm 0.09$ |
| 1     | $0.20 \pm 0.04$ | $0.20\pm0.04$   | $0.56 \pm 0.09$                            | $0.56 \pm 0.09$ | $0.56 \pm 0.09$                            | $0.56 \pm 0.09$ |
| 2     | $0.21\pm0.06$   | $0.21 \pm 0.06$ | $0.55 \pm 0.11$                            | $0.55 \pm 0.11$ | $0.55 \pm 0.11$                            | $0.55 \pm 0.11$ |
| 4     | $0.20\pm0.03$   | $0.20 \pm 0.03$ | $0.50\pm0.10$                              | $0.50\pm0.10$   | $0.50\pm0.10$                              | $0.50 \pm 0.10$ |
|       |                 |                 |  |                 |  |                 |

Table 10: Comparison to neural networks on the multi-curvature regression task.

| $\mathcal{P}$               | Product DT                  | Product RF      | $T_{\mu_0} \mathcal{P} \operatorname{MLP}$ | MLP             | $T_{\mu_0} \mathcal{P} \operatorname{GNN}$ | GNN             |
|-----------------------------|-----------------------------|-----------------|--|-----------------|--|-----------------|
| $\mathbb{E}^4$              | $0.27 \pm 0.04$             | $0.20 \pm 0.05$ | $0.56 \pm 0.09$                            | $0.56 \pm 0.09$ | $0.56 \pm 0.09$                            | $0.56 \pm 0.09$ |
| $\mathbb{H}^4$              | $0.24 \pm 0.06$             | $0.18\pm0.04$   | $0.54 \pm 0.10$                            | $0.54 \pm 0.10$ | $0.54 \pm 0.10$                            | $0.54 \pm 0.10$ |
| $\mathbb{H}^2\mathbb{E}^2$  | $0.26 \pm 0.05$             | $0.20\pm0.03$   | $0.54 \pm 0.12$                            | $0.54 \pm 0.12$ | $0.54\pm0.12$                              | $0.54\pm0.12$   |
| $(\mathbb{H}^2)^2$          | $0.24 \pm 0.06$             | $0.18 \pm 0.04$ | $0.52 \pm 0.11$                            | $0.52 \pm 0.11$ | $0.52 \pm 0.11$                            | $0.52 \pm 0.11$ |
| $\mathbb{H}^2 \mathbb{S}^2$ | $0.26 \pm 0.06$             | $0.20\pm0.04$   | $0.52 \pm 0.11$                            | $0.52 \pm 0.11$ | $0.52 \pm 0.11$                            | $0.52 \pm 0.11$ |
| $\mathbb{S}^4$              | $0.25 \pm 0.02$             | $0.19\pm0.03$   | $0.53 \pm 0.10$                            | $0.53 \pm 0.10$ | $0.53 \pm 0.10$                            | $0.53 \pm 0.10$ |
| $\mathbb{S}^2 \mathbb{E}^2$ | $0.26 \pm 0.08$             | $0.19 \pm 0.06$ | $0.51 \pm 0.10$                            | $0.51 \pm 0.10$ | $0.51 \pm 0.10$                            | $0.51 \pm 0.10$ |
| $(S^2)^2$                   | $\underline{0.28 \pm 0.06}$ | $0.20 \pm 0.04$ | $0.54 \pm 0.11$                            | $0.54 \pm 0.11$ | $0.54 \pm 0.11$                            | $0.54 \pm 0.11$ |

Table 11: Comparison to neural networks on graph node classification/regression tasks (Citeseer, Cora, Polblogs are classification tasks; CS PhDs is a regression task).

| _ |          |                 |                             |   |                 |  |                 |
|---|----------|-----------------|-----------------------------|---|-----------------|--|-----------------|
|   | Dataset  | Product DT      | Product RF                  | $T_{\mu_0}\mathcal{P} \operatorname{MLP}$ | MLP             | $T_{\mu_0} \mathcal{P} \operatorname{GNN}$ | GNN             |
|   | Citeseer | $0.23 \pm 0.04$ | $0.24 \pm 0.04$             | $0.23 \pm 0.03$                           | $0.24 \pm 0.02$ | $0.24 \pm 0.03$                            | $0.24 \pm 0.03$ |
|   | Cora     | $0.19 \pm 0.03$ | $0.22 \pm 0.04$             | $0.29 \pm 0.03$                           | $0.29 \pm 0.03$ | $0.29 \pm 0.03$                            | $0.29 \pm 0.03$ |
| ] | Polblogs | $0.89 \pm 0.19$ | $\underline{0.93 \pm 0.02}$ | $0.85\pm0.25$                             | $0.94 \pm 0.03$ | $0.88 \pm 0.17$                            | $0.82 \pm 0.38$ |
|   |          |                 |                             |   |                 |  |                 |

| 1458 | Table 12: Complexity comparison of machine learning models where: n: number of samples, d:        |
|------|---|
| 1459 | number of features, h: neurons per layer, L: number of layers, D: maximum tree depth, s: number   |
| 1460 | of support vectors. We include the complexity of computing pairwise distance, which are necessary |
| 1461 | for operating models like k-nearest neighbors and GNNs without topologies, as well.               |

|         |  |       | Tii            | me            | Space                   |        |
|---------|--|-------|----------------|---------------|-------------------------|--------|
|         | Model  | Phase | Worst          | Avg           | Worst                   | Avg    |
|         | Dists  |       | $n^2d$         | $n^2d$        | $n^2$                   | $n^2$  |
|         | MID  | Train | $ndh + Lnh^2$  | $ndh + Lnh^2$ | $nd + dh + L(h^2 + nh)$ | $h^2L$ |
|         | WILF   | Test  | $h^2L$         | $h^2L$        | $h^2L$                  | $h^2L$ |
| D       | arcontron  | Train | nd             | nd            | d                       | d      |
| 10      | erception  | Test  | d              | d             | d                       | d      |
|         | SVM  | Train | $n^3d$         | $n^2d$        | $n^2$                   | $n^2$  |
|         | 5 1 11   | Test  | sd             | sd            | sd                      | sd     |
|         | GNN  | Train | $n^2d$         | $n^2d$        | $n^2$                   | $n^2$  |
|         | UNIN   | Test  | $n^2$          | $n^2$         | $n^2$                   | $n^2$  |
|         | 1- NN  | Train | 1              | 1             | nd                      | nd     |
|         | K-ININ   | Test  | $nd + n\log n$ | $\log n$      | nd                      | nd     |
| Der     | rision Tree  | Train | Dnd            | Dnd           | $2^D$                   | $2^D$  |
| Du      |  | Test  | D              | D             | 1                       | 1      |
| Drodu   | otDT (vonillo)   | Train | Dnd            | Dnd           | $2^D$                   | $2^D$  |
| Flouud  | addi (vanna)   | Test  | d + D          | d + D         | d                       | d      |
| Draduat | $\mathbf{DT}\left(\begin{pmatrix} d \end{pmatrix} a \mathbf{n} \mathbf{l} \mathbf{t} \mathbf{s} \right)$ | Train | $Dnd^2$        | $Dnd^2$       | $2^D$                   | $2^D$  |
| Product | $D_1((2) \text{ spins})$   | Test  | $d^2 + D$      | $d^2 + D$     | $d^2$                   | $d^2$  |
|         |  |       |                |               |                         |        |

#### J **RUNTIMES AND COMPLEXITY**

#### We summarize complexities for models used in this paper, as well as the pairwise distance prepro-cessing necessary for operations such as computing nearest neighbors and creating reasonable graph edges for GNNs, in Table 12. Complexity estimates are adapted from Virgolin (2021).

To see that the training time complexity of ProductDT is O(Dnd), observe that we must first pre-process the data into angles, which takes O(nd) operations. From there, the angular comparison is a constant-time modification to the decision tree algorithm, so the complexity of ProductDT is O(nd + Dnd) = O(nd). For inference, asymptotic performance is slightly slower than decision trees because preprocessing an input requires O(d) operations. 

If using all  $\binom{d}{2}$  2-D projections, training time complexities are all multiplied by d, and the  $O(d^2)$ preprocessing step is added to test time complexities. 





## <sup>1512</sup> K INTERPRETABILITY AND VISUALIZATION

1513 1514

1515

1516

1517

1518 1519

1520

1521

1522

1530

1531

Alongside their demonstrated accuracy and efficiency, decision tree algorithms are attractive for their tractability and interpretability. In particular, given a trained decision tree T, it is possible to:

- 1. Predict its behavior on the entire space of possible inputs (equivalently:  $\mathcal{T}$  partitions  $\mathcal{P}$  in a tractable way).
- 2. Determine the importance of features (for classic decision trees) or feature pairs/components (for ProductDT) by observing how often and how early a feature(/pair/component) is used in the decision tree procedure. Heuristically, early-splitting features are more important.
  - 3. Visualize every node using a 2-dimensional projection of the input data and angle

# 1525 K.1 SUBMANIFOLD-LEVEL ATTRIBUTION EXPERIMENT

To determine whether our method could accurately distinguish between relevant and irrelevant submanifolds, we drew independent samples from Gaussian mixtures in  $\mathbb{H}^2$ ,  $\mathbb{E}^2$ , and  $\mathbb{S}^2$ , and yielding datasets  $(\mathbf{X}_{\mathbb{H}}, \mathbf{y}_{\mathbb{H}}), (\mathbf{X}_{\mathbb{R}}, \mathbf{y}_{\mathbb{R}}), (\mathbf{X}_{\mathbb{S}}, \mathbf{y}_{\mathbb{S}})$ . We then concatenated these embeddings together:

$$\mathbf{X}_{\mathcal{P}} = \mathbf{X}_{\mathbb{H}} \oplus \mathbf{X}_{\mathbb{R}} \oplus \mathbf{X}_{\mathbb{S}}.$$

(75)

We trained three separate decision tree models on  $\mathbf{X}_{\mathcal{P}}$ , using  $\mathbf{y}_{\mathbb{H}}, \mathbf{y}_{\mathbb{E}}$ , and  $\mathbf{y}_{\mathbb{S}}$  as labels. Because the labels and embeddings were drawn independently, it should be the case that only the component from the same manifold as the labels contains any relevant information, and the other two components are simply noise. Therefore, measuring the fraction of splits that fall in the "correct" manifold is a useful proxy for understanding tree models' ability to pick out signal that happens in individual component manifolds.

1538 Our results are summarized in Table 13. We found that both product space and ambient decision 1539 trees perform well at this task, which is to be expected.

We note that this analysis is unique to tree methods, where the split dimensions are part of the architecture; other methods, such as perceptrons, *k*-nearest neighbors, or neural networks are harder to query for feature(/component) importances. Therefore, we consider this simple experiment a useful demonstration of how decision tree learning can reveal aspects of structure in mixed-curvature datasets that other learning algorithms cannot reveal.

Table 13: Intepretability outcomes for Gaussian mixture. Percentages reflect the proportion of splits in the trained decision tree which fell in the non-spurious component manifold.

| Model     | $\mathbb{H}^2$ | $\mathbb{E}^2$ | $\mathbb{S}^2$ |
|-----------|----------------|----------------|----------------|
| oduct DT  | 100%           | 83%            | 86%            |
| mbient DT | 100%           | 83%            | 67%            |

1551 1552 1553

1545

1548 1549 1550

#### 1554 K.2 VISUALIZATION

A trained tree gives us all of the information we need to visualize the data and how it is split at every node, since each node looks at a 2-dimensional projection. We display three levels of a decision tree with a max depth of 3 in Figure 11. Note that, in this case, the decision tree also gives us relevant information about which 2-dimensional projections are worth looking at on the basis of their feature importances.

1560

- 1561
- 1502
- 1563 1564

1565



Figure 11: An example of a visualized decision tree for a Gaussian mixture in  $\mathcal{P} = \mathbb{S}^4 \mathbb{H}^4$ . Greyedout points are discarded "earlier" in the tree.