# Strain-Minimizing Hyperbolic Network Embeddings with Landmarks

Martin Keller-Ressel[1,*] and Stephanie Nargang[1,†]

[1]TU Dresden, Institute for Mathematical Stochastics, Dresden, 01062, Germany
[*]martin.keller-ressel@tu-dresden.de
[†]stephanie.nargang@tu-dresden.de

July 15, 2022

### Abstract

We introduce L-hydra (landmarked hyperbolic distance recovery and approximation), a method for embedding network- or distance-based data into hyperbolic space, which requires only the distance measurements to a few 'landmark nodes'. This landmark heuristic makes L-hydra applicable to large-scale graphs and improves upon previously introduced methods. As a mathematical justification, we show that a point configuration in $d$-dimensional hyperbolic space can be perfectly recovered (up to isometry) from distance measurements to just $d + 1$ landmarks. We also show that L-hydra solves a two-stage strain-minimization problem, similar to our previous (unlandmarked) method 'hydra'. Testing on real network data, we show that L-hydra is an order of magnitude faster than existing hyperbolic embedding methods and scales linearly in the number of nodes. While the embedding error of L-hydra is higher than the error of existing methods, we introduce an extension, L-hydra+, which outperforms existing methods in both runtime and embedding quality.

## 1 Introduction

Embeddings of networks and distance-based data into hyperbolic geometry have received substantial attention in the recent decade. Such embeddings have been used for link prediction [15, 16], visualization [19], and community detection [16, 14] in networks. In addition to providing insight into the tradeoff between popularity and similarity effects in network growth [15], they have interesting implications for routing, network navigability [11, 1] and efficient computation of shortest paths [21, 5]. In [10] we have introduced `hydra` (hyperbolic distance recovery and approximation) and `hydra+` as efficient methods to compute such embeddings. Similar to classic multidimensional scaling in Euclidean space, `hydra` uses an Eigendecomposition technique to minimize *strain*, a functional based on the hyperbolic Gram matrix of embedded points. Due to this technique, `hydra` provides an efficient alternative to methods like `Rigel` of [21] and `HyPy` of [5], which minimize *stress*, i.e., the least-squares embedding error, by numerical optimization. However, as `hydra` requires the full shortest path matrix of the input network, it can not be scaled to large and very large networks beyond the order of 100.000 nodes. `Rigel` and `HyPy` on the other side use a clever landmark heuristic (see also [6]), which replaces the full shortest path matrix by

only the shortest paths terminating at a small number of (randomly selected) landmark nodes. Here, we show that this landmark heuristic can also be applied to the method of strain minimization, and introduce `L-hydra` (Landmarked hyperbolic distance recovery and approximation) as a new hyperbolic embedding method for very large networks. On the theoretical side, we show that `L-hydra` perfectly recovers (up to isometry) any $n$-point configuration in $d$-dimensional hyperbolic space from its distances measured to just $d + 1$ landmarks. This result provides, for the first time, a theoretical basis to the use of landmark methods for hyperbolic network embeddings. In particular it shows that the required number of landmarks depends only on the intrinsic dimension $d$ of the network, but not on the number $n$ of vertices; hence the computational load of `L-hydra` scales only *linearly* with the number $n$ of vertices. Finally, we present in Section 4 numerical results for `L-hydra` and its extension `L-hydra+`, showing that they provide an improvement over `HyPy` in both runtime and embedding error.

## 2  Background

### 2.1  Hyperbolic space

Our embedding method is formulated in the mathematical framework of the $d$-dimensional hyperboloid model of hyperbolic geometry (cf. [17, 3]). For $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^{d+1}$ the *Lorentz product*, an indefinite inner product, is defined by

$$\boldsymbol{x} \circ \boldsymbol{y} := x_1 y_1 - (x_2 y_2 + \ldots + x_{d+1} y_{d+1}). \tag{1}$$

The real vector space $\mathbb{R}^{d+1}$ equipped with this inner product is called *Lorentz space* and denoted by $\mathbb{R}^{1,d}$. It contains, as subset, the *positive Lorentz space* $\mathbb{R}^{1,d}_+ = \{\boldsymbol{x} \in \mathbb{R}^{1,d} : x_1 > 0\}$. Within $\mathbb{R}^{1,d}_+$, the single-sheet hyperboloid $\mathcal{H}_d$ is given by

$$\mathcal{H}_d = \{\boldsymbol{x} \in \mathbb{R}^{1,d} : \boldsymbol{x} \circ \boldsymbol{x} = 1, x_1 > 0\}. \tag{2}$$

The *hyperboloid model* in dimension $d$ with curvature $-\kappa$, ($\kappa > 0$), consists of $\mathcal{H}_d$ endowed with the hyperbolic distance

$$\mathrm{d}_H^\kappa(\boldsymbol{x}, \boldsymbol{y}) = \frac{1}{\sqrt{\kappa}} \operatorname{arcosh}(\boldsymbol{x} \circ \boldsymbol{y}), \qquad \boldsymbol{x}, \boldsymbol{y} \in \mathcal{H}_d. \tag{3}$$

The hyperbolic distance $\mathrm{d}_H^\kappa$ is a distance on $\mathcal{H}_d$ in the usual mathematical sense; in particular it takes only positive values and satisfies the triangle inequality, cf. [17, §3.2]. In fact, equipped with the metric tensor $ds^2 = \frac{1}{\kappa}(d\boldsymbol{x} \circ d\boldsymbol{x})$, the hyperboloid $\mathcal{H}_d$ becomes a Riemannian manifold of constant sectional curvature $-\kappa$ and $\mathrm{d}_H^\kappa$ is exactly its geodesic distance. Note that the curvature parameter does not enter into the description of the hyperboloid $\mathcal{H}_d$, but only in the distance metric $\mathrm{d}_H^\kappa$. Just as Euclidean space is the canonical model of geometry with zero curvature, hyperbolic space is the canonical model of geometry with negative curvature. Aside from the hyperbolid model, other equivalent models of hyperbolic geometry exist, including the Poincaré ball model and the upper half-space model, see [17, §4.2, §4.6] and [3].

### 2.2  Embedding of distances and networks

To formulate the problem of embedding network or other data into hyperbolic space, let some objects $\boldsymbol{o}_1, \ldots, \boldsymbol{o}_n$ be given. Let $D = [d_{ij}] \in \mathbb{R}^{n \times n}_{\geqslant 0}$ be a symmetric matrix with zero diagonal, which represents the induced pairwise dissimilarities between the objects. The goal is to find a low dimensional coordinate representation $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n$ of the objects in hyperbolic space $\mathcal{H}_d$, such that the hyperbolic distances

between the coordinate representations approximate the given dissimilarities as closely as possible, i.e., such that

$$d_H^\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) \approx d_{ij}. \tag{4}$$

In Euclidean space, such approximations are well studied and can be calculated e.g. by multidimensional scaling (MDS), see also [2]. A survey of the hyperbolic setting, discussing in particular the properties of hyperbolic distance matrices, is given in [18].

An important special case is the *network embedding problem*, where a (unweighted, undirected) graph $G = (V, E)$ is given. In this setting, the objects $\boldsymbol{o}_1, \ldots, \boldsymbol{o}_n$ are given by the vertices $v_1, \ldots, v_n$ of $G$ and as dissimilarities $d_{ij}$ between two vertices $v_i$ and $v_j$ for $i, j \in V$ we consider the length of the shortest path between $i$ and $j$, i.e., $D = [d_{ij}]$ is the graph distance matrix of $G$. We assume that the graph is connected, i.e., that the shortest path distances $d_{ij}$ are finite for any pair of vertices $v_i$ and $v_j$.

While the volume of Euclidean space expands polynomially, exponential expansion can be observed in hyperbolic space, see e.g. [8]. Because of this property, hyperbolic space is expected to give a better representation for hierarchical or tree-like structures than Euclidean geometry, and therefore is a favorable target space for graph embeddings, see e.g [8, 11].

## 2.3 Landmark-based network embeddings

For large graphs or data sets, embedding methods encounter several limitations: First, the embedding method may become computationally too costly; second, the dissimilarity matrix $D$ may become too large to hold in memory; and third, the pre-computation of the dissimilarities $d_{ij}$ themselves may become infeasible. In the network embedding problem, for instance, computation of the full distance matrix of an unweighted, undirected graph requires $\mathcal{O}(|V|^3)$ computations with the Floyd-Warshall algorithm ([7]).*
To alleviate this problem, *landmark-based embedding methods* can be applied; see [6] for an application to MDS. For these methods, a comparatively small number of nodes, indexed by $L \subset V$ are designated as *landmarks*. These landmarks are embedded into the target space, based on their pairwise distances. The remaining non-landmark points (indexed by $N = V \setminus L$) are then embedded by 'triangulation', i.e., based on their distances with respect to the landmarks, but not to each other. Overall, only all shortest-paths sourced at landmark nodes have to be computed, reducing the cost of distance-computations to $\mathcal{O}(|L| \cdot |V|^2)$. This represents an asymptotic increase in efficiency, *if $|L|$ scales sublinear with $|V|$*, i.e. if the number of landmarks does not need to be scaled up proportional to $|V|$. In Theorem 3.1 we show that the number of landmarks required to recover point configurations in hyperbolic space $\mathcal{H}_d$ depends only on the dimension $d$, not on the number of points. In this way, we provide a theoretical basis for the application of landmark methods to hyperbolic embeddings.

## 2.4 Stress- vs. strain-based embeddings

Stress-based embeddings solve the embedding problem (4) by minimizing the squared *stress functional*

$$\text{Stress}(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)^2 = \sum_{i,j=1}^{n} (d_{ij} - d_H^\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j))^2. \tag{5}$$

Equivalent loss functions are the root-mean-square error

$$RMSE = \sqrt{\frac{1}{n(n-1)} \text{Stress}(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)^2} \tag{6}$$

---

*Other methods may be more efficient under additional assumptions on the graph structure, e.g., under bounds on the number of edges.

and the relative embedding error

$$REE = \text{Stress}(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)/\sqrt{\sum_{i,j=1}^{n} d_{ij}}, \tag{7}$$

which can be compared across different dissimilarity matrices $D$ and input sizes $n$. Due to the properties of the hyperbolic distance $\text{d}_H^\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j)$, minimizing (5) is a challenging non-convex optimization problem, which has to be solved by numerical minimization with no guarantee of convergence to the global minimum. Approaches by gradient descent and neural-network-based minimization are given in [19, 4]. Stress-minimization has been combined with a landmark-based approach and with advanced optimization methods by [21] and [5] under the names of `Rigel` and `HyPy` respectively. While `Rigel` uses the derivative-free Nelder–Mead simplex optimization method, `HyPy` provides an improvement by replacing the simplex optimization with iterative quasi-Newton minimization, for which efficient routines such as LBFGS [22] can be used and supplied with the analytic gradient of stress, given in [5, Eqs. (3.1),(3.2)]. As an alternative, in our previous work [10], we have introduced `hydra`, a strain-based hyperbolic embedding method, based on minimization of the squared *strain functional*

$$\text{Strain}(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)^2 := \sum_{i,j}(\cosh(\sqrt{\kappa}\, d_{ij}) - (\boldsymbol{x}_i \circ \boldsymbol{x}_j))^2. \tag{8}$$

The strain functional is obtained from the stress functional by the transformation of all distances by hyperbolic cosine. As shown in [10], `hydra` has two important theoretical properties. The algorithm recovers any configuration of points in $d$-dimensional hyperbolic space up to isometry and it is guaranteed to return the globally optimal solutions of the strain-minimization problem (8). In numerical experiments, `hydra` is faster than `HyPy` and `Rigel` by several orders of magnitudes, see [10]. While these are attractive properties, the strain functional (8) is not as interpretable as the stress functional (5), which aligns with the popular 'least-squares' paradigm of regression analysis and numerical approximation. Therefore, if we accept RMSE/REE as the target measure of embedding quality, there are two practical uses of the `hydra` method, cf. [10]:

- Use `hydra` as a stand-alone method. This gives a very fast embedding method (several orders of magnitude faster than `HyPy`/`Rigel`), but with a relative embedding error (REE) that is approx. 1.0 to 1.5 times larger than the REE of stress-based embeddings produced by `HyPy`/`Rigel`.

- Use the result of `hydra` as an initial condition for stress minimization. This gives an embedding method (called `hydra+` in [10]) that is typically 20% - 50% faster than `HyPy`/`Rigel` and results in a smaller REE in all numerical tests with improvements up to 40%.

We emphasize that these properties are limited by the fact that `hydra` is not a landmark-based method, and therefore can not be scaled up to the very large network instances considered for `HyPy` in [5].

## 2.5   Strain-minimization with landmarks

Both the `Rigel` algorithm of [21] and the `HyPy` algorithm of [5] combine the method of stress-minimization with the landmark heuristic. That is, the stress functional is first minimized on landmarks only, and then the stress between landmark and non-landmark-nodes is minimized. In addition, `HyPy` applies heuristics of gradually expanding the landmark set during optimization and of restarting the optimization procedure from multiple random initial point configurations. Due to the landmark approach, `Rigel`/`HyPy` are able to deal with very large instances of the network embedding problem with millions of nodes. The `hydra` method, as proposed in [10], on the other side, requires knowledge of the full dissimilarity matrix $D = [d_{ij}]$ and is therefore only applicable to small or medium sized networks. This is our motivation, to introduce

here a landmarked version of `hydra`, and to show that strain-minimization can be combined with the landmark heuristic while retaining all its attractive theoretical properties. In particular, we show in Theorem 3.1 that the necessary size of the landmark set only depends on the intrinsic dimension $d$ of the embedding space, but not on the total number of nodes (or other objects) to be embedded. This gives, for the first time, a sound theoretical basis for the application of the landmark heuristic for hyperbolic embedding problems. In terms of practical applications, we demonstrate in Section 4 that the landmarked `hydra` method `L-hydra` can be scaled up to all network embedding problems considered in [5] with up to almost 4 million nodes. The advantages of `hydra` as compared with HyPy/Rigel are effectively retained in the landmarked setting (see Section 4 for details):

- As a stand-alone method, `L-hydra` is a very fast embedding method (about ten times faster than `HyPy` with landmarks), but with a higher relative embedding error than the stress-based embedding produced by `HyPy`.

- Using the result of `L-hydra` as an initial condition for stress minimization (we call this method `L-hydra+`) gives an embedding method that is both faster than `HyPy` and produces embeddings with smaller error.

## 3   The landmarked hydra algorithm

### 3.1   Formulation of the embedding problem

Consider a set of objects $\boldsymbol{o}_1, \ldots, \boldsymbol{o}_n$, which are partitioned into landmarks $(\boldsymbol{o}_i)_{i \in L}$ and non-landmarks $(\boldsymbol{o}_i)_{i \in N}$. The pairwise dissimilarity $d_{ij}$ of $\boldsymbol{o}_i$ and $\boldsymbol{o}_j$ is only assumed to be known when $\boldsymbol{o}_i$ or $\boldsymbol{o}_j$ is a landmark. By permuting objects, the dissimilarity matrix $D$ can be arranged as $D = \begin{pmatrix} D_L & D_N^\top \\ D_N & D_R \end{pmatrix}$ where only $D_L$ – containing all pairwise dissimilarities between landmark nodes – and $D_N$ – containing all dissimilarities between one landmark and one non-landmark node – are known.

We first discuss the problem of exact recovery in hyperbolic space. That is, we assume that for a given embedding dimension $d$, the objects $\boldsymbol{o}_i$ can be perfectly described by points $\boldsymbol{x}_i$ in the hyperbolic manifold $\mathcal{H}_d$, i.e. the dissimilarity $d_{ij}$ between any pair of objects $(\boldsymbol{o}_i, \boldsymbol{o}_j)$ can be described exactly by the hyperbolic distance between the corresponding points $\boldsymbol{x}_i, \boldsymbol{x}_j \in \mathcal{H}_d$:

$$d_{ij} = \mathrm{d}_H^\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) = \frac{1}{\sqrt{\kappa}} \operatorname{arcosh}(\boldsymbol{x}_i \circ \boldsymbol{x}_j).$$

Transforming all distances by the hyperbolic cosine, we obtain

$$\cosh(\sqrt{\kappa}\, d_{ij}) = \boldsymbol{x}_i \circ \boldsymbol{x}_j. \tag{9}$$

We set $A = \cosh\left(\sqrt{\kappa}\, D\right)$, where the hyperbolic cosine is applied elementwise to the dissimilarity matrix $D$. Setting $l = |L|$ and $m = |N|$ we write

$$X = \begin{pmatrix} X_L \\ X_N \end{pmatrix} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_l, \boldsymbol{x}_{l+1}, \ldots, \boldsymbol{x}_{l+m})^\top \in \mathbb{R}^{(l+m)\times(d+1)}$$

for the coordinate matrix of some points $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_{l+m}$ in $\mathcal{H}_d$. Finally, let $J$ be the $(d+1) \times (d+1)$ diagonal matrix

$$J = \operatorname{diag}(1, -1, \ldots, -1),$$

cf. [17, §3.1]. Equation (9) for $i \in L \cup N$ and $j \in L$ can now be written in compact form as

$$\begin{pmatrix} A_L \\ A_N \end{pmatrix} = \begin{pmatrix} X_L \\ X_N \end{pmatrix} J\, X_L^\top. \tag{10}$$

Thus, the coordinate matrix $X = \begin{pmatrix} X_L \\ X_N \end{pmatrix}$ can be recovered from the known dissimilarities $D_L, D_N$ if we can solve (10) for $X_L$ and $X_N$. Due to its block structure, equation (10) can be split into two parts and solved in two stages, where the coordinates for landmarks must be determined from the first equation and the coordinates of non-landmarks from the second one:

$$\begin{cases} A_L = X_L J X_L^\top & \text{(11a)} \\ A_N = X_N J X_L^\top .. & \text{(11b)} \end{cases}$$

However, we can only hope to recover $(X_L, X_N)$ from these equations up to hyperbolic isometry, i.e., up to a distance-preserving transformation $\psi : \mathcal{H}_d \to \mathcal{H}_d$. From [17, Sec. 3.1,3.2], any hyperbolic isometry can be written as

$$\psi(\boldsymbol{x}) = T\boldsymbol{x}, \tag{12}$$

where $T \in \mathbb{R}^{(d+1) \times (d+1)}$ is an invertible matrix, which has $T_{11} > 0$ and satisfies

$$T^\top J T = T J T^\top = J. \tag{13}$$

Matrices with this property are called *positive Lorentz matrices* and both $\mathcal{H}_d$ and $\mathbb{R}_+^{1,d}$ are invariant under such transformations. Our observation can now be formalized as follows:

**Lemma 3.1.** *Let $T \in \mathbb{R}^{(d+1) \times (d+1)}$ be a positive Lorentz matrix and $(X_L, X_N)$ a solution of (11). Then $(\tilde{X}_L, \tilde{X}_N) = (X_L T, X_N T)$ is also a solution of (11). If $X_L$ and $X_N$ are hyperbolic coordinate matrices of points $(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_{l+m}) \in \mathcal{H}_d$, then $\tilde{X}_L$ and $\tilde{X}_N$ are also hyperbolic coordinate matrices of hyperbolically isometric points $(\tilde{\boldsymbol{x}}_1, \ldots, \tilde{\boldsymbol{x}}_{l+m}) \in \mathcal{H}_d$.*

*Proof.* The first statement follows directly from

$$\tilde{X}_L J \tilde{X}_L^\top = X_L T J T^\top X_L^\top = X_L J X_L^\top = A_L$$

and

$$\tilde{X}_N J \tilde{X}_L^\top = X_N T J T^\top X_L^\top = X_N J X_L^\top = A_N.$$

The second statement follows from the invariance of $\mathcal{H}_d$ under the positive Lorentz matrix $T$.  $\square$

## 3.2   Solving the embedding problem

In general, one cannot expect given dissimilarities $d_{ij}$ to be represented exactly by hyperbolic distances in $\mathcal{H}_d$, and therefore should not expect an exact solution to (11). For this reason, we replace (11) by its relaxation

$$\begin{cases} \hat{X}_L = \arg\min \left\{ \left\| A_L - X_L J X_L^\top \right\|_F : X_L \in \mathbb{R}^{l \times (d+1)} \right\} & \text{(14a)} \\ \hat{X}_N = \arg\min \left\{ \left\| A_N - X_N J X_L^\top \right\|_F : X_N \in \mathbb{R}^{m \times (d+1)} \right\}, & \text{(14b)} \end{cases}$$

where $\|.\|$ indicated the Frobenius norm on matrices. Note that the relaxation consists of two parts: First, equality has been replaced by minimality of the Frobenius distance. Second, $\hat{X}_L$ and $\hat{X}_N$ are no longer constrained to coordinate matrices of points in $\mathcal{H}_d$, but of points in the ambient Lorentz space $\mathbb{R}^{1,d}$.

The `L-hydra` algorithm consists of the stepwise solution of (14a):

- In step 1 and 2, the minimization problem (14a) is solved by a matrix Eigendecomposition. This is equivalent to the `hydra` embedding of [10], applied to the landmarks only.

- In step 3, the least-squares problem (14b) is solved using the Moore-Penrose inverse of $\hat{X}_L J$.

The details of each step are listed in Algorithm 1.

---

**Algorithm 1** L-hydra $(D_L, D_N, d, \kappa)$

---

**Input**    • Two matrices $D_L \in \mathbb{R}_{\geqslant 0}^{l \times l}$ and $D_N \in \mathbb{R}_{\geqslant 0}^{m \times l}$, of which the first represents the dissimilarities between pairs of landmarks and the second between pairs of landmarks and non-landmarks. The matrix $D_L$ must be symmetric with zero diagonal.

   • Embedding dimension $d \leq (l + m - 1)$

   • Parameter $\kappa > 0$; the negative of the hyperbolic curvature $-\kappa$.

**Step 1** Set

$$\begin{cases} A_L = \cosh\left(\sqrt{\kappa} D_L\right) & \text{(15a)} \\ A_N = \cosh\left(\sqrt{\kappa} D_N\right) & \text{(15b)} \end{cases}$$

with cosh applied elementwise, and compute the Eigendecomposition of the matrix $A_L$

$$A_L = Q \Lambda_L Q^\top, \tag{16}$$

where $\Lambda_L$ is the diagonal matrix of the Eigenvalues $\lambda_1 \geq \cdots \geq \lambda_l$ and the columns of $Q$ are the corresponding Eigenvectors $\boldsymbol{q}_1, \ldots, \boldsymbol{q}_l$.

**Step 2** Assuming that the last $d$ Eigenvalues are negative, allocate the $l \times (d + 1)$-matrix

$$\hat{X}_L := \left[ \sqrt{\lambda_1}\, \boldsymbol{q}_1 \quad \sqrt{(-\lambda_{l-d+1})}\, \boldsymbol{q}_{l-d+1} \quad \cdots \quad \sqrt{(-\lambda_l)}\, \boldsymbol{q}_l \right]. \tag{17}$$

If not all of the last $d$ Eigenvalues are negative, return **Null**. The algorithm must be rerun with smaller embedding dimension $d$ for non-null result.

**Step 3** Allocate the $m \times (d + 1)$-matrix

$$\begin{aligned} \hat{X}_N :&= A_N \hat{X}_L \operatorname{diag}\left( \frac{1}{\lambda_1} \quad -\frac{1}{\lambda_{l-d+1}} \quad \cdots \quad -\frac{1}{\lambda_l} \right) \\ &= A_N \left[ \frac{\boldsymbol{q}_1}{\sqrt{\lambda_1}} \quad -\frac{\boldsymbol{q}_{l-d+1}}{\sqrt{-\lambda_{l-d+1}}} \quad \cdots \quad -\frac{\boldsymbol{q}_l}{\sqrt{-\lambda_l}} \right]. \end{aligned} \tag{18}$$

**Return** Matrix $\hat{X} = \begin{pmatrix} \hat{X}_L \\ \hat{X}_N \end{pmatrix}$ whose rows are coordinates in positive Lorentz space $\mathbb{R}_+^{1,d}$.

---

## 3.3 Theoretical Properties

The key theoretical properties of the `L-hydra` algorithm are summarized in the following theorems.

**Theorem 3.1** (Exact Recovery). *Let $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_{m+l}$ be points in hyperbolic d-space $\mathcal{H}_d$, of which the first $l \geq d$ are designated as landmarks. Assume that the landmarks are not all contained in a single hyperplane of $\mathcal{H}_d$ and let $D = [d_{ij}] = [\mathrm{d}_H^\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j)]$ be the matrix of their hyperbolic distances with curvature $-\kappa$, partitioned as*

$$D = \begin{pmatrix} D_L & D_N^\top \\ D_N & D_R \end{pmatrix}.$$

*Then the algorithm $\mathtt{L\text{-}hydra}\,(D_L, D_N, d, \kappa)$ recovers all points $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_{l+m}$ up to isometry. That is, the rows $\hat{\boldsymbol{x}}_1, \ldots, \hat{\boldsymbol{x}}_{l+m}$ of the matrix $\hat{X}$ returned by $\mathtt{L\text{-}hydra}\,(D_L, D_N, d, \kappa)$ are points in $\mathcal{H}_d$ and satisfy*

$$\mathrm{d}_H(\hat{\boldsymbol{x}}_i, \hat{\boldsymbol{x}}_j) = \mathrm{d}_H(\boldsymbol{x}_i, \boldsymbol{x}_j) \qquad \text{for all } i, j = 1, \ldots, l + m. \tag{19}$$

**Theorem 3.2** (Optimal Approximation). *Suppose that $A_L = \cosh(\sqrt{\kappa} D_L)$ has at least $d \in \mathbb{N}$ strictly negative Eigenvalues. Then, the matrices $\hat{X}_L, \hat{X}_N$ returned by the algorithm $\mathtt{L\text{-}hydra}\,(D_L, D_N, d, \kappa)$ solve the minimization problem (14a). Moreover, the first columns of $\hat{X}_L$ and $\hat{X}_N$ are strictly positive; equivalently, all rows of $\hat{X}_L$ and $\hat{X}_N$ represent points in positive Lorentz space $\mathbb{R}_+^{1,d}$.*

The first part of the result, i.e., the optimality of the embedding of landmarks, follows from [10]. For convenience, we provide a self-contained proof of both parts:

*Proof.* Let $A_L = Q\Lambda_L Q^\top$ be the Eigendecomposition of $A_L = X_L J X_L^\top$, where $\Lambda$ is the diagonal matrix of the Eigenvalues $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_l$ of $A_L$ in descending order and $Q$ is unitary, i.e., $Q^\top Q = I$. By the unitary invariance of the Frobenius norm, solving the minimization problem (14a) is equivalent to solving

$$\hat{Y}_L = \arg\min\left\{ \left\| Y_L J Y_L^\top - \Lambda \right\|_F : Y_L \in \mathbb{R}^{l \times (d+1)} \right\}, \tag{20}$$

where $\hat{Y}_L$ is related to $\hat{X}_L$ by $X_L = Q\hat{Y}_L$. We claim that (20) is minimized by

$$\hat{Y}_L := \begin{bmatrix} \sqrt{\lambda_1}\, \boldsymbol{e}_1 & \sqrt{(-\lambda_{l-d+1})}\, \boldsymbol{e}_{l-d+1} & \cdots & \sqrt{(-\lambda_l)}\, \boldsymbol{e}_l \end{bmatrix}. \tag{21}$$

To see this, note that on the one hand,

$$\left\| \hat{Y}_L J \hat{Y}_L^\top - \Lambda \right\| = \sum_{i=2}^{l-d} \lambda_i^2.$$

On the other hand, for an arbitrary $Y_L \in \mathbb{R}^{l \times (d+1)}$, denote the Eigenvalues of $Y_L J Y_L^\top$ by $\eta_1 \geq \eta_2 \geq \cdots \geq \eta_l$. By Sylvester's law of inertia, owing to the structure of $J$, only $\eta_1$ is strictly positive, only the last $d$ Eigenvalues $\eta_{l-d+1}, \ldots, \eta_l$ are strictly negative, and all other Eigenvalues must be zero. Hence, using the Wielandt-Hofmann inequality, we obtain

$$\left\| Y_L J Y_L^\top - \Lambda \right\|_F \geq \sum_{i=1}^{l} (\eta_i - \lambda_i)^2 \geq \sum_{i=2}^{l-d} \lambda_i^2 = \left\| \hat{Y}_L J \hat{Y}_L^\top - \Lambda \right\|_F$$

for all $Y_L \in \mathbb{R}^{l \times (d+1)}$, showing the optimality of $\hat{Y}_L$ for (20). Transforming back to $\hat{X}_L = Q\hat{Y}_L$ yields (17) and shows that $\hat{X}_L$ solves (14a). It remains to solve (14b), which is a least-squares problem, whose solution can be written as

$$\hat{X}_N = A_N (J \hat{X}_L^\top)^+, \tag{22}$$

where $(J\hat{X}_L^\top)^+$ denotes the Moore-Penrose-Pseudoinverse; see [12, Theorem. (8.1) and Remark. (8.2)]. Since $J\hat{X}_L^\top$ has full row rank, we can write

$$(J\hat{X}_L^\top)^+ = \hat{X}_L J (J\hat{X}_L^\top \hat{X}_L J)^{-1} = \hat{X}_L (\hat{X}_L^\top \hat{X}_L)^{-1} J = \hat{X}_L \operatorname{diag}\left( \frac{1}{\lambda_1} \quad -\frac{1}{\lambda_{l-d+1}} \quad \cdots \quad -\frac{1}{\lambda_l} \right). \quad (23)$$

Inserting this into (22) yields (18), completing the main part of the proof. Finally, note that by the Perron-Frobenius theorem, the leading Eigenvector $\boldsymbol{q}_1$ of the positive matrix $A_L$ must also be positive. Together with (18), this shows that the first columns of both $\hat{X}_L$ and $\hat{X}_N$ are positive. $\quad\square$

*Proof of Thm. 3.1.* Denote by $X = \left( \begin{smallmatrix} X_L \\ X_N \end{smallmatrix} \right)$ the coordinate matrix of the original points $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_{l+m}$ and by $\hat{X} = \left( \begin{smallmatrix} \hat{X}_L \\ \hat{X}_N \end{smallmatrix} \right)$ the coordinate matrix of the points $\hat{\boldsymbol{x}}_1, \ldots, \hat{\boldsymbol{x}}_{l+m}$ returned by `L-hydra`. Since the points $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_l$ are not all contained in a single hyperplane of $\mathcal{H}_d$, it follows from [17, §3.2] that there are $d+1$ among them which are linearly independent as points in $\mathbb{R}^{d+1}$, and thus that $X_L$ has the full column rank $d+1$. Let $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_l$ be the Eigenvalues of $A_L = X_L J X_L^\top$. By the same argument as in the proof of Theorem 3.2 it follows that only $\lambda_1$ is strictly positive, only the $d$ last Eigenvalues $\lambda_{l-d+1}, \ldots, \lambda_l$ are strictly negative and all others are zero. From the proof of Theorem 3.2 we find that the residual of (14a) is

$$\left\| A_L - \hat{X}_L J \hat{X}_L^\top \right\|_F^2 = \sum_{i=2}^{l-d} \lambda_i^2 = 0,$$

i.e., the embedding of the landmarks is exact. This means that $X_L J X_L^\top = \hat{X}_L J \hat{X}_L^\top$, i.e., $\hat{X}_L$ and $X_L$ are hyperbolically isometric and by (12), there exists a positive Lorentz matrix $T$ such that

$$\hat{X}_L = X_L T.$$

For the Moore-Penrose-Pseudoinverse $(J X_L^\top)^+$ it follows that

$$(J\hat{X}_L^\top)^+ = (J X_L^\top)^+ T.$$

Moreover, since $X_L$ has the full column rank, $(J X_L^\top)^+$ is a right-inverse of $(J X_L^\top)$, and we have

$$\hat{X}_N = A_N (J\hat{X}_L^\top)^+ = A_N (J X_L^\top)^+ T = X_N J X_L^\top (J X_L^\top)^+ T = X_N T.$$

Together, it follows that $\hat{X} = XT$, i.e., $\hat{X}$ and $X$ are hyperbolically isometric and (19) must hold for all (both landmark and non-landmark) points. $\quad\square$

**Corollary 3.1** (Consistency of landmark and non-landmark embedding)**.** *The landmark embedding is consistent with the non-landmark embedding, meaning that re-embedding the $j$-th landmark point ($j \in \{1, \ldots, l\}$) as a non-landmark will not change its representation $\hat{\boldsymbol{x}}_j$, or equivalently*

$$\hat{X}_L = A_L \left[ \frac{\boldsymbol{q}_1}{\sqrt{\lambda_1}} \quad -\frac{\boldsymbol{q}_{l-d+1}}{\sqrt{-\lambda_{l-d+1}}} \quad \cdots \quad -\frac{\boldsymbol{q}_l}{\sqrt{-\lambda_l}} \right]. \quad (24)$$

*Remark* 3.1. Re-embedding of landmarks is equivalent to inserting the transformed landmark dissimilarity matrix $A_L$, instead of $A_N$, into (18). Hence, consistency can be expressed by equation (24).

*Proof.* For each eigenvector $\boldsymbol{q}_j$ with $j = l-d, \ldots, l$ of $A_L$ we have

$$\frac{A_L \boldsymbol{q}_j}{\sqrt{-\lambda_j}} = \frac{\lambda_j \boldsymbol{q}_j}{\sqrt{-\lambda_j}} = \sqrt{-\lambda_j} \boldsymbol{q}_j,$$

and – omitting the minus sign – also for $j = 1$. Applying this column-by-column and taking into account (17) we obtain (24). $\quad\square$

## 3.4 Practical implementation

After having analyzed the theoretical properties of `L-hydra`, we point out some more practical issues in the implementation of `L-hydra`:

**Reduced Eigendecomposition.** In (17) only the first and the last $d$ Eigenvalues and Eigenvectors of the matrix $A_L$ are needed. There are efficient numerical routines to perform such a reduced Eigendecomposition without computing the full Eigendecomposition of $A_L$.

**Optimizing Curvature.** The `L-hydra` algorithm treats the curvature parameter $\kappa$ as fixed input. In a practical implementation it is usually desirable to also optimize over $\kappa$ by running `L-hydra` for several different values of $\kappa$. In a similar way, curvature is optimized in `HyPy` by treating it as an additional free variable of the stress functional (5).

**Projection to $\mathcal{H}_d$.** `L-hydra` returns points $\hat{\boldsymbol{x}}_1, \ldots, \hat{\boldsymbol{x}}_{l+m}$ in positive Lorentz space $\mathbb{R}^{1,d}$, which are 'reasonably close' to $\mathcal{H}_d$, due to Theorem 3.1, but typically not exactly located on the hyperboloid $\mathcal{H}_d$. To obtain points on $\mathcal{H}_d$, a projection method must be applied, cf. [10]. Here, we project parallel to the $x_1$-axis, i.e. given $\boldsymbol{x} \in \mathbb{R}^{1,d}$ we set

$$\tilde{x}_1 := \sqrt{1 + x_2^2 + \cdots + x_{d+1}^2}$$

to obtain a projected point $\tilde{\boldsymbol{x}} = (\tilde{x}_1, x_2, \ldots, x_{d+1}) \in \mathcal{H}_d$ for each point returned by `L-hydra`.

**Analysis of Runtime.** Recall that the total number of points, the number of landmarks and the number of non-landmarks are denoted by $n$, $l$ and $m = n - l$ respectively. For Step 1 and 2 of hydra, we expect a runtime of $\mathcal{O}(l^\alpha)$ with $\alpha$ slightly above, but close, to 2, cf. [10, 9], using reduced Eigendecomposition. For Step 3 of hydra, we expect a runtime of $\mathcal{O}(lm)$ for the multiplication of a $m \times l$ and a $l \times (d+1)$ matrix, resulting in an overall runtime of $\mathcal{O}(l^{\alpha-1}n)$. On the basis of Theorem 3.1, the necessary size of the landmark set depends only on the intrinsic dimension $d$, but not on the total number of points to be embedded. This implies, that the runtime of `L-hydra` scales effectively *linear* as $\mathcal{O}(n)$ in the number of points to be embedded.

**The `L-hydra+` method.** If minimization of the stress functional (5) (or, equivalently, relative embedding error (7)) is the ultimate goal, then the result of `L-hydra` can be used as an initial condition for its numerical minimization. Effectively, this corresponds to a chaining of `L-hydra` and `HyPy`, where `L-hydra` replaces the random initial condition of `HyPy`.

## 4 Numerical Results

We test the practical performance of `L-hydra` and `L-hydra+` on five social networks from SNAP ([20]), with ~300,000 up to almost 4 million nodes and ~1 million up to ~117 million edges; see Table 1 for a brief description of the networks. The same networks were used in [5] to evaluate the performance of `HyPy`, which serves as a benchmark for our methods. All networks are unweighted and undirected. The $l = 100$ landmark nodes are chosen with probability proportional to the node degrees without replacement as proposed by [5], which combines the benefits of random sampling and the selection of highest-degree nodes. All shortest-path distances are computed using SNAP for Python ([13]). `L-hydra` and `L-hydra+` were implemented in Python as described in Algorithm 1 and Section 3.4, i.e., with curvature optimization, reduced Eigendecomposition and projection to $\mathcal{H}_d$, using parallelized routines for matrix operations. For `HyPy`, we used the Python code kindly supplied by the authors of [5], using the same parallelization as for `L-hydra` and `L-hydra+`. All calculations were performed on a Dual socket Intel server using an

| Network | Description | Source `http://snap.stanford.edu/data/...` | # Nodes | # Edges |
|---|---|---|---|---|
| Amazon | Network based on the 'Customers Who Bought This Item Also Bought' feature of the Amazon website | `.../com-Amazon.html` | $334,863$ | $925,872$ |
| DBLP | Collaboration network of the DBLP computer science bibliography | `.../com-DBLP.html` | $317,080$ | $1,049,866$ |
| YouTube | Friendships in the YouTube social network | `.../com-Youtube.html` | $1,134,890$ | $2,987,624$ |
| Live Journal | Friendship network of a free online blogging community | `.../com-LiveJournal.html` | $3,997,962$ | $34,681,189$ |
| Orkut | Friendships in an online social network | `.../com-Orkut.html` | $3,072,441$ | $117,185,083$ |

Table 1: Description of networks used for numerical experiments

Intel Xeon CPU E5-2680 v3 with 12 cores at 2.50GHz and with 64GB RAM.

The embedding results for all methods are shown in Figures 1 and 2. For each method and network, we calculate the network embedding for dimensions $d = 2, 3, \ldots, 10$. Following [5], we show the relative embedding error (7) in three variations: calculated over all pairs of *landmark* nodes; calculated over all pairs consisting of one landmark and one *non-landmark* node; and over 100,000 randomly chosen pairs of non-landmark nodes ('*validation error*'). Note that due to the size of the networks, calculation of the full embedding error over all node pairs is computationally infeasible, but the validation error should provide a good proxy. Since HyPy starts from a random initial condition, we show the $5\%-$ and $95\%-$quantiles in addition to its average result over 20 runs.

The runtimes for embedding into dimension $d = 2$ are shown in Figure 3; results for other dimensions are similar. Note that the runtime can be split into distance calculation – which is the same for all methods – and embedding calculation – which differs from method to method. We summarize our observations as follows:

**Consistency with [5].** The results obtained for HyPy are consistent with the results reported in [5]; note that we report REE while RMSE is reported in [5].

**Performance of L-hydra.** As expected, the strain-minimizing algorithm L-hydra typically achieves poorer results in terms of the relative embedding error (REE) than the stress-minimizing algorithms L-hydra+ and HyPy. Compared to the average performance of HyPy, the validation REE of L-hydra is on average (median over all networks and dimensions) 2.17 times larger. However, the total runtime of L-hydra (the bulk of which is spent on distance calculation) is between 5 and 10 times faster than the average runtime of HyPy. In essence, replacing HyPy by L-hydra trades a moderate increase in embedding error against a substantial decrease in computation time.

**Performance of L-hydra+.** The embedding result of L-hydra+ – which combines L-hydra with stress minimization – is consistently (over all networks, dimensions and error types) better than the average result of HyPy, and in the large majority of cases even better than the $5\%-$quantile of HyPy results. Averaged over all dimensions and networks, the validation REE of L-hydra+ is 12% smaller than the validation REE of HyPy. At the same time, L-hydra+ is faster than HyPy for all network embeddings except for YouTube, with a relative difference in average runtime from $-52.5\%$ (DBLP network) to $+5.7\%$ (YouTube network).

**Stability and Reproducibility.** Both L-hydra and L-hydra+ do not suffer from the noise introduced

by the random initial condition of `HyPy`, which – in particular for small dimension – causes a notable variation in embedding quality of `HyPy` and yields different embedding results upon each rerun.

**Runtime analysis.** In panels A, B of Figure 3 we show regression lines to estimate the exponent $\alpha$ in the conjectured complexity $\mathcal{O}(n^\alpha)$ for distance calculation and embedding. For distance calculation we obtain the estimate $\alpha_{\text{dist}} \approx 1.31$, showing that due to the sparsity of the networks, shortest path computations are more efficient than the upper bound $\alpha = 2$ obtained from the Floyd-Warshall method. For the embedding step of `L-hydra+` we estimate $\alpha_{\text{L-hydra+}} \approx 1.32$, while for `L-hydra` we obtain $\alpha_{\text{L-hydra+}} \approx 0.99$, confirming the linear scaling of the method.

# References

[1] Marian Boguna, Dmitri Krioukov, and Kimberly C Claffy. Navigability of complex networks. *Nature Physics*, 5(1):74, 2009.

[2] Ingwer Borg and Patrick JF Groenen. *Modern multidimensional scaling: Theory and applications.* Springer Science & Business Media, 2005.

[3] W. James Cannon, William J. Floyd, Richard Kenyon, and Walter R. Parry. Hyperbolic geometry. In Silvio Levy, editor, *Flavors of Geometry*, pages 59–115. MSRI Publications, 31 edition, 1997.

[4] Benjamin Paul Chamberlain, James Clough, and Marc Peter Deisenroth. Neural embeddings of graphs in hyperbolic space. *arXiv:1705.10359*, 2017.

[5] Kenny Chowdhary and Tamara G Kolda. An improved hyperbolic embedding algorithm. *Journal of Complex Networks*, 2017.

[6] Vin De Silva and Joshua B Tenenbaum. Sparse multidimensional scaling using landmark points. Technical report, Technical report, Stanford University, 2004.

[7] Robert W Floyd. Algorithm 97: shortest path. *Communications of the ACM*, 5(6):345, 1962.

[8] Tobias Friedrich. From graph theory to network science: The natural emergence of hyperbolicity. In *Symposium Theoretical Aspects of Computer Science (STACS)*, pages 5:1–5:9, 2019.

[9] Leslie Hogben. *Handbook of linear algebra.* CRC Press, 2006.

[10] Martin Keller-Ressel and Stephanie Nargang. Hydra: a method for strain-minimizing hyperbolic embedding of network- and distance-based data. *Journal of Complex Networks*, 8(1):cnaa002, 2020.

[11] Robert Kleinberg. Geographic routing using hyperbolic space. In *INFOCOM 2007. 26th IEEE International Conference on Computer Communications*, pages 1902–1909, 2007.

[12] Alan J. Laub. *Matrix Analysis For Scientists And Engineers.* Society for Industrial and Applied Mathematics, USA, 2004.

[13] Jure Leskovec and Rok Sosic. SNAP: A general purpose network analysis and graph mining library. *CoRR*, abs/1606.07550, 2016.

[14] Alessandro Muscoloni, Josephine Maria Thomas, Sara Ciucci, Ginestra Bianconi, and Carlo Vittorio Cannistraci. Machine learning meets complex networks via coalescent embedding in the hyperbolic space. *Nature Communications*, 8(1):1615, 2017.
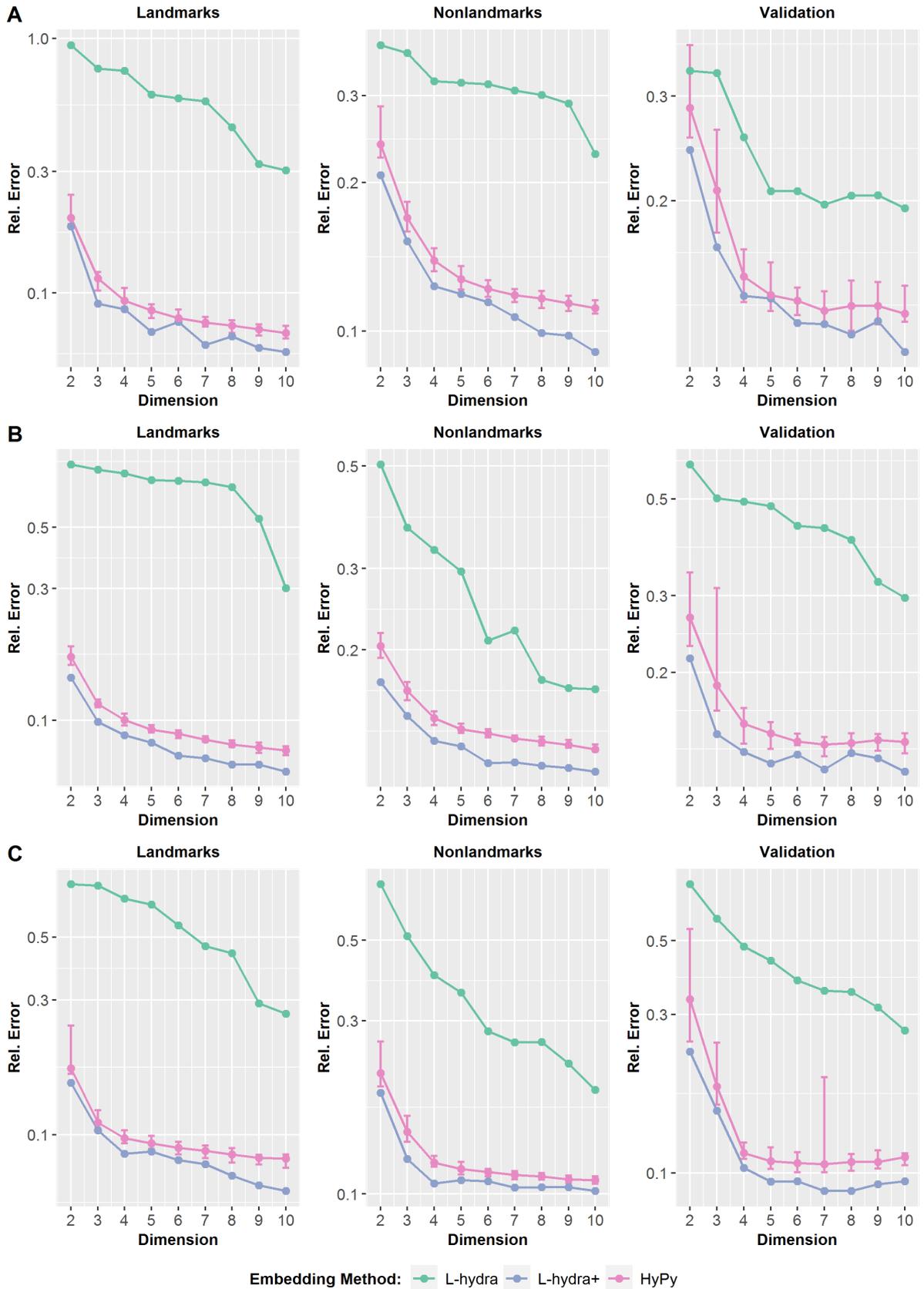
Figure 1: Comparison of the embedding performance of `HyPy`, `L-hydra` and `L-hydra+` on the real network data sets (A) Amazon, (B) DBLP and (C) YouTube. Embedding quality is measured by the relative embedding error (7) over pairs of landmark nodes (left column), pairs of landmark-non-landmark nodes (middle column) and over 100,000 randomly selected validation pairs of non-landmark nodes (right column). The same $l = 100$ landmarks are used for all three methods. For `L-hydra` a $5 - 95\%$ error bar is shown, corresponding to 20 runs with randomized initial condition.
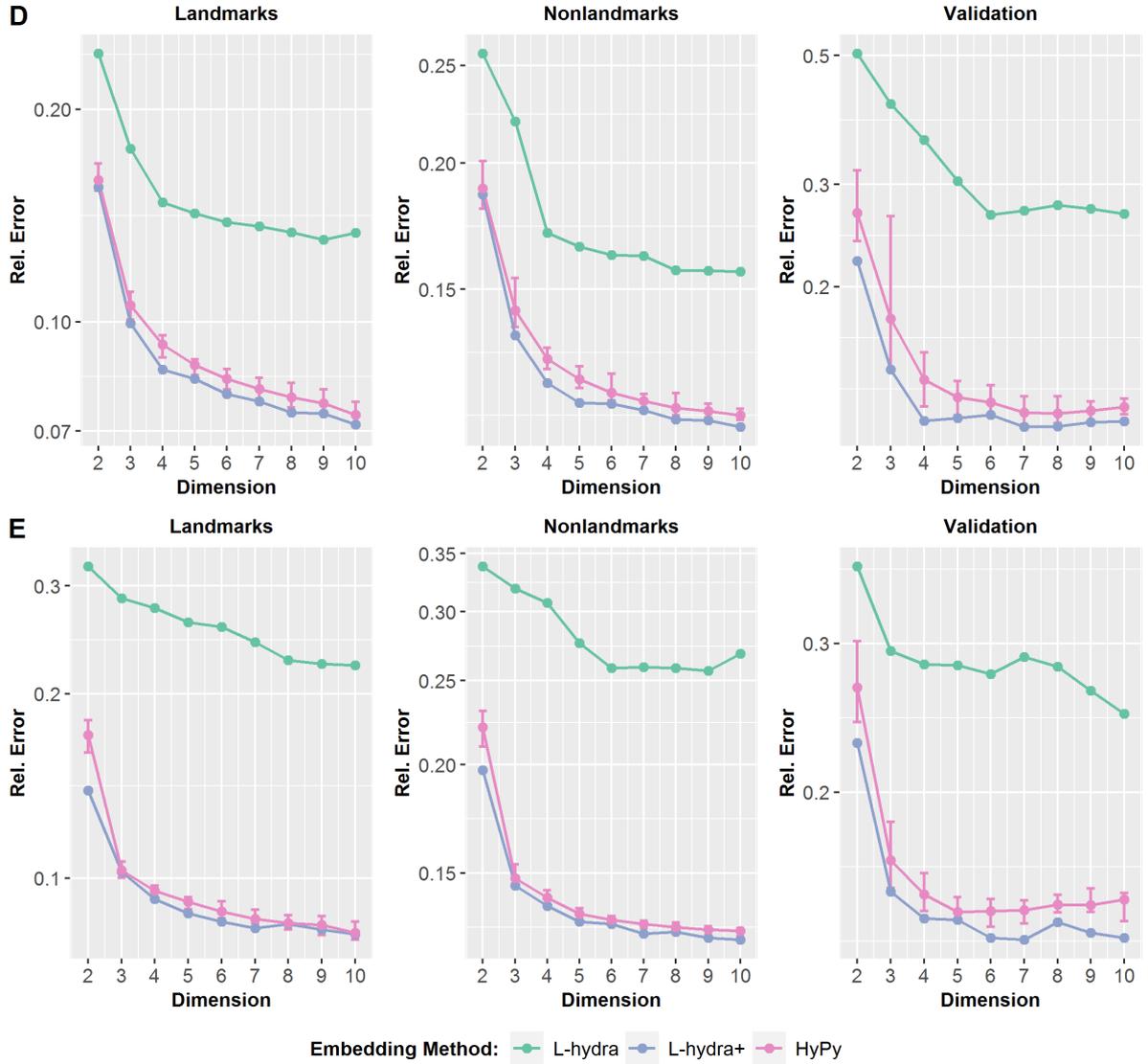
Figure 2: Comparison of the embedding performance of `HyPy`, `L-hydra` and `L-hydra+` on the real network data sets (D) LiveJournal and (E) Orkut. Embedding quality is measured by the relative embedding error (7) over pairs of landmark nodes (left column), pairs of landmark-non-landmark nodes (middle column) and over 100,000 randomly selected validation pairs of non-landmark nodes (right column). The same $l = 100$ landmarks are used for all three methods. For `L-hydra` a $5-95\%$ error bar is shown, corresponding to 20 runs with randomized initial condition.
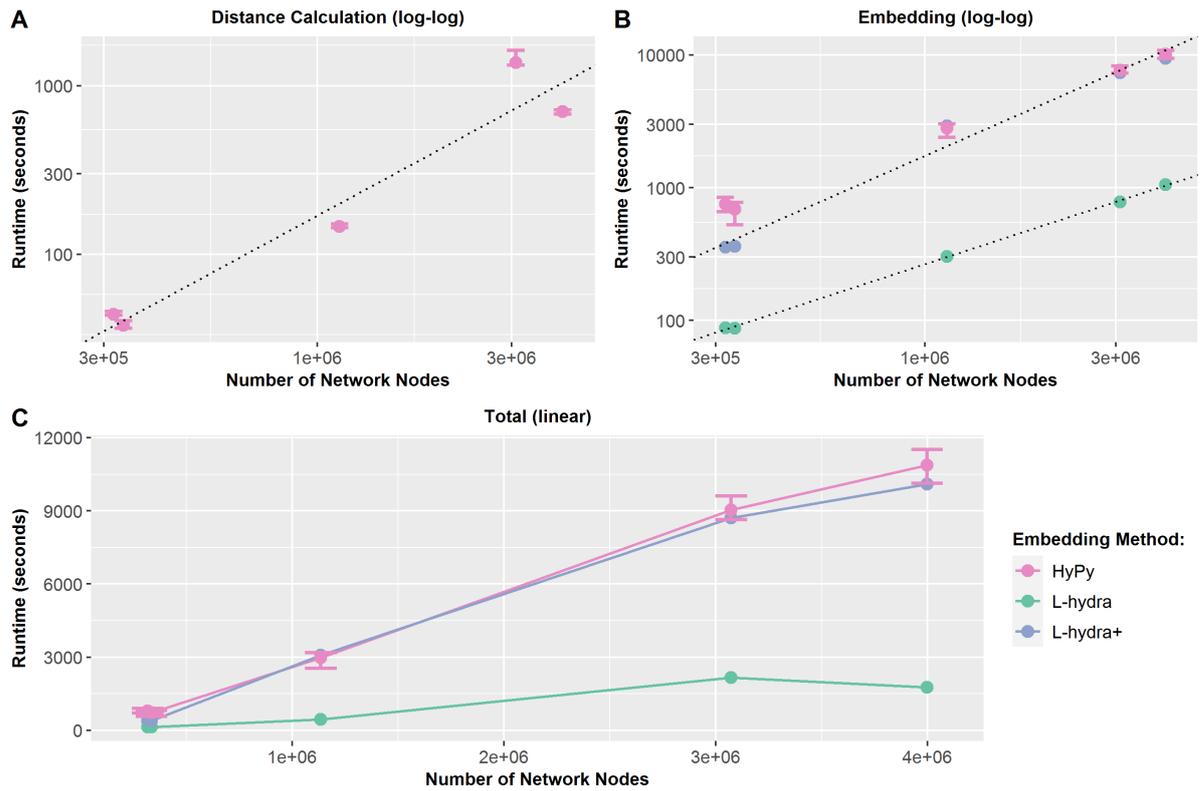
Figure 3: Computation time (in seconds) of the hyperbolic embedding methods `HyPy`, `L-hydra` and `L-hydra+` applied to the five networks listed in 1. For `HyPy`, average computation time and a $5 - 95\%$ error bar is shown, corresponding to 20 runs with randomized initial condition. Regression lines for `L-hydra` and `L-hydra+` are added in the doubly logarithmic plots (A) and (B).

[15] Fragkiskos Papadopoulos, Maksim Kitsak, M Ángeles Serrano, Marián Boguná, and Dmitri Krioukov. Popularity versus similarity in growing networks. *Nature*, 489(7417):537, 2012.

[16] Fragkiskos Papadopoulos, Constantinos Psomas, and Dmitri Krioukov. Network mapping by replaying hyperbolic growth. *IEEE/ACM Transactions on Networking (TON)*, 23(1):198–211, 2015.

[17] John Ratcliffe. *Foundations of hyperbolic manifolds*, volume 149. Springer Science & Business Media, 2006.

[18] Puoya Tabaghi and Ivan Dokmanić. Hyperbolic distance matrices. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1728–1738, 2020.

[19] Jörg A Walter. H-mds: a new approach for interactive visualization with multidimensional scaling in the hyperbolic space. *Information systems*, 29(4):273–292, 2004.

[20] Jaewon Yang and Jure Leskovec. Defining and evaluating network communities based on ground-truth. 2012.

[21] Xiaohan Zhao, Alessandra Sala, Haitao Zheng, and Ben Y Zhao. Fast and scalable analysis of massive social graphs. *arXiv preprint arXiv:1107.5114*, 2011.

[22] Ciyou Zhu, Richard H Byrd, Peihuang Lu, and Jorge Nocedal. Algorithm 778: L-BFGS-B: Fortran subroutines for large-scale bound-constrained optimization. *ACM Transactions on Mathematical Software (TOMS)*, 23(4):550–560, 1997.