# Vector-valued Distance and Gyrocalculus on the Space of Symmetric Positive Definite Matrices

**Federico López**[*]
HITS

**Beatrice Pozzetti**
Heidelberg University

**Steve Trettel**
Stanford University

**Michael Strube**
HITS

**Anna Wienhard**
Heidelberg University
HITS

## Abstract

We propose the use of the vector-valued distance to compute distances and extract geometric information from the manifold of symmetric positive definite matrices (SPD), and develop gyrovector calculus, constructing analogs of vector space operations in this curved space. We implement these operations and showcase their versatility in the tasks of knowledge graph completion, item recommendation, and question answering. In experiments, the SPD models outperform their equivalents in Euclidean and hyperbolic space. The vector-valued distance allows us to visualize embeddings, showing that the models learn to disentangle representations of positive samples from negative ones.

## 1 Introduction

Symmetric Positive Definite (SPD) matrices have been applied in many tasks in computer vision such as pedestrian detection [75, 79], action [36, 51, 60] or face recognition [41, 42], object [43, 90] and image set classification [86], visual tracking [87], and medical imaging analysis [5, 65] among others. They have been used to capture statistical notions (Gaussian distributions [68], covariance [78]), while respecting the Riemannian geometry of the underlying SPD manifold, which offers a convenient trade-off between structural richness and computational tractability [27]. Previous work has applied approximation methods that locally flatten the manifold by projecting it to its tangent space [22, 83], or by embedding the manifold into higher dimensional Hilbert spaces [35, 90].

These methods face problems such as distortion of the geometrical structure of the manifold and other known concerns with regard to high-dimensional spaces [29]. To overcome these issues, several distances on SPD manifolds have been proposed, such as the Affine Invariant metric [65], the Stein metric [71], the Bures–Wasserstein metric [13] or the Log-Euclidean metric [5, 6], with their respective geometric properties. However, the representational power of SPD is not fully exploited in many cases [6, 65]. At the same time, it is hard to translate operations into their non-Euclidean domain given the lack of closed-form expressions. There has been a growing need to generalize basic operations, such as addition, rotation, reflection or scalar multiplication, to their Riemannian geometric counterparts to leverage this structure in the context of Geometric Deep Learning [19].

SPD manifolds have a rich geometry that contains both Euclidean and hyperbolic subspaces. Thus, embeddings into SPD manifolds are beneficial, since they can accommodate hierarchical structure in data sets in the hyperbolic subspaces while at the same time represent Euclidean features. This makes them more versatile than using only hyperbolic or Euclidean spaces, and in fact, their different submanifold geometries can be used to identify and disentangle such substructures in graphs.

---

[*]Correspondence to federico.lopez@h-its.org

In this work, we introduce the vector-valued distance function to exploit the full representation power of SPD (§3.1). While in Euclidean or hyperbolic space the relative position between two points is completely captured by their distance (and this is the only invariant), in SPD this invariant is a vector, encoding much more information than a single scalar. This vector reflects the higher expressivity of SPD due to its richer geometry encompassing Euclidean as well as hyperbolic spaces. We develop algorithms using the vector-valued distance and showcase two main advantages: its versatility to implement universal models, and its use in explaining and visualizing what the model has learned.

Furthermore, we bridge the gap between Euclidean and SPD geometry by developing gyrocalculus in SPD (§4), which yields closed-form expressions of arithmetic operations, such as addition, scalar multiplication and matrix scaling. This provides means to translate previously implemented ideas in different metric spaces to their analog notions in SPD. These arithmetic operations are also useful to adapt neural network architectures to SPD manifolds.

We showcase this on knowledge graph completion, item recommendation, and question answering. In the experiments, the proposed SPD models outperform their equivalents in Euclidean and hyperbolic space (§6). These results reflect the superior expressivity of SPD, and show the versatility of the approach and ease of integration with downstream tasks.

The vector-valued distance allows us to develop a new tool for the analysis of the structural properties of the learned representations. With this tool, we visualize high-dimensional SPD embeddings, providing better explainability on what the models learn (§6.4). We show that the knowledge graph models are capable of disentangling and clustering positive triples from negative ones.

## 2   Related Work

Symmetric positive definite matrices are not new in the Machine Learning literature. They have been used in a plethora of applications [5, 29, 36, 40–43, 51, 60, 65, 68, 75, 78, 79, 86, 87, 90]), although not always respecting the intrinsic structure or the positive definiteness constraint [22, 31, 35, 83, 90]. The alternative has been to map manifold points onto a tangent space and employ Euclidean-based tools. Unfortunately, this mapping distorts the metric structure in regions far from the origin of the tangent space affecting the performance [44, 96].

Previous work has proposed alternatives to the basic neural building blocks respecting the geometry of the space. For example, transformation layers [29, 33, 40], alternate convolutional layers based on SPDs [94] and Riemannian means [23], or appended after the convolution [21], recurrent models [24], projections onto Euclidean spaces [50, 57] and batch normalization [20]. Our work follows this line, providing explicit formulas for translating Euclidean arithmetic notions into SPDs.

Our general view, using the vector-valued distance function, allows us to treat Riemannian and Finsler metrics on SPD in a unified framework. Finsler metrics have previously been applied in compressed sensing [30], information geometry [69], for clustering categorical distributions [63], and in robotics [67]. With regard to optimization, matrix backpropagation techniques have been explored [2, 17, 43], with some of them accounting for different Riemannian geometries [20, 40]. Nonetheless, we opt for tangent space optimization [26] by exploiting the explicit formulations of the exponential and logarithmic map.

## 3   The Space $\mathrm{SPD}_n$

The space $\mathrm{SPD}_n$ is a Riemannian manifold of non-positive curvature of $n(n + 1)/2$ dimensions. Points in $\mathrm{SPD}_n$ are positive definite real symmetric $n \times n$ matrices, with the identity matrix $I$ being a natural basepoint. The tangent space to any point of $\mathrm{SPD}_n$ can be identified with the vector space $S_n$ of all real symmetric $n \times n$ matrices. $\mathrm{SPD}_n$ contains $n$-dimensional Euclidean subspaces, $(n − 1)$-dimensional hyperbolic subspaces as well as products of $\lfloor \frac{n}{2} \rfloor$ hyperbolic planes (see Helgason [39] for an in-depth introduction).
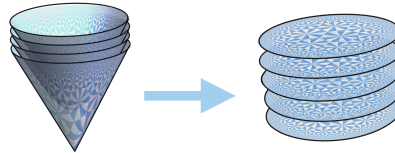


Figure 1: $\mathrm{SPD}_2$ is foliated by hyperboloids, each of which is a copy of the hyperbolic plane.
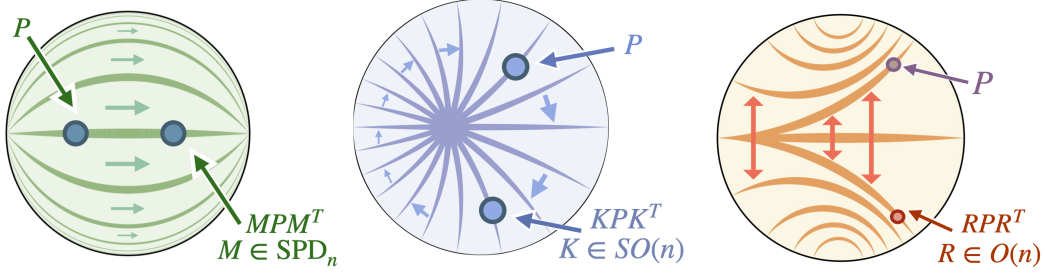
Figure 2: Some isometries of $\mathrm{SPD}_n$ have analogous Euclidean counterparts. Translation (left), rotation (center) and reflection (right).

In Figure 1 we visualize the smallest nontrivial example. $\mathrm{SPD}_2$ identifies with the inside of a cone in $\mathbb{R}^3$, cut out by requiring both eigenvalues of the matrix $\left(\begin{smallmatrix} x & y \\ y & z \end{smallmatrix}\right)$ to be positive. It carries the product geometry of the hyperbolic plane times a line.

**Exponential and logarithmic maps:** The exponential map, $\exp\colon S_n \to \mathrm{SPD}_n$, is a homeomorphism which links the Euclidean geometry of the tangent space $S_n$ and the curved geometry of $\mathrm{SPD}_n$. Its inverse is the logarithm map, $\log\colon \mathrm{SPD}_n \to S_n$. This pair of functions give diffeomorphisms that allows one to freely move between 'tangent space coordinates' or the original 'manifold coordinates'. We apply both maps based at $I \in \mathrm{SPD}_n$. The reason for this is that while mathematically any two points on $\mathrm{SPD}_n$ are equivalent, and we could obtain a different concrete expression for any other choice of basepoint $B \in \mathrm{SPD}_n$, the resulting formulas would be more complicated, and thus $I$ is the best choice from a computational point of view. We prove this in Appendix C.3.

**Symmetries:** The prototypical symmetries of $\mathrm{SPD}_n$ are parameterized by elements of $GL(n;\mathbb{R})$: any invertible matrix $M$ defines the symmetry $P \mapsto MPM^T$ acting on all points $P \in \mathrm{SPD}_n$. Thus many geometric transformations $\mathrm{SPD}_n$ can be completed using standard optimized matrix algorithms as opposed to custom-built procedures. See Appendix C.2. for a brief review of these symmetries.

Among these, we may find $\mathrm{SPD}_n$-generalizations of familiar symmetries of Euclidean geometry. When also $M$ is an element of $\mathrm{SPD}_n$, the symmetry $P \mapsto MPM^T$ is a generalization of an Euclidean *translation*, fixing no points of $\mathrm{SPD}_n$. When $M$ is an orthogonal matrix, the symmetry $P \mapsto MPM^T$ is conjugation by $M$, and thus fixes the basepoint $I = MIM^T = MM^{-1} = I$. We think of elements fixing the basepoint as being $\mathrm{SPD}_n$-*rotations* or $\mathrm{SPD}_n$-*reflections*, when the matrix $M$ is a familiar rotation or reflection (see Figure 2).

The Euclidean symmetry of *reflecting in a point* also has a natural generalization to $\mathrm{SPD}_n$. Euclidean reflection in the origin is given by $p \mapsto -p$; and its $\mathrm{SPD}_n$-analog, reflection in the basepoint $I$, is matrix inversion $P \mapsto P^{-1}$. The general $\mathrm{SPD}_n$-reflection in a point $Q \in \mathrm{SPD}_n$ is a conjugate of this by an $\mathrm{SPD}_n$ translation, given by $P \mapsto QP^{-1}Q$.

## 3.1 Vector-valued Distance Function

In Euclidean or hyperbolic spaces, the relative position between two points is completely determined by their distance, which is given by a scalar. For the space $\mathrm{SPD}_n$, it is determined by a vector.

**The VVD vector:** To assign this vector in SPD we introduce the vector-valued distance (VVD) function $d_{vv}\colon \mathrm{SPD}_n \times \mathrm{SPD}_n \to \mathbb{R}^n$. For two points $P, Q \in \mathrm{SPD}_n$, the VVD is defined as:

$$d_{vv}(P,Q) = \log(\lambda_1(P^{-1}Q), \ldots, \lambda_n(P^{-1}Q)) \tag{1}$$

where $\lambda_1(P^{-1}Q) \geq \ldots \geq \lambda_n(P^{-1}Q)$ are the eigenvalues of $P^{-1}Q$ sorted in descending order. This vector is an invariant of the relative position of two points up to isometry. This means that in $\mathrm{SPD}_n$, only if the VVD between two points $A$ and $B$ is the vector $v \in \mathbb{R}^n$, and the VVD between $P$ and $Q$ is also $v$, then there exists an isometry mapping $A$ to $P$ and $B$ to $Q$. Thus, we can recover completely the relative position of two points in $\mathrm{SPD}_n$ from this vector. For example, the Riemannian metric is obtained by using the standard $l_2$ norm on the VVD vector. This is: $d^R(P,Q) = ||d_{vv}(P,Q)||_2$. See [46] §2.6 and Appendix C.4 for a review of VVDs in symmetric spaces.

3

**Finsler metrics:** Any norm on $\mathbb{R}^n$ that is invariant under permutation of the entries induces a metric on $\mathrm{SPD}_n$. Moreover, $\mathrm{SPD}_n$ do not only support a Riemannian metric, but also Finsler metrics, a whole family of distances with the same symmetry group (group of isometries). These metrics are of special importance since distance minimizing geodesics are not necessarily unique in Finsler geometry. Two



VVD

Finsler Norms

0.8    1.9    1.5    1.7    1.1    1.3    2.0    Metric Distances

Figure 3: The vector-valued distance allows to reconstruct the Riemannian, or any Finsler distance.

different paths can have the same minimal length. This is particularly valuable when embedding graphs in $\mathrm{SPD}_n$, since in graphs there are generally several shortest paths. We obtain the Finsler metrics $\mathrm{F}_1$ or $\mathrm{F}_\infty$ by taking the respective $\ell_1$ or $\ell_\infty$ norms of the VVD in $\mathbb{R}^n$ (see Figure 3). See Planche [66] and Appendix C.6 for a review of the theory of Finsler metrics, [12] for the study of some Finsler metrics on $\mathrm{SPD}_n$, and [52] for applications of Finsler metrics on symmetric spaces in representation learning.
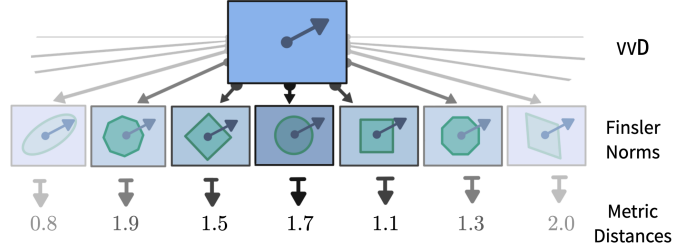
**Advantages of VVD:** The proposed metric learning framework based on the VVD offers several advantages. First, a single model can be run with different metrics, according to the chosen norm. The VVD contains the full information of the Riemannian distance and of all invariant Finsler distances, hence we can easily recover the Riemannian metric and extend the approach to many other alternatives (in Appendix C.7, we detail how the VVD generalizes other $\mathrm{SPD}_n$ metrics). Second, the VVD provides much more information than just the distance, and can be used to analyze the learned representation in $\mathrm{SPD}_n$, independent of the choice of the metric. Out of the VVD between two points, one can immediately read the regularity of the unique geodesics joining these two points. Geodesics in $\mathrm{SPD}_n$ have different regularity, which is related with the number of maximal Euclidean subspaces that contain the geodesic. The Riemannian or Finsler distances cannot distinguish the differences between these geodesics of different regularity, but the VVD function can. Third, the VVD function can be leveraged as a tool to visualize and analyze the learned high-dimensional representations (see §6.4).

## 4 Gyrocalculus

To build an analog of many Euclidean operators in $\mathrm{SPD}_n$, we require also a translation of operations internal to Euclidean geometry, chief among these being the vector space operations of addition and scalar multiplication. By means of tools introduced in pure mathematics literature[2], we describe a gyro-vector space structure on $\mathrm{SPD}_n$, which provides geometrically meaningful extensions of these vector space operations, extending successful applications of this framework in geometric deep learning to hyperbolic space [32, 54, 70]. These operations provide a template for translation, where one may attempt to replace $+, -, \times$ in formulas familiar from Euclidean spaces with the analogous operations $\oplus, \ominus, \otimes$ on $\mathrm{SPD}_n$. While straightforward, such translation requires some care, as gyro-addition is neither commutative nor associative. See Appendix D for a review of the underlying general theory of gyrogroups and additional guidelines for accurate formula translation.

**Addition and Subtraction:** Given a fixed choice $I$ of basepoint and two points $P, Q \in \mathrm{SPD}_n$, we define the gyroaddition of $P$ and $Q$ to be the point $P \oplus Q \in \mathrm{SPD}_n$ which is the image of $Q$ under the isometry which translates $I$ to $P$ along the geodesic connecting them. This directly generalizes the gyroaddition of hyperbolic space exploited by [7, 32, 70], via the geometric interpretation of Vermeer [84] (see Figure 4).

Fixing $P \in \mathrm{SPD}_n$, we may compute the value of $P \oplus Q$ for arbitrary $Q$ as the result of applying the $\mathrm{SPD}_n$-translation moving the basepoint to $P$, evaluated on $Q$. We see also the additive inverse of a point with respect to this operation must then be given by its geodesic reflection in $I$.

$$P \oplus Q = \sqrt{P}Q\sqrt{P} \qquad \ominus P = P^{-1} \tag{2}$$

---

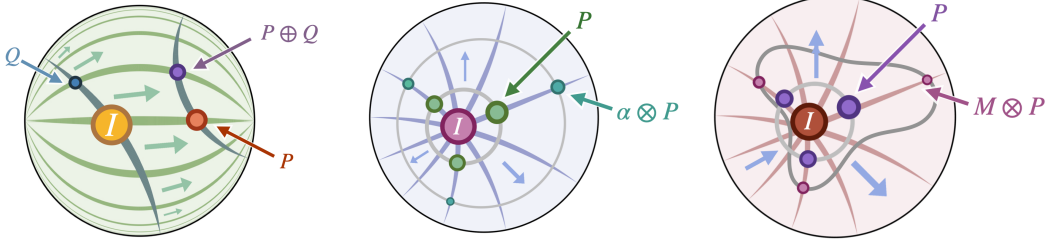[2]See [1, 37] for an abstract treatment of gyrocalculus, and [47] for the specific examples discussed here.

Figure 4: Gyro-addition (left), gyro-scalar multiplication (center) and matrix scaling (right).

As this operation encodes a symmetry of $\mathrm{SPD}_n$, it is possible to recast certain geometric statements purely in the gyrovector formalism. In particular, the vector-valued distance $d_{vv}(P, Q)$ may be computed as the logarithm of the eigenvalues of $\ominus P \oplus Q$ (see Appendix C.5).

**Scalar Multiplication and Matrix Scaling:** For a fixed basepoint $I$, we define the scalar multiplication of a point $P \in \mathrm{SPD}_n$ by a scalar $\alpha \in \mathbb{R}_+$ to be the point which lies at distance $\alpha d(I, P)$ from $I$ in the direction of $P$, where $d(\cdot, \cdot)$ is the metric distance on $\mathrm{SPD}_n$. Geometrically, this is a transfer of the vector-space scalar multiplication on the tangent space to $\mathrm{SPD}_n$:

$$\alpha \otimes P = P^\alpha = \exp(\alpha \log(P)), \tag{3}$$

where $\exp, \log$ are the matrix exponential and logarithm. We further generalize the notion of scalar multiplication to allow for different relative expansion rates in different directions. For a fixed basepoint $I$ and a point $P \in \mathrm{SPD}_n$, we can replace the scalar $\alpha$ from Eq. 3 with an arbitrary real symmetric matrix $A \in S_n$. We define this *matrix scaling* by:

$$A \otimes P = \exp(A \odot \log(P)) \tag{4}$$

where $A \odot X$ denotes the Hadamard product. We denote the matrix scaling with $\otimes$, extending the previous usage: for any $\alpha \in \mathbb{R}$, we have $[\alpha] \otimes P = \alpha \otimes P$ where $[\alpha]$ is the matrix with every entry $\alpha$.

## 5 Implementation

In this section we detail how we learn representations in $\mathrm{SPD}_n$, and implement different linear mappings so that they conform to the premises of each operator, yielding SPD neural layers.

**Embeddings in $\mathrm{SPD}_n$ and $S_n$:** We are interested in learning embeddings in $\mathrm{SPD}_n$. To do so we exploit the connection between $\mathrm{SPD}_n$ and its tangent space $S_n$ through the exponential and logarithmic maps. To learn a point $P \in \mathrm{SPD}_n$, we first model it as a symmetric matrix $U \in S_n$. We impose symmetry on $U$ by learning a triangular matrix $X \in \mathbb{R}^{n \times n}$ with $n(n+1)/2$ parameters, such that $U = X + X^T$. To obtain the matrix $P \in \mathrm{SPD}_n$, we employ the exponential map: $P = \exp(U)$. Modeling points on the tangent space offers advantages for optimization, explained in §5. For the matrix scaling $A \otimes P$, we impose symmetry on the factor matrix $A \in S_n$ in the same way that we learn the symmetric matrix $U$.

**Isometries: Rotations and Reflections:** Rotations in $n$ dimensions are described as collections of pairwise orthogonal 2-dimensional rotations in planes (with a leftover 1-dimensional "axis of rotation" in odd dimensions). We utilize this observation to efficiently build elements of $O(n)$ out of two-dimensional rotations in coordinate planes. More precisely, for any $\theta \in [0, 2\pi]$ and choice of sign $\{+, -\}$ we let $R^\pm(\theta)$ denote the 2-dimensional rotation (+) or reflection (−) as $R^\pm(\theta) = \left(\begin{smallmatrix} \cos\theta & \mp\sin\theta \\ \sin\theta & \pm\cos\theta \end{smallmatrix}\right)$.

Then for any pair $i < j$ in $1 \ldots n$, we denote by $R_{ij}^\pm(\theta)$ the transformation which applies $R^\pm(\theta)$ to the $x_i x_j$-plane of $\mathbb{R}^n$, and leaves all other coordinates fixed. For example, in $O(5)$ the element $R_{24}^+(\theta)$ (see on the right) denotes the transformation where we replace the entries $(ii, ij, ji, jj)$ of $I_n$ with the corresponding values of $R^+(\theta)$. More general, rotations and reflections are built by taking products of these basic transformations. Given a $n(n-1)/2$-dimensional vector of angles $\vec{\theta} = (\theta_{12}, \ldots, \theta_{ij}, \ldots, )$

$$R_{24}^+(\theta) = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & \cos\theta & 0 & -\sin\theta & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & \sin\theta & 0 & \cos\theta & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

and a choice of sign, we define the rotation and reflection corresponding to $\vec{\theta}$ by:

$$\mathrm{Rot}(\vec{\theta}) = \prod_{i<j} R^+_{ij}(\theta_{ij}) \qquad \mathrm{Refl}(\vec{\theta}) = \prod_{i<j} R^-_{ij}(\theta_{ij}) \tag{5}$$

where $\mathrm{Rot}(\vec{\theta}), \mathrm{Ref}(\vec{\theta}) \in \mathbb{R}^{n \times n}$ are the isometry matrices, and the vector of angles $\vec{\theta}$ can be regarded as a learnable parameter of the model. Finally, we denote the application of the transformation $M$ to the point $P \in \mathrm{SPD}_n$ by:

$$M \odot P = MPM^T \tag{6}$$

**Optimization:** For the proposed rotations and reflections, the learnable weights are vectors of angles $\vec{\theta} \in \mathbb{R}^{\frac{n(n-1)}{2}}$, which do not pose an optimization challenge. On the other hand, embeddings in SPD have to be optimized respecting the geometry of the manifold, but as already explained, we model them on the space of symmetric matrices $S_n$, and then we apply the exponential map. In this manner, we are able to perform tangent space optimization [26] using standard Euclidean techniques, and circumvent the need for Riemannian optimization [15, 11], which we found to be less numerically stable. Due to the geometry of $\mathrm{SPD}_n$ (see Appendix C.3), this is an exact procedure, which does not incur losses in representational power.

**Complexity:** The most frequently utilized operation when learning embeddings is the distance calculation, thus we analyze its complexity. In Appendix A.1 we detail the complexity of different operations. Calculating the distance between two points in $\mathrm{SPD}_n$ implies computing multiplications, inversions and diagonalizations of $n \times n$ matrices. We find that the cost of the distance computation with respect to the matrix dimensions is $\mathcal{O}(n^3)$. Although a matrix of rank $n$ implies $n(n+1)/2$ dimensions thus a large $n$ value is usually not required, the cost of many operations is polynomial instead of linear.

**Towards neural network architectures:** We employ the proposed mappings along with the gyro-vector operations as building blocks for SPD neural layers. This is showcased in the experiments presented below. Scalings, rotations and reflections can be seen as feature transformations. Moreover, gyro-addition allows us to define the equivalent of bias addition. Finally, although we do not employ non-linearities, our approach can be seamlessly integrated with the ReEig layer (adaptation of a ReLU layer for SPD) proposed in [40].

# 6 Experiments

In this section we employ the transformations developed on SPD to build neural models for knowledge graph completion, item recommendation and question answering. Task-specific models in different geometries have been developed in the three cases, hence we consider them adequate benchmarks for representation learning.[3]

## 6.1 Knowledge Graph Completion

Knowledge graphs (KG) represent heterogeneous knowledge in the shape of *(head, relation, tail)* triples, where *head* and *tail* are entities and *relation* represents a relationships among entities. KG exhibit an intricate and varying structure where entities can be connected by symmetric, anti-symmetric, or hierarchical relations. To capture these non-trivial patterns more expressive modelling techniques become necessary [25], thus we choose this application to showcase the capabilities of our transformations on SPD manifolds. Given an incomplete KG, the task is to predict which unknown links are valid.

**Problem formulation:** Let $\mathcal{G} = (\mathcal{E}, \mathcal{R}, \mathcal{T})$ be a knowledge graph where $\mathcal{E}$ is the set of entities, $\mathcal{R}$ is the set of relations and $\mathcal{T} \subset \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ is the set of triples stored in the graph. The usual approach is to learn a scoring function $\phi : \mathcal{E} \times \mathcal{R} \times \mathcal{E} \to \mathbb{R}$ that measures the likelihood of a triple to be true, with the goal of scoring all missing triples correctly. To do so, we propose to learn representations of entities as embeddings in $\mathrm{SPD}_n$, and relation-specific transformation in the manifold, such that the KG structure is preserved.

---

[3]Code available at `https://github.com/fedelopez77/gyrospd`

**Scaling model:** We follow the base hyperbolic model MuRP [8] and adapt it into $\mathrm{SPD}_n$ by means of the *matrix scaling*. Its scoring function has shown success in the task given that it combines multiplicative and additive components, which are fundamental to model different properties of KG relations [4]. We translate it into $\mathrm{SPD}_n$ as:

$$\phi(h, r, t) = -d((\mathbf{M}_r \otimes \mathbf{H}) \oplus \mathbf{R}, \mathbf{T})^2 + b_h + b_t \tag{7}$$

where $\mathbf{H}, \mathbf{T} \in \mathrm{SPD}_n$ are embeddings and $b_h, b_t \in \mathbb{R}$ are scalar biases for the head and tail entities respectively. $\mathbf{R} \in \mathrm{SPD}_n$ and $\mathbf{M}_r$ are matrices that depend on the relation. For $d(\cdot, \cdot)$, we experiment with the Riemannian and the Finsler One metric distances.

**Isometric model:** A possible alternative is to embed the relation-specific transformations as elements of the $O(n)$ group (*i.e.*, rotations and reflections). This technique has proven effective in different metric spaces [25, 88]. In this case, $\mathbf{M}_r$ is a rotation or reflection matrix as in Eq. 5, and the scoring function is defined as:

$$\phi(h, r, t) = -d((\mathbf{M}_r \odot \mathbf{H}) \oplus \mathbf{R}, \mathbf{T})^2 + b_h + b_t \tag{8}$$

**Datasets:** We employ two standard benchmarks, namely WN18RR [16, 28] and FB15k-237 [16, 76]. WN18RR is a subset of WordNet [59] containing 11 lexical relationships between $40,943$ word senses. FB15k-237 is a subset of Freebase [14], with $14,541$ entities and $237$ relationships.

**Training:** We follow the standard data augmentation protocol by adding inverse relations to the datasets [48]. We optimize the cross-entropy loss with uniform negative sampling defined in Equation 9, where

$$\mathcal{L} = \sum_{(h,r,t) \in \mathcal{T}} \log(1 + \exp(Y_t \phi(h, r, t))) \tag{9}$$

$\mathcal{T}$ is the set of training triples, and $Y_t = -1$ if $t$ is a factual triple or $Y_t = 1$ if $t$ is a negative sample. We employ the AdamW optimizer [55]. We conduct a grid search with matrices of dimension $n \times n$ where $n \in \{14, 20, 24\}$ (this is the equivalent of $\{105, 210, 300\}$ degrees of freedom respectively) to select optimal dimensions, learning rate and weight decay, using the validation set. More details and set of hyperparameters in Appendix B.1.

**Evaluation metrics:** At test time, we rank the correct tail or head entity against all possible entities using the scoring function, and use inverse relations for head prediction [48]. Following previous work, we compute two ranking-based metrics: mean reciprocal rank (MRR), which measures the mean of inverse ranks assigned to correct entities, and hits at K (H@K, $K \in \{1, 3, 10\}$), which measures the proportion of correct triples among the top K predicted triples. We follow the standard evaluation protocol of filtering out all true triples in the KG during evaluation [16].

**Baselines:** We compare our models with their respective equivalents in different metric spaces, which are also state-of-the-art models for the task. For the scaling model, these are MURE and MURP [8], which perform the scaling operation in Euclidean and hyperbolic space respectively. For the isometric models, we compare to ROTC [73], ROTE and ROTH, [25] (rotations in Complex, Euclidean and hyperbolic space respectively), and REFE and REFH [25] (reflections in Euclidean and hyperbolic space). Baseline results are taken from the original papers. We do not compare to previous work on SPD given that they lack the definition of an arithmetic operation in the space, thus a vis-a-vis comparison is not possible.

**Results:** We report the performance for all analyzed models, segregated by operation, in Table 1. On both dataset, the scaling model $\mathrm{SPD}_{\mathrm{Sca}}$ outperforms its direct competitors MuRE and MuRP, and this is specially notable in HR@10 for WN18RR: 59.0 for $\mathrm{SPD}_{\mathrm{Sca}}^{F_1}$ vs 55.4 and 56.6 respectively. SPD reflections are very effective on WN18RR as well. They outperform their Euclidean and hyperbolic counterparts RefE and RefH, in particular when equipped with the Finsler metric. Rotations on the SPD manifold, on the other hand, seem to be less effective. However, Euclidean and hyperbolic rotations require $500$ dimensions whereas the $\mathrm{SPD}_{\mathrm{Rot}}$ models are trained on matrices of rank 14 (equivalent to 105 dims). Moreover, the underperformance observed in some of the analyzed cases for rotations and reflections does not repeat in the following experiments (§6.2 & §6.3). Hence, we consider this is due to overfitting in some particular setups. Although we tried different regularization methods, we regard a sub-optimal configuration rather than a geometric reason to be the cause for the underperformance.

Regarding the choice of a distance metric, the Finsler One metric is better suited with respect to HR@3 and HR@10 when using scalings and reflections on WN18RR. For the FB15k-237 dataset, SPD models operating with the Riemannian metric outperform their Finsler counterparts. This

Table 1: Results for Knowledge graph completion.

| Operation | Model | WN18RR | | | | FB15k-237 | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | MRR | HR@1 | HR@3 | HR@10 | MRR | HR@1 | HR@3 | HR@10 |
| Scaling | MuRE | 47.5 | 43.6 | 48.7 | 55.4 | 33.6 | 24.5 | 37.0 | 52.1 |
| | MuRP | 48.1 | **44.0** | 49.5 | 56.6 | 33.5 | 24.3 | 36.7 | 51.8 |
| | $\mathrm{SPD}_{\mathrm{Sca}}^{R}$ | 48.1 | 43.1 | 50.1 | 57.6 | **34.5** | **25.1** | **38.0** | **53.5** |
| | $\mathrm{SPD}_{\mathrm{Sca}}^{F_1}$ | **48.4** | 42.6 | **51.0** | **59.0** | 32.9 | 23.6 | 36.3 | 51.5 |
| Rotations | RotC | 47.6 | 42.8 | 49.2 | 57.1 | 33.8 | 24.1 | 37.5 | 53.3 |
| | RotE | 49.4 | 44.6 | 51.2 | 58.5 | **34.6** | **25.1** | **38.1** | **53.8** |
| | RotH | **49.6** | **44.9** | **51.4** | 58.6 | 34.4 | 24.6 | 38.0 | 53.5 |
| | $\mathrm{SPD}_{\mathrm{Rot}}^{R}$ | 46.2 | 39.7 | 49.6 | 57.8 | 32.9 | 23.6 | 36.3 | 51.6 |
| | $\mathrm{SPD}_{\mathrm{Rot}}^{F_1}$ | 40.9 | 30.5 | 48.2 | 57.3 | 32.1 | 22.9 | 35.4 | 50.5 |
| Reflections | RefE | 47.3 | 43.0 | 48.5 | 56.1 | **35.1** | **25.6** | **39.0** | **54.1** |
| | RefH | 46.1 | 40.4 | 48.5 | 56.8 | 34.6 | 25.2 | 38.3 | 53.6 |
| | $\mathrm{SPD}_{\mathrm{Ref}}^{R}$ | 48.3 | 44.0 | 49.7 | 56.7 | 32.5 | 23.4 | 35.6 | 51.0 |
| | $\mathrm{SPD}_{\mathrm{Ref}}^{F_1}$ | **48.7** | **44.3** | **50.1** | **57.4** | 31.6 | 22.5 | 34.6 | 50.0 |

suggests that the Riemannian metric is capable of disentangling the large number of relationships in this dataset to a better extent.

In these experiments we have evaluated models applying equivalent operations and scoring functions in different geometries, thus they can be thought as a vis-a-vis comparison of the metric spaces. We observe that SPD models tie or outperform baselines in most instances. This showcases the improved representation capacity of the SPD manifold when compared to Euclidean and hyperbolic spaces. Moreover, it demonstrates the effectiveness of the proposed metrics and operations in this manifold.

## 6.2 Knowledge Graph Recommender Systems

Recommender systems (RS) model user preferences to provide personalized recommendations [93]. KG embedding methods have been widely adopted into the recommendation problem as an effective tool to model side information and enhance the performance [91, 34]. For instance, one reason for recommending a movie to a particular user is that the user has already watched many movies from the same genre or director [56]. Given multiple relations between users, items, and heterogeneous entities, the goal is to predict the user's next item purchase or preference.

**Model:** We model the recommendation problem as a link prediction task over users and items [49]. In addition, we aim to incorporate side information between users, items and other entities. Hence we apply our KG embedding method from §6.1 as is, to embed this multi-relational graph. We evaluate the capabilities of the approach by only measuring the performance over user-item interactions.

**Datasets:** To investigate the recommendation problem endowed with added relationships, we employ the Amazon dataset [58, 61] (branches "Software", "Luxury & Beauty" and "Prime Pantry"), with users' purchases of products, and the MindReader dataset [18] of movie recommendations. Both datasets provide additional relationships between users, items and entities such as product brands, or movie directors and actors. To generate evaluation splits, the penultimate and last item the user has interacted with are withheld as dev and test sets respectively.

**Training:** In this setup we also augment the data by adding inverse relations and optimize the loss from Equation 9. We set the size of the matrices to $10 \times 10$ dimensions (equivalent to 55 free parameters). More details about relationships and set of hyperparameters in Appendix B.2.

**Evaluation and metrics:** We follow the standard procedure of evaluating against 100 randomly selected samples the user has not interacted with [38, 53]. To evaluate the recommendation performance we focus on the *buys / likes* relation. For each user $u$ we rank the items $i_j$ according to the scoring function $\phi(u, buys, i_j)$. We adopt MRR and H@10, as ranking metrics for recommendations.

**Baselines:** We compare to TransE [16], RotC [73], MuRE and MuRP [8] trained with 55 dimensions.

**Results:** In Table 2 we observe that the SPD models tie or outperform the baselines in both MRR and HR@10 across all analyzed datasets. Rotations in both, Riemannian and Finsler metrics, are more effective in this task, achieving the best performance in 3 out of 4 cases, followed by the scaling models. Overall, this shows the capabilities of the systems to effectively represent user-item interactions enriched with relations between items and their attributes, thus learning to better model users' preferences. Furthermore, it displays the versatility of the approach to diverse data domains.

Table 2: Results for Knowledge graph-based recommender systems.

| | SOFTWARE | | LUXURY | | PANTRY | | MINDREADER | |
| Model | MRR | H@10 | MRR | H@10 | MRR | H@10 | MRR | H@10 |
|---|---|---|---|---|---|---|---|---|
| TRANSE | 28.5±0.1 | 47.2±0.5 | 35.6±0.1 | 52.3±0.1 | 16.6±0.0 | 35.3±0.1 | 19.1±0.4 | 37.6±0.1 |
| ROTC | 28.5±0.3 | 45.4±1.4 | 33.0±0.1 | 49.8±0.2 | 14.5±0.0 | 31.3±0.2 | 25.3±0.3 | 50.3±0.6 |
| MURE | 29.4±0.4 | 47.1±0.4 | 35.6±0.7 | 54.0±0.3 | 19.4±0.1 | 39.5±0.2 | 25.2±0.3 | 49.9±0.6 |
| MURP | 29.6±0.3 | 47.9±0.3 | **37.5±0.1** | **55.2±0.3** | 19.4±0.1 | 39.8±0.2 | 25.3±0.3 | 49.3±0.2 |
| $\text{SPD}_{\text{Sca}}^{R}$ | 29.4±0.4 | 48.1±0.8 | **37.5±0.2** | 55.1±0.2 | 19.5±0.0 | 39.6±0.3 | 25.4±0.1 | 49.8±0.3 |
| $\text{SPD}_{\text{Sca}}^{F_1}$ | 28.8±0.1 | 46.9±0.5 | 37.3±0.3 | 54.1±0.9 | 19.0±0.1 | 38.8±0.2 | **25.7±0.5** | 49.5±0.1 |
| $\text{SPD}_{\text{Rot}}^{R}$ | **30.3±0.2** | 48.6±0.9 | 37.2±0.1 | 54.8±0.4 | **20.0±0.1** | **40.3±0.1** | 25.3±0.0 | **50.5±0.3** |
| $\text{SPD}_{\text{Rot}}^{F_1}$ | 30.1±0.1 | **49.1±0.3** | 36.9±0.1 | 54.5±0.6 | 19.2±0.0 | 39.3±0.1 | **25.7±0.0** | 49.5±0.2 |
| $\text{SPD}_{\text{Ref}}^{R}$ | 29.6±0.2 | 48.0±0.5 | 37.3±0.2 | 55.0±0.2 | 19.3±0.0 | 39.7±0.3 | 25.3±0.0 | 49.1±0.1 |
| $\text{SPD}_{\text{Ref}}^{F_1}$ | 29.3±0.1 | 47.5±0.6 | 36.8±0.0 | 54.8±0.1 | 18.6±0.2 | 38.3±0.3 | 24.8±0.2 | 47.9±1.8 |

## 6.3 Question Answering

We evaluate our approach on the task of Question Answering (QA). In this manner we also showcase the capabilities of our methods to train word embeddings.

**Model:** We adapt the model from HyperQA [74] to SPD. We model word embeddings $t_i \in \text{SPD}_n$, and represent question/answers as a summation of the embeddings of their corresponding tokens. We apply a feature transformation $T(\cdot)$ followed by a bias addition, as an equivalent of a neural linear layer. $T(\cdot)$ can be a scaling, rotation or reflection. Finally we compute a distance-based similarity function between the resulting question/answer representations as defined in Equation 10, where $w_f, w_b \in \mathbb{R}$, $B \in \text{SPD}_n$ and the transformation $T$ are parameters of the model.

$$\text{sim}(q, a) = -w_f d(\mathbf{Q}, \mathbf{A}) + w_b,$$
$$\text{where } \mathbf{Q} = T(\bigoplus_{i=1}^{n} t_i^q) \oplus B \quad (10)$$

**Datasets:** We analyze two popular benchmarks for QA: TrecQA [85] (clean version) and WikiQA [89], filtering out questions with multiple answers from the dev and test sets.

**Training:** We optimize the cross-entropy loss from Eq. 9, where we replace $\phi$ for $\text{sim}(q, a)$ and for each question we use wrong answers as negative samples. We set the size of the matrices to $14 \times 14$ dimensions (equivalent to 105 free parameters). The set of hyperparameters can be found in Appendix B.3.

**Evaluation metrics:** At test time, for each question we rank its candidate answers according to Eq. 10. We adopt MRR and H@1 as evaluation metrics.

Table 3: Results for Question Answering.

| | TRECQA | | WIKIQA | |
| Model | MRR | H@1 | MRR | H@1 |
|---|---|---|---|---|
| Euclidean | 55.9±2.0 | 41.0±2.0 | 43.4±0.3 | 22.4±1.1 |
| Hyperbolic | 58.0±1.3 | 39.3±2.0 | 44.0±0.4 | 22.8±0.6 |
| $\text{SPD}_{\text{Sca}}^{R}$ | 55.4±0.1 | 37.1±0.1 | **45.5±0.5** | 24.4±1.1 |
| $\text{SPD}_{\text{Sca}}^{F_1}$ | 57.1±0.7 | 38.6±0.2 | 44.8±0.5 | 24.0±0.6 |
| $\text{SPD}_{\text{Rot}}^{R}$ | 58.7±1.5 | 41.4±2.9 | 44.6±0.6 | 23.6±0.6 |
| $\text{SPD}_{\text{Rot}}^{F_1}$ | 58.1±0.5 | **43.6±1.0** | 43.7±0.4 | 23.8±0.8 |
| $\text{SPD}_{\text{Ref}}^{R}$ | 57.3±0.3 | 40.7±1.1 | 43.9±0.7 | 23.4±2.0 |
| $\text{SPD}_{\text{Ref}}^{F_1}$ | **59.6±0.5** | 42.1±1.0 | 44.7±1.2 | **25.0±2.5** |

**Baselines:** We compare against Euclidean and hyperbolic spaces of 105 dimensions. For the Euclidean model we employ a linear layer as feature transformation. For the hyperbolic model, we operate on the tangent space and project the points into the Poincaré ball to compute distances.

**Results:** We present results in Table 3. In both datasets, we see that the word embeddings and transformations learned by the SPD models are able to place questions and answers representations in the space such that they outperform Euclidean and hyperbolic baselines. Finsler metrics seem to be very effective in this scenario, improving the performance of their Riemannian counterparts in many cases. Overall, this suggests that embeddings in SPD manifolds learn meaningful representations that can be exploited into downstream tasks. Moreover, we showcase how to employ different operations as feature transformations and bias additions, replicating the behavior of linear layers in classical deep learning architectures that can be seamlessly integrated with different distance metrics.

## 6.4 Analysis

One reason to embed data into Riemannian manifolds, such as SPD, is to use geometric properties of the manifold to analyze the structure of the data [52]. Visualizations in SPD are difficult due to their
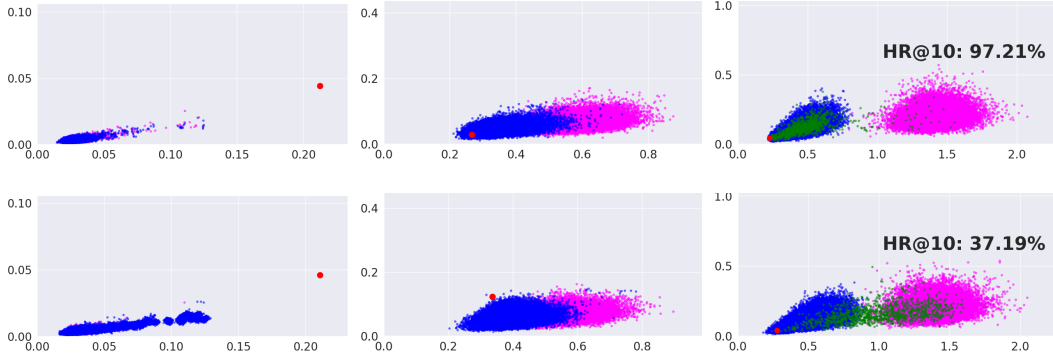
Figure 5: Train, negative and validation triples for relationships *'derivationally related form'* (top) and *'hypernym'* (bottom) for 5 (left), 50 (center) and 3000 (right) epochs for the $\mathrm{SPD}_{\mathrm{Sca}}^{F_1}$ model. The red dot corresponds to the relation addition $\mathbf{R}$.

high dimensionality. As a solution we use the vector-valued distance function to develop a new tool to visualize and analyze structural properties of the learned representations.

We adopt the vector $(n-1, n-3, \cdots, -n+3, -n+1)$, as the barycenter of the space in $\mathbb{R}^n$ where the VVD is contained. Then, we plot the norm of the VVD vector and its angle with respect to this barycenter. In Figure 5, we compute and plot the VVD corresponding to $d(\mathbf{M}_r \otimes \mathbf{H}, \mathbf{T})$ and $\mathbf{R}$ as defined in Eq. 7 for KG models trained on WN18RR. In early stages of the training, all points fall near the origin (left side of the plots). As training evolves, the model learns to separate true $(h, r, t)$ triples from corrupted ones (center part). When the training converges, the model is able to clearly disentangle and cluster positive and negative samples. We observe how the position of the validation triples (green points, not seen during training) directly correlates with the performance of each relation. Plots for more relations in Appendix B.1.

# 7 Conclusions

Riemannian geometry has gained attention due to its capacity to represent non-Euclidean data arising in several domains. In this work we introduce the vector-valued distance function, which allows to implement universal models (generalizing previous metrics on SPD), and can be leveraged to provide a geometric interpretation on what the models learn. Moreover, we bridge the gap between Euclidean and SPD geometry under the lens of the gyrovector theory, providing means to transfer standard arithmetic operations from the Euclidean setting to their analog notions in SPD. These tools enable practitioners to exploit the full representation power of SPD, and profit from the enhanced expressivity of this manifold. We propose and evaluate SPD models on three tasks and eight datasets, which showcases the versatility of the approach and ease of integration with downstream tasks. The results reflect the superior expressivity of SPD when compared to Euclidean or hyperbolic baselines.

This work is not without limitations. We consider the computational complexity of working with spaces of matrices to be the main drawback, since the cost of many operations is polynomial instead of linear. Nevertheless, a matrix of rank $n$ implies $n(n+1)/2$ dimensions thus a large $n$ value is usually not required.

## Acknowledgments and Disclosure of Funding

# References

[1] Abe, T. and Hatori, O. (2015). Generalized gyrovector spaces and a mazur–ulam theorem. *Publicationes Mathematicae Debrecen*, 87:393–413.

[2] Absil, P.-A., Mahony, R., and Sepulchre, R. (2009). *Optimization Algorithms on Matrix Manifolds*. Princeton University Press.

[3] Ai, Q., Azizi, V., Chen, X., and Zhang, Y. (2018). Learning heterogeneous knowledge base embeddings for explainable recommendation. *Algorithms*, 11(9).

[4] Allen, C., Balazevic, I., and Hospedales, T. (2021). Interpreting knowledge graph relation representation from word embeddings. In *International Conference on Learning Representations*.

[5] Arsigny, V., Fillard, P., Pennec, X., and Ayache, N. (2006a). Geometric means in a novel vector space structure on symmetric positive-definite matrices. *SIAM J. Matrix Anal. Appl.*, 29(1):328–347.

[6] Arsigny, V., Fillard, P., Pennec, X., and Ayache, N. (2006b). Log-euclidean metrics for fast and simple calculus on diffusion tensors. *Magnetic Resonance in Medicine*, 56(2):411–421.

[7] Bachmann, G., Becigneul, G., and Ganea, O.-E. (2020). Constant curvature graph convolutional networks. In *37th International Conference on Machine Learning (ICML)*.

[8] Balazevic, I., Allen, C., and Hospedales, T. (2019a). Multi-relational poincaré graph embeddings. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 32, pages 4463–4473. Curran Associates, Inc.

[9] Balazevic, I., Allen, C., and Hospedales, T. (2019b). TuckER: Tensor factorization for knowledge graph completion. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5185–5194, Hong Kong, China. Association for Computational Linguistics.

[10] Ballmann, W., Brin, M., and Eberlein, P. (1985). Structure of manifolds of nonpositive curvature. I. *Annals of Mathematics*, 122(1):171–203.

[11] Bécigneul, G. and Ganea, O. (2019). Riemannian adaptive optimization methods. In *7th International Conference on Learning Representations, ICLR*, New Orleans, LA, USA.

[12] Bhatia, R. (2003). On the exponential metric increasing property. *Linear Algebra and its Applications*, 375:211–220.

[13] Bhatia, R., Jain, T., and Lim, Y. (2019). On the Bures–Wasserstein distance between positive definite matrices. *Expositiones Mathematicae*, 37(2):165–191.

[14] Bollacker, K., Evans, C., Paritosh, P., Sturge, T., and Taylor, J. (2008). Freebase: A collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, SIGMOD '08, page 1247–1250, New York, NY, USA. Association for Computing Machinery.

[15] Bonnabel, S. (2011). Stochastic gradient descent on Riemannian manifolds. *IEEE Transactions on Automatic Control*, 58.

[16] Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., and Yakhnenko, O. (2013). Translating embeddings for modeling multi-relational data. In Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems*, volume 26, pages 2787–2795. Curran Associates, Inc.

[17] Boumal, N., Mishra, B., Absil, P.-A., and Sepulchre, R. (2014). Manopt, a matlab toolbox for optimization on manifolds. *Journal of Machine Learning Research*, 15(42):1455–1459.

[18] Brams, A. H., Jakobsen, A. L., Jendal, T. E., Lissandrini, M., Dolog, P., and Hose, K. (2020). Mindreader: Recommendation over knowledge graph entities with explicit user ratings. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, CIKM '20, page 2975–2982, New York, NY, USA. Association for Computing Machinery.

[19] Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., and Vandergheynst, P. (2017). Geometric deep learning: Going beyond Euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42.

[20] Brooks, D., Schwander, O., Barbaresco, F., Schneider, J.-Y., and Cord, M. (2019a). Riemannian batch normalization for SPD neural networks. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 32, pages 15489–15500. Curran Associates, Inc.

[21] Brooks, D. A., Schwander, O., Barbaresco, F., Schneider, J.-Y., and Cord, M. (2019b). Exploring complex time-series representations for riemannian machine learning of radar data. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3672–3676.

[22] Carreira, J., Caseiro, R., Batista, J., and Sminchisescu, C. (2012). Semantic segmentation with second-order pooling. In Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., and Schmid, C., editors, *Computer Vision – ECCV 2012*, pages 430–443, Berlin, Heidelberg. Springer Berlin Heidelberg.

[23] Chakraborty, R., Bouza, J., Manton, J., and Vemuri, B. C. (2020). Manifoldnet: A deep neural network for manifold-valued data with applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1.

[24] Chakraborty, R., Yang, C.-H., Zhen, X., Banerjee, M., Archer, D., Vaillancourt, D., Singh, V., and Vemuri, B. (2018). A statistical recurrent model on the manifold of symmetric positive definite matrices. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.

[25] Chami, I., Wolf, A., Juan, D.-C., Sala, F., Ravi, S., and Ré, C. (2020). Low-dimensional hyperbolic knowledge graph embeddings. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6901–6914, Online. Association for Computational Linguistics.

[26] Chami, I., Ying, Z., Ré, C., and Leskovec, J. (2019). Hyperbolic graph convolutional neural networks. In *Advances in Neural Information Processing Systems 32*, pages 4869–4880. Curran Associates, Inc.

[27] Cruceru, C., Becigneul, G., and Ganea, O.-E. (2020). Computationally tractable Riemannian manifolds for graph embeddings. In *37th International Conference on Machine Learning (ICML)*.

[28] Dettmers, T., Pasquale, M., Pontus, S., and Riedel, S. (2018). Convolutional 2d knowledge graph embeddings. In *Proceedings of the 32th AAAI Conference on Artificial Intelligence*, pages 1811–1818.

[29] Dong, Z., Jia, S., Zhang, C., Pei, M., and Wu, Y. (2017). Deep manifold learning of symmetric positive definite matrices with application to face recognition. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, AAAI'17, page 4009–4015. AAAI Press.

[30] Donoho, D. L. and Tsaig, Y. (2008). Fast solution of $\ell_1$-norm minimization problems when the solution may be sparse. *IEEE Trans. Information Theory*, 54(11):4789–4812.

[31] Feragen, A., Lauze, F., and Hauberg, S. (2015). Geodesic exponential kernels: When curvature and linearity conflict. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3032–3042.

[32] Ganea, O., Becigneul, G., and Hofmann, T. (2018). Hyperbolic neural networks. In *Advances in Neural Information Processing Systems 31*, pages 5345–5355. Curran Associates, Inc.

[33] Gao, Z., Wu, Y., Bu, X., Yu, T., Yuan, J., and Jia, Y. (2019). Learning a robust representation via a deep network on symmetric positive definite manifolds. *Pattern Recognition*, 92:1–12.

[34] Guo, Q., Zhuang, F., Qin, C., Zhu, H., Xie, X., Xiong, H., and He, Q. (2020). A survey on knowledge graph-based recommender systems. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1.

[35] Ha Quang, M., San Biagio, M., and Murino, V. (2014). Log-Hilbert-Schmidt metric between positive definite operators on Hilbert spaces. In Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.

[36] Harandi, M. T., Salzmann, M., and Hartley, R. (2014). From manifold to manifold: Geometry-aware dimensionality reduction for SPD matrices. In Fleet, D., Pajdla, T., Schiele, B., and Tuytelaars, T., editors, *Computer Vision – ECCV 2014*, pages 17–32, Cham. Springer International Publishing.

[37] Hatori, O. (2017). Examples and applications of generalized gyrovector spaces. *Results in Mathematics*, 71:295–317.

[38] He, X., Liao, L., Zhang, H., Nie, L., Hu, X., and Chua, T.-S. (2017). Neural collaborative filtering. In *Proceedings of the 26th International Conference on World Wide Web*, WWW '17, page 173–182, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.

[39] Helgason, S. (1978). *Differential geometry, Lie groups, and symmetric spaces*. Academic Press New York.

[40] Huang, Z. and Gool, L. V. (2017). A Riemannian network for SPD matrix learning. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, AAAI'17, page 2036–2042. AAAI Press.

[41] Huang, Z., Wang, R., Shan, S., and Chen, X. (2014). Learning Euclidean-to-Riemannian metric for point-to-set classification. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1677–1684.

[42] Huang, Z., Wang, R., Shan, S., Li, X., and Chen, X. (2015). Log-Euclidean metric learning on symmetric positive definite manifold with application to image set classification. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, ICML'15, page 720–729. JMLR.org.

[43] Ionescu, C., Vantzos, O., and Sminchisescu, C. (2015). Matrix backpropagation for deep networks with structured layers. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 2965–2973.

[44] Jayasumana, S., Hartley, R., Salzmann, M., Li, H., and Harandi, M. (2013). Kernel methods on the Riemannian manifold of symmetric positive definite matrices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

[45] Kapovich, M., Leeb, B., and Millson, J. (2009). Convex functions on symmetric spaces, side lengths of polygons and the stability inequalities for weighted configurations at infinity. *J. Differential Geom.*

[46] Kapovich, M., Leeb, B., and Porti, J. (2017). Anosov subgroups: dynamical and geometric characterizations. *European Journal of Mathematics*, 3(3):808–898.

[47] Kim, S. (2016). Gyrovector spaces on the open convex cone of positive definite matrices. *Mathematics Interdisciplinary Research*, 1(1):173–185.

[48] Lacroix, T., Usunier, N., and Obozinski, G. (2018). Canonical tensor decomposition for knowledge base completion. In Dy, J. and Krause, A., editors, *Proceedings of Machine Learning Research*, volume 80 of *Proceedings of Machine Learning Research*, pages 2863–2872, Stockholmsmässan, Stockholm Sweden. PMLR.

[49] Li, J., Zhang, L., Meng, F., and Li, F. (2014). Recommendation algorithm based on link prediction and domain knowledge in retail transactions. *Procedia Computer Science*, 31:875 – 881. 2nd International Conference on Information Technology and Quantitative Management, ITQM 2014.

[50] Li, P., Xie, J., Wang, Q., and Gao, Z. (2018). Towards faster training of global covariance pooling networks by iterative matrix square root normalization. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 947–955.

[51] Li, Y. and Lu, R. (2018). Locality preserving projection on SPD matrix lie group: algorithm and analysis. *Sci. China Inf. Sci.*, 61(9):092104:1–092104:15.

[52] López, F., Pozzetti, B., Trettel, S., Strube, M., and Wienhard, A. (2021a). Symmetric spaces for graph embeddings: A finsler-riemannian approach. In Meila, M. and Zhang, T., editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 7090–7101. PMLR.

[53] López, F., Scholz, M., Yung, J., Pellat, M., Strube, M., and Dixon, L. (2021b). Augmenting the user-item graph with textual similarity models.

[54] López, F. and Strube, M. (2020). A fully hyperbolic neural model for hierarchical multi-class classification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 460–475, Online. Association for Computational Linguistics.

[55] Loshchilov, I. and Hutter, F. (2019). Decoupled weight decay regularization. In *International Conference on Learning Representations*.

[56] Ma, W., Zhang, M., Cao, Y., Jin, W., Wang, C., Liu, Y., Ma, S., and Ren, X. (2019). Jointly learning explainable rules for recommendation with knowledge graph. In *The World Wide Web Conference*, WWW '19, page 1210–1221, New York, NY, USA. Association for Computing Machinery.

[57] Mao, Y., Wang, R., Shan, S., and Chen, X. (2019). Cosonet: Compact second-order network for video face recognition. In Jawahar, C. V., Li, H., Mori, G., and Schindler, K., editors, *Computer Vision – ACCV 2018*, pages 51–67, Cham. Springer International Publishing.

[58] McAuley, J. and Leskovec, J. (2013). Hidden factors and hidden topics: Understanding rating dimensions with review text. In *Proceedings of the 7th ACM Conference on Recommender Systems*, RecSys '13, page 165–172, New York, NY, USA. Association for Computing Machinery.

[59] Miller, G. A. (1992). WordNet: A lexical database for English. In *Speech and Natural Language: Proceedings of a Workshop Held at Harriman, New York, February 23-26, 1992*.

[60] Nguyen, X. S., Brun, L., Lezoray, O., and Bougleux, S. (2019). A neural network based on SPD manifold learning for skeleton-based hand gesture recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

[61] Ni, J., Li, J., and McAuley, J. (2019). Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 188–197, Hong Kong, China. Association for Computational Linguistics.

[62] Nickel, M. and Kiela, D. (2017). Poincaré embeddings for learning hierarchical representations. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems 30*, pages 6341–6350. Curran Associates, Inc.

[63] Nielsen, F. and Sun, K. (2019). *Clustering in Hilbert's Projective Geometry: The Case Studies of the Probability Simplex and the Elliptope of Correlation Matrices*, pages 297–331. Springer International Publishing, Cham.

[64] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R., editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.

[65] Pennec, X., Fillard, P., and Ayache, N. (2006). A Riemannian framework for tensor computing. *Int. J. Comput. Vision*, 66(1):41–66.

[66] Planche, P. (1995). *Géométrie de Finsler sur les espaces symétriques*. PhD thesis, Université de Genève, Geneve, Switzerland.

[67] Ratliff, N. D., Wyk, K. V., Xie, M., Li, A., and Rana, M. A. (2020). Generalized nonlinear and Finsler geometry for robotics. *CoRR*, abs/2010.14745.

[68] Said, S., Bombrun, L., Berthoumieu, Y., and Manton, J. H. (2017). Riemannian Gaussian distributions on the space of symmetric positive definite matrices. *IEEE Transactions on Information Theory*, 63(4):2153–2170.

[69] Shen, Z. (2006). Riemann-Finsler geometry with applications to information geometry. *Chinese Annals of Mathematics, Series B*, 27:73–94.

[70] Shimizu, R., Mukuta, Y., and Harada, T. (2021). Hyperbolic neural networks++. In *International Conference on Learning Representations*.

[71] Sra, S. (2012). A new metric on the manifold of kernel matrices with application to matrix geometric means. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc.

[72] Sra, S. (2015). Positive definite matrices and the S-divergence. *Proceedings of the American Mathematical Society*. Published electronically: October 22, 2015.

[73] Sun, Z., Deng, Z.-H., Nie, J.-Y., and Tang, J. (2019). Rotate: Knowledge graph embedding by relational rotation in complex space. In *International Conference on Learning Representations*.

[74] Tay, Y., Tuan, L. A., and Hui, S. C. (2018). Hyperbolic representation learning for fast and efficient neural question answering. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, WSDM '18, pages 583–591, New York, NY, USA. ACM.

[75] Tosato, D., Farenzena, M., Spera, M., Murino, V., and Cristani, M. (2010). Multi-class classification on Riemannian manifolds for video surveillance. In Daniilidis, K., Maragos, P., and Paragios, N., editors, *Computer Vision - ECCV 2010, 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5-11, 2010, Proceedings, Part II*, volume 6312 of *Lecture Notes in Computer Science*, pages 378–391. Springer.

[76] Toutanova, K. and Chen, D. (2015). Observed versus latent features for knowledge base and text inference. In *Proceedings of the 3rd Workshop on Continuous Vector Space Models and their Compositionality*, pages 57–66, Beijing, China. Association for Computational Linguistics.

[77] Trouillon, T., Welbl, J., Riedel, S., Gaussier, E., and Bouchard, G. (2016). Complex embeddings for simple link prediction. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, page 2071–2080. JMLR.org.

[78] Tuzel, O., Porikli, F., and Meer, P. (2006). Region covariance: A fast descriptor for detection and classification. In Leonardis, A., Bischof, H., and Pinz, A., editors, *Computer Vision – ECCV 2006*, pages 589–600, Berlin, Heidelberg. Springer Berlin Heidelberg.

[79] Tuzel, O., Porikli, F., and Meer, P. (2008). Pedestrian detection via classification on Riemannian manifolds. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30(10):1713–1727.

[80] Ungar, A. (2005). Gyrovector spaces and their differential geometry. *Nonlinear Functional Analysis and Applications*, 10.

[81] Ungar, A. A. (2008). *A Gyrovector Space Approach to Hyperbolic Geometry*. Morgan & Claypool.

[82] Ungar, A. A. (2018). *Beyond Pseudo-Rotations in Pseudo-Euclidean Spaces*. Mathematical Analysis and its Applications. Academic Press.

[83] Vemulapalli, R. and Jacobs, D. (2015). Riemannian metric learning for symmetric positive definite matrices. *ArXiv*, abs/1501.02393.

[84] Vermeer, J. (2005). A geometric interpretation of Ungar's addition and of gyration in the hyperbolic plane. *Topology and Its Applications: a journal devoted to general, geometric, set-theoretic and algebraic topology*, 152(3):226–242.

[85] Wang, M., Smith, N. A., and Mitamura, T. (2007). What is the Jeopardy model? A quasi-synchronous grammar for QA. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 22–32, Prague, Czech Republic. Association for Computational Linguistics.

[86] Wang, W., Wang, R., Huang, Z., Shan, S., and Chen, X. (2018). Discriminant analysis on Riemannian manifold of Gaussian distributions for face recognition with image sets. *IEEE Transactions on Image Processing*, 27(1):151–163.

[87] Wu, Y., Jia, Y., Li, P., Zhang, J., and Yuan, J. (2015). Manifold kernel sparse representation of symmetric positive-definite matrices and its applications. *IEEE Transactions on Image Processing*, 24(11):3729–3741.

[88] Yang, T., Sha, L., and Hong, P. (2020). *NagE: Non-Abelian Group Embedding for Knowledge Graphs*, page 1735–1742. Association for Computing Machinery, New York, NY, USA.

[89] Yang, Y., Yih, W.-t., and Meek, C. (2015). WikiQA: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2013–2018, Lisbon, Portugal. Association for Computational Linguistics.

[90] Yin, M., Guo, Y., Gao, J., He, Z., and Xie, S. (2016). Kernel sparse subspace clustering on symmetric positive definite manifolds. In *CVPR*, pages 5157–5164. IEEE Computer Society.

[91] Zhang, F., Yuan, N. J., Lian, D., Xie, X., and Ma, W.-Y. (2016). Collaborative knowledge base embedding for recommender systems. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, page 353–362, New York, NY, USA. Association for Computing Machinery.

[92] Zhang, S., Tay, Y., Yao, L., and Liu, Q. (2019a). Quaternion knowledge graph embeddings. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 32, pages 2735–2745. Curran Associates, Inc.

[93] Zhang, S., Yao, L., Sun, A., and Tay, Y. (2019b). Deep learning based recommender system: A survey and new perspectives. *ACM Comput. Surv.*, 52(1).

[94] Zhang, T., Zheng, W., Cui, Z., and Li, C. (2018a). Deep manifold-to-manifold transforming network. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 4098–4102.

[95] Zhang, Y., Ai, Q., Chen, X., and Wang, P. (2018b). Learning over knowledge-base embeddings for recommendation. *CoRR*, abs/1803.06540.

[96] Zhao, K., Wiliem, A., Chen, S., and Lovell, B. C. (2019). Convex class model on symmetric positive definite manifolds. *Image and Vision Computing*, 87:57–67.

# A    Implementation Details

## A.1    Computational Complexity of Operations

In this section we discuss the computational theoretical complexity of the different operations involved in the development of this work. We employ Big O notation[4]. Since in all cases operations are not nested, but are applied sequentially, the costs can be added resulting in a polynomial expression. Thus, by applying the properties of the notation, we disregard lower-order terms of the polynomial.

**Matrix Operations:**    For $n \times n$ matrices, the associated complexity of each operation is as follows:[5]

- Addition and subtraction: $\mathcal{O}(n^2)$
- Multiplication: $\mathcal{O}(n^{2.4})$
- Inversion: $\mathcal{O}(n^{2.4})$
- Diagonalization: $\mathcal{O}(n^3)$

**SPD Operations:**    For $n \times n$ SPD matrices, the associated complexity of each operation is as follows:

- Exp/Log map: $\mathcal{O}(n^3)$, due to diagonalizations.
- Gyro-Addition: $\mathcal{O}(n^{2.4})$, due to matrix multiplications
- Matrix Scaling: $\mathcal{O}(n^3)$, due to exp and log maps.
- Isometries: $\mathcal{O}(n^{2.4})$, due to matrix multiplications.

**Distance Calculation:**    The full computation of the distance algorithm in $\mathrm{SPD}_n$ involves matrix square root, inverses, multiplications, and diagonalizations. Since they are applied sequentially, without affecting the dimensionality of the matrices, we can take the highest value as the asymptotic cost of the algorithm, which is $\mathcal{O}(n^3)$.

## A.2    Tangent Space Optimization

Optimization in Riemannian manifolds normally requires Riemannian Stochastic Gradient Descent (RSGD) [15] or other Riemannian techniques [11]. We performed initial tests converting the Euclidean gradient into its Riemannian form, but found it to be less numerically stable and also slower than tangent space optimization [26]. With tangent space optimization, we can use standard Euclidean optimization techniques, and respect the geometry of the manifold. Note that tangent space optimization is an exact procedure, which does not incur losses in representational power. This is the case in $\mathrm{SPD}_n$ specifically because of a completeness property given by the choice of $I \in \mathrm{SPD}_n$ as the basepoint: there is always a global bijection between the tangent space and the manifold.

# B    Experimental Details

All models and experiments were implemented in PyTorch [64] with distributed data parallelism, for high performance on clusters of CPUs/GPUs.

**Hardware:**    All experiments were run on Intel Cascade Lake CPUs, with microprocessors Intel Xeon Gold 6230 (20 Cores, 40 Threads, 2.1 GHz, 28MB Cache, 125W TDP). Although the code supports GPUs, we did not utilize them due to higher availability of CPU's.

---

[4]https://en.wikipedia.org/wiki/Big_O_notation

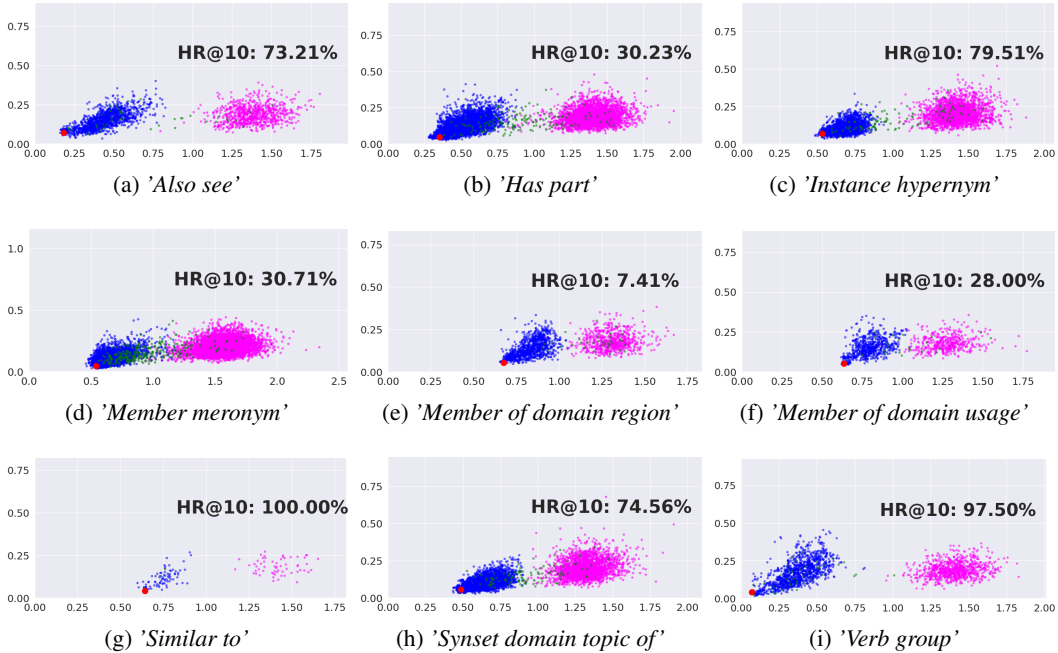[5]https://en.wikipedia.org/wiki/Computational_complexity_of_mathematical_operations

Figure 6: Train, negative and validation triples for WN18RR relationships of the $\text{SPD}_{\text{Sca}}^{F_1}$ model after convergence. The red dot corresponds to the relation addition **R**.

## B.1 Knowledge Graph Completion

**Setup:** We train for $5000$ epochs, with batch size of $4096$ and $10$ negative samples, reducing the learning rate by a factor of $2$ if the model does not improve the performance on the dev set after $50$ epochs, and early stopping based on the MRR if the model does not improve after $500$ epochs. We use the burn-in strategy [62] training with a $10$ times smaller learning rate for the first $10$ epochs. We experiment with matrices of dimension $n \times n$ where $n \in \{14, 20, 24\}$ (this is the equivalent of $\{105, 210, 300\}$ degrees of freedom respectively), learning rates from $\{1e{-}4, 5e{-}5, 1e{-}5\}$ and weight decays of $\{1e{-}2, 1e{-}3\}$.

**Datasets:** Stats about the datasets used in Knowledge graph experiments can be found in Table 4.

**Results:** In addition to the results provided in §6.1, in Table 5 we provide a comparison with other state-of-the-art models. We include ComplEx [77], Tucker [9], and Quaternion [92].

**Analisis:** In Figure 6 we add equivalent plots to the ones explained in §6.4 for other relations from WN18RR.

## B.2 Knowledge Graph Recommender Systems

**Setup:** We train for 3000 epochs, with batch size from $\{512, 1024\}$ and 10 negative samples, and early stopping based on the MRR if the model does not improve after 200 epochs. We use the burn-in strategy [62] training with a 10 times smaller learning rate for the first 10 epochs. We report average $\pm$ standard deviation of 3 runs. We experiment with matrices of dimension $10 \times 10$ (this is the

Table 4: Statistics for Knowledge Graph completion datasets.

| Dataset | Entities | Relations | Train | Dev | Test |
|---------|----------|-----------|-------|-----|------|
| WN18RR | 40943 | 11 | 86835 | 3034 | 3134 |
| FB15k-237 | 14541 | 237 | 272115 | 17535 | 20466 |

18

Table 5: Results for Knowledge graph completion.

| Space | Model | WN18RR | | | | FB15k-237 | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | MRR | HR@1 | HR@3 | HR@10 | MRR | HR@1 | HR@3 | HR@10 |
| $\mathbb{C}$ | ComplEx | 48.0 | 43.5 | 49.5 | 57.2 | 35.7 | 26.4 | 39.2 | 54.7 |
| | RotC | 47.6 | 42.8 | 49.2 | 57.1 | 33.8 | 24.1 | 37.5 | 53.3 |
| $\mathbb{Q}$ | Quaternion | 48.8 | 43.8 | 50.8 | 58.2 | **36.6** | **27.1** | **40.1** | **55.6** |
| $\mathbb{R}$ | Tucker | 47.0 | 44.3 | 48.2 | 52.6 | 35.8 | 26.6 | 39.4 | 54.4 |
| | MuRE | 47.5 | 43.6 | 48.7 | 55.4 | 33.6 | 24.5 | 37.0 | 52.1 |
| | RotE | 49.4 | 44.6 | 51.2 | 58.5 | 34.6 | 25.1 | 38.1 | 53.8 |
| | RefE | 47.3 | 43.0 | 48.5 | 56.1 | 35.1 | 25.6 | 39.0 | 54.1 |
| $\mathbb{H}$ | MuRP | 48.1 | 44.0 | 49.5 | 56.6 | 33.5 | 24.3 | 36.7 | 51.8 |
| | RotH | **49.6** | **44.9** | **51.4** | 58.6 | 34.4 | 24.6 | 38.0 | 53.5 |
| | RefH | 46.1 | 40.4 | 48.5 | 56.8 | 34.6 | 25.2 | 38.3 | 53.6 |
| SPD | $\text{SPD}_{\text{Sca}}^{R}$ | 48.1 | 43.1 | 50.1 | 57.6 | 34.5 | 25.1 | 38.0 | 53.5 |
| | $\text{SPD}_{\text{Sca}}^{F_1}$ | 48.4 | 42.6 | 51.0 | **59.0** | 32.9 | 23.6 | 36.3 | 51.5 |
| | $\text{SPD}_{\text{Rot}}^{R}$ | 46.2 | 39.7 | 49.6 | 57.8 | 32.7 | 23.4 | 36.1 | 51.4 |
| | $\text{SPD}_{\text{Rot}}^{F_1}$ | 38.6 | 25.6 | 48.7 | 56.8 | 31.4 | 22.3 | 34.7 | 49.8 |
| | $\text{SPD}_{\text{Ref}}^{R}$ | 48.3 | 44.0 | 49.7 | 56.7 | 32.5 | 23.4 | 35.6 | 51.0 |
| | $\text{SPD}_{\text{Ref}}^{F_1}$ | 48.7 | 44.3 | 50.1 | 57.4 | 30.7 | 21.7 | 33.7 | 48.8 |

equivalent of $\{55\}$ degrees of freedom), learning rates from $\{5e-4, 1e-4, 5e-5\}$ and weight decay of $1e-3$. Same grid search is applied to baselines.

**Datasets:** On the Amazon dataset we adopt the 5-core split for the branches "Software", "Luxury & Beauty" and "Industrial & Scientific", which form a diverse dataset in size and domain. We add relationships used in previous work [95, 3]. These are:

- *also_bought*: users who bought item A also bought item B.
- *also_view*: users who bought item A also viewed item B.
- *category*: the item belongs to one or more categories.
- *brand*: the item belongs to one brand.

On the MindReader dataset, we consider a user-item interaction when a user gave an explicit positive rating to the movie. The relationships added are:

- *directed_by*: the movie was directed by this person.
- *produced_by*: the movie was produced by this person/company.
- *from_decade*: the movie was released in this decade.
- *followed_by*: the movie was followed by this other movie.
- *has_genre*: the movie belongs to this genre.
- *has_subject*: the movie has this subject.
- *starring*: the movie was starred by this person.

Statistics of the datasets with the added relationships can be seen in Table 6. For dev/test we only consider users with 3 or more interactions.

Table 6: Statistics for KG Recommender datasets.

| Dataset | Users | Items | Other Entities | Train Relations | | Dev/Test |
|---|---|---|---|---|---|---|
| | | | | User-item | Others | |
| Software | 1826 | 802 | 689 | 8242 | 6078 | 1821 |
| Luxury Beauty | 3819 | 1581 | 2 | 20796 | 26044 | 3639 |
| Prime Pantry | 14180 | 4970 | 1100 | 102848 | 99118 | 14133 |
| MindReader | 961 | 2128 | 11775 | 11279 | 99486 | 953 |

19

Table 7: Statistics for Question Answering datasets.

|  | TRECQA | WIKIQA |
|---|---|---|
| Train Qs | 1227 | 2119 |
| Dev Qs | 65 | 127 |
| Test Qs | 68 | 244 |
| Train pairs | 53417 | 20361 |
| Dev Pairs | 1117 | 1131 |
| Test Pairs | 1442 | 2352 |

### B.3 Question Answering

**Setup:** We train for 300 epochs, with 2 negative samples and early stopping based on the MRR if the model does not improve after 20 epochs. We use the burn-in strategy [62] training with a 10 times smaller learning rate for the first 10 epochs. We report average $\pm$ standard deviation of 3 runs. We experiment with matrices of dimension $14 \times 14$ (equivalent of 105 degrees of freedom respectively), batch size from $\{512, 1024\}$, learning rate from $\{1e-4, 5e-5, 1e-5\}$ and weight decays of $1e-3$. Same grid search was applied to baselines.

**Datasets:** Stats about the datasets used for Question Answering experiments can be found in Table 7.

## C  Differential Geometry of $\mathrm{SPD}_n$

### C.1 Orthogonal Diagonalization

Every real symmetric matrix may be orthogonally diagonalized: For every point $P \in \mathrm{SPD}_n$ we may find a positive diagonal matrix $D$ and an orthogonal matrix $K$ such that $P = KDK^T$. This diagonalization has two practical consequences: it allows efficient computation of important $\mathrm{SPD}_n$ operations, and provides another means of generalizing Euclidean notions to $\mathrm{SPD}_n$.

With respect to computation, if $P \in \mathrm{SPD}_n$ has orthogonal diagonalization $P = KDK^T$, we may compute its square root and logarithm as $\sqrt{P} = K\sqrt{D}K^T$ and $\log(P) = K\log(D)K^T$ where $\sqrt{D} = \mathrm{diag}(\sqrt{d_1}, \ldots, \sqrt{d_n})$ and $\log(D) = \mathrm{diag}(\log d_1, \ldots, \log d_n)$ for $D = \mathrm{diag}(d_1, \ldots, d_n)$. Similarly, if a tangent vector $U \in S_n$ has orthogonal diagonalization $U = K\Lambda K^T$ (here $\Lambda = \mathrm{diag}(\lambda_1, \ldots, \lambda_n)$ not necessarily positive definite), the exponential map is computed as $\exp(U) = Ke^\Lambda K^T$, where $e^\Lambda = \mathrm{diag}(e^{\lambda_1}, \ldots e^{\lambda_n})$.

We verify this in the two lemmas below.

**Lemma 1.** *If $K \in O(n)$ and $X$ is any $n \times n$ matrix, then $\exp(KXK^T) = K\exp(X)K^T$.*

*Proof.* As $K$ is orthogonal, $K^T = K^{-1}$. Conjugation is an automorphism of the algebra of $n \times n$ matrices, and so applying this to any partial sum of the exponential $\exp(X) = \sum_{n=0}^\infty \frac{1}{n!}X^n$ yields

$$\sum_{n=0}^{N} \frac{1}{n!}(KXK^{-1})^n = K\left(\sum_{n=0}^{N} \frac{1}{n!}X^n\right)K^{-1}.$$

Taking the limit of this equality as $N \to \infty$ gives the claimed result. $\qquad\square$

**Lemma 2.** *If $D = \mathrm{diag}(d_1, \ldots, d_n)$ is a diagonal matrix, then $\exp(D) = \mathrm{diag}(e^{d_1}, \ldots, e^{d_n})$.*

*Proof.* The multiplication of diagonal matrices coincides with the elementwise product of their diagonal entries. Again applying this to any partial sum of the exponential of $D = \mathrm{diag}(d_1, \ldots, d_n)$ gives

$$\sum_{n=0}^{N} \frac{1}{n!} \mathrm{diag}(\ldots, d_i \ldots)^n = \mathrm{diag}\left(\ldots, \sum_{n=0}^{N} \frac{1}{n!}d_i^n, \ldots\right).$$

Taking the limit of this equality as $N \to \infty$ gives the claimed result. $\qquad\square$

## C.2 Metric and Isometries

The Riemannian metric on $\mathrm{SPD}_n$ is defined as follows: if $U, V \in S_n$ are tangent vectors based at $P \in \mathrm{SPD}_n$, their inner product is:

$$\langle U, V \rangle_P = \mathrm{tr}(P^{-1}UP^{-1}V).$$

Note that at the basepoint, this is just the standard matrix inner product $\langle U, V \rangle_I = \mathrm{tr}(UV^T)$ as $U, V$ are symmetric. We now verify the $GL(n, \mathbb{R})$ action given by $M$ acting as $P \mapsto MPM^T$ is an action by isometries of this metric.

**Lemma 3.** *The action $f(P) = MPM^T$ extends to tangent vectors $U$ based at $P$ without change in formula: $f_*(U) = MUM^T$*

*Proof.* Let $P \in \mathrm{SPD}_n$ and $U \in S_n$ be a tangent vector based at $P$. Then by definition, $U = P_0'$ is the derivative of some path $P_t$ of some path of matrices in $\mathrm{SPD}_n$ throguh $P_0 = P$. We compute the action of $P \to MPM^T$ on $U$ by taking the derivative of its action on the path:

$$\frac{d}{dt}\Big|_{t=0} MP_tM = MP_t'M\Big|_{t=0} = MUM^T$$

$\square$

**Proposition 1.** *For every $M \in GL(n; \mathbb{R})$ the transformation $M \mapsto MPM^T$ preserves the Riemannian metric on $\mathrm{SPD}_n$.*

*Proof.* Let $M \in GL(n; \mathbb{R})$ and choose arbitrary point $P \in \mathrm{SPD}_n$, and tangent vectors $U, V \in T_P\,\mathrm{SPD}_n$. We compute the pullback of the metric under the symmetry $f(P) = MPM^T$. Computing directly from the definition an the previous lemma,

$$\begin{aligned}
f^*\langle U, V \rangle_P &= \langle f_*U, f_*V \rangle_{f(P)} \\
&= \langle MUM^T, MVM^T \rangle_{MPM^T} \\
&= \mathrm{tr}\left( \left(MPM^T\right)^{-1} MUM^T \left(MPM^T\right)^{-1} MVM^T \right) \\
&= \mathrm{tr}\left( M^{-T}P^{-1}UP^{-1}VM^T \right) \\
&= \mathrm{tr}\left( P^{-1}UP^{-1}V \right) \\
&= \langle U, V \rangle_P,
\end{aligned}$$

where the penultimate equality uses that trace is invariant under conjugacy.

$\square$

This provides a vivid geometric interpretation of the previously discussed orthogonal diagonalization operation on $\mathrm{SPD}_n$.

**Corollary 1.** *Given any $P \in \mathrm{SPD}_n$, there exists a symmetry fixing $I$ which moves $P$ to a diagonal matrix.*

This subspace of diagonal matrices plays an essential role in working with $\mathrm{SPD}_n$. As we verify below, the intrinsic geometry of this subspace of diagonal matrices inherited from the Riemannian metric on $\mathrm{SPD}_n$ is flat.

**Proposition 2.** *Let $\mathcal{D} \subset \mathrm{SPD}_n$ be the set of diagonal matrices, and define $f: \mathbb{R}^n \to \mathcal{D}$ by $f(x_1, \ldots, x_n) = \mathrm{diag}(e^{x_1}, \ldots, e^{x_n})$. Then $f$ is an isometry from the Euclidean metric on $\mathbb{R}^n$ to the metric on $\mathcal{D}$ induced from $\mathrm{SPD}_n$.*

*Proof.* We pull back the metric on $\mathcal{D}$ by $f$, and see that on $\mathbb{R}^n$ this results in the standard Euclidean metric. Given a point $x \in \mathbb{R}^n$ with tangent vectors $y, z \in \mathbb{R}^n$, we compute this as

$$f^*\langle y, z \rangle_x = \langle f_*y, f_*z \rangle_{f(x)}$$

From the definition of $f$, we see that the pushforward of $y$ along $f$ is $\mathrm{diag}(\ldots, e^{x_i}y_i, \ldots)$ and similarly for $z$. Thus we may compute directly and see the result is the standard dot product on $\mathbb{R}^n$.

$$
\begin{aligned}
\langle \mathrm{diag}(e^{x_i}y_i), \mathrm{diag}(e^{x_i}z_i) \rangle_{\mathrm{diag}(e^{x_i})} &= \mathrm{tr}\left( \mathrm{diag}(e^{x_i}y_i)\,\mathrm{diag}(e^{-x_i})\,\mathrm{diag}(e^{x_i}y_i)\,\mathrm{diag}(e^{-x_i}) \right) \\
&= \mathrm{tr}\left( \mathrm{diag}(y_i z_i) \right) \\
&= \sum_{i=1}^{n} y_i z_i
\end{aligned}
$$

$\square$

This subspace $\mathcal{D}$ is in fact a *maximal flat* for $\mathrm{SPD}_n$, the largest dimensional totally geodesic Euclicean submanifold embedded in $\mathrm{SPD}_n$. For more information on the general theory of symmetric spaces from which the notion of maximal flats arises, see Helgason [39]. For our purposes, it is only important to note the following fact.

**Corollary 2.** *The set of diagonal matrices in $\mathrm{SPD}_n$ is an isometrically and totally geodesically embedded copy of euclidean $n$-space.*

### C.3 Exponential and Logarithmic Maps

The Riemannian exponential map gives a connection between the Euclidean geometry of the tangent space $S_n$ and the curved geometry of $\mathrm{SPD}_n$. It assigns the tangent vector $U$ to the point $Q = \exp(U)$ of $\mathrm{SPD}_n$ reached by traveling along the geodesic starting from the basepoint $I$ in direction $U$ for distance $\|U\|$.

As a consequence of non-positive curvature, $\exp$ is a diffeomorphism of $S_n$ onto $\mathrm{SPD}_n$, and so has an inverse: the Riemannian logarithm $\log\colon \mathrm{SPD}_n \to S_n$. See [10] for a review of the general theory of manifolds of non-positive curvature. Together, this pair of functions allows one to freely move between 'tangent space coordinates' or the original 'manifold coordinates' which we exploit to transfer Euclidean optimization schemes to $\mathrm{SPD}_n$ (see §5).

Secondly, the geometry of $\mathrm{SPD}_n$ is so tightly tied to the algebra of $n \times n$ matrices that the Riemannian exponential agrees exactly with the usual matrix exponential, and the Riemannian logarithm is the matrix logarithm (because of this, we do not distinguish the two notationally), as we verify in the proposition below. Both of these are readily computable via orthogonal diagonalization, as given in §C.1. This is in stark contrast to general Riemannian manifolds, where the exponential map may have no simple formula.

**Proposition 3.** *Let $\exp_{Riem}\colon S_n \to \mathrm{SPD}_n$ be the Riemannian exponential map based at $I \in \mathrm{SPD}_n$, and $\exp$ be the matrix exponential. Then $\exp_{Riem} = \exp$.*

*Proof.* Let $U \in S_n$ be a tangent vector to $\mathrm{SPD}_n$ at the basepoint $I$, and orthogonally diagonalize as $U = KDK^T$ for some $K \in O(n)$, $D = \mathrm{diag}(d_1, \ldots, d_n)$. As $D$ is tangent to the maximal flat $\mathcal{D}$ of diagonal matrices, the geodesic segment $\exp_{Riem}(tD)$ must be a geodesic in $\mathcal{D}$, which we know from Lemma 3 to be the coordinate-wise exponential of a straight line in $\mathbb{R}^n$. Precisely, this geodesic is $\mathrm{diag}(\ldots, e^{d_i t}, \ldots)$, and so the original geodesic with initial tangent $U = KDK^T$ is $\exp_{Riem}(tU) = K\,\mathrm{diag}(\ldots, e^{d_i t}, \ldots)K^T$ by Lemma 1. Specializing to $t = 1$, this gives the claim:

$$
\begin{aligned}
\exp_{Riem}(U) &= K \exp_{Riem}(D) K^T \\
&= K \,\mathrm{diag}(\ldots, e^{d_i}, \ldots) K^{-1} \\
&= K \exp(D) K^{-1} \\
&= \exp(KDK^{-1}) \\
&= \exp(U)
\end{aligned}
$$

$\square$

This easily transfers to an understanding of the Riemannian exponential at an arbitrary point $P \in \mathrm{SPD}_n$, if we identify the tangent space at $P$ with the symmetric matrices $S_n$ as well.

**Corollary 3.** *The exponential based at an arbitrary point $P \in \mathrm{SPD}_n$ is given by*

$$\exp_{Riem,P}(U) = \sqrt{P}\exp(\sqrt{P^{-1}}U\sqrt{P^{-1}})\sqrt{P}$$

*Proof.* Given $P \in \mathrm{SPD}_n$ and tangent vector $U \in T_P \mathrm{SPD}_n$ identified with the set $S_n$ of symmetric matrices, note that $X \mapsto \sqrt{P^{-1}}X\sqrt{P^{-1}}$ is a symmetry of $\mathrm{SPD}_n$ taking $P$ to $I$ and $U$ to $\sqrt{P^{-1}}U\sqrt{P}^{-1}$. Using the fact that we understand the Riemannian exponential at the basepoint, we see $\exp_{Riem}(\sqrt{P^{-1}}U\sqrt{P^{-1}}) = \exp(\sqrt{P^{-1}}U\sqrt{P^{-1}})$. It only remains to translate the result back to $P$, giving the claimed formula. $\qquad\square$

**Proposition 4.** *Let $\log_{Riem} \colon \mathrm{SPD}_n \to S_n$ be the Riemannian logarithm map based at $0 \in S_n$, and $\log$ be the matrix logarithm (note that while the matrix logaritm is multivalued in general, it is uniquely defined on $S_n$). Then $\log_{Riem} = \log$.*

*Proof.* Defined as the inverse of $\exp_{Riem}$, the Riemannian logarithm must satisfy

$$\log_{Riem} \circ \exp_{Riem} = \mathrm{id}_{S_n}$$

Let $U \in S_n$ and orthogonally diagonalize as $U = KDK^T$. Applying the Riemannian exponential, we see $\log_{Riem}(K\exp(D)K^T) = KDK^T$. Recalling from Lemma 3 the relation between isometries of $\mathrm{SPD}_n$ and their application on tangent vectors, we see that we may rewrite the left hand side as $\log_{Riem}(K\exp(D)K^T) = K\log_{Riem}(\exp(D))K^T$. Appropriately cancelling the factors of $K, K^T$ we arrive at the relationship

$$\log_{Riem}(\exp(D)) = D.$$

That is, restricted to the diagonal matrices, the Riemannian logarithm is an inverse of the matrix exponential, so Riemannian log equals matrix log. Re-absorbing the original factors of $K$ shows the same to be true for any positive definite symmetric matrix; thus $\log_{Riem} = \log$. $\qquad\square$

As for the exponential, conjugating by a symmetry moving $I$ to an arbitrary point $P$, we may describe the Riemannian logarithm at any point of $\mathrm{SPD}_n$.

**Corollary 4.** *The logarithm based at an arbitrary point $P \in \mathrm{SPD}_n$ is given by*

$$\log_{Riem,P}(Q) = \sqrt{P}\log(\sqrt{P^{-1}}Q\sqrt{P^{-1}})\sqrt{P}$$

### C.4 Vector-valued Distance

Here we collect useful observations about the vector-valued distance on $\mathrm{SPD}_n$, culminating in a proof of the fact that it is a complete invariant of pairs of points, as claimed in §3.

**Proposition 5.** *The vector-valued distance is well-defined: given any pair $P, Q \in \mathrm{SPD}_n$ of points and any two isometries taking $P, Q$ to the basepoint, a diagonal matrix respectively, the diagonal matrices differ at most by a permutation of their entries.*

*Proof.* We see heuristically that there is no remaining continuous degree of freedom by dimension count: the isometry group $GL(n; \mathbb{R})$ has dimension $n^2$, and we require $\dim(\mathrm{SPD}_n) = n(n+1)/2$ degrees of freedom to translate $P$ to the origin, and a further $\dim O(n) = n(n-1)/2$ degrees of freedom to diagonalize the image of $Q$ while fixing $I$. As $\dim GL(n; \mathbb{R}) = \dim \mathrm{SPD}_n + \dim O(n)$, there are no remaining continuous degrees of freedom. To see that the remaining ambiguity is precisely permutation of coordinates, note that conjugating a diagonal matrix by an orthogonal matrix results in another diagonal matrix only if the conjugating matrix is a permutation matrix. $\qquad\square$

**Proposition 6.** *If two points $P, Q \in \mathrm{SPD}_n$ have the same vector-valued distance from the basepoint $I$, then there is an isometry fixing $I$ taking $P$ to $Q$.*

*Proof.* For two matrices to have the same vector-valued distance from $I$ is equivalent to those two matrices having the same set of eigenvalues. Let $\lambda_1, \ldots, \lambda_n$ be a list of these eigenvalues with multiplicity, and construct two orthonormal bases $(v_i), (w_i)$ of $\mathbb{R}^n$ as follows. For each $i$, let $v_i$ be an eigenvector of $P$ with eigenvalue $\lambda_i$, and $w_i$ an eigenvector of $Q$ with eigenvalue $\lambda_i$ (in the case the eigenvalues are distinct, such bases are unique up to flipping the sign of each vector, but nontrivial choices must be made in the case of coincident eigenvalues). Given this pair of orthonormal bases, let $K \in O(n)$ be the orthogonal matrix which takes $(v_i)$ to $(w_i)$. It is then an easy observation of linear algebra to note that $Q = KPK^{-1}$, but recalling $K^T = K^{-1}$ we see this is interpreted in the geometry of $\mathrm{SPD}_n$ to say that there is an isometry $X \mapsto KXK^T$ fixing $I$ and taking $P$ to $Q$. $\qquad \square$

Combining Propositions 5 and 6, after translating appropriately to the basepoint yields the following cornerstone of the theory, showing the vector-valued distance to be the *best possible* invariant.

**Corollary 5.** *The vector-valued distance is a complete invariant of pairs of points. Two pairs of points $(P, Q)$ and $(P', Q')$ cam be mapped to one another by an isometry if and only if they have the same vector-valued distance.*

It's important to note that while the vector-valued distance is not *literally* a metric distance (it is vector valued, instead of positive-real-number valued, for one) it enjoys some properties analogous to traditional metric distances. For a brief review of some of these (the vector-valued triangle inequality, etc) see Kapovich, Leeb & Porti. [46], and Kapovich, Leeb & Millison [45].

One property distinguishing the vector-valued distance from traditional metrics is its assymmetry. We will wish to recall this relationship later on, and so prove it here for completeness.

**Lemma 4.** *For $P, Q \in \mathrm{SPD}_n$, the vector-valued distance satisfies*

$$d_{vv}(P, Q) = -d_{vv}(Q, P)$$

*with equality understood up to permutation of coordinates.*

*Proof.* The computation of $d_{vv}(P, Q)$ differs from that of $d_{vv}(Q, P)$ in the first step, where we reduce it to computing a function of the eigenvalues of $P^{-1}Q$ or $Q^{-1}P$ respectively. Noting these are inverses of one another, their eigenvalues are reciprocals we may perform the following calculation, where $\{\lambda_i(X)\}$ denotes the eigenvalues of $X$.

$$
\begin{aligned}
d_{vv}(Q, P) &= \log(\ldots, \lambda_i(Q^{-1}P), \ldots) \\
&= \log(\ldots, \lambda_i((P^{-1}Q)^{-1}), \ldots) \\
&= \log(\ldots, \lambda_i((P^{-1}Q)^{-1}), \ldots) \\
&= \log\left(\ldots, \frac{1}{\lambda_i((P^{-1}Q))}, \ldots\right) \\
&= -\log(\ldots, \lambda_i((P^{-1}Q)), \ldots) \\
&= -d_{vv}(P, Q)
\end{aligned}
$$

$\qquad \square$

### C.5 Riemannian Distance

This Riemannian metric allows the computation of the length of curves $\gamma \colon [0, 1] \to \mathrm{SPD}_n$ as

$$\mathrm{length}(\gamma) = \int_0^1 \sqrt{\langle \gamma'(t), \gamma'(t) \rangle_{\gamma(t)}} \, dt.$$

This in turn induces a distance function $d : \mathrm{SPD}_n \times \mathrm{SPD}_n \to \mathbb{R}$, by taking the infimum of the lengths of all paths joining two points:

$$d(P, Q) = \inf_{\substack{\gamma \colon [0,1] \to \mathrm{SPD}_n \\ \gamma(0) = P,\, \gamma(1) = Q}} \left\{ \mathrm{length}(\gamma) \right\}$$

While for general Riemannian manifolds such a distance function may be impossible to explicitly compute, the symmetries of $\mathrm{SPD}_n$ provide a readily computable formula.

**Proposition 7.** *The Riemannian distance from the basepoint $I$ to a point $P \in \mathrm{SPD}_n$ is given by $d(I, P) = \sqrt{\sum_{i=0}^n \log(\lambda_i(P))}$ where $\{\lambda_i(P)\}$ are the eigenvalues of of $P$.*

*Proof.* Let $P \in \mathrm{SPD}_n$ be arbitrary, and orthogonally diagonalize as $P = KDK^T$. As $K \in O(n)$, the isometry $X \mapsto KDK^T$ fixes $I$, so the distance $d(I, P)$ equals the distance $d(I, D)$. Note as this action of $K$ is by conjugacy, the diagonal entries $d_i$ of $D$ are precisely the eigenvalues of $P$. As $D$ lies in the totally geodesic Euclidean subspace $\mathcal{D}$, this distance is realized by the unique Euclidean geodesic connecting $I$ to $D$. Using Lemma 2, we may translate to familiar coordinates on $\mathbb{R}^n$ and notice this is the distance from the origin $0$ to the point $x = (\log(d_1), \ldots, \log(d_n))$. That is, $d(I, D) = \sqrt{\sum_i \log(d_i)^2}$ as claimed. $\qquad\square$

This immediately generalizes to the distance between a pair of arbitrary points, via conjugating by a symmetry moving one to the origin. However, with a little more work one may get a simpler expression for the general distance.

**Proposition 8.** *The Riemannian distance between two arbitrary points $P, Q \in \mathrm{SPD}_n$ is given by $d(P, Q) = \sqrt{\sum_i \log(\lambda_i(P^{-1}Q))}$ where $\{\lambda_i(P^{-1}Q)\}$ are the eigenvalues of of $P^{-1}Q$.*

*Proof.* If $P, Q$ are arbitrary points in $\mathrm{SPD}_n$, we may use an isometry to translate $P$ to the basepoint, while simultaneously moving $Q$ to $R = \sqrt{P^{-1}}Q\sqrt{P^{-1}}$. As isometries preserve distances, we have $d(P, Q) = d(I, R)$, and by Proposition 7, this distance is completely determined by the eigenvalues of $R$. As these are invariant under conjugacy, we replace $R$ with its conjugate by $\sqrt{P^{-1}}$ to get the matrix

$$
\begin{aligned}
R' &= \sqrt{P}R\sqrt{P}^{-1} \\
&= \sqrt{P^{-1}}\sqrt{P^{-1}}Q\sqrt{P^{-1}}\sqrt{P} \\
&= P^{-1}Q
\end{aligned}
$$

$\square$

### C.6 Finsler Distances

The Riemannian distance function on a manifold is completely determined by its Riemannian metric, a choice of inner product on the tangent bundle. Generalizing this, Finsler metrics are the class of distance functions which may be constructed from a smoothly varying choice of norm $\|\cdot\|_F$ on the tangent bundle (which need not be induced by an inner product). The basic theory proceeds in direct analogy to the Riemannian case: the length of a curve $\gamma \colon [0,1] \to \mathrm{SPD}_n$ with respect to a Finsler metric is still defined via integration of this norm along the path, and the distance between points by the infimum of this over all rectifiable curves joining them

$$
\mathrm{length}_F(\gamma) = \int_0^1 \|\gamma'\|_F \, dt, \qquad d_F(P, Q) = \inf_{\substack{\gamma \colon [0,1] \to \mathrm{SPD}_n \\ \gamma(0) = P, \, \gamma(1) = Q}} \left\{ \mathrm{length}_F(\gamma) \right\}
$$

The geometry of $\mathrm{SPD}_n$ allows the computaiton of all Finsler metrics directly from the vector-valued distance. As Riemannian metrics are in particular a special case of Finsler metrics, we begin by recasting our previous observations in this light. In §C.5 we derived a formula for the Riemannian distance function directly from the infintesimal Riemannian metric. But in light of Corollary 5, since the Riemannian distance is a function which depends only on its input points up to isometry, it must also be recoverable from the vector-valued distance. Indeed, looking at Proposition 7 we see there is a simple rephrasing to this effect:

**Corollary 6.** *The Riemannian distance from the basepoint $I$ to an arbitrary point $P \in \mathrm{SPD}_n$ is the Euclidean norm of the vector-valued distance from $I$ to $P$.*

One of the great advantages of higher rank symmetric spaces is the generalizations to which this rephrasing lends itself. Namely, the Euclidean metric is not special in this construction, and any sufficiently symmetric norm on $\mathbb{R}^n$ can induce a distance function on $\mathrm{SPD}_n$ in this way.

**Proposition 9.** *Let* $\|-\|$ *be any norm on* $\mathbb{R}^n$ *which is invariant under the permutation of coordinates. Then* $\|-\|$ *induces a distance function* $d$ *on* $\mathrm{SPD}_n$ *by*

$$d(P,Q) = \|d_{vv}(P,Q)\|$$

*Proof.* We first note this function is well-defined, as by Proposition 5 the vector-valued distance of $(P,Q)$ is well-defined up to permutation of coordinates, and our norm is invariant under this symmetry by hypothesis. To see that $d$ is in fact a distance function on $\mathrm{SPD}_n$, we now need to show it satisfies the axioms of a metric:

1. $d(P,Q) \geq 0,\ d(P,Q) = 0 \implies P = Q$

2. $d(P,Q) = d(Q,P)$

3. $d(P,R) \leq d(P,Q) + d(Q,R)$

To check the identity of indescernibles (1), note that $d$ is necessarily nonnegative as $\|-\|$ is, and if $d_(P,Q) = 0$ then the norm of $d_{vv}(P,Q)$ is zero, so the vector-valued distance itself is zero. But as this is a complete invariant and $d_{vv}(P,P) = 0$, this means $P = Q$.

Note that symmetry (2) is not automatic as the vector-valued distance itself is asymmetric. However recalling Lemma 4, we see that changing the order causes only a global negative sign, and the central symmetry of $\|-\|$, as a virtue of being a norm, gives equality.

The triangle inequality (3) is more subtle, and requires an understanding of the triangle inequality for the vector-valued distance. See the dissertation of Planche [66], Chapter 6 and the work of Kapovich, Leeb and Millson [45] for details. □

For our experiments, the most important such distances are induced by the $\ell_1$ and $\ell_\infty$ norms on $\mathbb{R}^n$. For completeness, the resulting distance functions are described below.

**Proposition 10.** *The distance function induced from the* $\ell_1$ *metric applied to the vector-valued distance can be computed as* $d_{F_1}(P,Q) = \sum_{i=1}^{n} |\log \lambda_i(P^{-1}Q)|$, *where* $\lambda_i(P^{-1}Q)$ *runs over the eigenvalues of* $P^{-1}Q$.

*Proof.* The vector-valued distance $d_{vv}(P,Q)$ is the vector of logarithms of the eigenvalues of $R = P^{-1}Q$, and its $\ell^1$ norm is the sum of their absolute values:

$$\|(\log(\lambda_1(R), \ldots, \lambda_n(R))\|_{\ell^1} = \sum_{i=1}^{n} |\log \lambda_i(R)|$$

where $\lambda_i(R)$ is the $i^{th}$ eigenvalaue of $R$, in decreasing order. □

A similar calculation yields the formula for the $F^\infty$ distance function.

**Proposition 11.** *The distance function induced from the* $\ell_\infty$ *metric applied to the vector-valued distance can be computed as* $d_{F_\infty}(P,Q) = \lambda_1(P^{-1}Q)$ *where* $\lambda_1(-)$ *returns the largest eigenvalue of the input matrix.*

### C.7 Relations with Other Metrics

Other distances previously used in the literature can be reconstructed from the vector-valued distance, by applying a suitable function:

The *Affine Invariant metric* of [65] is nothing but the usual Riemannian metric discussed in §C.5.

The *symmetric Stein divergence* [71], is given by

$$S(P,Q) := \log \det \frac{P+Q}{2} - \frac{1}{2} \log \det(PQ)$$

This can be computed from the vector-valued distance

$$d_{vv}(P,Q) = \log(\lambda_1(P^{-1}Q), \ldots, \lambda_n(P^{-1}Q))$$

by applying the function

$$\|v\|_S = \sum_{i=1}^{n} \log \frac{e^{-v_i/2} + e^{v_i/2}}{2}.$$

Indeed

$$
\begin{aligned}
S(P,Q) &= \log \det \frac{P+Q}{2} - \tfrac{1}{2} \log \det(PQ) \\
&= \log \det P \left( \frac{\mathrm{Id} + P^{-1}Q}{2} \right) - \log \det(P\sqrt{P^{-1}Q}) \\
&= \log \det \frac{\mathrm{Id} + P^{-1}Q}{2} - \log \det(\sqrt{P^{-1}Q}) \\
&= \sum_{i=1}^{n} \log \lambda_i \left( \frac{\mathrm{Id} + P^{-1}Q}{2\sqrt{P^{-1}Q}} \right) \\
&= \sum_{i=1}^{n} \log \left( \frac{\lambda_i(P^{-1}Q)^{-1/2} + \lambda_i(P^{-1}Q)^{1/2}}{2} \right)
\end{aligned}
$$

In particular we obtain, thanks to the vector-valued distance, a more direct proof of [72].

Instead the *Log-Euclidean* metric $d_{LE}$ [5, 6] is flat, and as such doesn't reflect the curved geometry of SPD. More precisely $d_{LE}$ is the pushforward, through the exponential map $\exp_{Riem} : S_n \to$ SPD of the Euclidean metric on $S_n$. As a result, for this choice $(\mathrm{SPD}_n, d_{LE})$ is *isometric* to the flat manifold $S_n$. Since the group $GL(n, \mathbb{R})$ does not act by isometries on $(\mathrm{SPD}_n; d_{LE})$, and the distance is therefore not related to the vector-valued distance nor can be computed from it.

Similarly the *Bures-Wasserstein* metric $d_{BW}$ inspired from quantum information theory [13] leads to a non-negatively curved manifold, and thus, again, has a different isometry group. More precisely

$$d_{BW}(P,Q) = \sqrt{\mathrm{tr}(P) + \mathrm{tr}(Q) - 2\sqrt{\mathrm{tr}(PQ)}}.$$

It is computed in [13, Page 15] that the group of isometries of $(\mathrm{SPD}_n, d_{BW})$ is reduced to $O(n)$. As a result, once again, $d_{BW}$ cannot be reconstructed from $d_{vv}$.

# D   Gyrocalculus

A primary difficulty of building analogs of Euclidean quantities in curved spaces is the lack of a vector space structure, making the translation of operations like vector addition or scalar multiplication difficult to immediately interpret. The need for these is already well-noted stumbling block in hyperbolic geometry, as any algorithm using the Euclidean addition of points cannot be implemented directly (for example considering the Poincare disk model, the sum of two points in the disk need not lie in the disk: and even when it does, the result is rarely geometrically meaningful). To combat this, means of interfacing with hyperbolic geometry using "vector-space-like" operations was developed by Ungar [81], which provides an analog of addition $\oplus \colon \mathbb{H}^n \times \mathbb{H}^n \to \mathbb{H}^n$ and of scalar multiplication $\otimes \colon \mathbb{R} \times \mathbb{H}^n \to \mathbb{H}^n$ called 'gyro-addition' and 'gyro-scalar multiplication' respectively. We give a brief introduction to this general theory below, see Ungar's treatment from the lens of differential geometry [80] for further information.

## D.1   Gyrogroups

Gyrogroups are a generalization of groups which encode algebraically some of the geometric properties of symmetric spaces. More precisely, a gyrogroup structure on a set $G$ is given by a binary operation $\oplus$, which is assumed to have an identity element $0 \in G$ and left inverses $\ominus g$ for each $g \in G$. Keeping with the conventions familiar from arithmetic, we write $a \ominus b$ to mean $a \oplus (\ominus b)$. The crucial difference from group theory is that $\oplus$ is *not* required to be associative. Instead, the additional structure of a *gyration operator* $\mathrm{gyr} \colon G \times G \to \mathrm{Aut}(G)$ captures the nonassociativity of $\oplus$ by

$$a \oplus (b \oplus c) = (a \oplus b) \oplus \mathrm{gyr}(a,b)c$$

.

For $(G, \oplus, \mathrm{gyr})$ to form a gyrogroup, an additional axiom is imposed on this gyration, namely that it satisfy the *left loop identity*, $\mathrm{gyr}(a, b) = \mathrm{gyr}(a \oplus b, b)$.

Gyrogroups generalize groups in the sense that every group $G$ is a gyrogroup with its usual binary operation as $\oplus$, and trivial gyration. As with standard groups, it is helpful to have at one's disposal a collection of elementary deductions from these axioms, which may significantly simplify further calculations.

**Proposition 12.** *The identity of a gyrogroup is unique, every left inverse is also a right inverse, and every element has a unique (left, and hence also right) inverse.*

See [80] §3 for a proof of this proposition, which uses only the axioms of a gyrogroup. It can be shown that when a gyrogroup structure exists on a set $G$, it is determined by the operation $\oplus$ alone, in the sense that for any $a, b, c$ we have

$$\mathrm{gyr}(a, b)c = (\ominus (a \oplus b)) \oplus (a \oplus (b \oplus c)) \tag{11}$$

We record also useful properties of the gyration operator following from this, which simplify calculation.

**Proposition 13.** *The following gyrations are trivial: the gyration of any element with zero* $\mathrm{gyr}(0, a) = \mathrm{gyr}(0, a) = \mathrm{gyr}(\ominus a, a) = \mathrm{id}_G$, *or with its inverse* $\mathrm{gyr}(\ominus a, a) = \mathrm{gyr}(a, \ominus a) = \mathrm{id}_G$. *A useful consequence of these is the* nested gyration identity*:*

$$\mathrm{gyr}(a, \ominus \mathrm{gyr}(a, b)b)\, \mathrm{gyr}(a, b) = \mathrm{id}_G$$

These are also proven in [80] §3 , and follow directly from the axioms of a gyrogroup.

Because of the additional complexity of $\oplus$ compared to the binary operation of a standard group, it is often useful in applications to introduce a second binary operation, the *gyrogroup co-operation* $\boxplus$ and its inverse $\boxminus$, defined by

$$a \boxplus b = a \oplus \mathrm{gyr}(a, \ominus b)b \qquad a \boxminus b = a \boxplus \ominus a$$

This operation provides a useful shorthand for solving equations in gyrogroups, which we discuss in D.2.1.

Finally, we give a means of computing the VVD in terms of these operations as claimed in §4.

**Proposition 14.** *The vector-valued distance from $P$ to $Q$ is the vector of logarithms of the eigenvalues of $(\ominus P) \oplus Q$.*

*Proof.* This is the matrix $(\ominus P) \oplus Q = P^{-1} \oplus Q = \sqrt{P^1}Q\sqrt{P^{-1}}$, which is conjugate to $P^{-1}Q$ (as in 8), and so has the same eigenvalues. But the logarithm of these eigenvalues is precisely the vector value distance as defined in §3.1. $\square$

### D.2 Gyro-vector Spaces

Though the operation $\oplus$ is not commutative in the usual sense, a gyrogroup $G$ is called *gyro-commutative* if it commutes *up to gyrations*: ie for every $a, b \in G$, $a \oplus b = \mathrm{gyr}(a, b)(b \oplus a)$. It is within this restricted class of gyro-commutative gyrogroups that a satisfactory analog of familiar vector space operations can be constructed [82].

A gyrovector space is a gyro-commutative gyrorgrup $(G, \oplus)$ together with a scalar multiplication $\otimes \colon \mathbb{R} \times G \to G$ such that 1 acts as the identity, and its interaction with standard multiplication, gyro-addition and gyration are constrained by

$$
\begin{aligned}
r_1 \otimes (r_2 \otimes a) &= r_1 r_2 \otimes a \\
(r_1 + r_2) \oplus a &= (r_1 \otimes a) \oplus (r_2 \otimes a) \\
r \otimes \mathrm{gyr}(a, b)c &= \mathrm{gyr}(a, b)(r \otimes c) \\
\mathrm{gyr}(r_1 \otimes a, r_2 \otimes a) &= I
\end{aligned}
\tag{12}
$$

Typically a gyrovector space is also assumed to be constructed within an ambient real inner product space, and there are additional compatibility relations between the operations of $(G, \oplus, \otimes)$ and the ambient vector space addition $(+)$ and norm $\|v\| = \sqrt{v \cdot v}$.

Gyro-vector spaces generalize vector spaces much as gyro-groups generalized groups: every vector space is a gyro-vector space with trivial gyration. As such, the formalism of gyro-vector spaces provides a convenient generalization where one may attempt to replace $+, -, \times$ in formulas familiar from Euclidean spaces with $\oplus, \ominus, \otimes$; being careful to recall that gyro-addition is neither commutative nor associative, and gyro-multiplication rarely distributes over $\oplus$.

### D.2.1 Solving Equations in Gyrogroups

As an example of the difficulties posed by this, if one requires the solution to the Euclidean equation $a + x = b$, it is equally correct to write $x = b - a$ or $x = -a + b$. But the translations $x = b \ominus a$ and $x = \ominus a \oplus b$ into a gyrogroup $G$ need not be equal, and generically only the latter solves the gyrovector equation $a \oplus x = b$.

To make this more systematic, note that working inwards respecting the order of operations, we are able to solve any equation in a gyrogroup if we compute a *left cancellation law*, *right cancellation law* and *invert scalar multiplication*.

**Proposition 15** (Left-Cancellation)**.** *Let $a, b$ be elements of a gyrogroup $G$. Then the relation $a \oplus x = b$ is satisfied by the unique value $x = (\ominus a) \oplus b$.*

*Proof.* Substituting the claimed expression for $x$, we verify by direct computation from the axioms of a gyrogroup, and the basic properties of Propositions 12.

$$
\begin{aligned}
a \oplus x &= a \oplus ((\ominus a) \oplus b) \\
&= (a \oplus \ominus a) \oplus \mathrm{gyr}(a, \ominus a) b \\
&= 0 \oplus \mathrm{gyr}(a, \ominus a) b \\
&= id_G(b) \\
&= b
\end{aligned}
$$

$\square$

**Proposition 16** (Right-Cancellation)**.** *Let $a, b$ be elements of a gyrogroup $G$. Then the relation $x \oplus a = b$ is satisfied by the unique value $x = b \boxminus a = a \ominus \mathrm{gyr}(a, \ominus b) b$, where $\boxminus$ is the additive inverse of the gyrogroup co-operation from Section D.1.*

*Proof.* To begin, we put the proposed solution $b \boxminus a$ in a more usable form:

$$
\begin{aligned}
b \boxminus a &= b \boxplus \ominus a \\
&= b \oplus \mathrm{gyr}(b, \ominus \ominus a) \ominus a \\
&= b \ominus \mathrm{gyr}(b, a) a
\end{aligned}
$$

We now verify the claim by subsituting the given value of $x$, and using the properties described in Propositions 12 and 13, (in particular, in the third step we expand $a$ using nested gyration)

$$
\begin{aligned}
x \oplus a &= (b \boxminus a) \oplus a \\
&= (b \ominus \mathrm{gyr}(b, a) a) \oplus id_G(a) \\
&= (b \ominus \mathrm{gyr}(b, a) a) \oplus (\mathrm{gyr}(b, \ominus \mathrm{gyr}(b, a) a) \, \mathrm{gyr}(b, a) a) \\
&= b \oplus (\ominus \mathrm{gyr}(b, a) a \oplus \mathrm{gyr}(b, a) a) \\
&= b \oplus 0 \\
&= b
\end{aligned}
$$

$\square$

**Proposition 17** (Inverting Scalar Multiplication)**.** *Let $r \in \mathbb{R}$ be any scalar, and $a$ an element of a gyrogroup $G$. Then the relation $r \otimes x = a$ is satisfied by the unique element $x = \left(\frac{1}{r}\right) \otimes a$.*

*Proof.* Substituting $x$ immediately yeidls the result given the axioms of gyro-scalar multiplication:

$$
\begin{aligned}
r \otimes &= r \otimes \left(\frac{1}{r} \otimes a\right) \\
&= \left(r \times \frac{1}{r}\right) \otimes a \\
&= 1 \otimes a \\
&= a
\end{aligned}
$$

$\square$

These three cancellation laws allow one work correctly with the gyro-translations of Euclidean vector space statements. Take for example the vector space expression $a + rx + b = c$ for vectors $a, b, c, x$ and scalar $r$. One possible gyro-vector space translation of this is $(a \oplus (r \otimes x)) \oplus b = c$ — and given this translation, we may work fully within the gyrovector space to solve for $x$ as follows:

$$
\begin{aligned}
(a \oplus (r \otimes x)) \oplus b &= c \\
a \oplus (r \otimes x) &= c \boxminus b \\
r \otimes x &= (\ominus a) \oplus (c \boxminus b) \\
x &= \frac{1}{r} \otimes ((\ominus a) \oplus (c \boxminus b))
\end{aligned}
$$