



Isometric Gaussian Process Latent Variable Model for Dissimilarity Data

Jorgensen, Martin; Hauberg, Soren

Published in:
Proceedings of the 38th International Conference on Machine Learning

Publication date:
2021

Document Version
Publisher's PDF, also known as Version of record

[Link back to DTU Orbit](#)

Citation (APA):
Jorgensen, M., & Hauberg, S. (2021). Isometric Gaussian Process Latent Variable Model for Dissimilarity Data. In *Proceedings of the 38th International Conference on Machine Learning* (Vol. 139). International Machine Learning Society (IMLS).

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Isometric Gaussian Process Latent Variable Model for Dissimilarity Data

Martin Jørgensen¹ Søren Hauberg²

Abstract

We present a probabilistic model where the latent variable respects both the distances and the topology of the modeled data. The model leverages the Riemannian geometry of the generated manifold to endow the latent space with a well-defined stochastic distance measure, which is modeled locally as Nakagami distributions. These stochastic distances are sought to be as similar as possible to observed distances along a neighborhood graph through a censoring process. The model is inferred by variational inference based on observations of pairwise distances. We demonstrate how the new model can encode invariances in the learned manifolds.

1. Introduction

Dimensionality reduction aims to compress data to a lower dimensional representation while preserving the underlying signal and suppressing noise. Contemporary nonlinear methods mostly call upon the *manifold assumption* (Bengio et al., 2013) stating that the observed data is distributed near a low-dimensional manifold embedded in the observation space. Beyond this unifying assumption, methods often differ by focusing on one of three key properties (Table 1).

Topology preservation. A *topological space* is a set of points whose *connectivity* is invariant to continuous deformations. For finite data, connectivity is commonly interpreted as a clustering structure, such that topology preserving methods do not form new clusters or break apart existing ones. For visualization purposes, the *uniform manifold approximation projection* (UMAP) (McInnes et al., 2018) appears to be the current state-of-the-art within this domain.

¹Department of Engineering Science, University of Oxford
²Department of Mathematics and Computer Science, Technical University of Denmark. Correspondence to: Martin Jørgensen <martinj@robots.ox.ac.uk>, Søren Hauberg <sohau@dtu.dk>.

	Probabilistic	Topology	Distance
PCA	(✓)	✗	(✓)
MDS	✗	✗	✓
IsoMap	✗	(✗)	✓
t-SNE	✗	(✓)	✓
UMAP	✗	✓	✓
GPLVM	✓	✗	✗
Iso-GPLVM (our)	✓	✓	✓

Table 1. A list of common dimensionality reduction methods and coarse overview of their features.

Distance preservation. Methods designed to find low-dimensional representation with pairwise distances that are similar to those of the observed data may generally be viewed as a variant of *multi-dimensional scaling* (MDS) (Ripley, 2007). Usually, this is achieved by a direct minimization of the *stress* defined as

$$\text{stress} = \sum_{i < j \leq N} (d_{ij} - \|z_i - z_j\|)^2, \quad (1)$$

where d_{ij} are the *dissimilarity* (or *distance*) of two data points x_i and x_j , and $\mathcal{Z} = \{z_i\}_{i=1}^N$ denote the low-dimensional representation in \mathbb{R}^q .

More advanced methods have been built on top of this idea. In particular, *IsoMap* (Tenenbaum et al., 2000) computes d_{ij} along a neighborhood graph using Dijkstra’s algorithm. This bears some resemblance to *t-SNE* (Maaten & Hinton, 2008) that uses the Kullback-Leibler divergence to match distribution in low-dimensional Euclidean spaces with the data in high dimensions.

Probabilistic models. A common trait for the mentioned methods is that they learn features in a mapping from high-dimensions to low, but not the reverse. This makes the methods mostly useful for visualization. *Generative models* (Kingma & Welling, 2014; Rezende et al., 2014; Lawrence, 2005; Goodfellow et al., 2014; Rezende & Mohamed, 2015) allow us to make new samples in high-dimensional space. Of particular relevance to us, is the *Gaussian process latent variable model* (GP-LVM) (Lawrence, 2005; Titsias & Lawrence, 2010) which learns a stochastic mapping $f : \mathbb{R}^q \rightarrow \mathbb{R}^D$ jointly with the latent representations z .

This is achieved by marginalizing the mapping under a Gaussian process prior (Rasmussen & Williams, 2006). The generative approach allows the methods to extend beyond visualization to e.g. missing data imputation, data augmentation and semi-supervised tasks (Mattei & Frellsen, 2019; Urtasun & Darrell, 2007).

In this paper, we learn a Riemannian manifold using Gaussian processes on which distances on the manifold match the *local* distances as is implied by the Riemannian assumption. Assuming the observed data lies on a Riemannian q -submanifold of \mathbb{R}^D with infinite injectivity radius, then our approach can learn a q -dimensional representation that is isometric to the original manifold. Similar statements only hold true for traditional manifold learning methods that embed into \mathbb{R}^q if the original manifold is flat. We learn global and local structure through a common technique from survival analysis, combined with a likelihood model based on the theory of Gaussian process arc-lengths. Lastly, we show how the GP approach allow us to marginalize the latent representation and produce a fully Bayesian non-parametric model. We envision how learning probabilistic models by pairwise dissimilarities easily allow for encoding invariances.

The data handled in this paper are *pairwise distances* between instances. This naturally gives a geometrical flavour to the approach since distances fall within the geometrical ontology. Note that this does not exclude tabular data — we only require a distance can be computed between points. Further, many modern datasets come in form of pairwise distances: proteins based on their distance on a phylogenetic tree, simple GPS data for place recognition, perception data from psychology, etc.

2. Background material

2.1. Gaussian Processes

A Gaussian process (GP) (Rasmussen & Williams, 2006) is a distribution over functions, $f : \mathbb{R}^q \rightarrow \mathbb{R}$, which satisfy that for any finite set of points $\{z_i\}_{i=1}^N$, in the domain \mathbb{R}^q , the output $\mathbf{f} = (f(z_1), \dots, f(z_N))$ have a joint Gaussian distribution. This Gaussian is fully determined by a mean function $\mu : \mathbb{R}^q \rightarrow \mathbb{R}$ and a covariance function $k : \mathbb{R}^q \times \mathbb{R}^q \rightarrow \mathbb{R}$, such that

$$p(\mathbf{f}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}), \quad (2)$$

where $\boldsymbol{\mu} = (\mu(z_1), \dots, \mu(z_N))$ and \mathbf{K} is the $N \times N$ -matrix with (i, j) -th entry $k(z_i, z_j)$.

GPs are well-suited for Bayesian non-parametric regression, since if we condition on data $\mathcal{D} = \{z, x\}$, where x denote the labels, then the posterior of $f(z^*)$, at a test location z^* , is given as

$$p(f(z^*)|\mathcal{D}) = \mathcal{N}(\boldsymbol{\mu}^*, \mathbf{K}^*), \quad (3)$$

where

$$\boldsymbol{\mu}^* = \mu(z^*) + k(z^*, z)^\top k(z, z)^{-1}x, \quad (4)$$

$$\mathbf{K}^* = k(z^*, z^*) - k(z^*, z)^\top k(z, z)^{-1}k(z, z^*) \quad (5)$$

We see that this posterior computation involves inversion of the $N \times N$ -matrix \mathbf{K} , which has complexity $\mathcal{O}(N^3)$. To overcome this computational burden in inference we consider variational sparse GP regression, which introduces M auxiliary points \mathbf{u} , that approximate the posterior of f with a variational distribution q . For a review of variational GP methods, we refer to Titsias (2009).

2.2. Riemannian Geometry

A *manifold* is a topological space, for which each point on it has a neighborhood that is homeomorphic to Euclidean space; that is, manifolds are locally linear spaces. Such manifolds can be embedded into spaces of higher dimension than the dimensionality of the associated Euclidean space; the manifold *itself* has the same dimension as the local Euclidean space. A q -dimensional manifold \mathcal{M} can, for our purposes thus, be seen as a surface embedded in \mathbb{R}^D . In order to make quantitative statements along the manifold we require it to be *Riemannian*.

Definition 1. A Riemannian manifold \mathcal{M} is a smooth q -manifold equipped with an inner product

$$\langle \cdot, \cdot \rangle_x : \mathcal{T}_x \mathcal{M} \times \mathcal{T}_x \mathcal{M} \rightarrow \mathbb{R}, \quad x \in \mathcal{M}, \quad (6)$$

that is smooth in x . Here $\mathcal{T}_x \mathcal{M}$ denotes the tangent space of \mathcal{M} evaluated at x .

The length of a curve is easily defined from the Riemannian inner product. If $c : [0, 1] \rightarrow \mathcal{M}$ is a smooth curve, its length is given by $s = \int_0^1 \|\dot{c}(t)\| dt$. On an embedded manifold $f(\mathcal{M})$ this becomes

$$s = \int_0^1 \|\dot{f}(c(t))\dot{c}(t)\| dt. \quad (7)$$

A metric on \mathcal{M} can then, for $x, y \in \mathcal{M}$, be defined as

$$d_{\mathcal{M}}(x, y) = \inf_{c \in C^1(\mathcal{M})} \{s | c(0) = x \text{ and } c(1) = y\}. \quad (8)$$

2.3. The Nakagami distribution

We consider random manifolds immersed by a GP. The length of a curve (7) on such a manifold is necessarily random as well. Fortunately, since this manifold is a Gaussian field, then curve lengths are well-approximated with the Nakagami m -distribution (Bewsher et al., 2017).

The Nakagami distribution (Nakagami, 1960) describes the length of an isotropic Gaussian vector, but Bewsher et al.

(2017) have meticulously demonstrated that this also provides a good approximation to the arc length of a GP. The Nakagami has density function

$$g(s) = \frac{2m^m}{\Gamma(m)\Omega^m} s^{2m-1} \exp\left(-\frac{m}{\Omega} s^2\right), \quad s \geq 0, \quad (9)$$

and it is parametrised by $m \geq 1/2$ and $\Omega > 0$; here Γ denotes the Gamma function. The parameters are interpretable by the equations

$$\Omega = \mathbb{E}[s^2] \quad \text{and} \quad m = \frac{\Omega^2}{\text{Var}(s^2)}, \quad (10)$$

which can be used to infer the parameters through samples, although it does involve a fourth moment.

3. Model and variational inference

With prerequisites settled, we now set up a Gaussian process latent variable model that is *locally* distance preserving and *globally* topology preserving. Notation-wise we let \mathcal{Z} denote the latent representation of a dataset $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^N$, $\mathbf{x}_i \in \mathbb{R}^D$, and let $f : \mathbf{z} \mapsto \mathbf{x}$ be the generative mapping.

3.1. Distance and topology preservation

The *manifold assumption* hypothesizes that high-dimensional data in \mathbb{R}^D lie near a manifold with small intrinsic dimension. A manifold suggests that, a neighborhood around any point is approximately homeomorphic to a linear space. So nearby points are approximately linear, but non-nearby points have distances *greater* than the linear approximation suggests.

We build a Gaussian process latent variable model (GP-LVM) (Lawrence, 2005) that is explicitly designed for distance and topology preservation. The vanilla GP-LVM takes on the Gaussian likelihood where observations \mathcal{X} are assumed i.i.d. when conditioned on a Gaussian process f . That is, $p(\mathcal{X}|f) = \prod_{i=1}^N p(\mathbf{x}_i|f(\mathbf{z}_i))$ and $p(\mathbf{x}_i|f(\mathbf{z}_i)) = \mathcal{N}(\mathbf{x}_i|f(\mathbf{z}_i), \sigma^2)$. In contrast, we consider a likelihood over pairwise distances between observations.

Neighborhood graph. To model locality, we condition our model on a graph embedding of the observed data \mathcal{X} . The graph is the ϵ -nearest neighbor embedded graph; that is, the undirected graph with vertices $V = \mathcal{X}$ and edges $E = \{e_{ij}\}$, where e_{ij} is in E , only if $d(\mathbf{x}_i, \mathbf{x}_j) < \epsilon$, for some metric d . Equivalently, $G = (V, E)$ can be represented by its adjacency matrix A_G with entries

$$a_{ij} = \mathbf{1}_{d(\mathbf{x}_i, \mathbf{x}_j) < \epsilon}. \quad (11)$$

In Sec. 3.4 we discuss how to choose ϵ informedly, but for now we view it as a hyperparameter.

Manifold distances. To arrive at a likelihood over pairwise distances, we first recall that the linear interpolation between \mathbf{z}_i and \mathbf{z}_j in the latent space has curve length

$$s_{ij} = \int_0^1 \|\mathbf{J}(\mathbf{c}(t))\dot{\mathbf{c}}(t)\| dt, \quad \mathbf{c}(t) = \mathbf{z}_i(1-t) + \mathbf{z}_j t, \quad (12)$$

where \mathbf{J} denotes the Jacobian of f , which is our generative manifold approximation.

As the manifold distance $d_{\mathcal{M}}$ is the length of the shortest connecting curve, then s_{ij} is by definition an upper bound on $d_{\mathcal{M}}$. However, as the manifold is locally homeomorphic to a Euclidean space, then we can expect s_{ij} to be a good approximation of the distance to nearby points, i.e.

$$d_{\mathcal{M}}(\mathbf{z}_i, \mathbf{z}_j) \approx s_{ij} \quad \text{for } \|\mathbf{x}_i - \mathbf{x}_j\| < \epsilon \quad (13)$$

$$d_{\mathcal{M}}(\mathbf{z}_i, \mathbf{z}_j) \leq s_{ij} \quad \text{otherwise.} \quad (14)$$

The behavior we seek is that local interpolation in latent space should mimic local interpolation in data space only if the points are close in data space. If they are far apart, they should *repel* each other in the sense that the linear interpolation in latent space should have *large* curve length.

Censoring. To encode this behavior in the likelihood, we introduce *censoring* (Lee & Wang, 2003) into our objective function. This method is usually applied to missing data in survival analysis, when the event of something happening is known to occur later than some time point.

We may think of censoring as modeling inequalities in data. The censored likelihood function for i.i.d. data t_i following distribution function G_θ , with density function g_θ , is defined

$$L(\{t_i\}_{i=1}^N | \theta, T) = \prod_{t_i < T} g_\theta(t_i) \prod_{t_i \geq T} (1 - G_\theta(T)), \quad (15)$$

where θ are the parameters of the distribution G and T is some ‘time point’, where the experiment ended. Carreira-Perpiñan (2010) remark that most neighborhood-embedding methods have loss functions with two terms: one attracting close point and one scattering term for far away connections. Censoring provides a *likelihood* with similar such terms. It may be viewed as a probabilistic version of the ideas in *maximum variance unfolding* (Weinberger & Saul, 2006).

Local distance likelihood. From earlier, we know that if the manifold $f(\mathcal{M})$ is a Gaussian field, then distances (12) are approximately Nakagami distributed. Thus, we write our likelihood as

$$L(\{\{e_{ij}\}_{i < j}\}_{i=1}^{N-1} | \theta, \epsilon) = \prod_{e_{ij} < \epsilon} g_\theta(e_{ij}) \prod_{e_{ij} \geq \epsilon} (1 - G_\theta(\epsilon)),$$

where G_θ is the distribution function of a Nakagami with parameters $\theta = \{m, \Omega\}$. The resulting log-likelihood is given in Eq. 16 within Fig. 2.

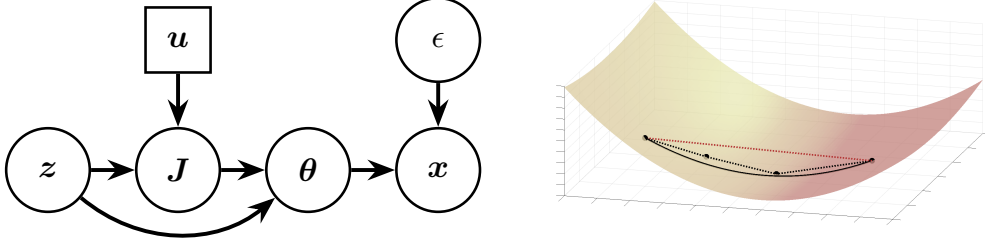


Figure 1. *Left:* A graphical representation of the model: x is the observational input, J is the Gaussian process manifold and θ are the parameters it yields based on latent embedding z . ϵ is a hyperparameter for the neighbor-graph embedding and u are variational parameters. *Right:* Illustration of the task: the dashed lines are Euclidean distances in three dimensions. The black ones are *neighbors* and their distance along the two-dimensional manifold should *match* the 3d-Euclidean distance. The red is not a neighbor-pair and the manifold distance should not match it.

$$l\left(\left\{\{e_{ij}\}_{i < j}\right\}_{i=1}^{N-1} \middle| \theta, \epsilon\right) = -\sum_{e_{ij} < \epsilon} \left(\log \Gamma(m_{ij}) + m_{ij} \log\left(\frac{\Omega_{ij}}{m_{ij}}\right) - (2m_{ij} - 1) \log(e_{ij}) + \frac{m_{ij} e_{ij}^2}{\Omega_{ij}} \right) \\ - \sum_{e_{ij} \geq \epsilon} \left(\log \Gamma(m_{ij}) - \log\left(\Gamma(m_{ij}) - \gamma\left(m_{ij}, \frac{m_{ij}}{\Omega_{ij}} e_{ij}^2\right)\right) \right), \quad (16)$$

Figure 2. The likelihood of our model. Here Γ and γ denotes the Gamma function and lower incomplete gamma function respectively and m_{ij} and Ω_{ij} are the Nakagami-parameters of Eq. 12.

Until now, we have introduced the log-likelihood based on an ϵ -NN graph, that preserves geometric features. Next we marginalize all other parameters to make a Bayesian model.

3.2. Marginalizing the representation

We have a loss function (16) that matches distances e_{ij} with parameters $\theta_{ij} = \{m_{ij}, \Omega_{ij}\}$. We now seek to first fit these parameters and marginalize them to obtain a full Bayesian approach. First, we will assume that conditioned on θ , we get the independent observations, i.e.

$$p(\mathcal{E}|\theta, \epsilon) = \prod_{1 \leq i < j \leq N} p(e_{ij}|\theta_{ij}, \epsilon) \quad (17)$$

$$= L\left(\left\{\{e_{ij}\}_{i < j}\right\}_{i=1}^{N-1} \middle| \theta, \epsilon\right), \quad (18)$$

as known from Eq. 3.1. We infer these parameters of the Nakagami by introducing a latent Gaussian field J and a latent representation z . This allows us to define curve length (12), which we assume is also Nakagami distributed. In practice, we draw¹ m samples of s_{ij} from Eq. 12, and estimate the mean and variance of their second moment. This gives estimates of m_{ij} and Ω_{ij} via Eq. 10.

Essentially, we match distances on the manifold J with the

¹We can approximate s by finely discretizing c and sum over the integrand.

observed distances \mathcal{E} . We marginalize this manifold

$$p(\mathcal{E}|z) = \int p(\mathcal{E}|\theta) p(\theta|J, z) p(J) d\theta dJ, \quad (19)$$

where

$$p(\theta|J, z) := \int p(\theta|s) p(s|J, z) ds, \quad (20)$$

$$\text{and } p(\theta|s) = \begin{cases} \delta_{\mathbb{E}, s^2}(\Omega) \\ \delta_{\Omega/\text{Var}(s^2)}(m), \end{cases} \quad (21)$$

and δ denotes the Dirac probability measure and $p(s|J, z)$ is the approximate Nakagami distribution (12). This means that s_{ij} and e_{ij} are both Nakagami variables that share the same parameters, which interpretively means the manifold distances s_{ij} match the embedding distances e_{ij} .

Further, we can pose a prior on z and marginalize this in Eq. 19. We infer everything variationally (Blei et al., 2017), and choose a variational distribution over the marginalized variables. We approximate the posterior $p(\theta, J, z, u|\mathcal{E})$ with

$$q(\theta, J, z, u) := q(\theta|J, z) q(J, u) q(z), \quad (22)$$

where u is an inducing variable (Titsias, 2009), and

$$q(\theta|J, z) = p(\theta|J, z), \quad q(J, u) = p(J|u) q(u) \quad (23)$$

$$\text{and } q(z) = \mathcal{N}(\mu_z, \mathbf{A}_z), \quad (24)$$

where μ_z is a vector of size N and A_z is a diagonal $N \times N$ -matrix. Further $q(u) = \mathcal{N}(\mu_u, S)$ is a full M -dimensional Gaussian.

This allow us to bound the log-likelihood (16), with the evidence lower bound (ELBO)

$$\log p(\mathcal{E}) = \log \int \frac{p(\mathcal{E}, \theta, \mathbf{J}, \mathbf{z}, \mathbf{u})}{q(\theta, \mathbf{J}, \mathbf{z}, \mathbf{u})} q(\theta, \mathbf{J}, \mathbf{z}, \mathbf{u}) d\theta d\mathbf{J} d\mathbf{u} d\mathbf{z} \quad (25)$$

$$\geq \mathbb{E}_\theta[l(\mathcal{E}|\theta)] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})) - \text{KL}(q(\mathbf{z})||p(\mathbf{z})), \quad (26)$$

where both KL-terms are analytically tractable, but the first term has to be approximated using Monte Carlo. The right hand side here is readily optimized with gradient descent type algorithms.

In summary, we have a latent representation \mathcal{Z} and a Riemannian manifold immersed as a GP \mathbf{J} . This implies that between any two points z_i and z_j , we can compute s_{ij} , which is approximately Nakagami. With censoring we can match s_{ij} with observation e_{ij} , if $e_{ij} < \epsilon$; else we push s_{ij} to have all its mass on $[\epsilon, \infty)$. It is optimized with variational inference by maximizing Eq. 26.

3.3. Invariances and geometric constraints

Why is it worth learning the manifold in a coordinate-free way? Invariances are easily encoded via dissimilarity pairs by introducing equivalence classes in saying $d(x_i, x_j) = 0$ if x_i and x_j are in the same equivalence class. Popular choices of such equivalence classes are rotations, translations and scaling. Many constraints one could wish to impose on models can be formulated as geometric constraints. It holds true also for GPLVM-based models as seen in Urtasun et al. (2008), who wish to encode topological information, and Zhang et al. (2010), who highlight invariant models' usefulness in causal inference. Geometric constraints can alternatively be encoded with GPs that take their output directly on a Riemannian manifold (Mallasto et al., 2018). Kato et al. (2020) try to enforce geometric constraints in Euclidean autoencoders by changing the optimisation, and Miolane & Holmes (2020) build Riemannian VAEs.

The geometry of latent variable models in general is an active field of study (Arvanitidis et al., 2018; Tosi et al., 2014), and Simard et al. (2012) and Kumar et al. (2017) argues that the tangent (Jacobian) space serves a convenient way to encode invariances. Recently, Borovitskiy et al. (2020) developed a framework for GPs defined on Riemannian manifold. Contrary to their method, we learn the manifold where they a priori determine it.

3.4. Topological Data Analysis and the influence of ϵ

The model is naturally affected by the hyperparameter ϵ . We argue that it can be chosen in a geometrically founded way using Topological Data Analysis (Carlsson, 2009). By constructing a *Rips diagram* (Fasy et al., 2014) one can find ϵ such that the ϵ -NN graph captures the right topology of data. It is beyond this paper to summarize the techniques; we refer readers to Chazal & Michel (2017).

To understand what ϵ means in broader terms we can study corner cases. If $\epsilon = \infty$ we would match *all* observed distances, which resembles MDS. If the covariance function of the marginalized \mathbf{J} is constant² the latent space is also preserved (scaled) Euclidean, hence iso-GPLVM may in this setting be viewed as a probabilistic MDS. This links well with how the GPLVM generalized the probabilistic PCA (Lawrence, 2005).

Although we shall not further discuss it in this paper, the Bayesian setup also suggests ϵ could potentially be marginalized. The argument why this is not as straightforward as one could hope is that the model has a pathological solution in the corner case $\epsilon = 0$. In this case, all points would repel each other, and a high likelihood can be obtained without a meaningful representation.

4. Experiments

We perform experiments first on a classical toy dataset and on the image datasets COIL20 and MNIST. We refer to the presented model as *Isometric Gaussian Process Latent Variable Model* (Iso-GPLVM). For comparisons we evaluate other models also based on dissimilarity data. In all cases we initialize Iso-GPLVM with IsoMap, as it is known GP-based methods are sensitive to initialization (Bitzer & Williams, 2010). We use the Adam-optimizer (Kingma & Ba, 2014) with a learning rate of $3 \cdot 10^{-3}$ and optimize sequentially $q(z)$ and $q(u)$ separately. We use $m = 100$ inducing points for $q(u)$ and an ARD-kernel as covariance function.

4.1. Swiss roll

The 'swiss roll' was introduced by Tenenbaum et al. (2000) to highlight the difficulties of non-linear manifold learning. The point cloud resides on a 2-dimensional manifold embedded in \mathbb{R}^3 and can be thought as a paper rolled around itself (see Fig 3A).

We find a 2-dimensional latent embedding by four methods: MDS, t-SNE, IsoMap and Iso-GPLVM. From Fig. 3 we observe the linear MDS is unable to capture the highly non-linear manifold. t-SNE captures some local structure, but the global outlook is far from the ground truth. We tried

²In this case the generating function f has a linear kernel.

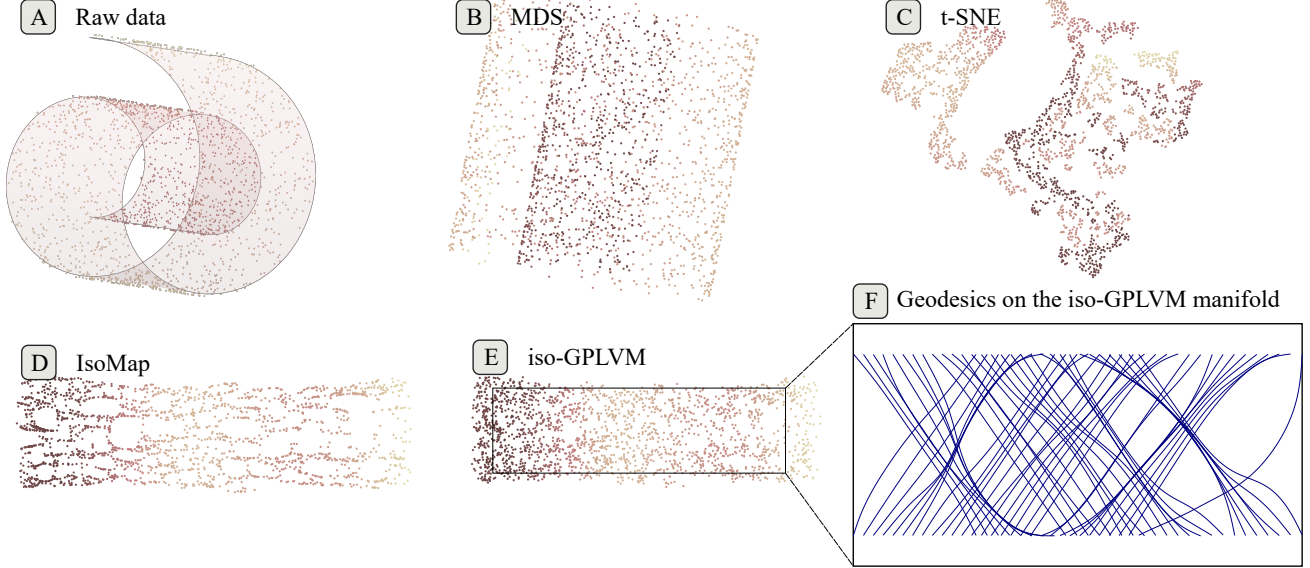


Figure 3. Data (A) and embeddings (B–E). All embeddings are shown with a unit aspect ratio to highlight that only IsoMap (D) and Iso-GPLVM (E) recover the elongated structure of the swiss roll. (F) shows some geodesics on the learned 2-dimensional manifold.

several tunings of the perplexity hyperparameter (60 in the plot), none successfully captured the structure. It is known that t-SNE is prone to create clusters, even if clusters are not a natural part of a dataset (Amid & Warmuth, 2018).

Naturally, as the dataset was constructed for the ‘geodesic’ approach of IsoMap, this captures both global and local structure. On closer inspection, we see the linear interpolations, stemming from Dijkstra’s algorithm, leaves some artificial ‘holes’ in the manifold. Hence, on a smaller scale it can be argued the topology of the manifold is captured imperfectly. The plot suggests Iso-GPLVM closes these holes and approximates the topology of an unfolded paper.

Figure 3F visualizes some geodesics and they appear roughly linear. There is some ‘gathering’ fix points which are due to the sparsity of the GP. These geodesics inform us that not only is the representation good, but the learned *geometry* is correct since the geodesics match those we know from Fig. 3A. We used $\epsilon = 0.4$.

4.2. COIL20

COIL20 (Nene et al., 1996) consists of greyscale images of 20 objects photographed from 72 different angles spanning a full rotation (see Figure 4 for some examples). This implies in total 1440 images — the version we use is of size 128×128 pixels, thus the original data resides in \mathbb{R}^{16384} .

First, we focus on only one object — a rotated rubber duck — to highlight the geodesic behaviour. Figure 4 shows the 2-dimensional embeddings and the geodesic curves on the learned manifold in latent space. We clearly observe the circular structure we expect from the rotated duck. On top

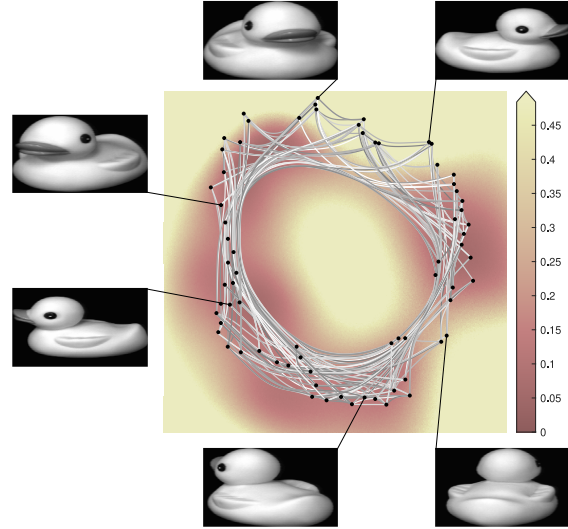


Figure 4. The 2-dimensional embeddings of the 72 images of a rubber duck. We observe from the geodesics (grey curves) how the latent manifold has learned the circular nature of the data.

of this the geodesics show the Riemannian geometry of the latent space: they move along the data manifold and avoid the space where no data is observed. The background color is the measure $\mathbb{E}[\sqrt{\det(J^T J)}]$, which provide a view of the Riemannian geometry of the latent space. Bishop et al. (1997) call this measure the *magnification factor*. Large values (light color) imply trajectories moving in this area are longer and likely also more uncertain (Hauberg, 2018).

IsoMap, t-SNE, UMAP and others, are also able to infer the circular embedding, but Iso-GPLVM is the only model

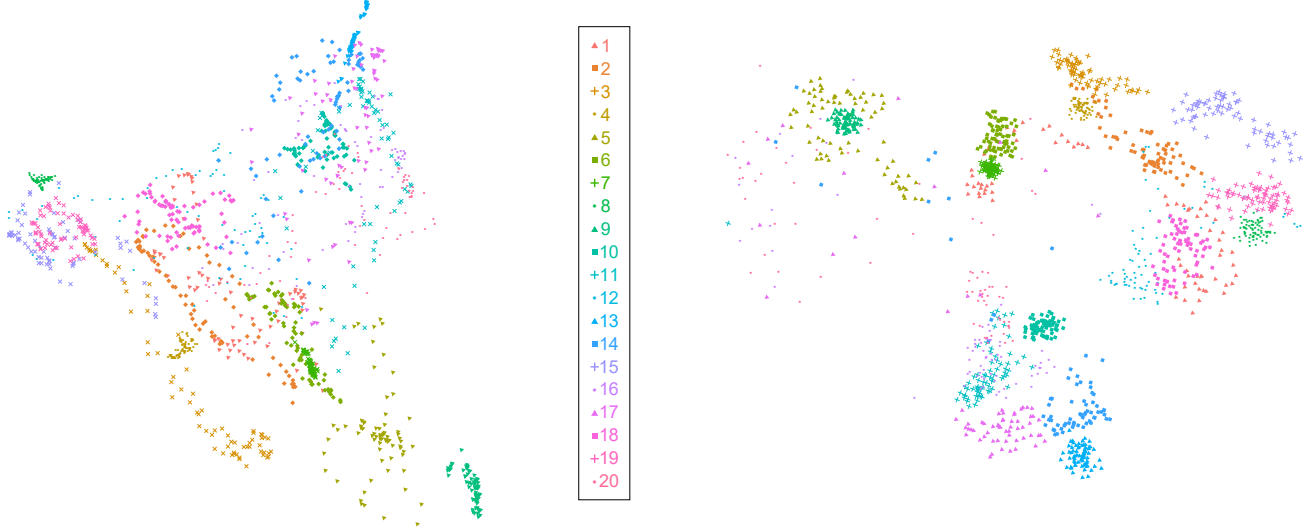


Figure 5. Embeddings of COIL20 objects. *Left*: IsoMap and *right*: Iso-GPLVM. We see that globally Iso-GPLVM can separate the objects (color- and shape coded), but is not able to find to all local structures.

to infer a *geometry* on latent space. For IsoMap the latent geometry is implicitly Euclidean through its loss (1), and t-SNE and UMAP do not allow for geodesic computations.

When considering all 20 objects at once a global element of separating the distinct objects is a key task to infer the topological structure. The embeddings for IsoMap and Iso-GPLVM are visible in Fig. 5. Here IsoMap struggle to clearly separate objects due to its implicit assumption of one connected manifold. Iso-GPLVM finds the *global* topological structure, but in no instances finds the local structure. So why is it unsuccessful here when successful in Fig. 4? When considering all 1440 images we only use 100 inducing points, and in this view it is unsurprising that the model has to use most capacity on the global structure. In Fig. 4 there is no sparsity required since there is only 72 images, and there is enough capacity to detect the hole in the manifold. This is a common problem for GP-based methods.

4.3. MNIST

Metrics. We evaluate our model on 5000 images from MNIST, and we foremost wish to highlight how invariances can be encoded with dissimilarity data. We consider fitting our model to data under three different distance measures. We consider the classical Euclidean distance measure

$$d(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|. \quad (27)$$

Further, we consider a metric that is invariant under image rotations

$$d_{\text{ROT}}(\mathbf{x}_i, \mathbf{x}_j) = \inf_{\theta \in [0, 2\pi)} \left\{ d(R_\theta(\mathbf{x}_i), \mathbf{x}_j) \right\}, \quad (28)$$

where R_θ rotates an image by θ radians. We note $d_{\text{ROT}}(\mathbf{x}_i, \mathbf{x}_j) \leq d(\mathbf{x}_i, \mathbf{x}_j)$ always. Finally, we introduce

a *lexicographic* metric (Rodriguez-Velazquez, 2018)

$$d_{\text{LEX}}(\mathbf{x}_i, \mathbf{x}_j) = \begin{cases} \epsilon, & \text{if } y_i \neq y_j \\ \min\{2r, d(\mathbf{x}_i, \mathbf{x}_j)\}, & \text{if } y_i = y_j \end{cases} \quad (29)$$

which in the censoring phase enforce images carrying different labels to repel each other. This is a handy way to encode a topology or clustering based on discrete variables, when such are available. For all metrics, we have normalized the data and have set $\epsilon = 7$.

Results. Figure 6(A—C) show the latent embeddings of the three metrics. The background color again indicates the magnification factor $\mathbb{E}[\sqrt{\det(J^T J)}]$. Panels A, D and E base their latent embedding on the Euclidean metric. We observe that IsoMap (D) and Iso-GPLVM (A) appear similar in shape, unsurprisingly as we initialize with IsoMap, but Iso-GPLVM finds a cleaner separation of the digits. Particularly, this is evident for the *six*, *three* and *eight* digits. The *fives* seem to group into several tighter cluster, and this behavior is found for t-SNE as well. Overall, from a clustering perspective, t-SNE visually is superior; but distances *between* clusters in (A) can be larger than the straight lines that connect them. This is evident from the lighter background color between cluster, say, *zeros* and *threes*. We note that IsoMap and t-SNE has no associated Riemannian metric and as such distances between any input cannot be computed.

The rotation invariant metric results in a latent embedding where different classes significantly overlap. Upon closer inspection we, however, note several interesting properties of the embedding. *Zero* digits are well separated from other classes as a rotated 0 does not resemble any other digits; the *one* digits form a cluster that is significantly more compact

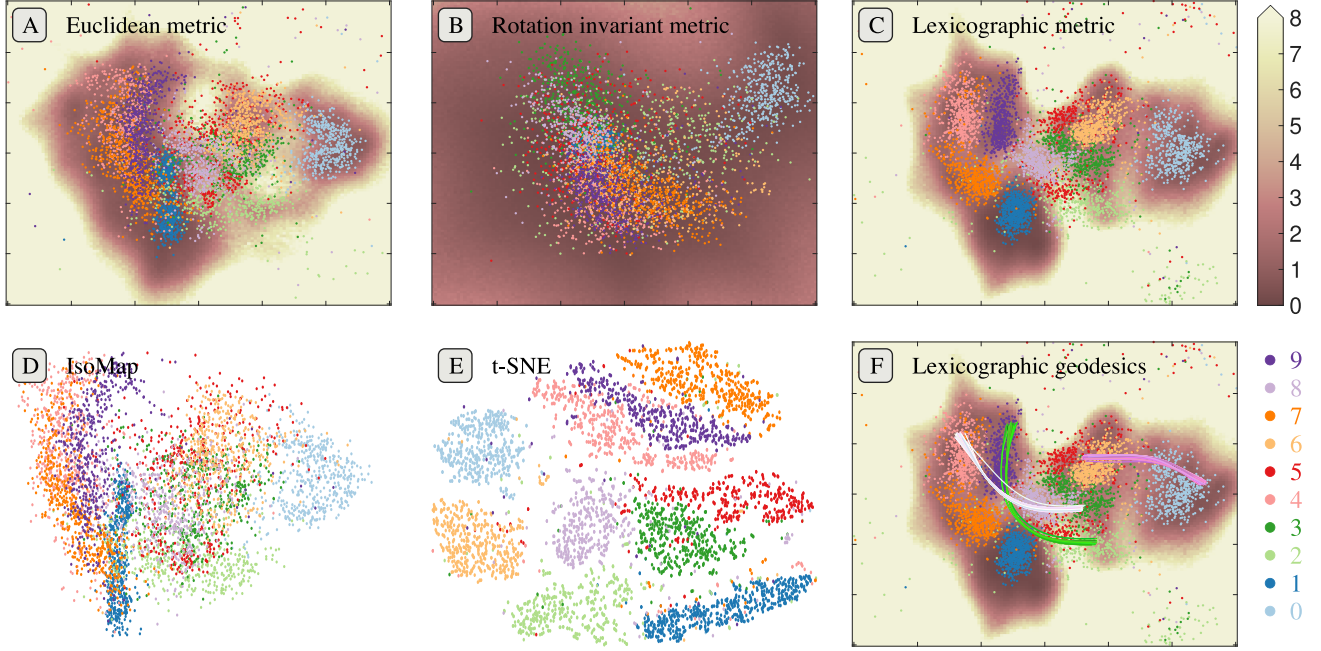


Figure 6. Embeddings of MNIST attained with our method under different metrics (A—C) and for baselines IsoMap (D) and t-SNE (E). The background color show the expected volume measure associated with the Riemannian metric $\mathbb{E}[\sqrt{\det(J^T J)}]$. A large measure generally indicate high uncertainty of the manifold. Panel F shows Riemannian geodesics under the lexicographic metric.

than other digits as there is limited variation left after rotations have been factored out; *two and five digits* significantly overlap, which is most likely due to 5 digits resembling 2 digits when rotated 180° ; similar observations hold for the *four, nine and six digits*; and a partial overlap between *three and eight digits* as is often observed. The overall darker background is due to the rotational invariant metric being shorter than the Euclidean counterpart.

In terms of clustering the lexicographic approach outshines the other metrics. This is expected as the metric use label information, but neatly illustrate how domain-specific metrics can be developed from weak or partial information. Most classes are well-separated except for a region in the middle of the plot. Note how this region has high uncertainty.

The Riemannian geometry of the latent space implies that geodesics (shortest paths) can be computed in our model. Figure 6F shows example geodesics under the lexicographic metric. Their highly non-linear appearance emphasizes the curvature of the learned manifold. The green geodesics has one endpoint in a cluster of nine digits and move along this cluster avoiding the uncertain area of eights and fives, as opposed to linearly interpolating through them.

5. Discussion

We introduced a model for non-linear dimensionality reduction from dissimilarity data. It is the first of its kind based on Gaussian processes. The non-linearity of the method stems both from the Gaussian processes, but also from the censoring in the likelihood. It unifies ideas from Gaussian processes, Riemannian geometry and neighborhood graph embeddings. Unlike traditional manifold learning methods that embed into \mathbb{R}^q , we embed into a q -dimensional Riemannian manifold through the learned metric. This allows us to learn latent representations that are isometric to the true underlying manifold.

The model does have limitations. Aesthetically the visualizations are not as satisfactory as e.g. t-SNE. However, the access to a geometrically founded GPLVM is of interest to many practitioners, since GPs are ubiquitous in many sub-disciplines of machine learning such as Bayesian optimization and reinforcement learning. Here, GPs are fundamental parts of decision-making pipelines, whereas t-SNE is a valuable visualization technique. The Nakagami distribution that approximates the arc lengths of Gaussian processes is prone to overestimate the variance (Bewsher et al., 2017) and better approximations would improve our method. Further, the model inherits problems of optimizing the latent variables and it has previously been noted that good performance in this regime is linked with good

initialization (Bitzer & Williams, 2010).

Our experiments highlight that Iso-GPLVM can learn the geometry of data and geometric constraints are easier encoded by learning a manifold contra doing GP regression. The uncertainty quantification associated with GPs follow through and further highlights the connection between uncertainty, geometry and topology. To the best of our knowledge, our model is the first of its kind that, locally, can asses the quality of the manifold approximation through the associated Riemannian measure.

Acknowledgements

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement no 757360). MJ and SH were supported in part by a research grant (15334) from VILLUM FONDEN. MJ is supported by the Carlsberg Foundation (CF20-0370). The majority of this work was done while MJ was affiliated with the Technical University of Denmark.

References

- Amid, E. and Warmuth, M. K. A more globally accurate dimensionality reduction method using triplets. *arXiv:1803.00854 [cs]*, March 2018.
- Arvanitidis, G., Hansen, L. K., and Hauberg, S. Latent space oddity: On the curvature of deep generative models. In *6th International Conference on Learning Representations, ICLR 2018*, 2018.
- Bengio, Y., Courville, A., and Vincent, P. Representation Learning: A Review and New Perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, August 2013.
- Bewsher, J., Tosi, A., Osborne, M., and Roberts, S. Distribution of gaussian process arc lengths. In *Artificial Intelligence and Statistics*, pp. 1412–1420, 2017.
- Bishop, C. M., Svens’ en, M., and Williams, C. K. Magnification factors for the som and gtm algorithms. In *Proceedings 1997 Workshop on Self-Organizing Maps*, 1997.
- Bitzer, S. and Williams, C. K. Kick-starting gplvm optimization via a connection to metric mds. In *NIPS 2010 Workshop on Challenges of Data Visualization*, 2010.
- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Borovitskiy, V., Terenin, A., Mostowsky, P., and Deisenroth (he/him), M. Matérn gaussian processes on riemannian manifolds. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 12426–12437. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/92bf5e6240737e0326ea59846a83e076-Paper.pdf>.
- Carlsson, G. Topology and data. *Bulletin of the American Mathematical Society*, 46(2):255–308, 2009.
- Carreira-Perpiñan, M. Á. The elastic embedding algorithm for dimensionality reduction. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pp. 167–174, 2010.
- Chazal, F. and Michel, B. An introduction to topological data analysis: fundamental and practical aspects for data scientists, 2017.
- Fasy, B. T., Lecci, F., Rinaldo, A., Wasserman, L., Balakrishnan, S., and Singh, A. Confidence sets for persistence diagrams. *Ann. Statist.*, 42(6):2301–2339, 12 2014. doi: 10.1214/14-AOS1252. URL <https://doi.org/10.1214/14-AOS1252>.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- Hauberg, S. Only bayes should learn a manifold (on the estimation of differential geometric structure from data). *arXiv preprint arXiv:1806.04994*, 2018.
- Kato, K., Zhou, J., Sasaki, T., and Nakagawa, A. Rate-distortion optimization guided autoencoder for isometric embedding in Euclidean latent space. In III, H. D. and Singh, A. (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5166–5176. PMLR, 13–18 Jul 2020.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kingma, D. P. and Welling, M. Auto-Encoding Variational Bayes. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR)*, 2014.
- Kumar, A., Sattigeri, P., and Fletcher, T. Semi-supervised learning with gans: Manifold invariance with improved inference. In *Advances in Neural Information Processing Systems*, pp. 5534–5544, 2017.

- Lawrence, N. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *Journal of machine learning research*, 6(Nov):1783–1816, 2005.
- Lee, E. T. and Wang, J. *Statistical Methods for Survival Data Analysis*, volume 476. John Wiley & Sons, 2003.
- Maaten, L. v. d. and Hinton, G. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.
- Mallasto, A., Hauberg, S., and Feragen, A. Probabilistic riemannian submanifold learning with wrapped gaussian process latent variable models. In *Proceedings of the 19th international Conference on Artificial Intelligence and Statistics (AISTATS)*, 2018.
- Mattei, P.-A. and Frellsen, J. Miwae: Deep generative modelling and imputation of incomplete data sets. In *International Conference on Machine Learning*, pp. 4413–4423, 2019.
- McInnes, L., Healy, J., and Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- Miolane, N. and Holmes, S. Learning weighted submanifolds with variational autoencoders and riemannian variational autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14503–14511, 2020.
- Nakagami, M. The m-distribution—a general formula of intensity distribution of rapid fading. In *Statistical Methods in Radio Wave Propagation*, pp. 3–36. Elsevier, 1960.
- Nene, S. A., Nayar, S. K., and Murase, H. Columbia university image library (coil-20). 1996. URL <http://www.cs.columbia.edu/CAVE/software/softlib/coil-20.php>.
- Rasmussen, C. and Williams, C. *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, USA, 2006.
- Rezende, D. J. and Mohamed, S. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic Backpropagation and Approximate Inference in Deep Generative Models. In *Proceedings of the 31st International Conference on Machine Learning (ICML)*, 2014.
- Ripley, B. D. *Pattern recognition and neural networks*. Cambridge university press, 2007.
- Rodriguez-Velazquez, J. A. Lexicographic metric spaces: Basic properties and the metric dimension, 2018.
- Simard, P. Y., LeCun, Y. A., Denker, J. S., and Victorri, B. Transformation invariance in pattern recognition—tangent distance and tangent propagation. In *Neural networks: tricks of the trade*, pp. 235–269. Springer, 2012.
- Tenenbaum, J. B., de Silva, V., and Langford, J. C. A Global Geometric Framework for Nonlinear Dimensionality Reduction. *Science*, 290(5500):2319, 2000.
- Titsias, M. Variational learning of inducing variables in sparse gaussian processes. In *Artificial Intelligence and Statistics*, pp. 567–574, 2009.
- Titsias, M. and Lawrence, N. D. Bayesian gaussian process latent variable model. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 844–851, 2010.
- Tosi, A., Hauberg, S., Vellido, A., and Lawrence, N. D. Metrics for Probabilistic Geometries. In *The Conference on Uncertainty in Artificial Intelligence (UAI)*, July 2014.
- Urtasun, R. and Darrell, T. Discriminative gaussian process latent variable model for classification. In *Proceedings of the 24th International Conference on Machine Learning*, pp. 927–934, 2007.
- Urtasun, R., Fleet, D. J., Geiger, A., Popović, J., Darrell, T. J., and Lawrence, N. D. Topologically-constrained latent variable models. In *Proceedings of the 25th International Conference on Machine Learning*, pp. 1080–1087, 2008.
- Weinberger, K. Q. and Saul, L. K. An introduction to non-linear dimensionality reduction by maximum variance unfolding. In *AAAI*, volume 6, pp. 1683–1686, 2006.
- Zhang, K., Schölkopf, B., and Janzing, D. Invariant gaussian process latent variable models and application in causal discovery. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pp. 717–724, 2010.