

INVERSE KERNEL DECOMPOSITION

Anonymous authors

Paper under double-blind review

ABSTRACT

The state-of-the-art dimensionality reduction approaches largely rely on complicated optimization procedures. On the other hand, closed-form approaches requiring merely eigen-decomposition do not have enough sophistication and nonlinearity. In this paper, we propose a novel nonlinear dimensionality reduction method—Inverse Kernel Decomposition (IKD)—based on an eigen-decomposition of the sample covariance matrix of data. The method is inspired by Gaussian process latent variable models (GPLVMs) and has comparable performance with GPLVMs. To deal with very noisy data with weak correlations, we propose two solutions—blockwise and geodesic—to make use of locally correlated data points and provide better and numerically more stable latent estimations. We use synthetic datasets and four real-world datasets to show that IKD is a better dimensionality reduction method than other eigen-decomposition-based methods, and achieves comparable performance against optimization-based methods with faster running speeds. Open-source IKD implementation in Python can be accessed at this <https://anonymous.4open.science/r/ikd-BABC>.

1 INTRODUCTION

Dimensionality reduction techniques have been widely studied in the machine learning field for many years, with massive applications in latent estimation (Wu et al., 2017; 2018), noise reduction (Sheybani & Javidi, 2009), cluster analysis (Bakrania et al., 2020), data visualization (Van der Maaten & Hinton, 2008a) and so forth. The most commonly used method is Principled Component Analysis (PCA), a linear dimensionality reduction approach. It is favored thanks to the easy use of a one-step eigen-decomposition. Its simple linear assumption, however, restricts its exploitation, especially in highly-nonlinear scenarios. On the other hand, nonlinear dimensionality reduction models, such as autoencoders (Kramer, 1991), variational autoencoders (VAE) (Kingma & Welling, 2013), t-SNE (Van der Maaten & Hinton, 2008b), UMAP (McInnes et al., 2020), and Gaussian process latent variable models (GPLVMs) (Lawrence, 2003; 2005) can achieve state-of-the-art (SOTA) performance in terms of finding (sub)optimal low-dimensional latent and rendering satisfactory downstream analyses (e.g., visualization, prediction, classification). However, all these nonlinear models involve intricate optimization which is time-consuming, easy to get stuck in bad local optima, and sensitive to initialization. In this paper, we propose a novel nonlinear eigen-decomposition-based dimensionality reduction approach that finds low-dimensional latent with a closed-form solution but intricate nonlinearity.

The proposed model is called Inverse Kernel Decomposition (IKD), inspired by GPLVMs. GPLVMs are probabilistic dimensionality reduction methods that use Gaussian Processes (GPs) to find a lower dimensional nonlinear embedding of high-dimensional data. They use a kernel function to form a nonlinear mapping from the embedded latent space to the data space, which is opposite to the use of kernel as in kernel PCA (Schölkopf et al., 1997). GPLVM and its many variants have been proposed in various domains (Bui & Turner, 2015; Wang et al., 2005; 2008; Urtasun et al., 2006; Wu et al., 2017) and proven to be powerful nonlinear dimensionality reduction and latent variable models. However, GPLVMs are highly nonlinear and non-convex due to the GP component, resulting in practical difficulties during optimization. IKD employs the same kernel function mapping the latent space to the data space capturing the nonlinearity, but solves the dimensionality reduction problem through eigen-decomposition. In the experiment section, we compare IKD against four eigen-decomposition-based and four optimization-based dimensionality reduction methods using synthetic datasets and four real-world datasets, and we can summarize four contributions of IKD:

- As an eigen-decomposition-based method, IKD achieves more reasonable latent representations than other eigen-decomposition-based methods with better classification accuracy in downstream classification tasks. The running time of IKD is on par with other eigen-decomposition-based methods.
- IKD is able to provide competitive performance against some SOTA optimization-based methods but at a much faster running speed.
- IKD promises a stable and unique optimal solution up to an affine transformation. In contrast, optimization-based methods do not guarantee unique optimal solutions and sometimes are not numerically stable due to the highly nonconvex optimization landscapes. Therefore, IKD can sometimes achieve better latent representations and classification performance than optimization-based methods like GPLVM and VAE.
- When the observation dimensionality is large (i.e. observation data is high-dimensional), a lot of methods have significant drawbacks. For example, t-SNE, UMAP, and VAE encounter the curse of dimensionality problem. The large dimensionality not only leads to longer running time but also hurts the dimensionality reduction performance. In contrast, IKD always obtains improved performance with an increasing observation dimensionality, and claims its absolute superiority when the observation dimensionality is very large.

Note that we are not claiming to propose the best dimensionality reduction approach that beats all other SOTAs. We propose an advanced eigen-decomposition-based method that (1) outperforms other eigen-decomposition-based methods in most of synthetic and real-world applications with the same scale of running speeds; and (2) reaches a comparable level against other optimization-based methods but with much faster running speed.

2 METHODOLOGY

2.1 GAUSSIAN PROCESS LATENT VARIABLE MODEL

Let $\mathbf{X} \in \mathbb{R}^{T \times N}$ be the observed data where T is the number of observations and N is the observation dimensionality of each data vector. Let $\mathbf{Z} \in \mathbb{R}^{T \times M}$ denote its associated latent variables where M is the latent dimensionality. Usually, we assume the latent space is lower-dimensional than the original observational space, leading to $M < N$. For each dimension of \mathbf{X} denoted as $\mathbf{X}_{:,n} \in \mathbb{R}^T$, $\forall n \in \{1, \dots, N\}$, GPLVM defines a mapping function that maps the latent to the observation which has a Gaussian process (GP) prior. Therefore, given the finite number of observations, we can write

$$\mathbf{X}_{:,n} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{K}), \quad \forall n \in \{1, 2, \dots, N\}. \quad (1)$$

where \mathbf{K} is a $T \times T$ covariance matrix generated by evaluating the kernel function k of GP at all pairs of rows in \mathbf{Z} . I.e., $k_{i,j} = k(\mathbf{z}_i, \mathbf{z}_j)$ where \mathbf{z}_i and \mathbf{z}_j are the i^{th} and j^{th} rows of \mathbf{Z} . The goal of GPLVM is to estimate the unknown latent variables \mathbf{Z} that is used for constructing the covariance matrix \mathbf{K} , from the observations \mathbf{X} . Note that we only consider noiseless GP for the derivation of IKD, but IKD can deal with noisy observations empirically.

2.2 INVERSE KERNEL DECOMPOSITION

In this section, we derive a novel nonlinear decomposition method, inverse kernel decomposition (IKD), inspired by GPVLM. Previous work has been solving GPLVM by maximizing the log-likelihood to obtain \mathbf{Z} from \mathbf{X} in a one-step fashion. Now let us break this process into two steps: (1) estimating \mathbf{K} from the observations \mathbf{X} , and (2) identifying the latent variables \mathbf{Z} from the estimated covariance matrix \mathbf{K} . The first step can be solved by estimating \mathbf{K} with the unbiased estimator, i.e., sample covariance $\mathbf{S} := \frac{1}{N-1} (\mathbf{X} - \bar{\mathbf{X}}\mathbf{1}^T) (\mathbf{X} - \bar{\mathbf{X}}\mathbf{1}^T)^T \approx \frac{1}{N-1} \mathbf{X}\mathbf{X}^T$, where $\bar{\mathbf{X}} = \frac{1}{N} \sum_{n=1}^N \mathbf{X}_{:,n}$ ought to be $\mathbf{0}$ since $\mathbf{X}_{:,n}$ are i.i.d. samples from a zero-mean Gaussian (Eq. 1).

Therefore, our main focus is to estimate the latent \mathbf{Z} given \mathbf{S} in the second step. In the following, we focus on the discussion of a commonly used stationary kernel, the squared exponential (SE) kernel. We will show in Sec. 2.3 that IKD can also work with various stationary kernels.

The SE kernel is defined as $k(\mathbf{z}_i, \mathbf{z}_j) = \sigma^2 \exp\left(-\frac{\|\mathbf{z}_i - \mathbf{z}_j\|^2}{2l^2}\right)$, where σ^2 is the marginal variance and l is the length-scale. Note that $\sigma^2 = k(\mathbf{z}_i, \mathbf{z}_i) = k_{i,i}$, $\forall i \in \{1, \dots, T\}$. Let $f(d)$ be the function mapping the scaled squared distance $d_{i,j} := \frac{\|\mathbf{z}_i - \mathbf{z}_j\|^2}{l^2}$ between latent $(\mathbf{z}_i, \mathbf{z}_j)$ to the scalar covariance $k_{i,j}$, i.e., $k_{i,j} = f(d_{i,j}) = \sigma^2 \exp\left(-\frac{d}{2}\right)$. Let us assume we know the true \mathbf{K} for now. Since $f(\cdot)$ is strictly monotonic, we can obtain $d_{i,j} = f^{-1}(k_{i,j}) = -2 \ln\left(\frac{k}{\sigma^2}\right)$. Writing $d_{i,j}$ in the matrix form $\mathbf{D} = (d_{i,j})_{T \times T}$, we have

$$\begin{aligned} \mathbf{D} &= \frac{1}{l^2} \begin{bmatrix} 0 & (\mathbf{z}_1 - \mathbf{z}_2)^\top (\mathbf{z}_1 - \mathbf{z}_2) & \cdots & (\mathbf{z}_1 - \mathbf{z}_T)^\top (\mathbf{z}_1 - \mathbf{z}_T) \\ (\mathbf{z}_2 - \mathbf{z}_1)^\top (\mathbf{z}_2 - \mathbf{z}_1) & 0 & \cdots & (\mathbf{z}_2 - \mathbf{z}_T)^\top (\mathbf{z}_2 - \mathbf{z}_T) \\ \vdots & \vdots & \ddots & \vdots \\ (\mathbf{z}_T - \mathbf{z}_1)^\top (\mathbf{z}_T - \mathbf{z}_1) & (\mathbf{z}_T - \mathbf{z}_2)^\top (\mathbf{z}_T - \mathbf{z}_2) & \cdots & 0 \end{bmatrix} \\ &= f^{-1}(\mathbf{K}) = (f^{-1}(k_{i,j}))_{T \times T}, \end{aligned} \quad (2)$$

where f^{-1} maps \mathbf{K} to \mathbf{D} element-wisely. We define $\tilde{\mathbf{z}} = \frac{\mathbf{z} - \mathbf{z}_1}{l}$ with $\tilde{\mathbf{z}}_1 = \mathbf{0}$. Now we have

$$\begin{aligned} d_{i,j} &= \frac{1}{l^2} (\mathbf{z}_i - \mathbf{z}_j)^\top (\mathbf{z}_i - \mathbf{z}_j) = \frac{1}{l^2} [(\mathbf{z}_i - \mathbf{z}_1) - (\mathbf{z}_j - \mathbf{z}_1)]^\top [(\mathbf{z}_i - \mathbf{z}_1) - (\mathbf{z}_j - \mathbf{z}_1)] \\ &= (\tilde{\mathbf{z}}_i - \tilde{\mathbf{z}}_j)^\top (\tilde{\mathbf{z}}_i - \tilde{\mathbf{z}}_j) = \tilde{\mathbf{z}}_i^\top \tilde{\mathbf{z}}_i + \tilde{\mathbf{z}}_j^\top \tilde{\mathbf{z}}_j - 2\tilde{\mathbf{z}}_i^\top \tilde{\mathbf{z}}_j. \end{aligned} \quad (3)$$

Since $\tilde{\mathbf{z}}_1 = \mathbf{0}$, we have $d_{1,j} = \tilde{\mathbf{z}}_1^\top \tilde{\mathbf{z}}_1 + \tilde{\mathbf{z}}_j^\top \tilde{\mathbf{z}}_j - 2\tilde{\mathbf{z}}_1^\top \tilde{\mathbf{z}}_j \implies \tilde{\mathbf{z}}_j^\top \tilde{\mathbf{z}}_j = d_{1,j}$, $\forall j \in \{1, \dots, T\}$. Therefore, we arrive at an expression of $\tilde{\mathbf{z}}_i^\top \tilde{\mathbf{z}}_j$ as $\tilde{\mathbf{z}}_i^\top \tilde{\mathbf{z}}_j = \frac{1}{2}(d_{i,1} + d_{1,j} - d_{i,j})$. Note that $d_{1,i} = d_{i,1}$ because of the symmetric property. Denote $\tilde{\mathbf{Z}} = [\tilde{\mathbf{z}}_1, \tilde{\mathbf{z}}_2, \dots, \tilde{\mathbf{z}}_T]^\top = [\mathbf{0}, \tilde{\mathbf{z}}_2, \dots, \tilde{\mathbf{z}}_T]^\top \in \mathbb{R}^{T \times M}$, we could write the matrix form $\tilde{\mathbf{Z}}\tilde{\mathbf{Z}}^\top = (\tilde{\mathbf{z}}_i^\top \tilde{\mathbf{z}}_j)_{T \times T}$ as

$$\begin{aligned} \tilde{\mathbf{Z}}\tilde{\mathbf{Z}}^\top &= \begin{bmatrix} 0 & \cdots & 0 & \cdots & 0 \\ 0 & d_{2,1} & \frac{1}{2}(d_{2,1} + d_{1,3} - d_{2,3}) & \cdots & \frac{1}{2}(d_{2,1} + d_{1,T} - d_{2,T}) \\ 0 & \frac{1}{2}(d_{3,1} + d_{1,2} - d_{3,2}) & d_{3,1} & \cdots & \frac{1}{2}(d_{3,1} + d_{1,T} - d_{3,T}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \frac{1}{2}(d_{T,1} + d_{1,2} - d_{T,2}) & \frac{1}{2}(d_{T,1} + d_{1,3} - d_{T,3}) & \cdots & d_{T,1} \end{bmatrix} \\ &=: g(\mathbf{D}) = g(f^{-1}(\mathbf{K})), \end{aligned} \quad (4)$$

which is a rank- M symmetric positive semi-definite matrix given $M < T$. g is the function mapping \mathbf{D} to $\tilde{\mathbf{Z}}\tilde{\mathbf{Z}}^\top$. Then, Eq. 4 has the unique ‘‘reduced’’ eigen-decomposition

$$g(f^{-1}(\mathbf{K})) = \mathbf{U}\mathbf{A}\mathbf{U}^\top = \left(\sqrt{\lambda_1}\mathbf{U}_{:,1}, \dots, \sqrt{\lambda_M}\mathbf{U}_{:,M}\right) \left(\sqrt{\lambda_1}\mathbf{U}_{:,1}, \dots, \sqrt{\lambda_M}\mathbf{U}_{:,M}\right)^\top =: \tilde{\mathbf{U}}\tilde{\mathbf{U}}^\top, \quad (5)$$

where $\mathbf{U}_{:,m} = [0, u_{2,m}, u_{3,m}, \dots, u_{T,m}]^\top \in \mathbb{R}^T$ is the m th column of $\mathbf{U} \in \mathbb{R}^{T \times M}$ and $\mathbf{A} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_M)$ with $\lambda_1 > \lambda_2 > \dots > \lambda_M > \lambda_{M+1} = \dots = \lambda_T = 0$. Note that the unique ‘‘reduced’’ singular value decomposition of $\tilde{\mathbf{Z}}$ is

$$\tilde{\mathbf{Z}} = \mathbf{U}\mathbf{A}^{\frac{1}{2}}\mathbf{V}^\top = \tilde{\mathbf{U}}\mathbf{V}^\top \implies \mathbf{z}_t = l\mathbf{V}\tilde{\mathbf{U}}_{t,:}^\top + \mathbf{z}_1, \quad \forall t \in \{1, \dots, T\}, \quad (6)$$

where \mathbf{z}_1 represents the reference translation, the length-scale l is a scaling factor, and \mathbf{V} is an orthogonal matrix that is responsible for the corresponding rotation and reflection. Since \mathbf{Z} and $\tilde{\mathbf{U}}$ span the same column space such that

$$\begin{aligned} k(\mathbf{z}_i, \mathbf{z}_j) &= \sigma^2 \exp\left(-\frac{\|\mathbf{z}_i - \mathbf{z}_j\|^2}{2l^2}\right) = \sigma^2 \exp\left(-\frac{\|l\mathbf{V}\tilde{\mathbf{U}}_{i,:}^\top - l\mathbf{V}\tilde{\mathbf{U}}_{j,:}^\top\|^2}{2l^2}\right) \\ &= \sigma^2 \exp\left(-\frac{\|\tilde{\mathbf{U}}_{i,:} - \tilde{\mathbf{U}}_{j,:}\|^2}{2}\right) = k_{l=1}(\tilde{\mathbf{U}}_{i,:}, \tilde{\mathbf{U}}_{j,:}). \end{aligned} \quad (7)$$

$\tilde{\mathbf{U}}$ contains all of the low-dimensional information of \mathbf{Z} . Therefore, we consider $\tilde{\mathbf{U}}$ as an estimator of \mathbf{Z} . Eq. 5 and Eq. 6 summarize the inverse relationship from a GP kernel covariance matrix to the latent variable.

Algorithm 1 Inverse kernel decomposition

-
- 1: **function** IKD($\mathbf{X} \in \mathbb{R}^{T \times N}, f$)
 - 2: $\mathbf{S} \leftarrow \frac{1}{N-1} (\mathbf{X} - \bar{\mathbf{X}}\mathbf{1}^\top) (\mathbf{X} - \bar{\mathbf{X}}\mathbf{1}^\top)^\top$ ▷ \mathbf{S} serves as an estimator of the covariance \mathbf{K}
 - 3: $\sigma^2 \leftarrow \frac{1}{T} \sum_{i=1}^T s_{i,i}$ ▷ estimate σ^2 through a statistic of the diagonal of \mathbf{S}
 - 4: $\hat{\mathbf{D}} = \hat{d}_{i,j T \times T} \leftarrow f^{-1}(\mathbf{S})$ ▷ $\hat{\mathbf{D}}$ serves as an estimation of \mathbf{D} (Eq. 2)
 - 5: $\mathbf{U}, \mathbf{A} \leftarrow$ eigen-decomposition of $g(\hat{\mathbf{D}})$ ▷ Eq. 5
 - 6: Form the optimal latent solution $\tilde{\mathbf{U}}$ using \mathbf{U} and \mathbf{A} ▷ Eq. 5
 - 7: **return** $\tilde{\mathbf{Z}} = \tilde{\mathbf{U}}$
 - 8: **end function**
-

Table 1: Stationary kernels that can be applied to IKD.

kernel	f	f^{-1}
squared exponential	$f(d) = \sigma^2 \exp(-\frac{d}{2})$	$f^{-1}(k) = -2 \ln(\frac{k}{\sigma^2})$
rational quadratic	$f(d) = \sigma^2 (1 + \frac{d}{2\alpha})^{-\alpha}$	$f^{-1}(k) = 2\alpha \left[\left(\frac{k}{\sigma^2}\right)^{-\frac{1}{\alpha}} - 1 \right]$
γ -exponential	$f(d) = \sigma^2 \exp(-d^{\frac{2}{\gamma}})$	$f^{-1}(k) = \left(-\ln \frac{k}{\sigma^2}\right)^{\frac{\gamma}{2}}$
Matérn	$f(d) = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\sqrt{2\nu}\sqrt{d}\right)^\nu K_\nu\left(\sqrt{2\nu}\sqrt{d}\right)$	no closed-form but solvable with root-finding algorithms

To date, we are able to find the exact estimation $\tilde{\mathbf{U}}$ given the true GP covariance kernel \mathbf{K} that is constructed from \mathbf{Z} (Eq. 7). In practice, we only have the sample covariance estimator \mathbf{S} , and neither rank- M nor positive semi-definite is guaranteed for $g(f^{-1}(\mathbf{S}))$. Therefore, we try to find its optimal rank- M positive definite approximation, i.e.

$$\underset{\tilde{\mathbf{U}} \in \mathbb{R}^{T \times M}}{\text{minimize}} \left\| g(f^{-1}(\mathbf{S})) - \tilde{\mathbf{U}}\tilde{\mathbf{U}}^\top \right\|. \quad (8)$$

Dax et al. (2014) shows that $\tilde{\mathbf{U}} = (\sqrt{\lambda_1}\mathbf{U}_{:,1}, \dots, \sqrt{\lambda_M}\mathbf{U}_{:,M})$ is the optimal solution for any unitarily invariant matrix norm $\|\cdot\|$, where $\lambda_1, \dots, \lambda_M$ are the first M largest positive eigenvalues of $g(f^{-1}(\mathbf{S}))$ and $\mathbf{U}_{:,1}, \dots, \mathbf{U}_{:,M}$ are the corresponding eigenvectors. We summarize the IKD algorithm in Alg. 1.

2.3 IKD WITH GENERAL STATIONARY KERNELS

Apart from the SE kernel, IKD also works for most commonly used stationary kernels, as long as the kernel function $f(d)$ is invertible (i.e., $f(d)$ is strictly monotonic over $[0, \infty)$) and we can find a unique non-negative solution for $d = f^{-1}(k)$. We summarize the kernels in Tab. 1.

For the SE kernel, we can generalize it to the ARD kernel $k(\mathbf{z}_i, \mathbf{z}_j) = \sigma^2 \exp\left(-\frac{1}{2} \sum_{d=1}^D \frac{1}{l_d^2} (z_{i,d} - z_{j,d})^2\right)$ and the Gaussian kernel $k(\mathbf{z}_i, \mathbf{z}_j) = \sigma^2 \exp\left(-\frac{1}{2} (\mathbf{z}_i - \mathbf{z}_j)^\top \mathbf{L}^{-1} (\mathbf{z}_i - \mathbf{z}_j)\right)$, where an extra affine transformation $\mathbf{L}^{\frac{1}{2}}$ is needed, rather than a constant scaling l .

For the Matérn kernel parameterized by ν , $K_\nu(\cdot)$ is the modified Bessel function of the second kind. Although it is complicated to obtain a closed-form of $f^{-1}(\cdot)$ for the Matérn kernel, $f^{-1}(\cdot)$ always exists since $f(\cdot)$ is strictly monotonically decreasing over $[0, \infty)$ for all $\nu > 0$. Note that for the commonly used $\nu = p + \frac{1}{2}$, $p \in \mathbb{N}$, it is easy to derive $f'(\cdot)$, e.g., when $\nu = \frac{3}{2}$, $f(d) = \sigma^2 \left(1 + \frac{\sqrt{3}d}{l}\right) \exp\left(-\frac{\sqrt{3}d}{l}\right)$, and $f'(d) = -\frac{3d\sigma^2}{l^2} \exp\left(-\frac{\sqrt{3}d}{l}\right)$. In such cases, higher-order root-finding algorithms (e.g., Newton’s method) can be used to solve $d = f^{-1}(k)$.

2.4 ERROR ANALYSIS OF IKD

IKD performs eigen-decomposition on $g(f^{-1}(\mathbf{S}))$ (Eq. 5), which uses the sample covariance \mathbf{S} as an empirical estimator of \mathbf{K} . In practice, sample covariance values $s_{i,j}$ in \mathbf{S} can be very noisy due to the noise in the data and an insufficient observation dimensionality N . There can be non-positive and

close-to-zero positive covariance values preventing from calculating $\hat{d}_{i,j} = f^{-1}(s_{i,j})$ accurately. Non-positive $s_{i,j}$ falls out of the input range of f^{-1} , i.e., $(0, \sigma^2]$. For close-to-zero positive $s_{i,j}$, the error between the estimation $\hat{d}_{i,j} = f^{-1}(s_{i,j})$ and the ground truth $d_{i,j} = f^{-1}(k_{i,j})$ can be large and sensitive to $s_{i,j}$. A sketch analysis of the error for the SE kernel is via the Taylor expansion of f^{-1} at $s_{i,j}$:

$$\begin{aligned} d_{i,j} &= f^{-1}(k_{i,j}) = -2 \ln \frac{s_{i,j} + (k_{i,j} - s_{i,j})}{\sigma^2} = -2 \ln \frac{s_{i,j}}{\sigma^2} - 2 \frac{k_{i,j} - s_{i,j}}{s_{i,j}} + O((k_{i,j} - s_{i,j})^2) \\ &= f^{-1}(s_{i,j}) + \frac{O(k_{i,j} - s_{i,j})}{s_{i,j}} = \hat{d}_{i,j} + \frac{O(k_{i,j} - s_{i,j})}{s_{i,j}}. \end{aligned} \quad (9)$$

We define the estimation error as $|d_{i,j} - \hat{d}_{i,j}| = \frac{O(|k_{i,j} - s_{i,j}|)}{s_{i,j}}$. For large $s_{i,j}$, the error is small; but for small $s_{i,j}$, the error is very sensitive to the covariance error $|k_{i,j} - s_{i,j}|$. To resolve the issue, there are two solutions:

- **Blockwise solution** We first throw away bad $s_{i,j}$ values by thresholding the sample covariance with a value s_0 , leading to a thresholded covariance matrix $\tilde{\mathbf{S}} = (s_{i,j} \cdot \mathbb{1}[s_{i,j} > s_0])_{T \times T}$. $\tilde{\mathbf{S}}$ is not a fully connected graph due to the zero values. We can not directly apply IKD to $\tilde{\mathbf{S}}$ for latent estimation. Then, we can use, for example the Bron–Kerbosch algorithm (Bron & Kerbosch, 1973), to find maximal cliques in $\tilde{\mathbf{S}}$. Consequently, each clique is a fully connected subgraph (block) of $\tilde{\mathbf{S}}$, which can be decomposed using IKD. After obtaining the latent for each clique, we merge all of the estimated latent variables according to the shared points between every clique pair. As long as the number of shared points between two cliques is greater than M , we are able to find the unique optimal rigid transformation that aligns every two cliques correctly. Although the complexity of the Bron-Kerbosch algorithm for finding maximal cliques is $\mathcal{O}(3^T)$, we can terminate the algorithm as long as the union of the existing cliques is the whole dataset. In other words, we only need up to T maximal cliques, so the clique finding time can be bounded by $\mathcal{O}(T^2)$. Then solving first M eigen-decomposition algorithm for up to T cliques takes $\mathcal{O}(T \times (MT^2))$. Therefore, the complexity of the entire procedure can be bounded by $\mathcal{O}(MT^3)$.

- **Geodesic solution** Since small values $s_{i,j} < s_0$ have significantly bad effects on eigen-decomposition, we can replace $s_{i,j}$, whose value is smaller than s_0 , with the geodesic covariance $s_{i,j} \leftarrow \max_{(t_1, t_2, \dots, t')} s_{i,t_1} \cdot s_{t_1, t_2} \cdot \dots \cdot s_{t', j}$, where $i \rightarrow t_1 \rightarrow \dots \rightarrow t' \rightarrow j$ is the geodesic path from i to j found by the Dijkstra algorithm (Dijkstra et al., 1959). The complexity of this approach is bounded by the complexity of the Dijkstra algorithm, which is $\mathcal{O}(T^2 \log(T))$. Since the complexity of the geodesic approach is smaller than that of the blockwise approach when T is larger (greater than 1000 in the following experiments), we choose the geodesic instead of the blockwise.

2.5 REFERENCE POINT SELECTION

In Eq. 3, we choose \mathbf{z}_1 as the reference point to calculate $d_{i,j}$. But the reference point can be any one in $\{\mathbf{z}_t\}_{t=1}^T$. If we choose \mathbf{z}_r , for an arbitrary index $r \in \{1, \dots, T\}$, to be the reference point, then similar to Eq. 4, the r^{th} row and the r^{th} column of $g(\mathbf{D})$ are 0s, and the remaining elements are $g(\mathbf{D})_{i,j} = \frac{1}{2}(d_{i,r} + d_{r,j} - d_{i,j}) \neq 0, \forall i \neq r, j \neq r$. Note that every $g(\mathbf{D})_{i,j}$ includes an element from $\{d_{r,i}\}_{i=1}^T$. Thus the quality of $\{d_{r,i}\}_{i=1}^T$ is vital for latent estimation. Based on the analysis in Eq. 9, we know that in practice we want to choose a good reference index r so that $\{\hat{d}_{r,i}\}_{i=1}^T$ are relatively small (i.e., $\{s_{r,i}\}_{i=1}^T$ are large, which means the r^{th} data point is highly correlated with the rest of the data points). Note that multidimensional scaling (MDS) (Kruskal & Wish, 1978) solves a similar eigen-decomposition problem, i.e., finding coordinates \mathbf{Z} from the distance matrix \mathbf{D} . It employs a centering idea which is equal to using the average of all latent variables $\frac{1}{T} \sum_{t=1}^T \mathbf{z}_t$ as the reference point. We choose the best reference point instead since we want to reduce the estimation error in Eq. 9 as much as we can so that the objective function in Eq. 8 can be minimized as much as possible. Mathematically, we obtain r such that $\|\hat{\mathbf{d}}_r\|_\infty \leq \|\hat{\mathbf{d}}_i\|_\infty, \forall i \in \{1, 2, \dots, T\}$ and $\hat{\mathbf{d}}_i$ is the i^{th} row of $\hat{\mathbf{D}} = (\hat{d}_{i,j})_{T \times T}$, i.e., $r = \arg \min_i \|\hat{\mathbf{d}}_i\|_\infty = \arg \min_i \left\{ \max_j \{\hat{d}_{i,j}\} \right\}$.

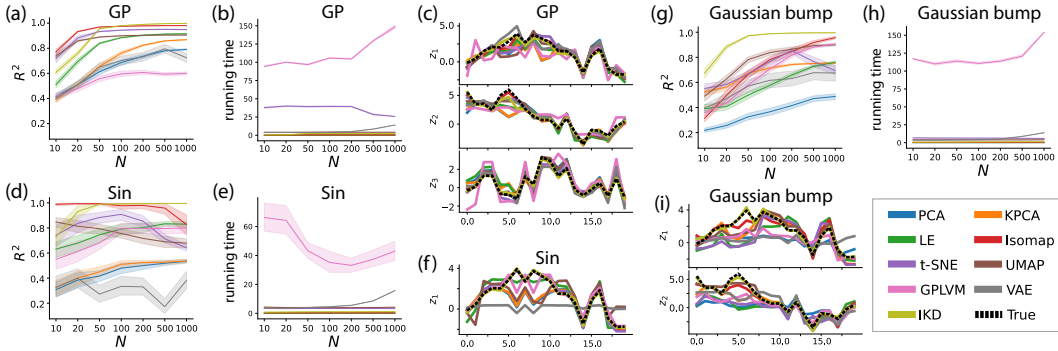


Figure 1: R^2 values (a,d,g) and running time (b,e,h) with respect to N with GP, sinusoidal, and Gaussian bump mapping functions. (c,f,i): Latent recovery visualization of an example trial with $N = 100$, for the first 20 points with GP, sinusoidal, and Gaussian bump mapping functions.

3 EXPERIMENTS

In this section, we evaluate IKD on three synthetic datasets, where we know the true latent representations, and four real-world datasets. We compare IKD with PCA, kernel PCA (KPCA), Laplacian eigenmaps (LE) (Belkin & Niyogi, 2003), Isomap (Tenenbaum et al., 2000), t-SNE, UMAP, GPLVM, and VAE. The first four are eigen-decomposition-based methods; the last four are optimization-based methods. For kernel PCA, we try different kernels (polynomial, SE, sigmoid, and cosine) and present the best one. For IKD, we use the SE kernel.

3.1 SYNTHETIC DATA

We first test all the methods on three synthetic datasets. All the following experiments are based on 50 independent repeats (trials). For each trial, we generate the true latent variables from

$$\mathbf{Z}_{m,1:T} \sim \mathcal{N}\left(\mathbf{0}, \left(6e^{-\frac{|i-j|}{5}}\right)_{T \times T}\right), \forall m \in \{1, \dots, M\}, \quad (10)$$

where M is the latent dimensionality, varying across different datasets. Then, we generate the noiseless data from GP, sinusoidal, and Gaussian bump mapping functions respectively. Afterwards, i.i.d. Gaussian noise is added to form the final noisy observations \mathbf{X} . We evaluate the performance using the R^2 metric. When computing R^2 values, we first align the estimated latent with the ground truth through an affine transformation; then compute R^2 for each latent dimension, and finally take an average across all latent dimensions $m \in \{1, 2, \dots, M\}$. The reasons for choosing the affine transformation are: (1) rigid transformation could lead to very negative R^2 values for those non-IKD methods (e.g., PCA), not shown here; and (2) affine transformation is the commonly used one for latent estimation and alignment.

GP mapping function We start our experiments from the GP mapping function. In each trial, we generate a 3D latent $\mathbf{Z} \in \mathbb{R}^{1000 \times 3}$ (i.e., $M = 3$) according to Eq. 10, and generate $\mathbf{X} \in \mathbb{R}^{1000 \times N}$ according to Eq. 1 with $\sigma^2 = 1$ and $l = 3$. Then Gaussian noise is added:

$$x_{t,n} \leftarrow x_{t,n} + \varepsilon_{t,n}, \forall (t, n) \in \{1, \dots, 1000\} \times \{1, \dots, N\}, \quad (11)$$

where noise $\varepsilon_{t,n} \sim \mathcal{N}(0, 0.05^2)$. Note that this generating process is consistent with the generating process of GPLVM. Thus it is well aligned with the model assumptions of IKD, deemed as a data-matching example. Fig. 1(a) shows that for $N = 10$ and $N = 20$, Isomap is the best; but when $N > 50$, IKD becomes the best and its R^2 converges to 1 as N increases. The latent recovery visualization of an example trial under $N = 100$ for the first 20 points (Fig. 1(c)) shows that Isomap and IKD match the true latent the best.

Sinusoidal mapping function In each trial, we generate a 1D latent (i.e., $M = 1$) according to Eq. 10, and generate the noisy observations $\mathbf{X} \in \mathbb{R}^{1000 \times N}$ as

$$\mathbf{x}_t = \sin(\boldsymbol{\Omega} \mathbf{z}_t + \varphi) + \varepsilon_t, \forall t \in \{1, \dots, 1000\}, \quad (12)$$

where $\Omega = (\omega_{n,m})_{N \times M}$ with $\omega_{n,m} \sim \mathcal{U}(-1, 1)$, $\varphi = [\varphi_1, \dots, \varphi_N]^T$ with $\varphi_n \sim \mathcal{U}(-\pi, \pi)$, and noise $\varepsilon_t \sim \mathcal{N}(\mathbf{0}, 0.1^2 \mathbf{I})$. The result in Fig. 1(d) indicates that even the observed data is not from a GP (data-mismatching), IKD is still able to discover the latent structure consistently better than others, except for $N = 10$ and $N = 20$ where Isomap is the best. When the observation dimensionality $N > 50$, the R^2 value of IKD approaches to 1, while Isomap, t-SNE, and UMAP all have decreasing performance due to the curse of dimensionality. The latent recovery visualization of an example trial under $N = 100$ for the first 20 points (Fig. 1(f)) shows that Isomap and IKD match the true latent the best.

Gaussian bump mapping function In each trial, we generate a 2D latent (i.e., $M = 2$) according to Eq. 10, and generate $\mathbf{X} \in \mathbb{R}^{1000 \times N}$ as

$$x_{t,n} = 20 \exp(-\|z_t - c_n\|_2^2) + \varepsilon_{t,n}, \forall (t, n) \in \{1, \dots, 1000\} \times \{1, \dots, N\}, \quad (13)$$

where $c_n \in \mathbb{R}^2$ is the center of the n^{th} Gaussian bump randomly selected from 10,000 grid points uniformly distributed in $[-6, 6]^2$, with noise $\varepsilon_{t,n} \sim \mathcal{N}(0, 0.05^2)$. This is another data-mismatching example. Fig. 1(g) shows that in this case, IKD is the best one among all methods for all observation dimensionality N . Fig. 1(i) also shows that only IKD matches the true latent accurately.

The running time of IKD in the three synthetic datasets above is on par with other eigen-decomposition-based methods (Fig.1(b,e,h)), and much less than optimization-based methods. In particular, GPLVM always takes a very long time; t-SNE requires more running time on datasets whose latent dimensionality is greater than 2; and VAE is more time-consuming when the observation dimensionality is large.

Note that we only vary the dimensionality N not the number of observations T . When increasing T , the running time of all methods will increase polynomially. For optimization-based methods, stochastic optimization can be employed to scale to large-scale datasets. For fair comparison, extra scaling techniques should be incorporated to deal with large-scale eigen-decomposition, which falls out of the scope of this paper.

In general, IKD performs the best for all three mapping functions especially when the observation dimensionality N is large. It is also very effective in capturing details in addition to recovering the general latent structure correctly. Same as other eigen-decomposition-based methods, IKD takes less time in solving the presented problems compared with optimization-based methods.

Varying dimensionality, kernels and latent structures We test the effectiveness of IKD on the observation data generated from the GP mapping function described in Eq. 1, for different observation dimensionality $N \in \{100, 200, 500, 1000, 2000, 5000, 10000\}$, different generating kernels (Tab. 1), and three different latent structures (hard, medium, and easy shown in Fig. 2(a) and Fig. 5 in Appendix according to their difficulty levels). We only compare IKD with PCA and GPLVM here. PCA is the most commonly used linear method, so it serves as a baseline. GPLVM is the traditional optimization-based model for data generated from the GP mapping function. From Fig. 2(b), we can tell that IKD is always the best for the most commonly used SE kernel. Compared with PCA and GPLVM, IKD is highly effective especially when N is large, where the R^2 values of IKD are very close to 1. Fig. 2(c) shows the latent recovery visualization from one example trial of the medium dataset when $N = 1000$. The kernel of the generating model is SE. We can tell that IKD matches the ground truth the best. For the third dimension, particularly, only the estimated latent from IKD reflects the linear increasing trend of the ground truth correctly. In terms of complexity, GPLVM is time-consuming compared with IKD (Fig. 2(d)). These results indicate that we can use IKD to recover the latent for data generated from the GP mapping function faster and more accurately.

Kernel mismatch In real-world applications, we do not have information about the kernel. Here, we conduct a kernel mismatching experiment to test if IKD is still able to recover an acceptable latent without knowing the generating kernel. Specifically, one kernel is used for generating the true covariance matrix \mathbf{K} based on the true latent \mathbf{Z} ; and another kernel is used for latent estimation. Fig. 6 (in Appendix) shows that IKD outperforms GPLVM in most cases under kernel mismatch conditions, implying a good generalization performance over different generating kernels.

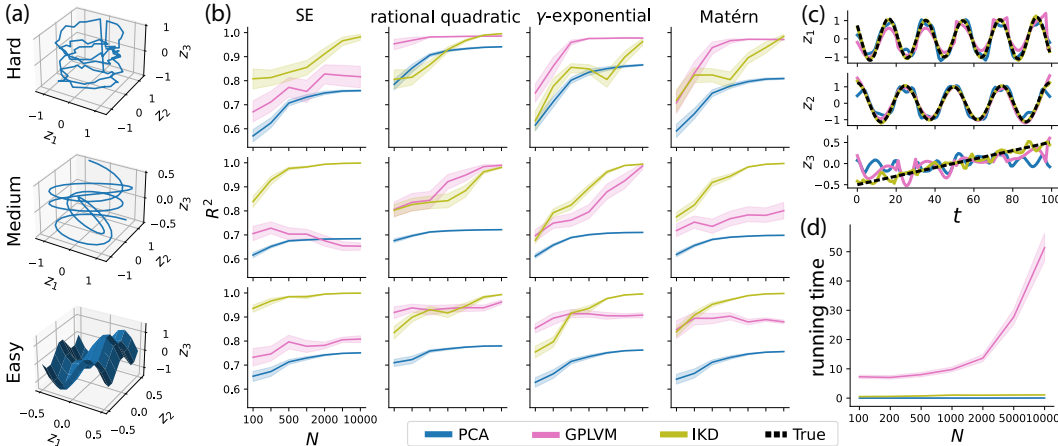


Figure 2: (a): Three latent datasets with different difficulty levels. (b): R^2 values, with respect to N , of PCA, GPLVM, and IKD for different datasets and kernels. (c): Latent recovery visualization of an example trial of the medium dataset with the SE kernel and $N = 1000$. (d): Average running time in seconds (across different kernels, different datasets, and 50 independent trials) of the three methods w.r.t. observation dimensionality N . We take averages across different kernels and different datasets because all of them share similar running time results.

3.2 REAL-WORLD DATA

We compare IKD against alternatives on four real-world datasets:

- Single-cell qPCR (PRC) (Guo et al., 2010): Normalized measurements of 48 genes of a single cell at 10 different stages. There are 437 data points in total, resulting in $\mathbf{X} \in \mathbb{R}^{437 \times 48}$.
- Hand written digits (digits) (Dua & Graff, 2017): It consists 1797 grayscale images of hand written digits. Each one is an 8×8 image, resulting in $\mathbf{X} \in \mathbb{R}^{1797 \times 64}$.
- COIL-20 (Nene et al., 1996): It consists 1440 grayscale photos. For each one of the 20 objects in total, 72 photos were taken from different angles. Each one is a 128×128 image, resulting in $\mathbf{X} \in \mathbb{R}^{1440 \times 16384}$.
- Fashion MNIST (F-MNIST) (Xiao et al., 2017): It consists 70000 grayscale images of 10 fashion items (clothing, bags, etc). We use a subset of it, resulting in $\mathbf{X} \in \mathbb{R}^{3000 \times 784}$.

Since there is no true latent to compare against, we first estimate the latent in a $\{2, 3, 5, 10\}$ -dimensional latent space and then use the k -nearest neighbor (k -NN) classifier to evaluate the performance of each dimensionality reduction method. Specifically, we apply 5-fold cross-validation k -NN ($k \in \{5, 10, 20\}$) on the estimated $\{2, 3, 5, 10\}$ -dimensional latent to evaluate the performance of each method on each dataset. The k -NN classification results of different methods under different latent dimensionality M , different datasets, and different choices of k are shown in Fig. 3(a).

Comparing IKD with other eigen-decomposition-based methods (PCA, KPCA, LE, Isomap), we can conclude that IKD is almost always the best one on all four datasets, except that when $M \in \{3, 5, 10\}$ in the digits dataset, Isomap is better than IKD. When comparing IKD with GPLVM, we find the performances of GPLVM on PRC, digits, and F-MNIST datasets are slightly better than IKD while GPLVM takes too much running time. Specifically, IKD is significantly better than GPLVM on the COIL-20 dataset but only slightly worse than GPLVM on the other three datasets. VAE only performs well on the most complicated dataset—F-MNIST, and is much worse than IKD in the other three datasets. Although IKD is worse than the remaining two optimization-based methods (t-SNE and UMAP), the performance of IKD is the best on the COIL-20 dataset. The reason is that the observation dimensionality is very large ($N = 16384$) in the COIL-20 dataset, and IKD is very effective for high-dimensional data as shown in the synthetic results (in Fig. 1(a)). A 2D visualization of the digits dataset is shown in Fig. 4. 2D visualizations of the other three datasets are shown in Fig. 7, 8, and 9 in the Appendix. Qualitatively, we can see that IKD consistently finds more separate clusters compared with all other eigen-decomposition-based methods

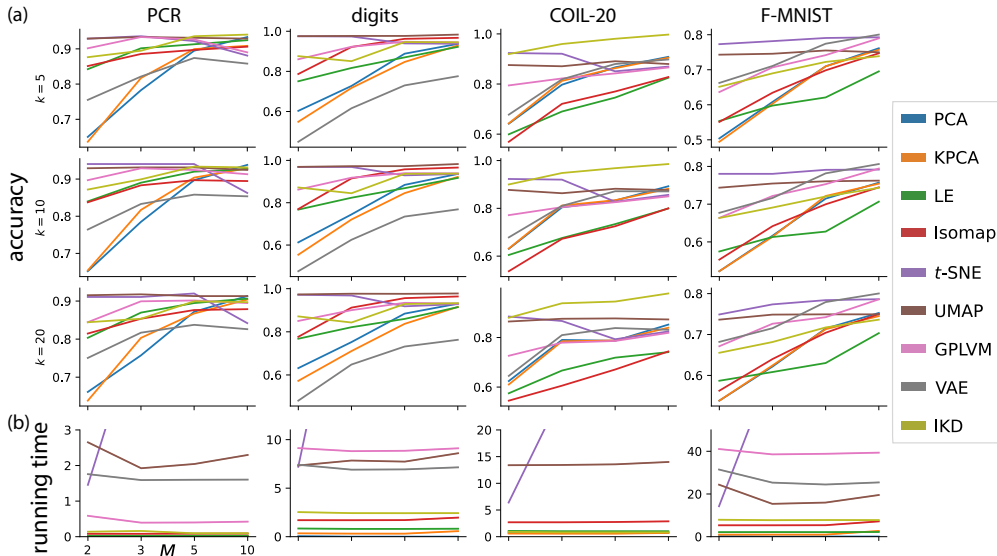


Figure 3: (a) k -NN 5-fold cross validation and (b) running time, on different methods, different latent dimensionality M , different datasets, and different choices of k .

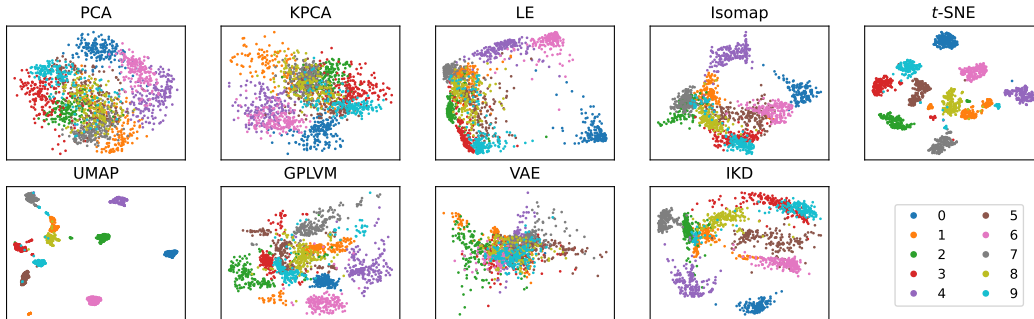


Figure 4: Visualization of the dimensionality reduction results of different methods on the digits dataset.

and two optimization-based methods (GPLVM and VAE) across all four datasets. Therefore, even though IKD is an eigen-decomposition-based method, its performance is significantly better than other eigen-decomposition methods on all four real-world datasets, sometimes as good as the best optimization-based method.

In terms of running time (Fig. 3(b)), IKD is on par with Isomap, and these eigen-decomposition-based methods are significantly faster than those four optimization-based methods. Note that if the desired latent dimensionality $M > 2$, the running time of t-SNE is barely acceptable. For the high dimensional COIL-20, the running time values of VAE and GPLVM are extremely high, getting out of the upper limit of the corresponding axes.

In summary, IKD, as an eigen-decomposition-based method, consumes short running time, but is able to obtain dimensionality reduction results better than other eigen-decomposition-based methods. When facing high-dimensional observation data, IKD can perform significantly better than all other methods in a very short time.

Note that although eigen-decomposition-based methods perform relatively worse than optimization-based methods, the benefit of fast running provides good initialization for sophisticated nonlinear optimization problems, mitigating the numerical instability and multi-modal issues commonly observed in methods such as GPLVM and VAE.

REFERENCES

- Mayur R Bakrania, I Jonathan Rae, Andrew P Walsh, Daniel Verscharen, and Andy W Smith. Using dimensionality reduction and clustering techniques to classify space plasma regimes. *Frontiers in Astronomy and Space Sciences*, pp. 80, 2020.
- Mikhail Belkin and Partha Niyogi. Laplacian Eigenmaps for Dimensionality Reduction and Data Representation. *Neural Computation*, 15(6):1373–1396, June 2003. ISSN 0899-7667. doi: 10.1162/089976603321780317. Conference Name: Neural Computation.
- Coen Bron and Joep Kerbosch. Algorithm 457: finding all cliques of an undirected graph. *Communications of the ACM*, 16(9):575–577, September 1973. ISSN 0001-0782, 1557-7317. doi: 10.1145/362342.362367. URL <https://dl.acm.org/doi/10.1145/362342.362367>.
- Thang D Bui and Richard E Turner. Stochastic variational inference for Gaussian process latent variable models using back constraints. In *Black Box Learning and Inference NIPS workshop*, 2015.
- Achiya Dax et al. Low-rank positive approximants of symmetric matrices. *Advances in Linear Algebra & Matrix Theory*, 4(03):172, 2014.
- Edsger W Dijkstra et al. A note on two problems in connexion with graphs. *Numerische mathematik*, 1(1):269–271, 1959.
- Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Guoji Guo, Mikael Huss, Guo Qing Tong, Chaoyang Wang, Li Li Sun, Neil D Clarke, and Paul Robson. Resolution of cell fate decisions revealed by single-cell gene expression analysis from zygote to blastocyst. *Developmental cell*, 18(4):675–685, 2010.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Mark A. Kramer. Nonlinear principal component analysis using autoassociative neural networks. *AIChE Journal*, 37(2):233–243, 1991. ISSN 1547-5905. doi: 10.1002/aic.690370209. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/aic.690370209>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/aic.690370209>.
- Joseph B Kruskal and Myron Wish. *Multidimensional scaling*. Number 11. Sage, 1978.
- Neil Lawrence. Gaussian Process Latent Variable Models for Visualisation of High Dimensional Data. In *Advances in Neural Information Processing Systems*, volume 16. MIT Press, 2003. URL <https://proceedings.neurips.cc/paper/2003/hash/9657c1fffd38824e5ab0472e022e577e-Abstract.html>.
- Neil Lawrence. Probabilistic Non-linear Principal Component Analysis with Gaussian Process Latent Variable Models. *Journal of machine learning research*, 6(11), 2005.
- Leland McInnes, John Healy, and James Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv:1802.03426 [cs, stat]*, September 2020. URL <http://arxiv.org/abs/1802.03426>. arXiv: 1802.03426.
- Sameer A Nene, Shree K Nayar, Hiroshi Murase, et al. Columbia object image library (coil-100). 1996.
- Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Kernel principal component analysis. In Wulfram Gerstner, Alain Germond, Martin Hasler, and Jean-Daniel Nicoud (eds.), *Artificial Neural Networks — ICANN’97*, Lecture Notes in Computer Science, pp. 583–588, Berlin, Heidelberg, 1997. Springer. ISBN 978-3-540-69620-9. doi: 10.1007/BFb0020217.
- Ehsan Sheybani and Giti Javidi. Dimensionality reduction and noise removal in wireless sensor network datasets. In *2009 Second International Conference on Computer and Electrical Engineering*, volume 2, pp. 674–677. IEEE, 2009.

- Joshua B Tenenbaum, Vin de Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- R. Urtasun, D.J. Fleet, and P. Fua. 3D People Tracking with Gaussian Process Dynamical Models. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 1, pp. 238–245, June 2006. doi: 10.1109/CVPR.2006.15. ISSN: 1063-6919.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008a.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11), 2008b.
- Jack Wang, Aaron Hertzmann, and David J Fleet. Gaussian Process Dynamical Models. In *Advances in Neural Information Processing Systems*, volume 18. MIT Press, 2005. URL <https://proceedings.neurips.cc/paper/2005/hash/ccd45007df44dd0f12098f486e7e8a0f-Abstract.html>.
- Jack M. Wang, David J. Fleet, and Aaron Hertzmann. Gaussian Process Dynamical Models for Human Motion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(2):283–298, 2008. ISSN 1939-3539. doi: 10.1109/TPAMI.2007.1167. Conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence.
- Anqi Wu, Nicholas A. Roy, Stephen Keeley, and Jonathan W Pillow. Gaussian process based nonlinear latent structure discovery in multivariate spike train data. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/b3b4d2dbedc99fe843fd3dedb02f086f-Abstract.html>.
- Anqi Wu, Stan Pashkovski, Sandeep R Datta, and Jonathan W Pillow. Learning a latent manifold of odor representations from neural responses in piriform cortex. *Advances in Neural Information Processing Systems*, 31, 2018.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.