

# S-Isomap++: Multi Manifold Learning from Streaming Data

Suchismit Mahapatra

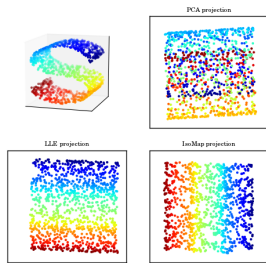
Department of Computer Science



# Motivation

## Massive amounts of data

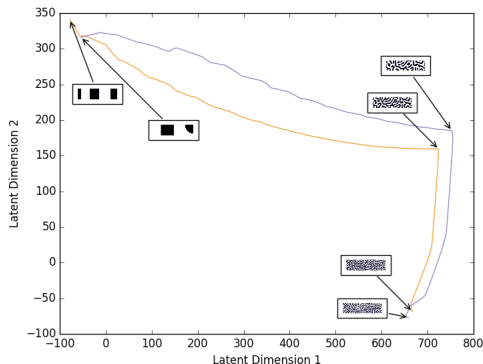
- Natural data tends to be generated by systems (physical or non-physical) that have **very few** degrees of underlying freedom.
- Real-world data is typically a result of **complex non-linear processes**, but can often be described by a **low-dimensional manifold**.



[Credit: Raymond Fu]

# Motivation

## Nonlinear Process Dynamics



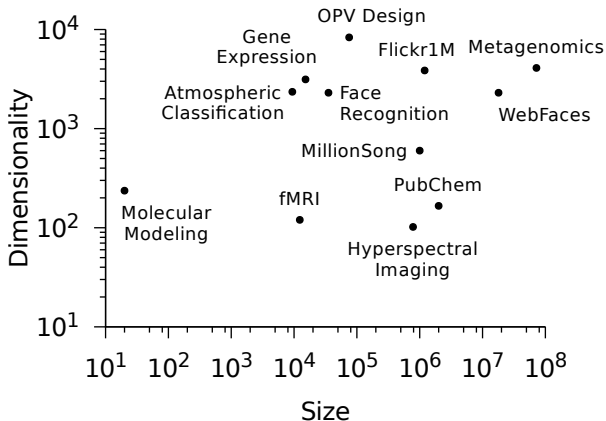
Morphological parametric trajectories for a nonlinear process.

[Click here for [simulation of all parametric trajectories](#)][Click here for [simulation of Manifold](#)]



# Motivation

## Massive amounts of data



Topology of high-dimensional, massive datasets

# Learning efficiently

## Common Approaches

- Smoothness
  - Try to learn functions that are **smooth**.
  - Examples - Spline based techniques, Kernel methods,  $L_2$ -regularization, etc.
- Sparsity
  - Represent in terms of **sparse/few** basis functions.
  - Examples - Lasso, Compressive Sensing, Wavelets
- Geometry
  - Data distribution is **not uniform**, try to **exploit geometry**.
  - Examples - Laplacian based techniques, Manifold learning

Even more **relevant** in high-dimensional spaces.

# Manifold Learning

## Assumptions

- Distribution of data **not uniform**.
- Data **lives on/near** some low-dimensional manifold, typically **embedded** in high dimensions and **separated** by **low-density regions**.
- Typically used as a generic **non-linear, non-parametric technique** to **approximate** probability distributions in high-dimensional spaces.

# Manifold

## Properties

### Definition

A manifold  $\mathcal{M}$  is a metric space with the following property: if  $x \in \mathcal{M}$ , then there exists some neighborhood  $\mathcal{U}$  of  $x$  and  $\exists n$  such that  $\mathcal{U}$  is homeomorphic to  $\mathbb{R}^n$ .

# Manifold

## Properties

### Definition

A manifold  $\mathcal{M}$  is a metric space with the following property: if  $x \in \mathcal{M}$ , then there exists some neighborhood  $\mathcal{U}$  of  $x$  and  $\exists n$  such that  $\mathcal{U}$  is homeomorphic to  $\mathbb{R}^n$ .

- Global structure can be more complicated.
- Usually embedded in high dimensional spaces, but the intrinsic dimensionality is typically low due to fewer degrees of freedom.
- Examples
  - Collection of news articles
  - Image data sets
  - State space of MDP's



# Manifold

Caltech 101 Dataset



[Credit: <https://lvdmaaten.github.io/tsne/>]

# Nonlinear Spectral Dimension Reduction

## Formulation

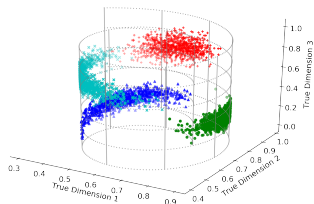
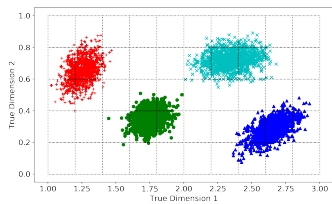
### Definition

Given  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T$ , where  $\forall \mathbf{x}_i \in \mathbb{R}^D$ , the task is to **find** a corresponding **low-dimensional** representation,  $\mathbf{y}_i \in \mathbb{R}^d$ , for each  $\mathbf{x}_i$ , where  $d \ll D$ .

- We assume there **exists**  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$  that **maps** each data sample  $\mathbf{y}_i \in \mathbb{R}^d$  to  $\mathbf{x}_i \in \mathbb{R}^D$ .
- The goal is to **learn the inverse mapping**,  $\phi^{-1}$ , that can be used to map high-dimensional  $\mathbf{x}_i$  to low-dimensional  $\mathbf{y}_i$ , i.e.  $\mathbf{y}_i = \phi^{-1}(\mathbf{x}_i)$ .

# Nonlinear Spectral Dimension Reduction

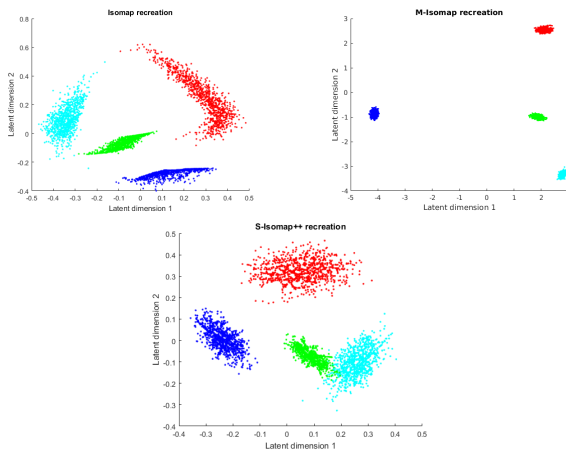
## Illustration



Typical real world scenario wherein we need to **learn the inverse mapping**,  $\phi^{-1}$ , to be able to uncover the intrinsic low-dimensional representation from high-dimensional data.

# Nonlinear Spectral Dimension Reduction

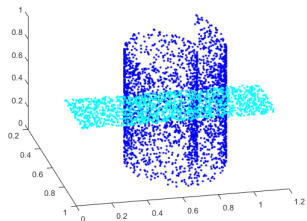
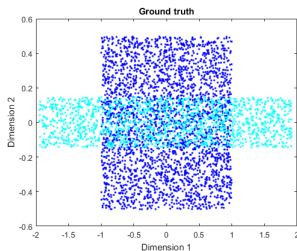
## Illustration



How well different algorithms could **recreate the latent ground truth** used to generate the high-dimensional data.

# Nonlinear Spectral Dimension Reduction

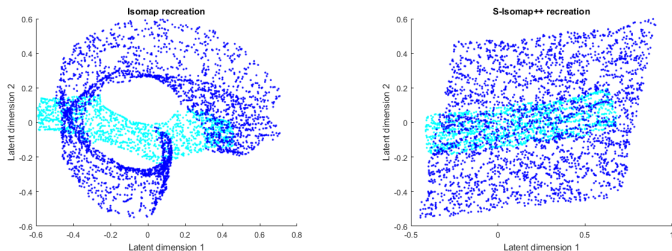
## Illustration



Multiple manifolds typically involve **dissimilar mappings**  $\{\phi_i\}_{i=1,2,\dots,p}$  projecting the intrinsic low-dimensional representation to higher dimensional real-world data.

# Nonlinear Spectral Dimension Reduction

## Illustration



In an ideal scenario, when manifolds are **densely sampled** and **sufficiently separated**, existing NLSDR methods can uncover individual manifolds. But **intersecting** manifolds are still a challenge.

# S-Isomap++ algorithm

## Introduction

The algorithm takes in as input, the batch and streaming data sets,  $\mathcal{B}$  and  $\mathcal{S}$  respectively and can be divided into two main phases:

- Batch processing phase
  - **Cluster** samples in  $\mathcal{B}$  into  $p$  clusters.
  - **Learn** individual manifolds corresponding to each cluster, and **map** samples from each cluster to its low-dimensional representation.
  - **Map** low-dimensional samples from individual manifolds into a global space.
- Stream mapping phase
  - **Map** each sample  $s$  from  $\mathcal{S}$  onto each of the  $p$  manifolds by **matching their inner products** to the computed geodesic distances with the  $k$  nearest neighbors, to determine which manifold  $s$  belongs to.

# S-Isomap++

## Batch Processing phase

- 1:  $\mathcal{C}_{i=1,2,\dots,p} \leftarrow \text{Find\_Clusters}(\mathcal{B}, \epsilon)$
- 2:  $\xi_S \leftarrow \emptyset$
- 3: **for**  $1 \leq i \leq p$  **do**
- 4:    $\mathcal{LDE}_i \leftarrow \text{Isomap}(\mathcal{C}_i)$
- 5: **end for**
- 6:  $\xi_S \leftarrow \bigcup_{i=1}^p \bigcup_{j=i+1}^p \text{NN}(\mathcal{C}_i, \mathcal{C}_j, \mathbf{k}) \cup \text{FN}(\mathcal{C}_i, \mathcal{C}_j, \mathbf{l})$
- 7:  $\mathcal{GE}_S \leftarrow \text{MDS}(\xi_S)$
- 8: **for**  $1 \leq j \leq p$  **do**
- 9:    $\mathcal{I} \leftarrow \xi_S \cap \mathcal{C}_j$
- 10:    $\mathcal{A} \leftarrow \begin{bmatrix} \mathcal{LDE}_j^{\mathcal{I}} \\ \mathbf{e}^T \end{bmatrix}$
- 11:    $\mathcal{R}_j, \mathbf{t}_j \leftarrow \mathcal{GE}_{\mathcal{I},S} \times \mathcal{A}^T (\mathcal{A}\mathcal{A}^T + \lambda \mathbf{I})^{-1}$
- 12: **end for**



# S-Isomap++

## Tangent Manifold Clustering

- Multiscale SVD (M-SVD) allows us to **estimate** the **intrinsic dimension** of noisy, high-dimensional point clouds.
- M-SVD estimates the intrinsic dimension by **computing singular values**  $\sigma_{i \in \{1, 2, \dots, D\}}^{x, r}$  of  $\mathcal{B}(x, r)$ ,  $\forall x \in \mathcal{M}$ , at different **scales**  $r > 0$ .

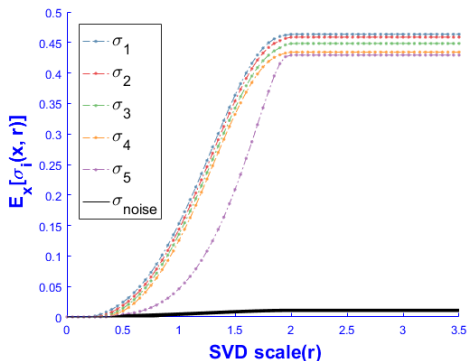
# S-Isomap++

## Tangent Manifold Clustering

- Multiscale SVD (M-SVD) allows us to **estimate** the **intrinsic dimension** of noisy, high-dimensional point clouds.
- M-SVD estimates the intrinsic dimension by **computing singular values**  $\sigma_{i \in \{1,2,\dots,D\}}^{x,r}$  of  $\mathcal{B}(x,r)$ ,  $\forall x \in \mathcal{M}$ , at different **scales**  $r > 0$ .
- Small  $r$  leads to **not enough samples** in  $\mathcal{B}(x,r)$ .
- Large  $r$  leads to **curvature** making the process **over estimate** the intrinsic dimension.
- True  $\{\sigma_i^{x,r}\}$  separate from the noise  $\{\sigma_i^{x,r}\}$  at the **right scale**, due to their **different rates of growth** and the intrinsic dimension of  $\mathcal{M}$  gets revealed.

# S-Isomap++

## Tangent Manifold Clustering



How  $\{\sigma_i^{X,r}\}$  behave over different scales when M-SVD is done on a noisy  $\mathbb{R}^5$  sphere embedded in  $\mathbb{R}^{100}$  ambient space. Notice how the noise dimensions decay out, leaving only the primary components at the appropriate scale.

# S-Isomap++

## Tangent Manifold Clustering

- Executing M-SVD on the local neighborhood of  $\forall \mathbf{x}_i \in \mathcal{B}$ , allows us to determine basis vectors,  $\mathbf{t}_{i1}, \mathbf{t}_{i2}, \dots, \mathbf{t}_{id'}$ , which define the tangent plane,  $\mathcal{T}_i$ .
- To determine the similarity between tangent planes  $\mathcal{T}_i$  and  $\mathcal{T}_j$ , we tried the following techniques, including two novel approaches :
  - **Gunawan's approach :**  

$$\phi(\mathcal{T}_i, \mathcal{T}_j) = \cos \theta = |\det(\mathcal{N})|, \text{ where } \mathcal{N}_{x,y} = \mathcal{T}_{ix}^T \mathcal{T}_{jy}$$
  - **$L_1$ -norm based metric :**  

$$\phi(\mathcal{T}_i, \mathcal{T}_j) = \frac{1}{k} \sum_{l=1}^k |\mathbf{t}_{il}^T \mathbf{t}_{jl}|$$
  - **$L_2$ -norm based metric :**  

$$\phi(\mathcal{T}_i, \mathcal{T}_j) = \sqrt{\frac{1}{k} \sum_{l=1}^k (\mathbf{t}_{il}^T \mathbf{t}_{jl})^2}$$

# S-Isomap++

## Tangent Manifold Clustering

- **Incremental** in nature.
- **Initially all points**  $\forall \mathbf{x}_i \in \mathcal{B}$  are **unlabelled**.
- An **unlabelled random** point  $\mathbf{x}_k$  is picked and is **labelled** as  $l_k$ , the next available label index.
- Subsequently, **similarity** of  $\mathbf{x}_k$  with all unlabelled  $x \in \mathcal{N}(\mathbf{x}_k)$  is evaluated. If similarity **exceeds certain threshold** i.e.  $\cos \theta \geq \epsilon_{thres}$ , points in  $\mathcal{N}(\mathbf{x}_k)$  also get **labelled** as  $l_k$ .
- **Repeat** above, till all points are labelled.

# S-Isomap++

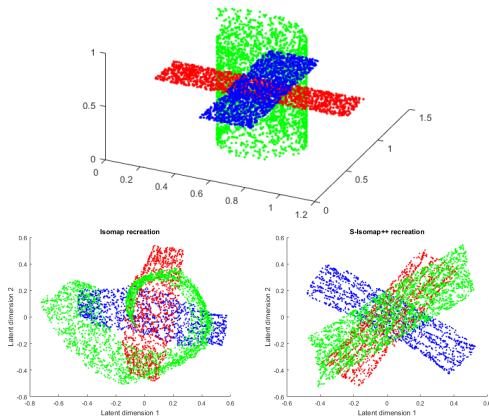
## Stream Mapping phase

- 1: **for**  $s \in \mathcal{S}$  **do**
- 2:   **for**  $1 \leq i \leq p$  **do**
- 3:      $y_s^i \leftarrow \text{S-Isomap}(s, C_i)$
- 4:      $\mathcal{G}\mathcal{E}_s^i \leftarrow \mathcal{R}_i y_s^i + t_i$
- 5:   **end for**
- 6: **end for**
- 7:  $\text{index} \leftarrow \underset{i}{\text{argmin}} \left| y_s^i - \mu(C_i, \mathcal{R}_i, t_i) \right|$
- 8:  $\mathcal{Y}_S \leftarrow \mathcal{Y}_S \cup y_s^{\text{index}}$
- 9: **return**  $\mathcal{Y}_S$

S-Isomap( $\cdot$ ) **maps** points  $s \in \mathcal{S}$  by **matching their inner products** with  $LDE_{C_i}$  to the computed geodesic distances with the  $k$  nearest neighbors of  $s$ .

# S-Isomap++

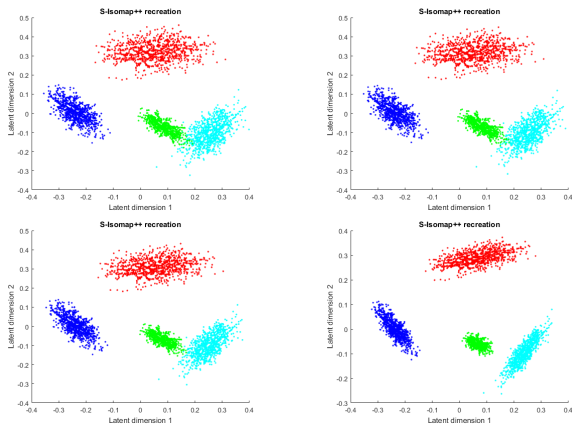
Multiple planes through swiss-roll



Top: Actual manifolds in  $\mathbb{R}^3$  space, clustered for demonstration, Bottom Left: Recreation by Isomap/M-Isomap, Bottom Row: Recreation by S-Isomap++.

# S-Isomap++

Effect of varying parameter  $\lambda$

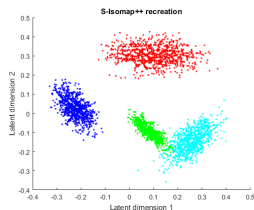
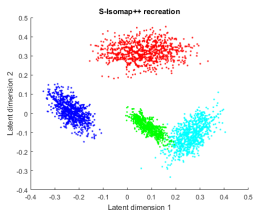
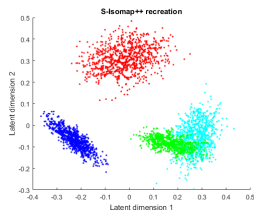
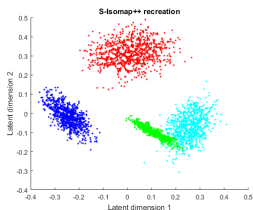


Top Left:  $\lambda = 0.01$ , Top Right:  $\lambda = 0.02$ , Bottom Left:  $\lambda = 0.04$ , Bottom Right:  $\lambda = 0.16$



# S-Isomap++

## Effect of varying parameter $k$



Top Left:  $k = 8$ , Top Right:  $k = 16$ , Bottom Left:  $k = 24$ , Bottom Right:  $k = 32$

# S-Isomap++

## Additional results

<i>Method</i>	<i>L-1</i>	<i>L-2</i>	<i>Gunawan</i>
Sphere-Sphere	<b>0.825</b>	0.619	0.5
Sphere-Plane	<b>0.759</b>	0.602	0.5
Swiss Roll-Plane	<b>0.838</b>	0.621	0.5

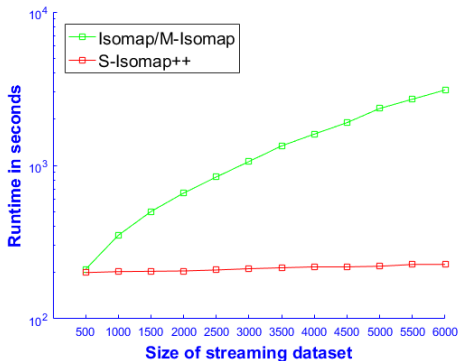
Accuracy scores for the different tangent manifold clustering approaches.

<i>digit '0'</i>	<b>0.0296</b>	<i>digit '3'</i>	<b>0.0364</b>	<i>digit '6'</i>	<b>0.0476</b>
<i>digit '1'</i>	<b>0.0806</b>	<i>digit '4'</i>	<b>0.0586</b>	<i>digit '8'</i>	<b>0.0712</b>
<i>digit '2'</i>	<b>0.0499</b>	<i>digit '5'</i>	<b>0.0449</b>	<i>digit '9'</i>	<b>0.0498</b>

Procrustes error values for different digits of MNIST, computed by comparing the original with 3-D recreation via S-Isomap++.

# S-Isomap++

## Scalability



The results are in log scale and demonstrate the scalability of our proposed algorithm.

## Summary & Future work

- The proposed algorithm allows for **scalable** non-linear dimensionality reduction of **streaming high-dimensional data**.
- By allowing for the samples to belong to **multiple** manifolds, or sampled **non-uniformly** in a single manifold, our approach can be applied to a **wide variety** of **practical** settings.
- The ability to **cluster** data lying on **multiple intersecting** manifolds is **significant** since it allows us to **automatically** identify the number of **underlying** manifolds.
- Our algorithm **assumes** that all manifolds are represented in the batch data set. This means that a novel manifold which might **appear** subsequently in the stream  $\mathcal{S}$ , does not get learned. We plan to **resolve** this limitation in our **future work**.