

# A Fast Manifold Learning Algorithm for Dimensionality Reduction

Yu Liang<sup>1</sup>, Furao Shen<sup>1,\*</sup>, Jinxi Zhao<sup>1</sup>, Yi Yang<sup>1</sup>

<sup>1</sup>National Key Laboratory for Novel Software Technology, and Department  
of Computer Science and Technology, Nanjing University, China

Email: mg1533023@smail.nju.edu.cn, frshen@nju.edu.cn, jxzhao@nju.edu.cn, yangyi868@gmail.com

**Abstract**—This paper proposes a new manifold learning method called “Soinnmanifold”. Traditional manifold learning method needs a lot of computation and appropriate priori parameters. This has somewhat restricted the domains in which manifold learning can potentially be applied. However, with the high-dimensional inputs, our method can generate a low-dimensional manifold in the high-dimensional space and determine the intrinsic dimension automatically. Then we will use this manifold to do dimensionality reduction quickly. Experiments demonstrate that our method can get promising results with less time and memory.

**Index Terms**—Manifold learning, Intrinsic Dimension, Dimensionality Reduction

## I. INTRODUCTION

With the coming of the era “Big Data”, the growing dimension of the data has led to a lot of problems such as curse of dimensionality. Therefore, how to represent the high-dimensional data in a lower-dimensional space becomes more and more important. As a consequence, dimensionality reduction has attracted many researchers in recent years.

Up to now, many classical algorithms have been proposed. For dimensionality reduction, Principal Component Analysis (PCA) can be defined as the orthogonal projection of the data onto a lower dimensional linear space, known as the principal subspace, such that the variance of the projected data is maximized [1]. Besides, MDS (Multidimensional Scaling) is to minimize the difference between corresponding distances in high and mapped low dimensional space [2]. PCA and MDS are both linear methods, they can not work well on nonlinear data. Hence, kernel method was proposed. Kernel method is based on the idea that linearly inseparable problem can be linearly separable in high dimensional space [3]. What's more, scientists have proposed a new method which is called manifold learning. Manifold learning holds the idea that the original high dimensional data may lie on a low dimensional manifold.

Isomap [4] is the first manifold learning method which uses geodesic distance to measure the similarity between data in high dimensional space [5]. Locally Linear Embedding (LLE) is another powerful algorithm in manifold learning. LLE supposes that each point can be represented by a linear combination of its several neighbor points. It computes the local combination weight matrix  $W$  in the original space

firstly. Then it minimizes the reconstructed error in the new space, in which each point is reconstructed by  $W$ . [6].

However, LLE adopts k-nearest neighbors method to generate the manifold. Choosing a very small neighborhood is not a satisfactory solution, which may fragment the manifold into a large number of disconnected regions [7]. Choosing a too large neighborhood may cause “short-circuit”. Besides, target dimension also has an effect on the performance of LLE. A too small dimension can not preserve the structure and relationships in the high-dimensional space after data are mapped into a low-dimensional space.

In this paper, we propose a novel method called “Soinnmanifold”. It can generate the manifold with the given input automatically. With the manifold, we can find the intrinsic dimension which is implicit in the input data without any parameters. For the reason that the manifold consists of a few landmarks, we can do dimensionality reduction on the manifold, which can save enormous amounts of time and memory. Afterwards, input data can be represented in the low-dimension space which is the result we want to get. Compared with other algorithms such as PCA, LLE, Isomap, experiments of Soinnmanifold demonstrate its promising results on swiss-roll data set, swissshole data set and Adult data set.

## II. DIMENSIONALITY REDUCTION

Our method is composed of three components: manifold generation, intrinsic dimension estimation and dimensionality reduction. We defer the these components and implementation details to subsequent sections.

### A. Manifold Generation

We adopt a self-organizing neural network (SOINN) algorithm to generate the manifold [8] [9]. SOINN can generate the appropriate manifold without the parameter  $k$ , which is the number of the neighbors.

#### 1) Initialization:

- Initialize the node set  $A$  to contain two nodes,  $x_1$  and  $x_2$ , which are chosen randomly from all the nodes.
- Initialize the connection set  $C$ ,  $C \subset A \times A$ , to the empty set. The output manifold is composed of  $A$  and  $C$ .
- Initialize the local accumulated error  $E_{x_1} = 0$  and  $E_{x_2} = 0$ .
- Initialize the local accumulated number of nodes  $M_{x_1} = 0$  and  $M_{x_2} = 0$ .

Corresponding author is Furao Shen.



- Initialize the similarity threshold  $T_{s_1}$  and  $T_{s_2}$  to an infinite number.

2) *Input and Find Winner*: Input a new node  $x_i$ . Search the nearest node  $s_1$  and the second nearest node  $s_2$  by

$$s_1 = \underset{s_c \in A}{\operatorname{argmin}} \|x_i - s_c\| \quad (1)$$

$$s_2 = \underset{s_c \in A \setminus \{s_1\}}{\operatorname{argmin}} \|x_i - s_c\| \quad (2)$$

3) *Nodes Update*: If the distance between  $x_i$  and  $s_1$  are greater than  $T_{s_1}$  and the distance between  $x_i$  and  $s_2$  are greater than  $T_{s_2}$ , it means that the node  $x_i$  is far away from all the nodes in A, hence we add the new node  $x_i$  to A. Then go to step 2 to process next node.

If a connection between  $s_1$  and  $s_2$  does not exist already, create it and set the age of the connection  $age(s_1, s_2)$  to 0, then add it to connection set C.

For all edges emanating from  $s_1$ :

$$age(s_1, s_i) = age(s_1, s_i) + 1 \quad (3)$$

In which  $s_i$  is next to  $s_1$ . If  $age(s_1, s_i)$  is too large, it means that the edge between  $s_i$  and  $s_1$  is created long time ago, hence we remove edges with an age greater than a predefined threshold  $age_{dead}$  and delete it from C.

For the winner  $s_1$ , add the Euclidian distance between the node  $x_i$  and the winner  $s_1$  to local accumulated error  $E_{s_1}$ , add 1 to the local accumulated number of signals  $M_{s_1}$ . Then, adapt the  $s_1$  and its direct topological neighbors to the node  $x_i$  by a certain fraction  $\varepsilon$ .

4) *Thresholds Update*: Adjust the thresholds  $T_{s_1}/T_{s_2}$  to the maximum distance between  $s_1/s_2$  and its neighbors.

5) *Noise Elimination*: If the number of input generated so far is an integer multiple of parameter  $\lambda$ , for all  $s_i$  in A, if  $s_i$  has no or only one neighbor and  $M_{s_1}$  is less than an predefined threshold, remove  $s_i$  from the node set A.

6) *Manifold Generation*: After the process above, we can get the manifold consisting of the node set A and the connection set C. However, sometimes the number of edges may be too small, which would fragment the manifold into some disconnected components. Therefore, for all nodes in A, if the number of neighbors  $k_i$  of the node  $s_i$  is less than the average number of neighbors  $k$ , then add the connection between the node  $s_i$  and its nearest node which are not adjacent to it until  $k_i$  becomes  $k$ .

Then, the number of the neighbors will be appropriate, which means that the manifold will not be fragmented into a large number of disconnected regions and will not cause "short-circuit" problem.

## B. Intrinsic Dimension Estimation

For the reason that the input data lies on a low-dimensional manifold and we have got the appropriate manifold after the above process, we can estimate the intrinsic dimension of the input data. The manifold is composed of a lot of local linear space, therefore we can estimate the dimension of the local linear space, then calculate the average value of the dimension, which is the intrinsic dimension of the input data. We adopt the locally principal component analysis algorithm to do the dimension estimation.

### Algorithm 1 Dimension Estimation

1. Input: the node set A and the connection set C.
2. For all the node  $s_i \in A$ , find the node set N, which consists of the neighbors of  $s_i$ , put the  $s_i$  into N, suppose the size of N be M.
3. Compute the mean  $\bar{s} = \frac{\sum_{i=1}^M n_i}{M}$ .
4. Compute the covariance matrix  $X = Cov(s) = \frac{1}{N} \sum_{i=1}^M (x_i - \bar{x})(x_i - \bar{x})^T$  and find the eigenvalues of X  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ .
5.  $k_i = \operatorname{argmin}_k \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} \geq 0.9$ ,  $k_i$  is the intrinsic dimension around the node  $s_i$
6.  $k = \frac{\sum_{s_i \in A} k_i}{|A|}$ .  $k$  is the intrinsic dimension of the input data.

Because of the small number of the neighbors, calculating the local dimension of the node  $s_i$  can be run in O(1) time. Therefore all the algorithm can be run in O(N) time, that is acceptable.

## C. Dimensionality Reduction

Now, we have the manifold of the input data and the target dimension, hence we can do the dimensionality reduction. For the reason that the input data lies on the manifold, we can do dimensionality reduction on the manifold made up of the node set A and the connection C. Because the number of nodes is much less than the original input, we can finish the dimensionality reduction with a high speed.

In this paper, we adopt Locally Linear Embedding (LLE) algorithm to do dimensionality reduction. LLE calculates the linear coefficients that reconstruct each data point from its neighbors [6]. Reconstruction errors are measured by the cost function:

$$\varepsilon(W) = \sum_i |s_i - W_{ij}s_j|^2 \quad (4)$$

If the  $s_j$  is not adjacent to  $s_i$ ,  $W_{ij} = 0$ . Calculate all the weight  $w_{ij}$  by minimizing the Eq 4. Then compute the low-dimensional embedding node that can be best reconstructed by minimizing the cost function:

$$\phi(M) = \sum_i |m_i - W_{ij}m_j|^2 \quad (5)$$