

# Linear Dimensionality Reduction Method Based on Topological Properties<sup>\*</sup>

Yuqin Yao<sup>a</sup>, Hua Meng<sup>a</sup>, Yang Gao<sup>b</sup>, Zhiguo Long<sup>b,c,\*</sup>, Tianrui Li<sup>b,c</sup>

<sup>a</sup>*School of Mathematics, Southwest Jiaotong University, Chengdu, 611756, China*

<sup>b</sup>*School of Computing and Artificial Intelligence, Southwest Jiaotong University, Chengdu, 611756, China*

<sup>c</sup>*Manufacturing Industry Chains Collaboration and Information Support Technology Key Laboratory of Sichuan Province, China*

---

## Abstract

Dimensionality reduction is an important data preprocessing technique that has been extensively studied in machine learning and data mining. Locality Preserving Projection (LPP) is a widely used linear unsupervised dimensionality reduction method, which maps high-dimensional data into low-dimensional subspace through linear transformation. Although various variants of LPP have been proposed to tackle different drawbacks of LPP, it is identified in this article that LPP does not possess the important topological property of translation invariance, that is, the linear transformation given by LPP is strongly related to the relative position between the data and the origin of the coordinate system. In this article, we theoretically analyze the reason why this drawback exists in LPP and propose to resolve it by introducing a kind of centralization to the model. Moreover, as topological properties are prominent information to characterize the structure of the data, this article proposes a further improvement of LPP to maintain topological connectivity of data after dimensionality reduction. Experiments on multiple synthetic and real-world datasets show that the new model incorporating topological properties outperforms not only the original LPP model but also several other classic linear or non-linear dimensionality reduction methods.

*Keywords:* Dimensionality reduction, linear projection, LPP, translation invariance, connectivity

---

## 1. Introduction

2 The era of big data brings us the problems of complexity, diversity and high dimension-  
3 ality in data, e.g., [1, 2, 3]. If this kind of data are directly used to obtain information,  
4 there will be not only troubles of irrelevant attributes, but also increased computational

---

\*Accepted version. Formal version available at: <http://dx.doi.org/10.1016/j.ins.2022.12.098>

\*Corresponding author.

*Email address:* zhiguolong@swjtu.edu.cn (Zhiguo Long)

5 complexity and reduced performance. Dimensionality reduction can remove irrelevant infor-  
6 mation and reduce the complexity of data, and thus has become an important preprocessing  
7 step in machine learning and data mining. Traditional dimensionality reduction methods  
8 can be divided into two categories: feature extraction [4, 5] and feature selection [6, 7, 8]  
9 methods. Feature extraction methods aim to produce new features by mapping the original  
10 high-dimensional data into low-dimensional space through algebraic transformation. A good  
11 dimensionality reduction method should retain essential characteristics of data as much as  
12 possible, and remove redundant information to reveal the underlying structure and pattern  
13 of data. In addition, dimensionality reduction is also helpful for data visualization. For  
14 applications, dimensionality reduction has achieved great success in face recognition [9, 10],  
15 handwritten numeral recognition [11], signature verification [12], disease diagnosis [13] and  
16 stock selection [14].

17 In the past, numerous feature extraction based dimensionality reduction techniques  
18 have been studied from different aspects. There are supervised, unsupervised, and semi-  
19 supervised ones according to the availability of label information, and there are linear and  
20 non-linear ones depending on the corresponding mapping types. For example, Principal  
21 Component Analysis (PCA) [15] is a classic unsupervised linear technique; Linear Dis-  
22 criminant Analysis (LDA) [16] is a widely used supervised linear technique, and there are  
23 many other supervised linear techniques (e.g., [17]); Locally Linear Embedding (LLE) [18],  
24 t-Distributed Stochastic Neighbor Embedding (t-SNE) [19], Isometric Feature Mapping  
25 (Isomap) [20] and Laplacian Embedding (LE) [21], and many more (e.g., [22, 23]) are on the  
26 other hand non-linear techniques. There are also dimensionality reduction methods based  
27 on nonnegative matrix factorization, e.g., [24, 25].

28 Locality Preserving Projection (LPP) [26] is a well-known feature extraction based un-  
29 supervised dimensionality reduction approach, and it is a linear approximation of LE. LPP  
30 aims to project the original data through a linear transformation while retaining nearest  
31 neighbor connections. Several variants have been proposed to deal with different draw-  
32 backs of LPP. For example, supervised information is exploited in some variants, including  
33 CLPP [27] and CdLPP [28], which only consider similarities between data points within the  
34 same class, and LPDP [29] and DLPP [30], which incorporate inter- and intra-class informa-  
35 tion;  $L_2$  norm is replaced by  $L_1$  norm to achieve better robustness, e.g., 2D-DLPP-L1 [31]  
36 and ILPP-L1 [32], where the latter also proposes to preserve similarities between points and  
37 their neighbors in addition to dissimilarities; similarity measure calculation is also improved,  
38 e.g., LAPP [33] adaptively measures similarities in new representations obtained by applying  
39 LPP iteratively.

40 Although LPP has been shown to perform well on various datasets, it is surprising to  
41 notice that its selected directions for projection are very sensitive to the coordinate system  
42 for data points. In fact, by only changing the position of the data points in the coordinate  
43 system, i.e. translating the data points as a whole, the projection directions given by LPP  
44 might change significantly, though the relative positions of data points are not changed (cf.  
45 Figure 1). This means LPP does not possess an important topological property that is  
46 *translation invariance*.

47 In fact, taking topological properties into account for data analysis is a trending topic

48 known as topological data analysis (TDA) [34]. One of the main differences of TDA ap-  
49 proaches from traditional statistical learning ones is that, the former concern more about  
50 local or global structural information, such as compactness, connectivity, or algebraic prop-  
51 erties like persistent homology [35], whereas the latter are more interested in distribution  
52 characteristics of data. Having seen its merits, researchers applied TDA to various fields  
53 and tasks, e.g., image processing [36], classification [37], and bioinformatics [38].

54 This article repairs the topological property of translation invariance by introducing a  
55 kind of *centralization* to the LPP model. Moreover, as topological properties concern im-  
56 portant information of data (e.g., the number of connected components in data is important  
57 structural information), we propose a novel characterization of the latent topological struc-  
58 ture of data, viz. *topological connectivity*, which reflects connections of different parts of data,  
59 and empirically show that it is useful to retain this kind of topological information when  
60 performing dimensionality reduction. Existing works [18, 21] usually consider maintaining  
61 connectivity between points and their neighbors, including LPP, LE, and LLE. In other  
62 words, they only capture local structures, but ignore higher level structures like connected  
63 components (cf. Figure 4). However, these higher level structures can reveal important in-  
64 formation of data, e.g., similarity and separation between data points in these structures  
65 should also be maintained after dimensionality reduction. Therefore, the approach proposed  
66 here not only tries to capture local structures, but also explores these higher level ones by  
67 constructing topological connectivity of data.

68 There are LPP variants that consider maintaining certain kinds of structures of data.  
69 For example, 2D-DLPP-L1 [31] proposes to maintain relative positions of image pixels by  
70 using matrices instead of vectors for finding projection subspaces. It considers structural  
71 information within each sample instead of among data points, which is different from the  
72 approach in this article. SSTNTL [39] makes use of supervised information to remove  
73 connections between points in different classes, and then reduces dimensionalities of data by  
74 using the new connection graph. It constructs topology on data points to help characterize  
75 similarities between points, but does not care about topological structural information like  
76 connected components.

77 The contributions of the article are as follows:

- 78 • We identify, theoretically analyze, and resolve the problem of the original LPP model  
79 that it is *highly sensitive to geometric translation*, i.e., the projection directions change  
80 significantly when data are moved around in the coordinate system.
- 81 • We propose to retain topological connectivity of data in dimensionality reduction, by  
82 exploring connectivity information in variant scales with novel connectivity measures  
83 for data, and devise an improved projection model for such purpose.
- 84 • We demonstrate the effectiveness of the improved model and its superiority over LPP  
85 and several other classic dimensionality reduction methods on multiple synthetic and  
86 real-world datasets.

87 The rest of this paper is arranged as follows. Section 2 briefly introduces the original LPP  
88 model. Section 3 analyzes the identified problems of LPP in details and proposes an improved

89 model (ConLPP) and its corresponding algorithms. Section 4 evaluates the performance of  
 90 the new method on both synthetic and real-world datasets. Section 5 concludes the paper.

## 91 2. Preliminaries

92 Given a dataset  $X = (x_1, x_2, \dots, x_m)$ , where  $x_i \in \mathbb{R}^n$  is a column vector with  $n$   
 93 feature values. A linear projection can be described by a matrix  $A_{n \times d}$  which maps  $X$  to  
 94  $Y = (y_1, y_2, \dots, y_m)$ , where  $y_i \in \mathbb{R}^d$  and  $y_i = A^T x_i$ .

### 95 2.1. Locality Preserving Projection

Locality Preserving Projection (LPP) is an unsupervised dimensionality reduction model,  
 which tries to preserve the neighborhood structure of data. It is a linear approximation of  
 Laplacian Eigenmaps (LE), and can achieve better results than LE on various types of  
 data. The goal of LPP is to select several projection directions  $a_i \in \mathbb{R}^n$  ( $1 \leq i \leq d$ ) to  
 form a projection matrix  $A = (a_1, \dots, a_d)$ , so that the projected data  $A^T X$  can satisfy  
 that nearest neighboring points are still neighbors. In order to achieve this goal, LPP  
 constructs  $\operatorname{argmin}_a \sum_{i,j} (a^T x_i - a^T x_j)^2 W_{ij}$  as the optimization goal, where  $W_{ij}$  is a value  
 characterizing the similarity between  $x_i$  and  $x_j$ . LPP adds the constraint  $a^T X D X^T a = 1$   
 to avoid the influence of scaling on the projection directions, where  $D$  is a diagonal matrix  
 s.t.  $D_{ii} = \sum_j W_{ij}$ . So the LPP model can be expressed as

$$\operatorname{argmin}_a \sum_{i,j} (a^T x_i - a^T x_j)^2 W_{ij}, \text{ s.t. } a^T X D X^T a = 1. \quad (\text{Model 1})$$

96 The following lemma shows that the objective function of Model 1 can be written in matrix-  
 97 vector form.

98 **Lemma 1** ([26]).  $\sum_{i,j} (a^T x_i - a^T x_j)^2 W_{ij} = a^T X L X^T a$ , where  $L = D - W$  is a Laplacian  
 99 matrix and  $D$  is a diagonal matrix s.t.  $D_{ii} = \sum_j W_{ij}$ .

Model 1 is then equivalently transformed to Model 2 below:

$$\operatorname{argmin}_a a^T X L X^T a \quad \text{s.t.} \quad a^T X D X^T a = 1. \quad (\text{Model 2})$$

100 The steps to solve Model 2 of LPP are as follows.

101 1. Construct the adjacency matrix:

102 If  $x_i$  and  $x_j$  are *neighbors*, then  $W_{ij} = \exp(-\frac{\|x_i - x_j\|^2}{t})$  is used to represent the similarity  
 103 between  $x_i$  and  $x_j$ ; otherwise  $W_{ij}$  is set to 0. There are two common ways to determine  
 104 neighbors:

- 105 (a)  $\varepsilon$ -neighborhood. For a given parameter  $\varepsilon > 0$ , if  $\|x_i - x_j\|^2 < \varepsilon$ , then  $x_i$  and  $x_j$   
 106 are neighbors.
- 107 (b)  $k$ -nearest neighborhood. For a given positive integer  $k$ , if  $x_i$  is among the  $k$   
 108 nearest neighbors of  $x_j$ , or  $x_j$  is among the  $k$  nearest neighbors of  $x_i$ , then  $x_i$  and  
 109  $x_j$  are neighbors.

110 2. Find the  $d$  optimal projection directions of Model 2 by calculating the eigenvectors of  
 111  $(X D X^T)^{-1} X L X^T$  corresponding to the  $d$  smallest eigenvalues.

112 *2.2. An Equivalent Model of LPP*

It can be easily seen that Model 2 of LPP is also equivalent to the following model.

$$\operatorname{argmin}_a \frac{a^T X L X^T a}{a^T X D X^T a} \quad \text{s.t.} \quad a^T X D X^T a = 1. \quad (\text{Model 3})$$

In order to solve Model 3, we can first solve the following Model 3\* and then adjust the length of  $a$  to satisfy the constraint  $a^T X D X^T a = 1$ . That is, for  $a$  being a solution to Model 3\* and  $a^T X D X^T a = h$ , let  $a' = \frac{1}{\sqrt{h}}a$ , then  $a'^T X D X^T a' = 1$  and is a solution to Model 3. The correctness of this process is ensured by Proposition 2.

$$\operatorname{argmin}_a \frac{a^T X L X^T a}{a^T X D X^T a}. \quad (\text{Model 3}^*)$$

113 **Proposition 2.** *If  $c$  is a solution to Model 3\*, then there is a solution  $b$  to Model 3 s.t. the*  
 114 *direction of  $c$  is the same as  $b$ , and vice versa.*

115 **PROOF.** Note that the value of  $\frac{c^T X L X^T c}{c^T X D X^T c}$  is only related to the direction of  $c$ , not its length,  
 116 i.e.,  $\frac{c^T X L X^T c}{c^T X D X^T c} = \frac{(\lambda c)^T X L X^T (\lambda c)}{(\lambda c)^T X D X^T (\lambda c)}$  for any  $\lambda \neq 0$ . Therefore, for a solution  $c$  of Model 3\*, we can  
 117 always find some  $\lambda$  s.t.  $(\lambda c)^T X D X^T (\lambda c) = 1$ , and thus  $b = \lambda c$  is also a solution of Model 3,  
 118 whereas  $c$  and  $b$  have the same direction.

119 On the other hand, if  $b$  is a solution to Model 3, then  $b$  is also a solution to Model 3\*.  
 120 This is because, if  $\frac{b^T X L X^T b}{b^T X D X^T b}$  is not minimal, then  $\exists c$  s.t.  $\frac{c^T X L X^T c}{c^T X D X^T c} < \frac{b^T X L X^T b}{b^T X D X^T b}$  and hence  $\exists \lambda$   
 121 s.t.  $(\lambda c)^T X D X^T (\lambda c) = 1$  and  $\frac{(\lambda c)^T X L X^T (\lambda c)}{(\lambda c)^T X D X^T (\lambda c)} = \frac{c^T X L X^T c}{c^T X D X^T c} < \frac{b^T X L X^T b}{b^T X D X^T b}$ . This means  $\lambda c$  is a  
 122 solution to Model 3 and with a smaller objective function value, which contradicts to that  
 123  $b$  is a solution to Model 3.  $\square$

124 Minimizing  $\frac{a^T X L X^T a}{a^T X D X^T a}$  is to balance the two goals: minimizing  $a^T X L X^T a$  and maximizing  
 125  $a^T X D X^T a$ . Minimizing  $a^T X L X^T a$  is to make the transformed points as close as possible,  
 126 i.e., to maintain neighbors. In next section, we will see that maximizing  $a^T X D X^T a$  is to  
 127 maximize a weighted sum of squares of the distances between the transformed points and  
 128 the origin of the coordinate system. This makes the projection directions found by LPP very  
 129 sensitive to the location of data in the coordinate system, instead of the relative position  
 130 between data points, which is unfavorable in many scenarios. For example, for self-driving  
 131 cars, as the car moves, a static obstacle will be at dynamic positions from the point of  
 132 view of the car, and if LPP is applied for data preprocessing, then this obstacle might not  
 133 be well recognized during the movement, since the projection direction keeps changing and  
 134 the resulting representation might lose important information for some of the projection  
 135 directions. We propose to repair this problem and to consider more topological properties  
 136 to improve the original LPP model.

137 **3. ConLPP Algorithm**

138 We improve LPP from two aspects, i.e., repairing translation invariance and introducing  
 139 topological connectivity. In the following, we discuss them in details.

140 *3.1. Translation Invariance*

141 *3.1.1. Analysis*

142 In the objective function of Model 3 of LPP, the term  $a^T XDX^T a$  is strongly affected by  
 143 the specific coordinates of data points. This means that for different coordinates of the data  
 144 points, even when the relative positions of them are the same, LPP might produce signif-  
 145 icantly different projection directions, i.e., LPP is very sensitive to geometric translations.  
 146 Take Figure 1 as an example. In this figure, there are five groups of data and each group  
 147 has two clusters, where the distributions of the data are exactly the same between groups  
 148 (different groups are simply geometric translations of each other), but their positions in  
 149 the coordinate system are different. Arrows represent projection directions found by LPP,  
 150 and the points on an arrow represent the data after projection on that direction. It can  
 151 be seen from the figure that the projection directions vary significantly across these five  
 152 groups although the distributions of points in different groups are the same. Notably, the  
 153 change of directions will significantly affect the resulting representation after dimensionality  
 154 reduction: some directions might destroy the structure within a group, resulting in a large  
 155 overlap between data points in the two clusters within a group.

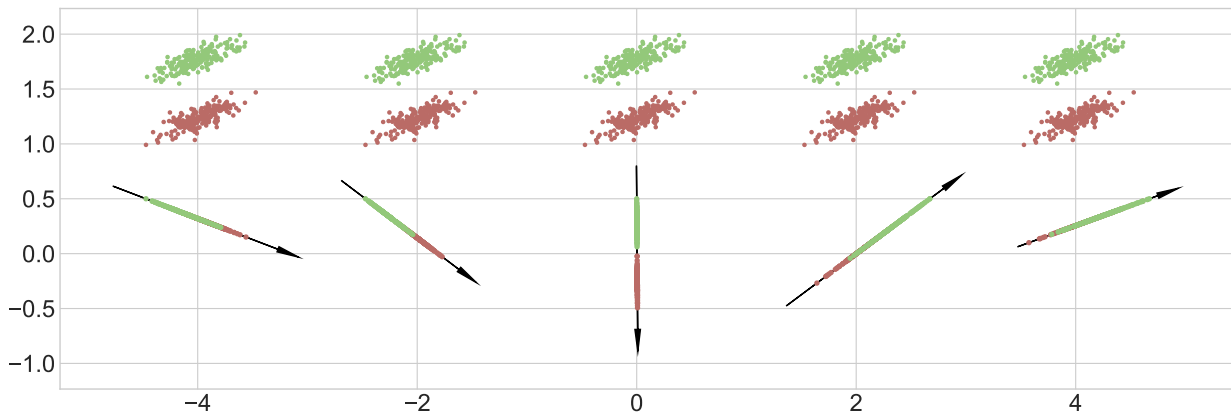


Figure 1: LPP projection results for geometric translations of the same group of data.

156 In order to analyze why LPP is seriously affected by coordinates, we make the following  
 157 observation.

158 **Proposition 3.**  $a^T XDX^T a = \sum (a^T x_i)^2 D_{ii}$ .

PROOF. Recall that  $X = (x_1, x_2, \dots, x_m)$ , then

$$\begin{aligned}
\sum_i (a^T x_i)^2 D_{ii} &= \sum_i (a^T x_i) (a^T x_i)^T D_{ii} \\
&= \sum_i a^T (x_i D_{ii} x_i^T) a \\
&= a^T \left( \sum_i x_i D_{ii} x_i^T \right) a \\
&= a^T \left( (x_1 D_{11}, x_2 D_{22}, \dots, x_m D_{mm}) (x_1, x_2, \dots, x_m)^T \right) a \\
&= a^T X D X^T a.
\end{aligned} \tag{1}$$

159 □

160 Note that  $(a^T x_i)^2$  is the square of the distance from the projected coordinates to the  
161 origin. Therefore, by Proposition 3, maximizing the term  $a^T X D X^T a$  is equivalent to max-  
162 imize the weighted sum of squares of the distance from the projected coordinates to the  
163 origin, where the weights are  $D_{ii}$  ( $i = 1, \dots, m$ ). This term in Model 3 of LPP tends to keep  
164 the data away from the origin after projection. For example, for data with higher density,  
165  $D_{ii}$  will be larger and this model will be more likely to keep the data points away from the  
166 origin after projection. This is why  $\max a^T X D X^T a$  is strongly affected by coordinates.

### 167 3.1.2. Improvement for Translation Invariance

In order to reduce the influence of the positions of the data in the coordinate system for LPP, we propose to *centralize* the coordinates of data. Let  $x'_i = x_i - \mu$  and  $X' = (x'_1, x'_2, \dots, x'_m)$ , where  $\mu = \sum_i x_i / m$ . It is easy to check that  $\|x_i - x_j\| = \|x'_i - x'_j\|$  for any  $i$  and  $j$ . Then, the nearest neighbors of data points do not change, and if we consider Model 3 on  $X'$ , we have

$$W'_{ij} = \exp\left(\frac{-\|x'_i - x'_j\|^2}{\sigma^2}\right) = \exp\left(\frac{-\|x_i - x_j\|^2}{\sigma^2}\right) = W_{ij}, \tag{2}$$

$$D'_{ii} = \sum_j W'_{ij} = \sum_j W_{ij} = D_{ii}, \tag{3}$$

$$L' = D' - W' = D - W = L. \tag{4}$$

Therefore, an improved model of Model 3 for LPP can be proposed to repair translation invariance as follows.

$$\arg \min_a \frac{a^T X' L X'^T a}{a^T X' D X'^T a} \quad \text{s.t.} \quad a^T X' D X'^T a = 1. \quad (\text{Model 4})$$

168 The relation between Model 3 of LPP and Model 4 can be seen from the following  
169 proposition.

170 **Proposition 4.**  $a^T X' L X'^T a$  is equal to  $a^T X L X^T a$ , but when  $\mu \neq 0$ ,  $a^T X' D X'^T a$  is not  
 171 always equal to  $a^T X D X^T a$ .

PROOF. From Lemma 1, we know that  $a^T X L X^T a = \sum_{i,j} (a^T x_i - a^T x_j)^2 W_{ij}$ , and similarly  
 $a^T X' L X'^T a = \sum_{i,j} (a^T x'_i - a^T x'_j)^2 W'_{ij}$ . Note that

$$\begin{aligned} \sum_{i,j} (a^T x'_i - a^T x'_j)^2 W'_{ij} &= \sum_{i,j} (a^T (x_i - \mu) - a^T (x_j - \mu))^2 W_{ij} \\ &= \sum_{i,j} (a^T x_i - a^T x_j)^2 W_{ij}. \end{aligned} \quad (5)$$

172 Therefore,  $\sum a^T X' L X'^T a = \sum a^T X L X^T a$ . On the other hand, by Proposition 3,  $a^T X D X^T a =$   
 173  $\sum (a^T x_i)^2 D_{ii}$ , and similarly  $a^T X' D X'^T a = \sum_i (a^T (x_i - \mu))^2 D_{ii}$ . It is easy to see that gener-  
 174 ally when  $\mu \neq 0$  then  $\sum (a^T x_i)^2 D_{ii}$  is not equal to  $\sum_i (a^T (x_i - \mu))^2 D_{ii}$ .  
 175 □

176 In fact, from the proof, we can see that the term  $a^T X D X^T a$  in Model 3 of LPP corresponds  
 177 to the “weighted distance” of data points to the origin, whereas  $a^T X' D X'^T a$  of Model  
 178 4 corresponds to the “weighted distance” of data points to the center of them, which is  
 179 invariant w.r.t. geometric translations of data. In addition, it is also easy to see that  
 180  $a^T X' L X'^T a = \sum_{i,j} (a^T x'_i - a^T x'_j)^2 W'_{ij}$  is also invariant w.r.t. geometric translations of data.  
 181 Then we have the following conclusion.

182 **Theorem 5.** For any dataset  $X$ , the projection directions found by Model 4 are independent  
 183 of geometric translations of the data.

184 Figure 2 illustrates the projection of the same data as in Figure 1 with Model 4. We  
 185 can see that the projection directions are not affected by the relative position of data to the  
 186 origin.

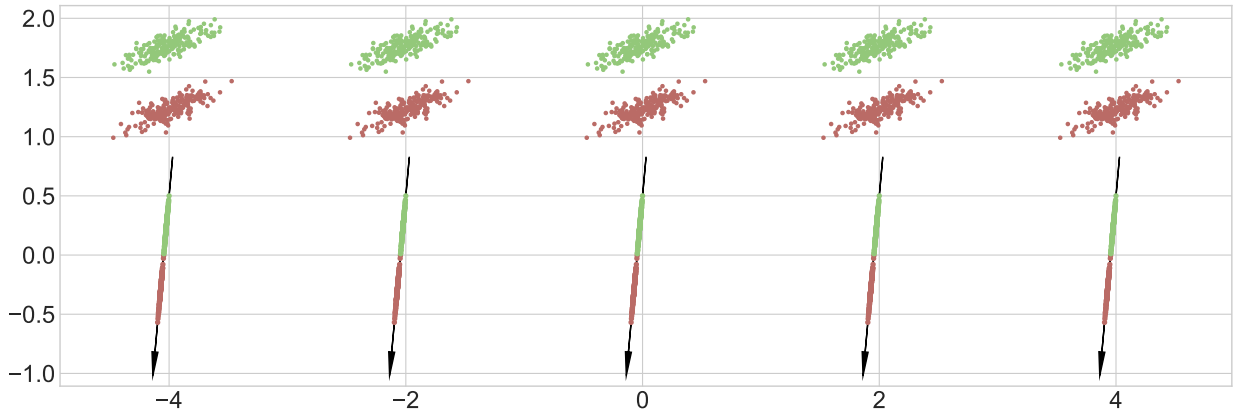


Figure 2: Projection results of the improved Model 4 for geometric translations of the same group of data.



Note that  $a^T X' D X^T a = \sum_i (a^T (x_i - \mu))^2 D_{ii}$ . In particular, when  $D_{ii} = 1$ ,  $a^T X' D X^T a$  degenerates to

$$\begin{aligned} \sum_i (a^T (x_i - \mu))^2 &= \sum_i (a^T x_i - a^T \mu) (a^T x_i - a^T \mu)^T \\ &= a^T \left( \sum_i (x_i - \mu) (x_i - \mu)^T \right) a \end{aligned} \quad (6)$$

187 In this situation,  $\operatorname{argmax} \sum_i (a^T (x_i - \mu))^2$  is exactly the goal of PCA, that is, to keep the data  
 188 as separated as possible. Therefore,  $\operatorname{argmax} \sum_i (a^T (x_i - \mu))^2 D_{ii}$  can be regarded as a weighted  
 189 PCA model, s.t. the points with higher density will have higher weights. Therefore, Model  
 190 4 is equivalent to finding the projection directions that *locally* make the nearest neighbors  
 191 of data as close as possible, and *globally* maintain the weighted separation of data as much  
 192 as possible. The chosen directions balance these two optimization objectives.

### 193 3.2. Topological Connectivity

194 Although Model 4 repairs the translation invariance property of LPP, it ignores the inter-  
 195 mediate structures that are between local nearest neighborhoods and global data separation.  
 196 In fact, data usually have several clusters, and clusters may be connected or separated. How-  
 197 ever, the original LPP and the revised Model 4 do not care about these structures, and this  
 198 may result in poor projection results when the data structure is relatively complex. For  
 199 example, consider Figure 3. The left of the figure shows the projection of the original data  
 200 by LPP, and the right shows the projection by Model 4. As both LPP and Model 4 only  
 201 consider to keep local neighboring points close and to maintain global separation of the  
 202 whole data, but ignore intermediate structures, the dimensionality reduction results of them  
 203 have obvious overlaps between data from different clusters. The main reason is that the  
 204 two smaller clusters were taken as a whole for global separation, and the two models did  
 205 not consider separating them as the main optimization goal. We argue that this kind of  
 206 intermediate separation information can be captured by *topological connectivity*.

#### 207 3.2.1. Topological Space and Its Connectivity

208 Connectivity is a fundamental concept in topology. Before discussing topological con-  
 209 nectivity of a dataset, we briefly introduce the concepts related to topological connectivity,  
 210 and more details can be found in [40].

211 For a nonempty set  $X$ , let  $T \subseteq 2^X$  be a family of subsets of  $X$ . If  $T$  contains  $\emptyset$  and  
 212  $X$ , and is closed with respect to set union operation and finite intersection operation, then  
 213  $T$  is called a *topology* on  $X$ . The elements in  $T$  are called *open sets*, and  $(X, T)$  is called  
 214 a *topological space*. A *closed set* in  $(X, T)$  is a subset of  $X$  whose complement is an open  
 215 set. Suppose  $(X, T)$  is a topological space and  $Y \subseteq X$ . Then the *subspace topology*  $T_Y$  is  
 216 a topology on  $Y$  such that  $V \in T_Y$  iff there is a  $U \in T$  and  $V = U \cap Y$ . For example, the  
 217  $n$ -dimensional Euclidean space  $\mathbb{R}^n$ , with the topology in which all the *open balls* are open  
 218 sets, is a topological space, called the *Euclidean topological space*. Here, an open ball  $B(x, \delta)$   
 219 in  $\mathbb{R}^n$  is  $B(x, \delta) = \{y \mid d(x, y) < \delta\}$ . A subspace topology of the Euclidean topological space

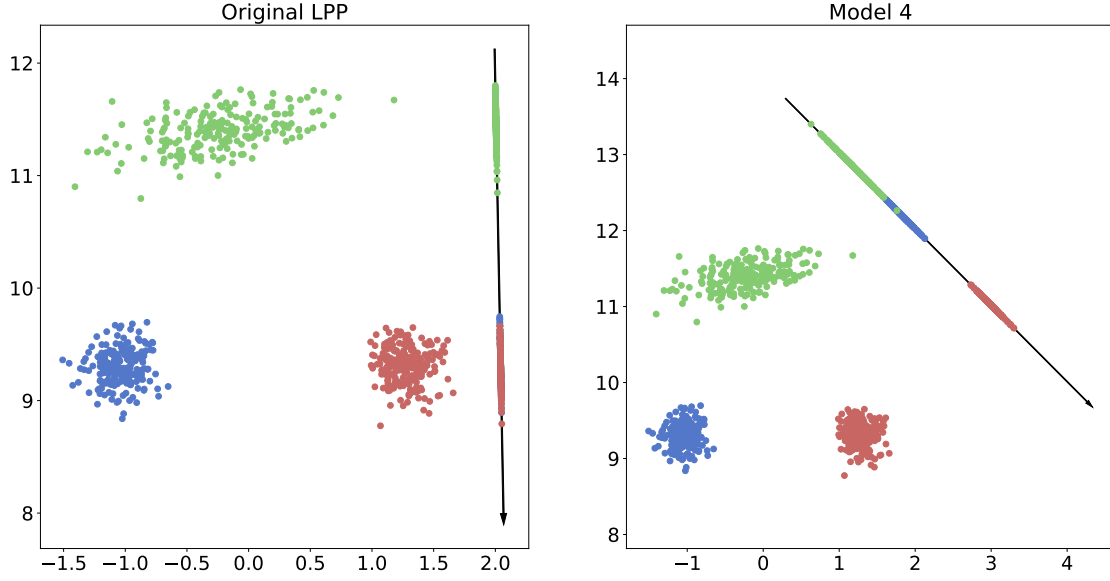


Figure 3: Illustration of the problems of the original LPP model and the improved Model 4.

220  $\mathbb{R}^2$  could be  $T_Y$  for  $Y = \{x \mid \|x\| \leq 2\}$ . Topology can be used to characterize the relative  
 221 distance between data points. For instance, if  $x \in \mathbb{R}^n$ , and  $y \in B(x, 1/5), z \notin B(x, 1/3)$  can  
 222 characterize that  $y$  is closer to  $x$  than  $z$  is.

223 Topology can also characterize connectivity information of data. A topological space is  
 224 said to be *connected* if and only if there is no nonempty set  $U$  in it such that both  $U$  and  
 225 its complement are open sets. Therefore, if a topological space is not connected, then it can  
 226 be divided into disjoint parts, i.e., *connected components*.

227 **Definition 1.** Suppose  $(X, T)$  is a topological space,  $Y \subseteq X$ , and  $Y \neq \emptyset$ . Then  $Y$  is a  
 228 *connected component* of  $(X, T)$  iff  $(Y, T_Y)$  is connected and for all  $Y'$  s.t.  $Y \subset Y' \subseteq X$ ,  
 229  $(Y', T_{Y'})$  is not connected.

230 If  $(X, T)$  is connected, then  $X$  itself is the unique connected component of  $(X, T)$ ; if  
 231  $(X, T)$  is not connected, then  $X$  can be decomposed into several connected components.  
 232 When two points  $x, y \in X$  are in different connected components, we say  $x$  and  $y$  are *strongly*  
 233 *separated*. Take Figure 4(a) as an example. The two geometries shown in it together form a  
 234 topological subspace of  $\mathbb{R}^2$ , which is not connected but contains two connected components.

235  
 236 In practice, *path connectivity* is usually used to characterize connectivity of a set in a  
 237 topological space.

238 **Definition 2 (Path connectivity).** Suppose  $(X, T)$  is a topological space, and  $x, y \in X$ .  
 239 Then if there is a continuous map  $f$  from  $[0, 1]$  to  $X$ , such that  $f(0) = x$  and  $f(1) = y$ , then  
 240  $f$  is a *path* from  $x$  to  $y$ . In this case,  $x$  and  $y$  are said to be *path connected*. A topological  
 241 space  $(X, T)$  is *path connected* if any two distinct  $x, y \in X$  are path connected.

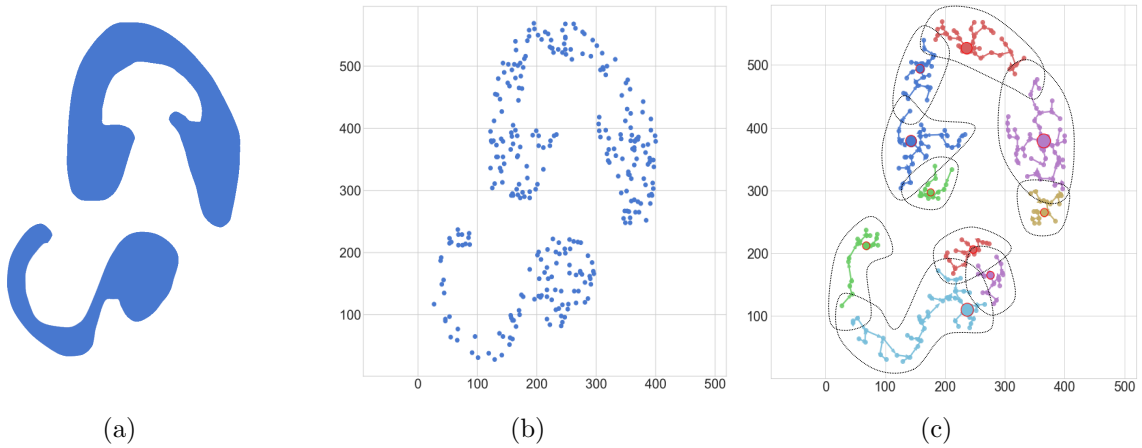


Figure 4: Illustration of connected components in a topological subspace and connectivity of data points.

242 In Euclidean space, a path between two points  $x$  and  $y$  is a continuous curve connecting  
 243  $x$  and  $y$ . Note that a path connected topological space is also connected, but the other  
 244 way around generally is not true. However, for most of the practical cases, a connected  
 245 topological space is also path connected, e.g., for a connected component in Figure 4(a),  
 246 it is both connected and path connected. We also have the following properties for path  
 247 connectivity.

248 **Lemma 6.** *Suppose  $(X, T)$  is a topological space, and  $x_1, x_2, x_3 \in X$ .*

- 249 • *If  $x_1, x_2$  and  $x_2, x_3$  are path connected then  $x_1, x_3$  are path connected.*
- 250 • *If  $x_1$  and  $x_2$  are path connected, then  $x_1$  and  $x_2$  are in the same connected component.*

251 To maintain the structure of data during dimensionality reduction, connectivity infor-  
 252 mation of data is an important aspect to consider. Topological connectivity provides a  
 253 promising candidate to model such information. However, as we usually do not have all of  
 254 the data points from a topological subspace, but only some samples of them, and sometimes  
 255 even the whole data subspace itself may be a discrete set. For example, in Figure 4(b), we  
 256 only have some sample data points from the topological subspace in Figure 4(a), and thus  
 257 we have to use only these sample data points to discover the structural information of the  
 258 whole topological subspace.

259 In the following, we will discuss how to model connectivity for sample data points with  
 260 the idea of topological connectivity. The general idea is to approximate connectivity of  
 261 a sample space by exploring a kind of “path connectivity” of data points and constructing  
 262 connected components with such information. Once connectivity and separation information  
 263 is available, we can compute similarity and separation measures for data points and thus  
 264 incorporate more topological information into the dimensionality reduction model.

265 *3.2.2. Connectivity of Data*

266 Suppose  $X = (x_1, x_2, \dots, x_m)$  contains the sample data points. For every  $x_i \in X$  we  
 267 denote by  $N_k(x_i)$  the nearest  $k$  neighbors of  $x_i$  and we require  $x_i \in N_k(x_i)$ . Intuitively, a  
 268 connected component in a topological space consists of two kinds of points: inner points and  
 269 boundary points. To distinguish between these points, we exploit the concept of density, as  
 270 inner points usually have higher density than boundary points.

**Definition 3 (*Density of points*).** For each  $x_i \in X$ , we define the *density* of  $x_i$  as follows:

$$\rho(x_i) = \sum_{x_j \in N_k(x_i) \setminus \{x_i\}} \exp\left(-\frac{d_{ij}^2}{\sigma^2}\right), \quad (7)$$

271 where  $d_{ij}$  is the Euclidean distance between  $x_i$  and  $x_j$ .

272 The density of data points reflects the tightness of the distribution of data points in the local  
 273 area. The inner-most point of a connected component tends to have the highest density,  
 274 and this point is called a *core point*.

275 **Definition 4 (*Leader point and core point*).** Denote by  $R(x_i)$  the *leader point* of  $x_i \in$   
 276  $X$ . Then  $R(x_i)$  can be defined as follows. If the density of  $x_i$  is greater than that of any  
 277  $x \in N_k(x_i) \setminus \{x_i\}$  then  $R(x_i) = x_i$ , otherwise  $R(x_i) = x_j$ , where  $x_j$  is the data point with  
 278 higher density than  $x_i$  in  $N_k(x_i)$  that is closest to  $x_i$  (if there are multiple, then we select  
 279 the first met one). We call  $x_i$  a *core point* if  $R(x_i) = x_i$ .

280 The leading relationship reflects connectivity to some extent. Generally, a point with  
 281 higher density is considered to be more representative for data. A point is likely to be con-  
 282 nected to some of its nearest neighbors, and the most possible connection happens between  
 283 this point and its nearest neighbor that has higher density (i.e. more representative). There-  
 284 fore, we consider there is a path between each point and its leader point, and they are in  
 285 the same connected component of a topological subspace. Thus, a point and its leader point  
 286 should be in the same connected component. In order to formally characterize connected  
 287 component with the help of the leading relationship between points, we need to define the  
 288 *leader set* of a point.

289 **Definition 5 (*Leader set*).** For each  $x \in X$ ,  $\mathcal{S}(x)$  is the *leader set* of  $x$ , which is recursively  
 290 defined as follows:

- 291 •  $R(x) \in \mathcal{S}(x)$ ;
- 292 • If  $y \in \mathcal{S}(x)$ , then  $R(y) \in \mathcal{S}(x)$ .

293 For convenience, we also call a point  $y$  in  $\mathcal{S}(x)$  a *leader* of  $x$ . It is easy to see that if  $z$  is  
 294 a core point then  $\mathcal{S}(z) = \{z\}$ . In the following, we also identify the set of the points whose  
 295 leader set contains the core point  $z$ , as a *density branch*, inspired by [41], and identify the  
 296 nearest  $k$  neighbors of each point in a density branch as an *expanded density branch*.

297 **Definition 6 (Density branch and expanded density branch).** For each core point  
 298  $z$ , we call the set of points that have  $z$  as a leader, i.e.,  $D(z) = \{x \mid z \in \mathcal{S}(x)\}$ , the  
 299 *density branch* of  $z$ . We also define the *expanded density branch* of a core point  $z$  as  
 300  $E(z) = \cup_{x \in D(z)} N_k(x)$ .

301 It can be seen from the definition that there is a 1-1 correspondence between density  
 302 branches and core points, i.e., a density branch contains exactly one core point and each  
 303 core point has a density branch.

304 Different density branches may belong to the same connected component and this hap-  
 305 pens if they are closely *adjacent*, that is, they have enough *shared nearest neighbors*. The  
 306 shared nearest neighbors of two core points  $z_1$  and  $z_2$  is the set  $\text{SNN}(z_1, z_2) = E(z_1) \cap E(z_2)$ .

307  
 308 **Definition 7 (Connectivity).** Given two core points  $z_1$  and  $z_2$ , the density branches  $D(z_1)$   
 309 and  $D(z_2)$  are *connected* if  $|\text{SNN}(z_1, z_2)| > \tau \times \min\{|E(z_1)|, |E(z_2)|\}$ , where  $\tau$  is some given  
 310 threshold parameter. Furthermore, if  $D(z_1)$  is connected with  $D(z_2)$  and  $D(z_2)$  is connected  
 311 with  $D(z_3)$ , then we also say  $D(z_1)$  is connected with  $D(z_3)$ . For any two point  $x_i \in D(z_1)$   
 312 and  $x_j \in D(z_2)$ ,  $x_i$  and  $x_j$  are said to be *connected* if  $D(z_1)$  is connected with  $D(z_2)$ .

313 Here, the idea is to consider a point  $x$  is connected with its leaders and thus the points in  
 314 each density branch are connected with each other, and two density branches are connected if  
 315 they have enough shared neighbors ( $\tau$  is set as 0.05 in experiments). Figure 4(c) illustrates  
 316 this idea. In this figure, each point is connected with its leaders and the points with a  
 317 red circle are core points. Each core point determines a density branch in which each  
 318 pair of points have a path between them. Expanded density branches are illustrated as  
 319 regions bounded by dashed lines. For two density branches, if they share enough neighbors  
 320 and thus have strong connection with each other, then we consider them as connected.  
 321 Although there is no path between points in two connected density branches, we think they  
 322 have a high possibility to be connected in the original topological subspace. In this way, we  
 323 obtain *connected components* of data points, and these components approximate those in  
 324 the original topological subspace.

325 **Definition 8 (Connected component).** A *connected component*  $C$  for a dataset  $X$  is the  
 326 maximal set in which each pair of points are connected.

327 Each density branch has a unique core point, which is the leader of every point in this  
 328 density branch. On the other hand, a connected component can contain more than one core  
 329 points. In fact, a connected component is exactly the union of all density branches that  
 330 are connected with each other. That is,  $C = \bigcup_{i=1}^l D(z_i)$ , where for any  $1 \leq i, j \leq l$ ,  $D(z_i)$   
 331 and  $D(z_j)$  are connected, but for any core point  $z \notin \{z_1, z_2, \dots, z_l\}$ ,  $D(z_i)$  and  $D(z)$  are not  
 332 connected. Roughly speaking, density branches characterize local structural information,  
 333 while connected components captures higher level structural information.

334 It is easy to see that the connectivity relation between data points is reflexive, transitive,  
 335 and symmetric. Therefore, the connectivity relation induces a partition of  $X$  into disjoint  
 336 connected components, as stated in the following proposition.

337 **Proposition 7.** Given  $k$  and  $\tau$ , each dataset  $X$  is divided into disjoint connected compo-  
 338 nents.

339 For example, as shown in Figure 5, given the values of  $k$  and  $\tau$ , for each  $x_i$ , we can  
 340 calculate  $N_k(x_i)$  and  $\rho(x_i)$  by definition. Then we can calculate  $R(x_i)$  and find all the core  
 341 points of the dataset. After that, for each core point  $z_t$  we can obtain its density branch  
 342  $D(z_t)$  and expanded density branch  $E(z_t)$ . According to Definition 8, the connected density  
 343 branches are combined to obtain connected components  $\{C_k\}$ . The corresponding process is  
 344 shown in Algorithm 1. Note that if the data contain outliers, then there are some connected  
 345 components that contain only few samples. Thus, to remove outliers, only the connected  
 components with more than 2 samples are retained by the algorithm.

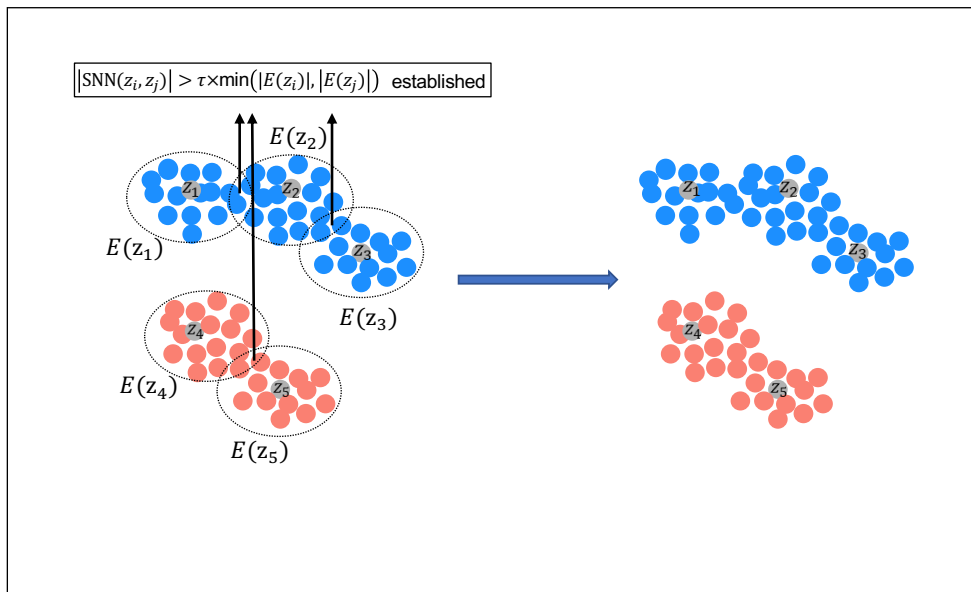


Figure 5: Illustration of connectivity and constructing connected components of data points.

346 Note that the number of connected components is related to the number of nearest  
 347 neighbors  $k$  and the given threshold  $\tau$ . Given  $\tau$ , if  $k = 1$ , then each data point itself forms  
 348 a connected component; if  $k = m$  ( $m$  is the total number of data points), then all the data  
 349 points together form a unique connected component. The smaller  $k$  is, the more connected  
 350 components there will be; the greater  $k$  is, the fewer connected components there will be,  
 351 and more density branches will be merged together. Moreover, if  $k$  is small and two points  
 352 are in the same density branch, then they should be considered as very similar. Conversely,  
 353 if  $k$  is large and two points belong to different connected components, then they should be  
 354 considered to be obviously different. Therefore, later we will take the results of different  $k$   
 355 into consideration to better capture the structure of data.

### 357 3.2.3. Similarity Matrix and Separation Matrix

358 With the connected components information obtained by Algorithm 1, we propose two  
 359 measures to characterize the similarity and difference between points, respectively.

---

**Algorithm 1:** ExploreStructure
 

---

**Input:** Dataset  $X$ ; number of nearest neighbors  $k$ ; threshold parameter  $\tau$ .

**Output:** Core points; density branches  $\{D(z_i)\}$ ; connected components  $\{C_i\}$

- 1 Calculate  $N_k(x)$ ,  $\rho(x)$  and leader point  $R(x)$  for each  $x \in X$ ;
  - 2 Get the core points:  $\text{Core} = \{z_1, z_2, \dots, z_t\}$ ;
  - 3 Calculate the density branches:  $D(z_1), D(z_2), \dots, D(z_t)$ ;
  - 4 Calculate the expanded density branches:  $E(z_1), E(z_2), \dots, E(z_t)$ ;
  - 5 **foreach**  $(z_i, z_j) \in \text{Core} \times \text{Core}$  **and**  $z_i \neq z_j$  **do**
  - 6  $G[i, j] \leftarrow 0$ ;
  - 7  $\text{SNN}(z_i, z_j) \leftarrow E(z_i) \cap E(z_j)$ ;
  - 8 **if**  $|\text{SNN}(z_i, z_j)| > \tau \times \min(|E(z_i)|, |E(z_j)|)$  **then**
  - 9  $G[i, j] \leftarrow 1$ ;
  - 10 Get connected components  $\{C_i\}$  from  $G$ ;
  - 11 **return**  $\text{Core}$ ,  $\{D(z_i)\}$ ,  $\{C_i\}$ .
- 

If two samples  $x_i$  and  $x_j$  are in the same density branch then they are very similar, and we denote a *similarity* matrix by  $\text{Sim}^{(k)}$  for a given  $k$  to describe the similarity of those samples, such that

$$\text{Sim}_{ij}^{(k)} = \begin{cases} \exp\left(-\frac{\|x_i - x_j\|^2}{\sigma^2}\right) & \text{if } i \neq j, \text{ and } \exists z_t, x_i, x_j \in D(z_t); \\ 0 & \text{otherwise,} \end{cases} \quad (8)$$

360 where the parameter  $\sigma$  is set as 1% of the square of the largest distance between data points.  
 361 In fact, the performance of the proposed approach is not sensitive to this parameter, so we  
 362 will just apply this setting to all of the experiments.

363 On the other hand, if  $x_i$  and  $x_j$  are in different connected components then they should  
 364 be considered as dissimilar and dissimilarity will be characterized by a *separation* matrix.  
 365 The separation matrix is over core points only, instead of all the data points, as core points  
 366 are the most representative ones for separation. For core points, the separation matrix  $\text{Sep}^{(k)}$   
 367 is defined as follows.

1. When there is only one connected component found by Algorithm 1, the separation matrix  $\text{Sep}^{(k)}$  is defined as:

$$\text{Sep}^{(k)} = \frac{1}{|\text{Core}|^2} \sum_{(z_i, z_j) \in \text{Core} \times \text{Core}} (z_i - z_j)(z_i - z_j)^T. \quad (9)$$

2. When there are more than one connected components found by Algorithm 1, then we can calculate the separation matrix  $\text{Sep}^{(k)}$  over pairs of core points from different connected components. Denote by  $P^{(k)}$  the set of all the pairs of  $z_i, z_j$  from different connected components, then

$$\text{Sep}^{(k)} = \frac{1}{|P^{(k)}|} \sum_{(z_i, z_j) \in P^{(k)}} (z_i - z_j)(z_i - z_j)^T \quad (10)$$

The corresponding algorithm to calculate  $\text{Sim}^{(k)}$  and  $\text{Sep}^{(k)}$  is shown in Algorithm 2.

---

**Algorithm 2:** ComputeMatrices

---

**Input:** Dataset  $X$ ; number of nearest neighbors  $k$ ; threshold parameter  $\tau$ .

**Output:**  $\text{Sim}^{(k)}$ ,  $\text{Sep}^{(k)}$

```

1 Core,  $\{D(z_i)\}, \{C_i\} \leftarrow \text{ExploreStructure}(X, k, \tau)$ 
2  $\text{Sim}^{(k)} \leftarrow (0)_{m \times m}$ ;
3 foreach  $z \in \text{Core}$  do
4   | foreach  $(x_u, x_v) \in D(z) \times D(z)$  and  $x_u \neq x_v$  do
5   |   |  $\text{Sim}_{uv}^{(k)} \leftarrow \exp(-\frac{\|x_u - x_v\|^2}{\sigma^2})$ ;
6  $M \leftarrow \emptyset$ ;
7 if  $|\{C_i\}| > 1$  then
8   | foreach  $(z_u, z_v) \in \text{Core} \times \text{Core}$  do
9   |   |  $C_u \leftarrow$  the connected component containing  $z_u$ ;
10  |   |  $C_v \leftarrow$  the connected component containing  $z_v$ ;
11  |   | if  $C_u \neq C_v$  then
12  |   |   |  $M \leftarrow M \cup \{(z_u, z_v)\}$ ;
13 else
14 |  $M \leftarrow \text{Core} \times \text{Core}$ ;
15  $\text{Sep}^{(k)} = \frac{1}{|M|} \sum_{(z_u, z_v) \in M} (z_u - z_v)(z_u - z_v)^T$ ;
16 return  $\text{Sim}^{(k)}$ ,  $\text{Sep}^{(k)}$ .

```

---

368

369 *3.3. Dimensionality Reduction Algorithm*

370 As we have mentioned before, the value of  $k$  will influence the number of connected  
371 components and thus influence the structure of data revealed by these components. For  
372 example, we increase the value of  $k$  from  $k_1$  to  $k_2$  and the results are shown in Figure 6. When  
373  $k = k_1$ , there are 7 connected components; when  $k$  is increased to  $k_2$ , the two connected  
374 components  $C_{11}$  and  $C_{12}$  are merged, and the same happens to  $C_{15}$  and  $C_{16}$ ; when  $k$  is  
375 increased even larger to  $k_3$ , more components are merged and a higher level structure of the  
376 data is then revealed.

377 Inspired by this observation, we vary the value of  $k$  to explore the structure of data in  
378 different scales. Several similarity matrices and separation matrices will be obtained through  
379 the process, and we propose to combine these matrices by a weighted sum. As the similarity  
380 with a smaller  $k$  will be more significant than that with a larger  $k$ , we assign larger weight  
381 to  $\text{Sim}^{(k)}$  for a smaller value of  $k$ . On the other hand, we assign larger weight to  $\text{Sep}^{(k)}$  for  
382 a larger value of  $k$ , since the separation with a larger  $k$  will be more significant.

383 In particular, the similarity matrix and the separation matrix will be defined as follows.



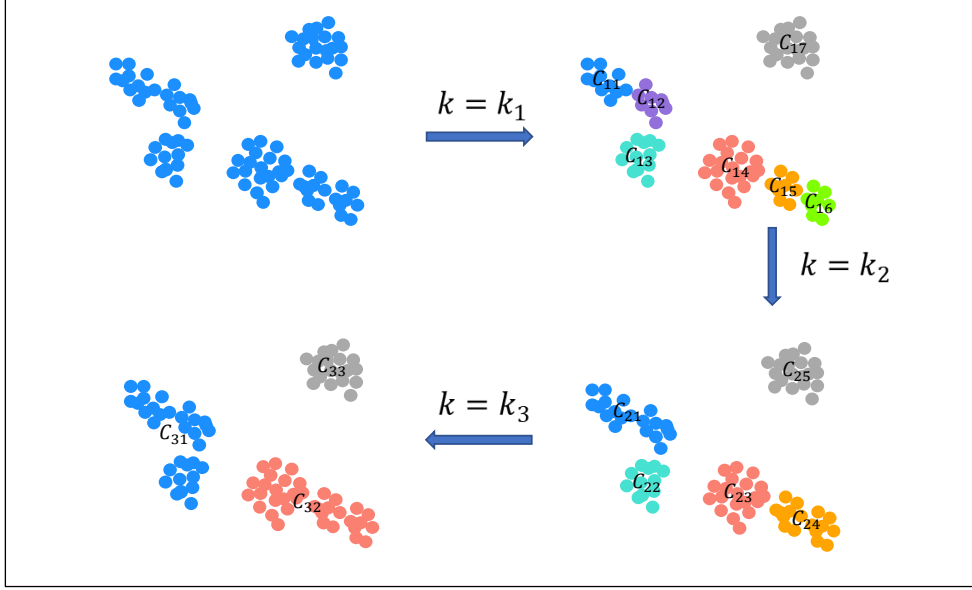


Figure 6: Illustration of the effect of different values of  $k$  on constructing connected components.

$$\text{Sim} = \sum_{i=1}^l \frac{\exp(\frac{1}{k_i})}{\sum_{i=1}^l \exp(\frac{1}{k_i})} \text{Sim}^{(k_i)}; \quad (11)$$

$$\text{Sep} = \sum_{i=1}^l \frac{\exp(k_i)}{\sum_{i=1}^l \exp(k_i)} \text{Sep}^{(k_i)} \quad (12)$$

384 The original model of LPP is thus improved by combining the idea in Model 4 and the  
 385 connectivity measures, as below.

$$\text{argmin}_a \frac{a^T X(L + L^*)X^T a}{a^T (XD^*X^T + \text{Sep}) a}, \quad \text{s.t.} \quad a^T (XD^*X^T + \text{Sep})a = 1. \quad (\text{ConLPP})$$

386 where  $X = (x_1 - \mu, \dots, x_m - \mu)$ ,  $L^* = D^* - \text{Sim}$ , and  $D^*$  is a diagonal matrix and  $D_{ii}^* =$   
 387  $\sum_{j=1}^{|X|} \text{Sim}_{ij}$ . Note that the constraint here is just to change the length of the projection  
 388 vectors and thus to scale the coordinates after projection (cf. Proposition 2). One can also  
 389 remove this constraint and use unit projection vectors instead to have orthogonal projections,  
 390 or any other lengths of interest.

391 The above optimization problem is a typical Rayleigh quotient problem and can be solved  
 392 by calculating eigenvectors. Let  $S_1 = X(L + L^*)X^T$  and  $S_2 = XD^*X^T + \text{Sep}$ . Then the  
 393 solution of the above optimization problem is the eigenvector corresponding to the smallest  
 394 eigenvalue of  $S_2^{-1}S_1$ , and this eigenvector is the projection direction for dimensionality re-  
 395 duction. When  $d$  directions are needed, the eigenvectors of the  $d$  smallest eigenvalues are  
 396 then used. The corresponding algorithm is shown in Algorithm 3.

---

**Algorithm 3: ConLPP**

---

**Input:** Dataset  $X = (x_1, \dots, x_m)$ ; the range of the numbers of nearest neighbors  $[k_0, k_1]$ ; target dimensionality  $d$ ; threshold parameter  $\tau$ .

**Output:** Projection directions  $A = (a_1, \dots, a_d)$ .

- 1  $X \leftarrow X - \mu$ , where  $\mu$  is the average vector of  $X$ ;
  - 2 **foreach**  $k \in [k_0, k_1]$  **do**
  - 3     |  $\text{Sim}^{(k)}, \text{Sep}^{(k)} \leftarrow \text{ComputeMatrices}(X, k, \tau)$ ;
  - 4     | Calculate Sim and Sep by Equations 11 and 12;
  - 5     | Calculate  $L$ , which is consistent with that of LPP;
  - 6     | Calculate  $D^* = (D_{ii}^*)_{m \times m}$ , where  $D_{ii}^* = \sum_{j=1}^{|X|} \text{Sim}_{ij}$ ;
  - 7     | Calculate  $L^* = D^* - \text{Sim}$ ;
  - 8     | Calculate  $S_1 = X(L + L^*)X^T$ ;
  - 9     | Calculate  $S_2 = XD^*X^T + \text{Sep}$ ;
  - 10    | Calculate the eigenvalues and eigenvectors of  $S_2^{-1}S_1$ , and then sort them in ascending order;
  - 11    | Select the eigenvectors corresponding to the first  $d$  non-zero minimum eigenvalues:  $\lambda_1, \lambda_2, \dots, \lambda_d$  to make up  $A = (\frac{a_1}{\sqrt{a_1^T S_2 a_1}}, \dots, \frac{a_d}{\sqrt{a_d^T S_2 a_d}})$ ;
  - 12 **return**  $A$ .
- 

### 3.4. Time Complexity Analysis and Comparison

397 Suppose the number of data points in a dataset is  $m$ , the original dimension of data is  
398  $n$ , the number of core points is  $t$ , and the number of nearest neighbors  $k$  is in range  $[k_0, k_1]$ ,  
399 where the length of the range is  $l$ .

400 In Algorithm 1, to get  $N_k(x)$ ,  $\rho(x)$ ,  $R(x)$ , and core points in lines 1 and 2, one can make  
401 use of a KD-tree, and the time complexity is  $\mathcal{O}((n+k)m \log m)$ . To obtain density branches,  
402 one can descendingly sort the points by their densities  $\rho(x)$ , and go through the points one  
403 by one to allocate them to the density branch that their leader point belongs to (if the leader  
404 point of a point is itself, then this point is a core point and it will be allocated to a new  
405 density branch). So line 3 of Algorithm 1 takes  $\mathcal{O}(m \log m + m)$  time. Line 4 expands the  
406 density branches by including the  $k$ -nearest neighbors of each point in the density branches,  
407 so its time complexity is  $\mathcal{O}(km)$ . For two expanded density branches, one can obtain their  
408 intersection by scanning through one of them and checking if each point of it is also in the  
409 other one, which takes  $\mathcal{O}(m)$  time in the worst case. In total, lines 5-9 take  $\mathcal{O}(t^2 m)$  time.  
410 Line 10 can obtain connected components with a depth first search in time  $\mathcal{O}(t^2)$ . Therefore,  
411 the time complexity of Algorithm 1 is  $\mathcal{O}((\log m + k + t^2)m)$ .

412 For Algorithm 2, line 1 exploits Algorithm 1 and has time complexity  $\mathcal{O}((\log m + k + t^2)m)$ .  
413 For lines 3-5, it takes  $\mathcal{O}(n)$  time to calculate each  $\text{Sim}_{uv}^{(k)}$ , and each core point will take  
414  $\mathcal{O}(m_s^2)$  steps to obtain all  $\text{Sim}_{uv}^{(k)}$ , where  $m_s$  is the number of points in density branch  $z_s$  and  
415  $\sum_{s=1}^t m_s = m$ . So the total time for lines 3-5 is  $\mathcal{O}(nm^2)$ . Lines 6-14 take  $\mathcal{O}(t^2)$  time, and  
416 line 15 takes  $\mathcal{O}(n^2 t^2)$  time. So in total, Algorithm 2 takes  $\mathcal{O}((\log m + k + t^2)m + nm^2 + n^2 t^2)$ .

418 The main algorithm ConLPP, i.e. Algorithm 3, basically has similar steps to LPP, with  
419 additional steps for centralizing  $X$  and obtaining Sim and Sep in lines 1-4, where these  
420 additional steps take  $\mathcal{O}(mn + l((\log m + k + t^2)m + nm^2 + n^2t^2))$  time. Lines 5-7 take  $\mathcal{O}(m^2)$   
421 time to obtain  $L$ ,  $D^*$ , and  $L^*$ . Lines 8-9 take  $\mathcal{O}(nm^2 + n^2m)$  time. Lines 10-11 take  $\mathcal{O}(n^3)$   
422 to obtain the final projection directions. In summary, the time complexity of ConLPP is  
423  $\mathcal{O}(mn + l((\log m + k + t^2)m + nm^2 + n^2t^2) + m^2 + nm^2 + n^2m + n^3)$ . As  $k$  and  $l$  is usually very  
424 small compared to  $m$  and  $n$ , the time complexity can be simplified to  $\mathcal{O}((m + n^2)t^2 + nm^2 +$   
425  $n^2m + n^3)$ , or  $\mathcal{O}(nm^2)$  if  $n \ll m$ . This is comparable to LPP, which has the time complexity  
426 of  $\mathcal{O}(nm^2 + n^2m + n^3)$ , and is also comparable to many other dimensionality reduction  
427 methods, like t-SNE, LE, and LLE, whose time complexity is  $\mathcal{O}(nm^2)$ . For example, on  
428 a dataset of 2000 points and 649 features, for a target dimensionality of 5, ConLPP takes  
429 3.98s, and LPP takes 0.54s, whereas LLE takes 0.23s, LE takes 45s, and t-SNE takes 178s.

## 430 4. Empirical Evaluations

431 In order to evaluate the proposed model, we conduct experiments on both synthetic and  
432 real-world datasets, against seven dimensionality reduction baseline methods, i.e., LLE [18],  
433 PCA [15], LDA [16], LE [21], t-SNE [19], LPP [26], ILPP-L1 [32], and LAPP [33], where  
434 LLE, LE, and t-SNE are non-linear methods, and LPP, ILPP-L1, and LAPP are LPP based  
435 methods. LDA is supervised and is used to demonstrate that ConLPP can capture inherent  
436 structures of data. Note that t-SNE, LLE, PCA, and LDA here are from the scikit-learn  
437 library and use default parameters; LE, ILPP-L1, and LAPP are implemented according to  
438 the corresponding papers. The range of  $k$  for ConLPP is set to [5, 15].

### 439 4.1. Datasets

440 There are one synthetic dataset and 14 datasets used in the experiments. The synthetic  
441 dataset is used to show that ConLPP model can repair the problem of LPP for translation  
442 invariance and can further maintain data connectivity and separation after dimensionality  
443 reduction. Real-world datasets are used to illustrate the advantages of ConLPP over other  
444 algorithms.

445 The synthetic dataset is visually illustrated in Figure 7. It is sampled from a multivariate  
446 normal distribution, and contains three clusters of points, each of which has 200 points. The  
447 details of the 14 real-world datasets are shown in Table 1. Among them, ORL, USPS, COIL-  
448 20, FashionMNIST, and CIFAR-10 are image datasets, where the last two are testing subsets  
449 of the original datasets.

450 All datasets are standardized, so that the mean value of a dataset is 0 and the standard  
451 deviation is 1. Note that, in this case, the original LPP and Model 4 are equivalent because  
452  $\mu = 0$ .

### 453 4.2. Results on Synthetic Datasets

454 As illustrated in Figure 2, Model 4 reduces the influence of the position of the origin  
455 on LPP and fixes translation invariance. However, from Figure 3, both LPP and Model 4

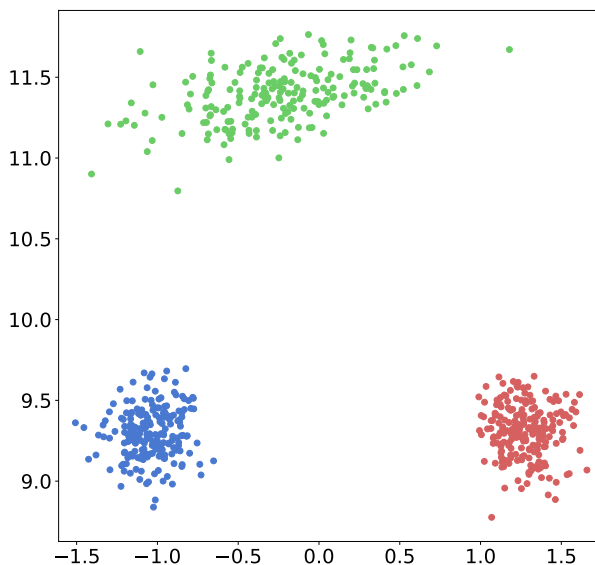


Figure 7: Visualization of the synthetic dataset.

456 might project clusters that are originally separated into overlapping ones, because they do  
 457 not consider higher level structural information.

458 Figure 8 shows the results of LPP (equivalently, Model 4, as data were standardized) and  
 459 ConLPP into one-dimensional space on the synthetic dataset after standardization. It can  
 460 be seen that LPP makes the originally separated clusters overlapping after dimensionality  
 461 reduction, while ConLPP can maintain the separation state of these clusters after projec-  
 462 tion. The main reason is that the data contain three connected components and ConLPP  
 463 maintains this structure by introducing characterizations of connectivity information into  
 464 the optimization model.

#### 465 4.3. Visualization on Real-World Datasets

466 Visualization can give a direct and intuitive impression on the performance of maintaining  
 467 the structure of complex high-dimensional data [42]. Figure 9 is the visualization of the  
 468 projection of the points for digits 0 to 9 in the dataset “digital” into two-dimensional plane  
 469 by various algorithms. By comparing ConLPP and LPP, we can find that ConLPP better  
 470 maintains the separation between different clusters and the closeness of data from the same  
 471 cluster, e.g., LPP mixes “2”, “3”, and “8” clusters while ConLPP does not. Note that  
 472 LDA is a supervised dimensionality reduction method, whereas ConLPP is unsupervised  
 473 but has a quite similar results as LDA. This means ConLPP can indeed capture inherent  
 474 structure of the data. The result of t-SNE for visualization is good, but it has several  
 475 disadvantages which can be seen from later experiments: its performance worsens for higher  
 476 target dimensionalities, and it cannot give the projection function which means it cannot  
 477 deal with unseen data points. For the other methods, the advantage of ConLPP is obvious.

Table 1: Details of real-world datasets.

Dataset	#instances	#features	#classes
indianliver	583	10	2
congressEW	435	16	2
vote	435	16	2
sonarEW	208	60	2
ORL	400	1024	40
digital	2000	649	10
pengcolonEW	62	2000	2
segmentation	2100	19	7
mfeat-fou	2000	76	10
mfeat-kar	2000	64	10
USPS	9298	256	10
COIL-20	1440	1024	20
FashionMNIST	10000	784	10
CIFAR-10	10000	3072	10

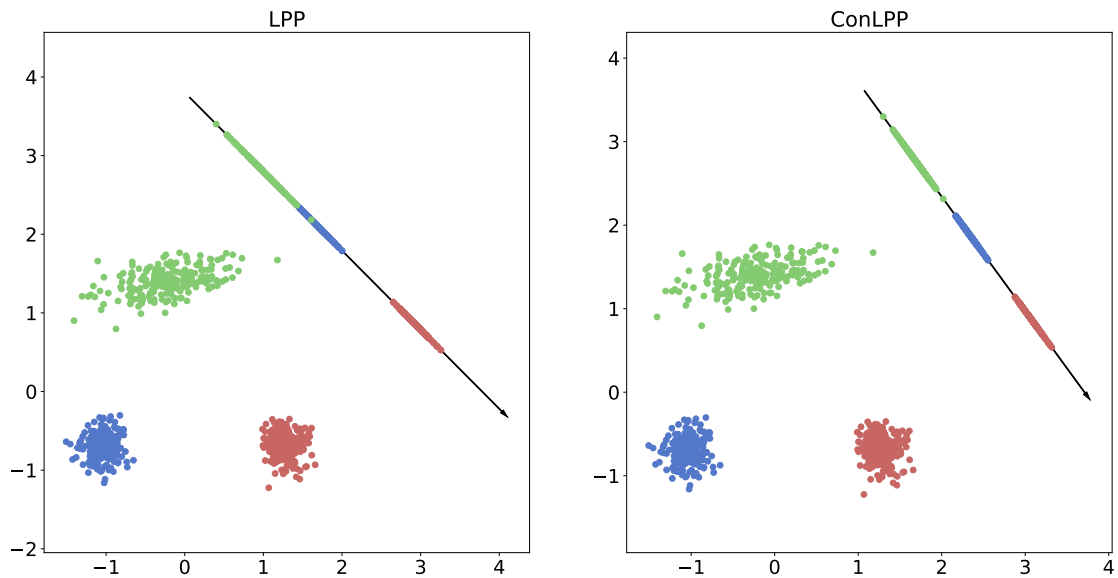


Figure 8: Projection results of the standardized synthetic dataset by LPP and ConLPP, respectively.

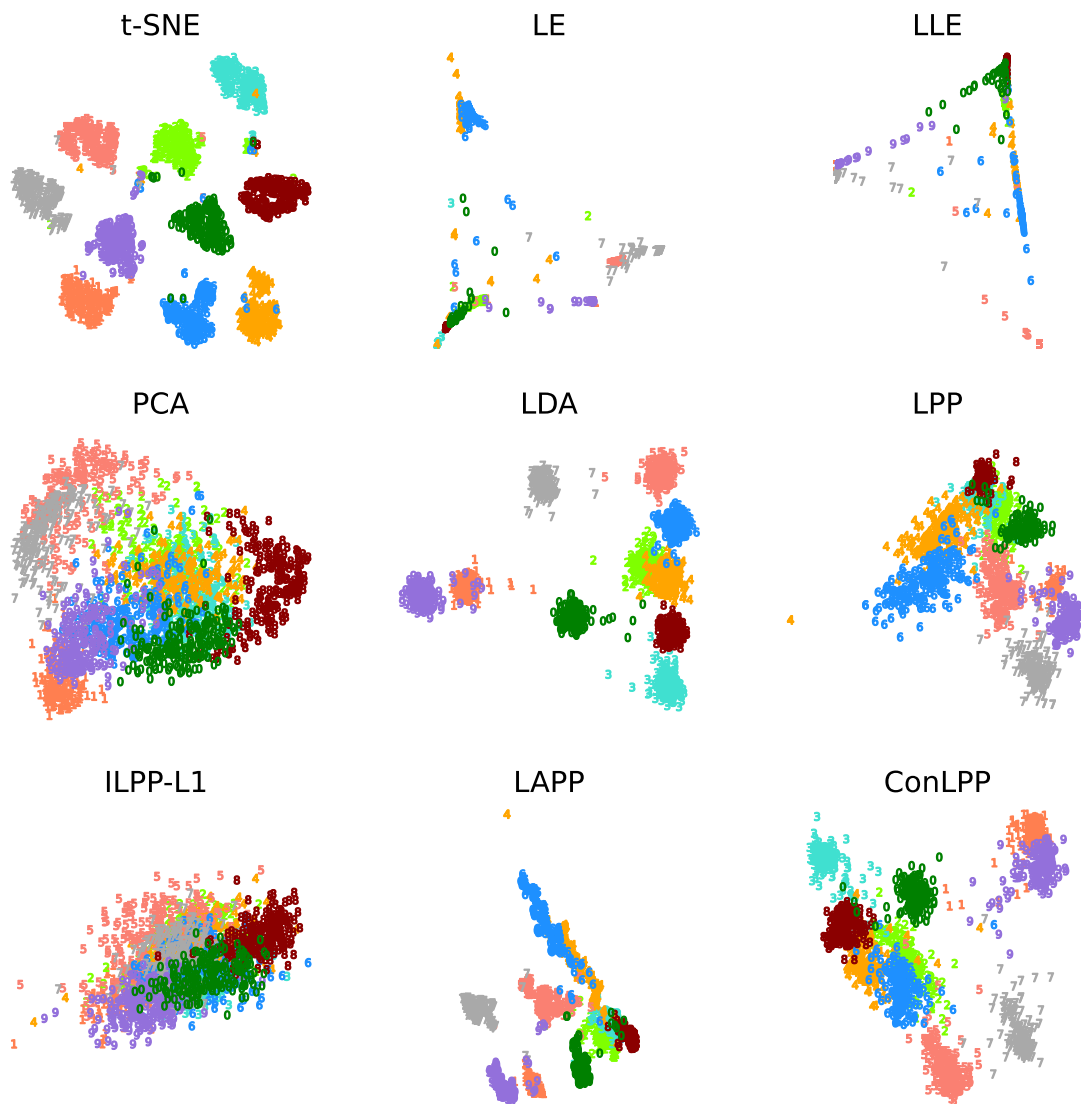


Figure 9: Visualization results of dimensionality reduction methods for mapping data points representing “0”-“9” in the dataset “digital” into two-dimensional space.

Table 2: Results on the real-world datasets in terms of 1NN classification accuracy.

Dataset	t-SNE	LE	LLE	PCA	LPP	ILPP-L1	LAPP	ConLPP
indianliver	66.21(1)	67.40(9)	64.14(8)	65.54(4)	66.22(7)	68.45(4)	66.57(2)	<b>68.46(4)</b>
congressEW	92.86(2)	91.48(14)	89.65(7)	92.41(13)	91.73(11)	<b>94.71(10)</b>	93.09(15)	93.33(5)
vote	93.09(2)	93.09(7)	89.42(8)	91.94(14)	91.26(8)	94.47(13)	94.03(8)	<b>94.70(6)</b>
sonarEW	62.92(2)	65.93(6)	69.70(19)	70.18(4)	63.07(5)	69.23(11)	68.90(14)	<b>71.08(19)</b>
ORL	92.24(2)	77.73(18)	86.48(19)	<b>93.00(19)</b>	90.74(19)	73.50(18)	60.75(13)	92.49(9)
digital	97.80(2)	98.10(12)	97.70(13)	97.75(19)	97.05(18)	87.30(19)	98.35(19)	<b>98.45(12)</b>
pengcolonEW	71.03(2)	71.23(10)	72.82(18)	71.23(4)	72.82(5)	66.27(1)	71.43(13)	<b>80.75(1)</b>
segmentation	93.19(13)	88.71(11)	88.62(11)	93.48(13)	93.57(10)	94.00(16)	<b>94.10(7)</b>	93.76(10)
mfeat-fou	79.55(3)	74.65(16)	76.20(16)	75.15(19)	74.65(15)	72.75(19)	<b>81.45(17)</b>	79.90(17)
mfeat-kar	96.00(4)	94.30(19)	94.9(16)	94.75(18)	91.20(14)	87.65(19)	96.15(15)	<b>96.45(15)</b>
USPS	<b>96.91(2)</b>	94.49(19)	96.10(16)	96.20(19)	89.85(19)	N/A	83.54(17)	94.67(19)
COIL-20	97.22(4)	89.03(15)	92.5(18)	95.97(18)	93.75(15)	85.77(18)	84.86(15)	<b>99.72(11)</b>
FashionMNIST	<b>81.40(9)</b>	70.58(15)	74.75(19)	79.40(19)	76.60(19)	N/A	73.14(19)	77.71(19)
CIFAR-10	30.3(14)	21.76(15)	22.69(17)	<b>30.50(19)</b>	22.25(18)	N/A	20.31(19)	19.45(19)

#### 4.4. Evaluation with 1NN on Results Obtained After Dimensionality Reduction

In this subsection, we first perform dimensionality reduction on the whole dataset, for each of the 14 real-world datasets, and then use 1NN [43] to conduct classification on the dimensionality reduction results. The motivation is that the performance of 1NN can reflect the quality of dimensionality reduction results, illustrating whether points from the same class are relatively close and points from different classes are relatively far away from each other.

For each algorithm, the best dimensionality for dimensionality reduction on each dataset is determined by first performing dimensionality reduction with target dimensionality from 2 to 19 and calculating the average accuracy of 1NN in 10-fold cross validation for the transformed data, and then selecting the dimensionality with the highest accuracy as the best one. Once the best dimensionality is determined, the corresponding average accuracy and standard deviation over ten folds is recorded. For ILPP-L1, because its results might vary largely for different runs, we perform the aforementioned process ten times for each dimensionality and use the highest score in these ten runs as the score of that dimensionality.

The results are shown in Table 2, where the numbers in braces are the best dimensionalities for corresponding algorithms and “N/A” means the corresponding algorithm cannot finish in the time limit of 10 hours.

From Table 2, we can see that on all of the 14 datasets, the average classification accuracy of the transformed data by ConLPP is significantly higher than LPP, and is the best in 7 out of 14 cases and close to the best in most of the other cases. Note that we select the best result of ILPP-L1 in 10 runs, so it is not surprising it is better than ConLPP on the dataset “congressEW”, where ConLPP outperforms all the others except ILPP-L1; for the datasets “segmentation” and “mfeat-fou” where LAPP is the best, we can see that ConLPP is very close to the best, and ConLPP actually outperforms LAPP on most of the other datasets and the advantage is very significantly on several datasets, e.g., ORL, pengcolonEW, USPS, and COIL-20.

505 We also illustrate the results of each algorithm for different target dimensionalities on  
 506 several datasets in Figure 10. From Figure 10, we observe that ConLPP is dominating for  
 507 almost all of the dimensionalities. For ORI and vote, ConLPP shows more stability than  
 508 the others across different dimensionalities. For mfeat-fou and mfeat-kar, although ConLPP  
 509 is not the best in lower dimensionalities, it quickly becomes evidently better than or close  
 510 to the others in higher dimensionalities. These demonstrate the effectiveness of the idea of  
 511 ConLPP for identifying topological structures.

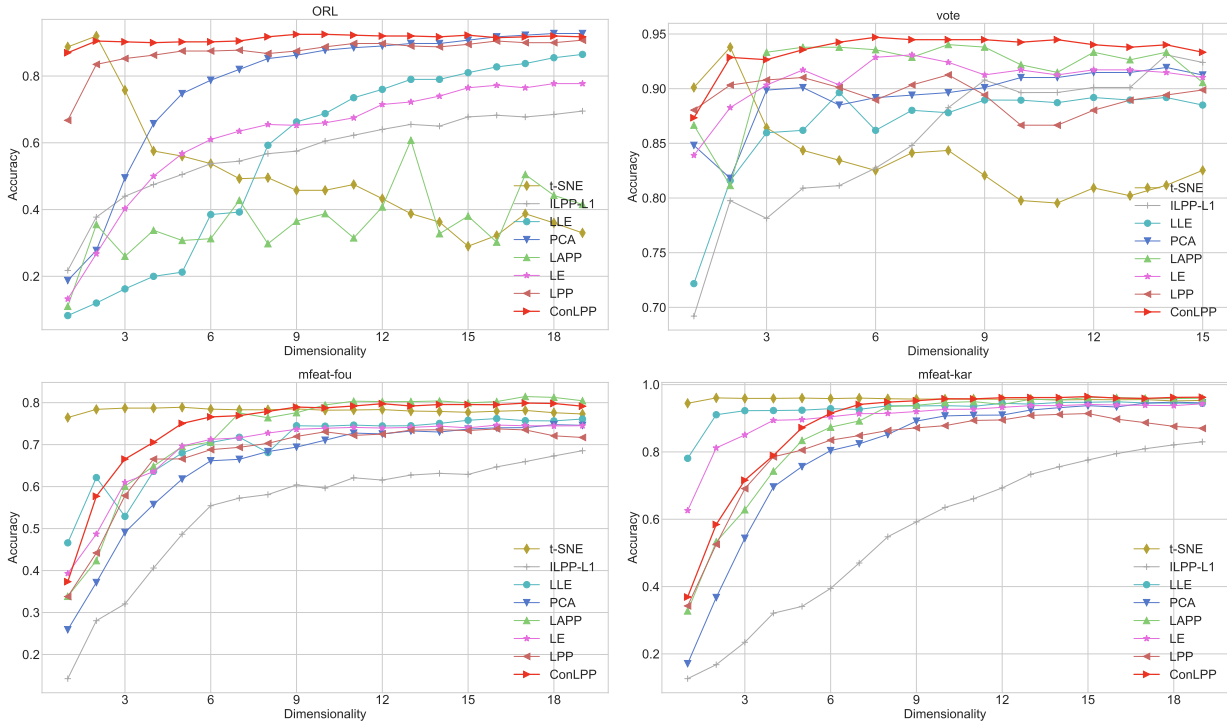


Figure 10: Results on four real-world datasets w.r.t. different target dimensionalities in terms of 1NN classification accuracy.

#### 512 4.5. Evaluation of Generalization Ability

513 In this subsection, instead of finding projection directions based on the whole dataset,  
 514 we first partition the dataset into training set and testing set, and then find projection  
 515 directions on training set and apply them on testing set. The transformed testing set is then  
 516 used to perform 1NN classification with 10-fold cross validation. The partition of training set  
 517 and testing set is also performed 10 times, and the results are averaged. The motivation of  
 518 the experiments here is to evaluate the generalization ability of the proposed dimensionality  
 519 reduction model.

520 Particularly, ConLPP is compared against four linear ones of the baseline methods, i.e.,  
 521 PCA, LPP, ILPP-L1, and LAPP, on the dataset “digital” of handwritten numbers. The  
 522 other methods used in this article are not applicable as they have to be based on the whole



Table 3: Accuracy results on the dataset “digital” for 1NN classification on transformed testing sets by different dimensionality reduction methods.

Method	20%	40%	60%	80%
PCA	<b>97.38(19)</b>	<b>97.28(19)</b>	96.72(19)	96.05(19)
LPP	42.06(15)	64.25(9)	93.23(11)	95.43(12)
ILPP-L1	85.49(19)	82.92(19)	81.49(19)	80.27(19)
LAPP	28.31(17)	92.00(15)	97.13(19)	97.49(14)
ConLPP	39.18(8)	87.17(10)	<b>97.17(11)</b>	<b>97.65(11)</b>

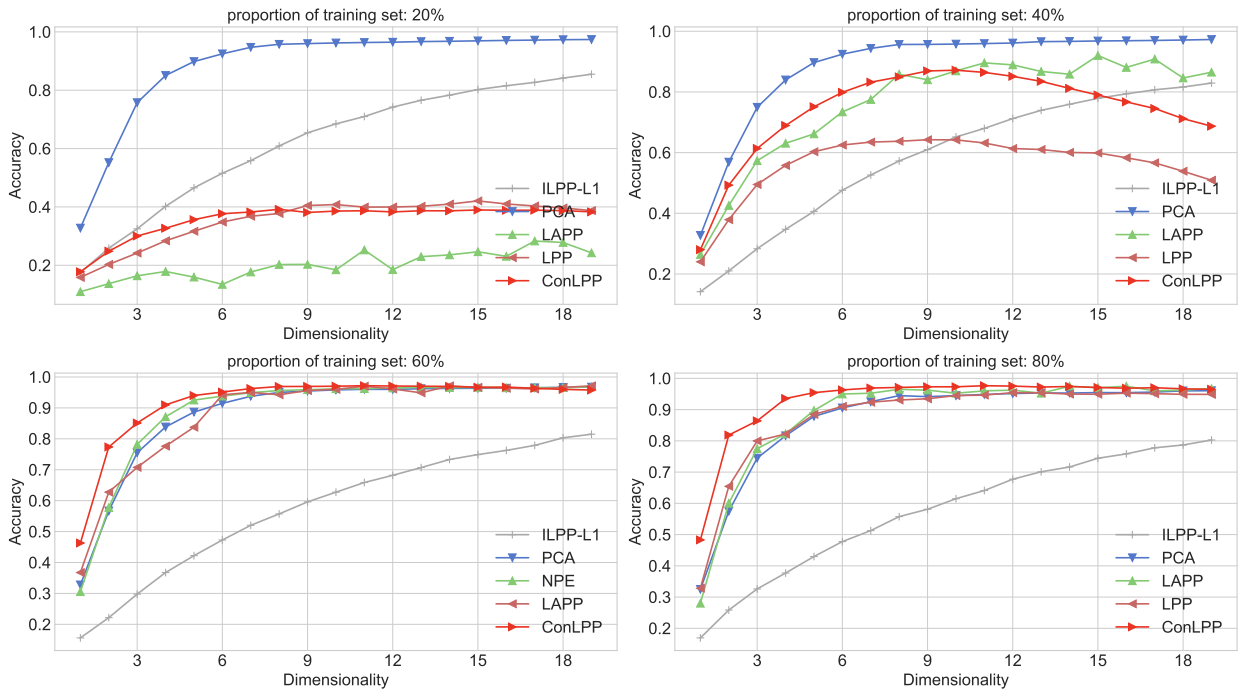


Figure 11: Results on the dataset “digital” w.r.t. different target dimensionalities in terms of 1NN classification accuracy on different portions of transformed testing sets.

523 dataset. The proportion of training set is set to be 20%, 40%, 60% and 80%, respectively.  
 524 Table 3 and Figure 11 show the results.

525 From Table 3, it is easy to see that ConLPP can have better results with larger proportion  
 526 of training samples and it performs better than LPP, ILPP-L1, LAPP in these cases. It is  
 527 worth noting that the performance of ConLPP is not ideal when the proportion of training  
 528 samples is very small (20%). This is expected, as in this case the training samples cannot  
 529 accurately reflect the structure of the whole dataset, and the confidence of the structure  
 530 found by ConLPP will be low in this case. This phenomenon becomes clearer in Figure  
 531 11. With relatively larger proportion of training samples, ConLPP can better identify the  
 532 structure of the dataset. Note that for the 40% case, the performance of both LPP and  
 533 ConLPP becomes worse when the target dimensionality is over 10. The reason is that the

534 constraint in the model scales the projection coordinates differently in different dimension,  
 535 which stretches the projected data and affects the performance of 1NN on the transformed  
 536 data. This can be avoided by discarding the constraint but using unit projection vectors.  
 537 In summary, ConLPP can have better generalization ability than LPP, ILPP-L1, LAPP,  
 538 and PCA with relatively larger proportion of training samples (e.g., over 60%), and it is an  
 539 effective model for dimensionality reduction.

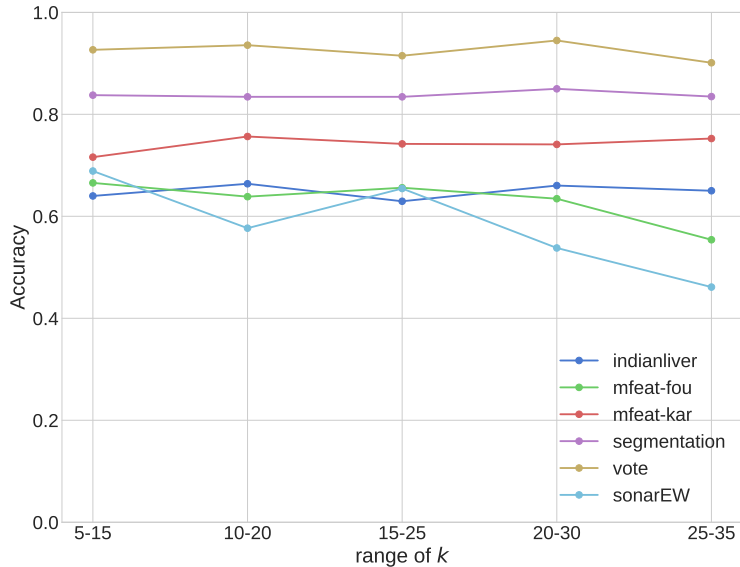


Figure 12: Illustration of the robustness of ConLPP w.r.t. different ranges of  $k$ .

#### 540 4.6. Parameter Analysis

541 Sensitivity to parameters could limit the applicability of a dimensionality reduction  
 542 method in practice. In the proposed ConLPP model, we would like to see how the choice  
 543 of the ranges of  $k$  will affect the accuracy of ConLPP. Figure 12 shows the results of 1NN  
 544 classification on the transformed data of the datasets indianliver, mfeat-fou, mfeat-kar, seg-  
 545 mentation and vote, w.r.t. different ranges of  $k$ , which are  $[5, 15]$ ,  $[10, 20]$ ,  $[15, 25]$ ,  $[20, 30]$ ,  
 546 and  $[25, 35]$ , respectively. The results illustrate that the performance of ConLPP is relatively  
 547 stable for different choices of parameters, which is favorable for practical applications.

### 548 5. Conclusion

549 This paper proposed a novel model for dimensionality reduction by improving the well-  
 550 known LPP model to maintain topological properties including translation invariance and  
 551 topological connectivity. We demonstrated and analysed the translation sensitivity problem  
 552 of LPP, which, to the best of our knowledge, was not noticed in the literature. We proposed  
 553 a new model to repair this problem and theoretically proved its effectiveness. We also  
 554 proposed to take more topological properties into consideration and devised a novel model  
 555 that considers topological connectivity. Extensive experimental results on synthetic and

556 real-world datasets demonstrated the effectiveness and superiority of the new model over  
557 the original LPP model and several other widely used dimensionality reduction models. In  
558 the future, we plan to incorporate more topological properties to the model so that it can  
559 discover more complex structures.

## 560 Acknowledgment

561 This work was supported by the National Natural Science Foundation of China [grant  
562 number 61806170]; the Humanities and Social Sciences Fund of Ministry of Education [grant  
563 number 18XJC72040001]; and the National Key Research and Development Program of  
564 China [grant number 2019YFB1706104].

## 565 References

- 566 [1] S. Yan, D. Xu, B. Zhang, H. Zhang, Q. Yang, S. Lin, Graph embedding and extensions: A general  
567 framework for dimensionality reduction, *IEEE Transactions on Pattern Analysis and Machine Intelli-*  
568 *gence* 29 (1) (2007) 40–51.
- 569 [2] R. He, N. Xiong, L. T. Yang, J. H. Park, Using multi-modal semantic association rules to fuse keywords  
570 and visual features automatically for web image retrieval, *Information Fusion* 12 (3) (2011) 223–230.
- 571 [3] R. Zhou, X. Wang, J. Wan, N. Xiong, EDM-Fuzzy: An Euclidean distance based multiscale fuzzy  
572 entropy technology for diagnosing faults of industrial systems, *IEEE Transactions on Industrial Infor-*  
573 *matics* 17 (6) (2021) 4046–4054.
- 574 [4] H. Li, T. Jiang, K. Zhang, Efficient and robust feature extraction by maximum margin criterion, *IEEE*  
575 *Transactions on Neural Networks* 17 (1) (2006) 157–165.
- 576 [5] M. Unser, M. Eden, Multiresolution feature extraction and selection for texture segmentation, *IEEE*  
577 *Transactions on Pattern Analysis and Machine Intelligence* 11 (7) (1989) 717–728.
- 578 [6] J. Yan, B. Zhang, N. Liu, S. Yan, Q. Cheng, W. Fan, Q. Yang, W. Xi, Z. Chen, Effective and effi-  
579 cient dimensionality reduction for large-scale and streaming data preprocessing, *IEEE Transactions on*  
580 *Knowledge and Data Engineering* 18 (3) (2006) 320–333.
- 581 [7] A. Blum, P. Langley, Selection of relevant features and examples in machine learning, *Artificial Intel-*  
582 *ligence* 97 (1-2) (1997) 245–271.
- 583 [8] D. Koller, M. Sahami, Toward optimal feature selection, in: *International Conference on Machine*  
584 *Learning*, 1996, pp. 284–292.
- 585 [9] X. He, D. Cai, S. Yan, H. Zhang, Neighborhood preserving embedding, in: *IEEE International Confer-*  
586 *ence on Computer Vision*, 2005, pp. 1208–1213.
- 587 [10] G. D. C. Cavalcanti, T. I. Ren, J. F. Pereira, Weighted modular image principal component analysis  
588 for face recognition, *Expert Systems with Applications* 40 (12) (2013) 4971–4977.
- 589 [11] E. Silva Jr, A. Oliveira, W. Santos, C. Mello, C. Zanchettin, G. D. C. Cavalcanti, Feature selection and  
590 model design through GA applied to handwritten digit recognition from historical document images,  
591 in: *International Conference on Frontiers in Handwritten Recognition*, 2008, pp. 562–567.
- 592 [12] B. Li, D. Zhang, K. Wang, Online signature verification based on null component analysis and principal  
593 component analysis, *Pattern Analysis and Applications* 8 (4) (2006) 345–356.
- 594 [13] H. Gunduz, An efficient dimensionality reduction method using filter-based feature selection and varia-  
595 tional autoencoders on Parkinson’s disease classification, *Biomedical Signal Processing and Control* 66  
596 (2021) 102452.
- 597 [14] J. Han, Z. Ge, Effect of dimensionality reduction on stock selection with cluster analysis in different  
598 market situations, *Expert Systems with Applications* 147 (2020) 113226.
- 599 [15] I. T. Jolliffe, *Principal Component Analysis*, Springer Series in Statistics, Springer, 1986.
- 600 [16] R. A. Johnson, D. W. Wichern, *Applied Multivariate Statistical Analysis*, Vol. 6, Pearson London,  
601 2014.

- 602 [17] S. Chen, C. H. Q. Ding, B. Luo, Linear regression based projections for dimensionality reduction,  
603 *Information Sciences* 467 (2018) 74–86.
- 604 [18] S. T. Roweis, L. K. Saul, Nonlinear dimensionality reduction by locally linear embedding, *Science*  
605 290 (5500) (2000) 2323–2326.
- 606 [19] L. van der Maaten, G. Hinton, Visualizing data using t-SNE, *Journal of Machine Learning Research*  
607 9 (86) (2008) 2579–2605.
- 608 [20] M. Balasubramanian, E. L. Schwartz, The Isomap algorithm and topological stability, *Science*  
609 295 (5552) (2002) 7–7.
- 610 [21] M. Belkin, P. Niyogi, Laplacian eigenmaps and spectral techniques for embedding and clustering, in:  
611 *Advances in Neural Information Processing Systems*, 2001, pp. 585–591.
- 612 [22] Y. Li, Y. Chai, H. Zhou, H. Yin, A novel dimension reduction and dictionary learning framework for  
613 high-dimensional data classification, *Pattern Recognition* 112 (2021) 107793.
- 614 [23] J. Chen, L. Liao, W. Zhang, L. Du, Mixture factor analysis with distance metric constraint for dimen-  
615 sionality reduction, *Pattern Recognition* 121 (2022) 108156.
- 616 [24] X. Peng, D. Xu, D. Chen, Robust distribution-based nonnegative matrix factorizations for dimension-  
617 ality reduction, *Information Sciences* 552 (2021) 244–260.
- 618 [25] Z. Chen, A. Fu, R. H. Deng, X. Liu, Y. Yang, Y. Zhang, Secure and verifiable outsourced data dimension  
619 reduction on dynamic data, *Information Sciences* 573 (2021) 182–193.
- 620 [26] X. He, P. Niyogi, Locality preserving projections, in: *Advances in Neural Information Processing*  
621 *Systems*, 2003, pp. 153–160.
- 622 [27] J. Li, J. Pan, S. Chu, Kernel class-wise locality preserving projection, *Information Sciences* 178 (7)  
623 (2008) 1825–1835.
- 624 [28] E. R. Silva, G. D. C. Cavalcanti, T. I. Ren, Class-wise feature extraction technique for multimodal  
625 data, *Neurocomputing* 214 (2016) 1001–1010.
- 626 [29] J. Gui, W. Jia, L. Zhu, S. Wang, D. Huang, Locality preserving discriminant projections for face and  
627 palmprint recognition, *Neurocomputing* 73 (13-15) (2010) 2696–2707.
- 628 [30] W. Yu, X. Teng, C. Liu, Face recognition using discriminant locality preserving projections, *Image and*  
629 *Vision Computing* 24 (3) (2006) 239–248.
- 630 [31] Y. Tang, Z. Zhang, Y. Zhang, F. Li, Robust L1-norm matrixed locality preserving projection for  
631 discriminative subspace learning, in: *International Joint Conference on Neural Networks*, 2016, pp.  
632 4199–4204.
- 633 [32] W. Yu, R. Wang, F. Nie, F. Wang, Q. Yu, X. Yang, An improved locality preserving projection with  
634  $l_1$ -norm minimization for dimensionality reduction, *Neurocomputing* 316 (2018) 322–331.
- 635 [33] A. Wang, S. Zhao, J. Liu, J. Yang, L. Liu, G. Chen, Locality adaptive preserving projections for linear  
636 dimensionality reduction, *Expert Systems with Applications* 151 (2020) 113352.
- 637 [34] L. Wasserman, Topological data analysis, *Annual Review of Statistics and Its Application* 5 (2018)  
638 501–532.
- 639 [35] I. Chevyrev, V. Nanda, H. Oberhauser, Persistence paths and signature features in topological data  
640 analysis, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42 (1) (2020) 192–202.
- 641 [36] H. Elhamdadi, S. Canavan, P. Rosen, AffectiveTDA: Using topological data analysis to improve analysis  
642 and explainability in affective computing, *IEEE Transactions on Visualization and Computer Graphics*  
643 28 (1) (2022) 769–779.
- 644 [37] W. Guo, K. Manohar, S. L. Brunton, A. G. Banerjee, Sparse-TDA: Sparse realization of topological  
645 data analysis for multi-way classification, *IEEE Transactions on Knowledge and Data Engineering*  
646 30 (7) (2018) 1403–1408.
- 647 [38] T. Feng, J. I. Dávila, Y. Liu, S. Lin, S. Huang, C. Wang, Semi-supervised topological analysis for  
648 elucidating hidden structures in high-dimensional transcriptome datasets, *IEEE/ACM Transactions on*  
649 *Computational Biology and Bioinformatics* 18 (4) (2021) 1620–1631.
- 650 [39] A. J. Ma, P. C. Yuen, W. W. W. Zou, J. Lai, Supervised spatio-temporal neighborhood topology  
651 learning for action recognition, *IEEE Transactions on Circuits and Systems for Video Technology*  
652 23 (8) (2013) 1447–1460.

- 653 [40] J. R. Munkres, *Topology*, Pearson, 2000.
- 654 [41] D. Cheng, S. Zhang, J. Huang, Dense members of local cores-based density peaks clustering algorithm,  
655 *Knowledge-Based Systems* 193 (2020) 105454.
- 656 [42] I. S. Dhillon, D. S. Modha, W. S. Spangler, Class visualization of high-dimensional data with applica-  
657 tions, *Computational Statistics and Data Analysis* 41 (1) (2002) 59–90.
- 658 [43] N. Kwak, Principal component analysis based on L1-norm maximization, *IEEE Transactions on Pattern*  
659 *Analysis and Machine Intelligence* 30 (9) (2008) 1672–1680.