

An Improved Isomap Algorithm Based on LP-Centers

W.W. Chen, W.J. Mao

State Key Laboratory of Industrial Control Technology
Institute of Cyber-Systems and Control
Zhejiang University
Hangzhou, P.R.China

Abstract—The isometric feature mapping (Isomap) is one of the widely used low-dimensional embedding methods. However, Isomap has three bottlenecks: sensitivity to outliers, high computational complexity and topological instability caused by short-circuit edges. To address above problems, by introducing the concepts of linear patch (LP) and k-nearest center (k-NC), a LP-center based Isomap (LPC-Isomap) approach is proposed in this paper. Firstly, before constructing a neighborhood graph, outliers are eliminated by LPC-Isomap. Then LP-centers are selected reasonably as landmark points by LPC-Isomap to reduce computational complexity. We compute the intersections between k-NC of each pair of points. If an intersection is empty, the edge connecting the corresponding pair of points is considered as a short-circuit edge. Lastly, the useful behavior of LPC-Isomap is verified through numerical examples with several data sets.

Keywords—Isomap; linear patch (LP); k-nearest center (k-NC); LPC-Isomap

I. INTRODUCTION

Isomap is one of representative algorithms [1-3] in manifold learning community, which extends classical multidimensional scaling (MDS) by exploiting geodesic distances instead of Euclidean distances. In contrast to previous algorithms for nonlinear dimensionality reduction, Isomap computes a globally optimal solution and has lower residual variance [1]. Despite of the superiority of global optimization, there also exist some problems about Isomap.

The first problem of Isomap is its sensitivity to outliers, which may collapse the total manifold structure [4]. A lot of algorithms have been developed to detect outliers [5-7]. However, they also caused mis-deletion of normal points when neighborhood size is large. In order to overcome this situation, BORDER algorithm [8] is adopted in this paper, which detects the outliers with their reverse k-nearest neighbor (RkNN) number.

The second problem is the high computational complexity of Isomap. Isomap with landmark points (L-Isomap) was developed to solve this bottleneck [9]. In L-Isomap algorithm, some points are selected to be landmark points and it only calculates the distances from each data point to each landmark point. But it is a hard work to pick out the appropriate landmark points. Fast Isomap based on minimum set coverage is a feasible road to solve the problem [10]. Unfortunately, it only applies to the situation of low curvature. The centers of the linear patches (LP-centers) are taken as the landmark points to solve this problem.

The third problem is the topological instability caused by short-circuit edges which may devastate global manifold

structure. An easy way to avoid this short-circuit edge is to decrease the neighborhood size. Several methods of selecting neighbourhood size were presented [11], [12], and an adaptive neighbourhood selection algorithm proposed [13] has achieved a better result. In addition, some algorithms select neighborhood according to the dissimilarity costs between each pair of points [14], [15]. Another feasible method is to find out all the short-circuit edges and remove them from the neighborhood graph. Obviously, the latter is better because of the free choice of the neighborhood size. To realize the latter idea, we compute the LP-centers, and find out the k-nearest center (k-NC) for each data point. An edge will be regarded as a short-circuit edge if there is no intersection between the k-NC of the two points it connects.

The rest of this paper is organized as follows. Next section introduces relevant definitions proposed in this paper. Section 3 develops our new algorithm, LPC-Isomap. Numerical experiments are shown in Section 4, in order to demonstrate the better performance of LPC-Isomap. Conclusions are summarized in Section 5.

II. RELEVANT CONCEPTS

Before introducing the LPC-Isomap algorithm, the relevant definitions and formulas must be put forward. The first is the definition of a linear patch:

Definition 1 (d-Dimensional Linear Patch): Given a D -dimensional X , and a user defined threshold ε_0 , X is a d -dimensional linear patch ($d \ll D$), denoted as $dLP(X)$, if we project X to d -dimensional linear space by linear analysis, e.g. principal component analysis (PCA), with a projection error smaller than ε_0 .

The projection error ε of data set X is measured by the form:

$$\varepsilon = \frac{\sum_{i=d+1}^D \lambda_i}{\sum_{i=1}^D \lambda_i} \quad (1)$$

where $X = [x_1, x_2, \dots, x_N] \in \mathcal{R}^{D \times N}$, λ_i is the eigenvalues of the covariance matrix of X , and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$.

A linear patch LP can be considered as a class, which employs Gaussian distribution as its discriminant function [16], denoted as $N(x)$. We can split a D -dimensional sample data S into several linear patches LP_1, LP_2, \dots , when take appropriate values for projection error threshold ε_0 and destination dimensionality d . A center of LP_i , denoted as μ_i , refers to the D -dimensional mean vector of LP_i . For $\forall x \in S$, the k-nearest center (k-NC) of x can be defined as follows:

Definition 2 (k-Nearest Center): Given a set of linear patches LP_i with centers μ_i and subordinate degree functions $N_i(\bullet)$, and a data point x , k nearest centers of x are a set of centers μ_j , where $N_j(x)$ is the first k largest values of $N_i(x)$ ($1 \leq j \leq k$).

III. MAIN RESULTS

In this section, LPC-Isomap will be proposed based on the concepts of section 2. BORDER algorithm is originally designed to find boundary points and outliers are extraordinary boundary points, being away from densely distributed data. The data points whose RkNN number equals a user defined threshold η or less are output as outliers. Generally, η is set to 1.

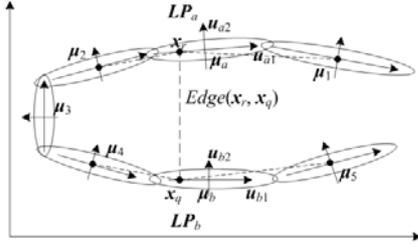


FIGURE 1. THE ELLIPSES SHOW THE CONTOURS OF LINEAR PATCHES, AND M_i ($i=A, B, 1, 2, \dots$) ARE LP-CENTERS. U_{a1} AND U_{b1} ($i=1, 2$) ARE ORTHONORMAL EIGENVECTORS OF COVARIANCE MATRICES OF LP_A AND LP_B , RESPECTIVELY. $EDGE(x_r, x_q)$ IS A SHORT-CIRCUIT EDGE CONNECTING x_r AND x_q . K-NC SIZE k_c IS 3. M_A, M_1 AND M_2 ARE K-NC OF x_r . M_B, M_4 AND M_5 ARE K-NC OF x_q

After BORDER, we will divide the sample data set S in to several linear patches LP_1, LP_2, \dots , with the user defined projection error threshold ε_0 and destination dimensionality d . Firstly, we calculate the projection error ε_i of each point $x_i \in S$ and its k -NN, and search the point x_m with minimum projection error x_m satisfying $\varepsilon_m < \varepsilon_0$. Then x_m and its k -NN are marked as a new linear patch LP_j , and points not in LP_j nearest to the LP-center μ_j are added to LP_j continuously until projection error ε^j of $LP_j \geq \varepsilon_0$. The above procedures are repeated until all points are added to the linear patches or there is no point x satisfying $\varepsilon < \varepsilon_0$.

Considering the flexibility of our proposed algorithm, we design an overlapping factor α ($0 \leq \alpha \leq 1$). According the properties of Gaussian distribution, the contour of a linear patch LP is a hyper-ellipsoid, with the center at μ and axes oriented along the eigenvectors of covariance matrix of LP , and with scaling factors in the directions of the axes given by corresponding eigenvalues. When we shrink the contour of LP_j , the point x_i belongs to overlapping region if it satisfies:

$$M_j(x_i) > (1 - \alpha) \max_{y \in LP_j} M_j(y) \quad (2)$$

where $M_j(x)$ is the Mahalanobis distance [16]. The points not covered by linear patches are considered as outliers.

For $\forall x_r \in LP_a, x_q \in LP_b$, if $Edge(x_r, x_q)$ is a short-circuit edge connecting x_r and x_q , $Edge(x_r, x_q)$ passes through a region containing few data points and the variances of LP_a and LP_b along the direction of $Edge(x_r, x_q)$ are very small. Therefore, both $N_a(x_q)$ and $N_b(x_r)$ are so small that the inequality

$N_j(x_r) > N_b(x_r)$ or $N_j(x_q) > N_a(x_q)$ holds for most of the linear patches LP_j in sample data. Therefore, there is no intersection between k -nearest center (k -NC) of x_r and x_q . An edge will be marked as a short-circuit edge if there is no intersection between k -NC of the two points it connects (see fig. 1).

The principle of selecting landmark points is that the information of intrinsic manifold structure can be recovered as much as possible according to the landmark points. In other words, more landmark points are required in the region with relatively large curvature. More linear patches will be generated in the place of large curvature because of containing fewer points. Therefore, the LP-centers can be exploited as landmark points. For convenience, we select a point which is nearest to the LP-center and in the same linear patch as a landmark point instead.

TABLE I. LPC-ISOMAP ALGORITHM

Input: D -dimensional data set S composed of data points x_i ; Projection error threshold ε_0 , Overlapping factor α and k -NC size k_c ; Neighborhood size k , Destination dimensionality d .
Output: A set of outliers OT and A set of landmarks LM ; A set of short-circuit edges SC and An embedding result.
1. $OT = \{\text{Outliers identified by BORDER}\}$, $S = S - OT$;
2. For $\forall x_i \in S$ do
2.1. Calculate U nearest neighbors $KNN_U(x_i)$ of x_i ($U = d + 1$);
2.2. Obtain the projection error ε_i of $\{x_i\} \cup KNN_U(x_i)$ when project it to \mathcal{R}^d ;
3. The set of overlapping region $OR = \emptyset$;
4. While $S \neq \emptyset$ do
4.1. Find the point x_m with minimum ε_m ;
4.2. If $\varepsilon_m > \varepsilon_0$ then break;
4.3. Mark the set $\{x_m\} \cup KNN_U(x_m)$ as a new linear block LP_j ;
4.4. While projection error ε_j^p of $LP_j < \varepsilon_0$
4.4.1. Find the point $x_n \in S \cup OR - LP_j$ nearest to center μ_j ;
4.4.2. Add x_n to LP_j and recalculate ε_j^p ;
4.5. Remove the last added point x_l from LP_j and $S = S - LP_j$;
4.6. $OR = OR \cup \{x \in LP_j M_j(x) > (1 - \alpha) \max_{y \in LP_j} M_j(y)\}$;
5. $OT = OT \cup S$, $LM = \cup_j \{\text{The point in } LP_j \text{ nearest to } \mu_j\}$;
6. For $\forall x_i \in \cup_j LP_j$ do
6.1. Calculate k_c nearest centers of x_i , denoted as $KNC_{k_c}(x_i)$;
7. $SC = \emptyset$;
8. For $\forall x_r, x_q \in LP_j$ do
8.1. If $KNC_{k_c}(x_r) \cap KNC_{k_c}(x_q) = \emptyset$;
8.1.1. $SC = SC \cup \{Edge(x_r, x_q)\}$;
9. Construct the neighborhood graph G with neighborhood size k ;
10. $G = G - SC$, then compute shortest path graph between LM and $\cup_j LP_j$;
11. Make a d -dimensional embedding by means of MDS.

Based on above analysis, we summarize the proposed LPC-Isomap algorithm in Table 1. Furthermore, an incremental PCA algorithm can be utilized to reduce the computational complexity of recalculating ε_j^p in step 4.4.2 of LPC-Isomap [17]. When the dimensionality D of data set is very high and the data set size N is insufficient, we only calculate the first p largest

eigenvalues and the corresponding orthonormal eigenvectors, where $d < p \ll \min(N, D)$. Eqn (1) is transformed into:

$$\varepsilon = \sum_{i=d+1}^p \lambda_i / \sum_{i=1}^p \lambda_i \quad (3)$$

The Mahalanobis distance $M(x)$ of LP is reformulated as:

$$M(x) = \sqrt{\|x - \mu\|^2 / \sigma + \sum_{i=1}^p (1/\lambda_i - 1/\sigma) (\mathbf{u}_i^T (x - \mu))^2} \quad (4)$$

where $\sigma \in \mathcal{R}^+$, and $\lambda_i (1 \leq i \leq p)$ is the positive eigenvalues of the covariance matrix satisfying $\lambda_1 \geq \dots \geq \lambda_p \geq \sigma > 0$, with corresponding orthonormal eigenvectors $\mathbf{u}_i (1 \leq i \leq p)$.

IV. NUMERICAL EXPERIMENTS

The first experiment is parameter setting of threshold ε_0 of projection error and overlapping factor α . ε_0 determines the approximate scope of center quantity where as α just trims the quantity. Table 2 shows the ability to identify outliers. Different values of k_c determine different thresholds of short-circuit edges. We can record the change process of residual variance of the embedding result by increasing k_c (neighborhood size is large enough), and a marked rise of residual variance indicates the appearance of short-circuit edges.

TABLE II. ABILITY OF LPC-ISOMAP TO IDENTIFY OUTLIERS IN NOISY SWISS-ROLL DATA (%)

Size of data set	BORDER		$\varepsilon_0 = 0.03$ $\alpha = 0.5$		$\varepsilon_0 = 0.003$ $\alpha = 0.5$		$\varepsilon_0 = 0.003$ $\alpha = 0.8$	
	FRR	FAR	FRR	FAR	FRR	FAR	FRR	FAR
1000 + 10 noise	50	0.35	93.8	0.54	32.7	3.1	18	0.92
3000 + 10 noise	28	0.24	100	0.20	46	0.80	47.7	0.15
3000 + 30 noise	32	0.22	98.7	0.22	31	0.93	42.6	0.16

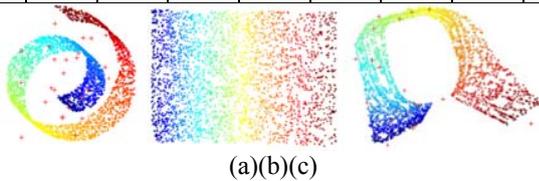


FIGURE II. (A) NOISY SWISS-ROLL DATA. (B) RESULT OF LPC-ISOMAP. (C) RESULT OF ISOMAP

Then we verify the behavior of LPC-Isomap based on noisy Swiss-Roll data. The size of the data set is 3030, including 30 noisy points (see fig.2a). The parameters are set to $\varepsilon_0 = 0.003$, $\alpha = 0.7$, and $k_c = 5$, the neighborhood size k of k-NN is set to 10. As shown in fig. 2b, most outliers are eliminated and the intrinsic manifold structure of sample data set is worked out by LPC-Isomap, whereas Isomap fails to learn the intrinsic structure in fig.2c. To verify effectiveness of deleting the short-circuit edges, we increase the neighborhood size k from 5 to 100 based on the Swiss-Roll data of 3000 points without noise (see fig.3). Table 3 shows the comparison on computational

complexity between LPC-Isomap and Isomap based on Swiss-Roll data.

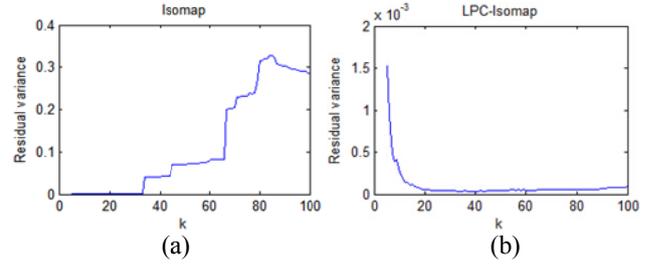


FIGURE III. THE CHANGE OF RESIDUAL VARIANCE WITH NEIGHBORHOOD SIZE k BEING INCREASED FROM 5 TO 100

TABLE III. RUNNING TIME OF ISOMAP AND LPC-ISOMAP IN SWISS-ROLL DATA (S)

Data set size	1000	1500	2000	2500	3000	3500
Isomap	2.611	5.525	9.285	13.887	20.330	30.131
LPC-Isomap	1.187	1.970	2.884	4.104	5.458	7.012

V. CONCLUSION

We have utilized BORDER algorithm and linear patches to eliminate the outliers from noisy data set and select landmark points. Moreover, most of the short-circuit edges have been removed from the neighborhood graph by analyzing intersection of k-nearest center (k-NC) between any two points. These three contributions led to the LPC-Isomap. Experiment results based on Swiss-Roll data and noisy sculpture face data illustrate that our new algorithm achieves better precision and speed than that of Isomap.

ACKNOWLEDGEMENTS

The work was supported by the Zhejiang Provincial Natural Science Foundation of China (No. LR12F03002) and the National Natural Science Foundation of China (No. 61473252).

REFERENCES

- [1] J. B. Tenenbaum, V. de Silva, and J. C. Langford, *A global geometric framework for nonlinear dimensionality reduction*, Science, vol. 290, pp. 2319–2323, 2000.
- [2] S. T. Roweis and L. K. Saul, *Nonlinear dimensionality reduction by locally linear embedding*, Science, vol. 290, pp. 2323–2326, 2000.
- [3] R. Wang, S. Shan, X. Chen, J. Chen, and W. Gao, *Maximal linear embedding for dimensionality reduction*, Pattern Analysis and Machine Intelligence, vol. 33, no. 9, pp. 1776–1792, Sept 2011.
- [4] M. Balasubramanian, E. L. Schwartz, J. B. Tenenbaum, V. de Silva, and J. C. Langford, *The Isomap algorithm and topological stability*, Science, vol. 295, no. 5500, p. 7a, January 2002.
- [5] Z. G. Tang, J. Yang, and B. R. Yang, *A new outlier detection algorithm based on manifold learning*, in Control and Decision Conference (CCDC), 2010 Chinese, May 2010, pp. 452–457.
- [6] P. Feng, L. Li, H. Zhang, and G. Qian, *Outliers detection in ICA*, in Communications, Circuits and Systems (ICCCAS), 2013 International Conference on, vol. 2, Nov 2013, pp. 328–330.
- [7] B. Liu, Y. Xiao, P. Yu, Z. Hao, and L. Cao, *An efficient approach for outlier detection with imperfect data labels*, Knowledge and Data Engineering, IEEE Transactions on, vol. 26, no. 7, pp. 1602–1616, July 2014.
- [8] C. Y. Xia, W. Hsu, M. L. Lee, and B. C. Oo, *BORDER: Efficient computation of boundary points*, IEEE Trans. Knowl. Data Eng., vol. 18,

- no. 3, pp. 289–303, 2006.
- [9] J. B. Tenenbaum and V. de Silva, *Global versus local methods in nonlinear dimensionality reduction*, Advances in Neural Information Processing Systems, vol. 15, pp. 721–728, 2003.
- [10] Y. K. Lei, Y. M. Xu, *Fast Isomap based on minimum set coverage*, Lecture Notes in Artificial Intelligence, vol. 6216, pp. 173–179, 2010.
- [11] L. Feng, C. K. Gao, T. Sun*, and H. Wu, *A neighborhood selection algorithm for manifold learning*, International Conference On Computer Design And Applications, vol. 2, pp. 339–342, 2010.
- [12] Q. S. Jiang, Y. P. Lu, and Z. K. Hong, *A neighborhood parameter optimization method of lle based on topology preservation*, in Electronic and Mechanical Engineering and Information Technology (EMEIT), 2011 International Conference on, vol. 8, August 2011, pp. 4231–4234.
- [13] Y. B. Zhan, J. P. Yin, and X. W. Liu, *Adaptive neighborhood selection based on local linearity for nonlinear dimensionality reduction*, Lecture Notes in Computer Science Volume, vol. 5821, pp. 337–348, 2009.
- [14] T. Wang, T. Xia, and X. M. Hu, *Minimax distance metric-based neighborhood selection algorithm for Isomap*, in Education Technology and Computer, 2010 2nd International Conference on, vol. 3, 2010, pp. 342–346.
- [15] V. Premachandran and R. Kakarala, *Consensus of k-nns for robust neighborhood selection on graph-based manifolds*, in Computer Vision and Pattern Recognition, 2013 IEEE Conference on, 2013, pp.1594–1601.
- [16] C. M. Bishop, Ed., *Pattern recognition and machine learning*. New York: Springer, 2006.
- [17] D. Huang, Z. Yi, and X. R. Pu, *A new incremental PCA algorithm with application to visual learning and recognition*, Neural Process Letters, vol. 30, pp. 171–185, 2009.