

基于图的半监督降维算法

杨格兰¹ 金辉霞² 孟令中³ 朱幸辉⁴

(湖南城市学院信息科学与工程学院 益阳 413000)¹ (湖南城市学院通信与电子工程学院 益阳 413000)²

(中国科学院软件研究所基础软件测评实验室 北京 100190)³

(湖南农业大学信息科学工程学院 长沙 410128)⁴

摘要 非线性降维和半监督学习都是近年来机器学习的热点。将半监督的方法运用到非线性降维中,提出了基于图的半监督降维的算法。该算法用等式融合的方法推出了标记传播算法的另一种表达形式,用标记传播的结果作为初始的数据映射,然后在图谱张成的线性空间中寻找最逼近初始映射的数据作为最后的半监督降维的结果。实验表明,所提算法可以获得平滑的数据映射,更接近于理想的降维效果。与标记传播算法、图谱逼近算法、无监督的降维算法的比较也体现出本算法的优越性。

关键词 半监督学习,流形学习,标记传播,图谱理论

中图分类号 TP301 **文献标识码** A

Graph-based Semi-supervised Dimensionality Reduction Algorithm

YANG Ge-lan¹ JIN Hui-xia² MENG Ling-zhong³ ZHU Xing-hui⁴

(School of Information Science and Engineering, Hunan City University, Yiyang 413000, China)¹

(School of Communication and Electronic Engineering, Hunan City University, Yiyang 413000, China)²

(Laboratory of Fundamental Software Testing and Evaluation, Institute of Software, Chinese Academy of Sciences, Beijing 100190, China)³

(Institute of Information Engineering, Hunan Agricultural University, Changsha 410128, China)⁴

Abstract Nonlinear dimensionality reduction and semi-supervised learning are both hot issues in machine learning area. Based on semi-supervised method, the article solved nonlinear dimensionality reduction problem to make up for the shortfall of ordinary methods. By using integration of equalities, a novel expression of label propagation algorithm was proposed. We used the label propagation result as the initial value mapping, and then found the best approximation to it in the graph spectral space. The experiment shows that our semi-supervised dimensionality reduction method can achieve smooth data mapping that is closer to the ideal effect.

Keywords Semi-supervised learning, Manifold learning, Label propagation, Spectral graph theory

1 引言

随着信息技术的飞速发展,信息数据具有一些新的特点,如维数高、属性强相关、样本数巨大但标签样本很少等,这使得传统的机器学习方法难以直接用来分析来源于真实世界的高维和非线性数据。故产生了很多解决高维数据问题的方法。降维是从高维观测空间通过线性或非线性映射投影到一个低维特征空间,发现隐藏在高维观测数据中有意义的低维结构,来研究数据特性。传统的线性降维算法有主分量分析(PCA)、独立分量分析(ICA)、线性判别分析(LDA)等,这些算法产生简单的变换函数,运算简便,尤其对线性结构降维效果好^[1]。但是现实中的高维数据大多是非线性的,用线性方法很难发现高维数据的内在性质,于是近年来涌现了许多基于流形的非线性降维方法,如等度规映射(ISOMAP)^[2]、局部线性嵌入(LLE)^[3]、拉普拉斯特征映射(Laplacian Eigen-

maps)^[4]、图册化流形(Charting a Manifold)^[5]、基于局部切空间(LTSA)^[6]等,这些算法能较好地解决数据处理中的非线性高维数据问题,保持原始数据的拓扑结构不变。

由于图和流形有很多相近的性质,最重要的一点就是都可以嵌入到 Euclid 空间,因此很多研究人员都使用图的方法逼近流形,并利用图的理论求解低维嵌入。流谱图理论在流形学习中已经得到了广泛的应用,很多流形降维的算法都是以谱理论为基础的。图谱分解方法更加促进了非线性降维的发展,文献[10]用对齐不同流形内在关系的方法,对两组流形数据之间的降维问题进行了探讨,阐明了怎样根据一些附加的关于这些数据集之间的映射信息来对齐它们的低维表达,很好地解决了两个流形降维后在同一个坐标系中的排列问题。但是在流形降维算法实现中,提出用标记传播(label propagation)做半监督降维的算法,存在的问题是定义的状态转移矩阵是能量消耗型,降维易导致标记点在外围,未标记点

到稿日期:2013-05-11 返修日期:2013-07-29 本文受国家科技支撑计划课题(2012BAD35B07),湖南省教育厅优秀青年项目(12BQ23)资助。
杨格兰(1975—),男,硕士,副教授,主要研究方向为机器学习、模式识别, E-mail: glyang@mail.ustc.edu.cn; 孟令中(1981—),男,博士,主要研究方向为软件可信性、软件测试。

和标记点的映射没有平滑地衔接起来,与实际降维效果有些差距^[11-14]。文献[8]提出了反映标签逼近和流形学习权重变化的正则化形式。由于标签逼近的二次型损失函数的约束性太强,该正则化形式直接的最优解并不能提供很好的学习性能,而是用图谱张成的特征空间去逼近标签信息,得到最优解。然而谱分解是个非常耗时的过程,特别是当样本的数目非常大的时候,大矩阵特征向量的迭代求解需要消耗大量时间和内存,这样效率不高,而且降维后效果也不够好。

针对以上问题,本文以标记传播算法和图谱理论为基础,以标记传播算法获取样本的初始映射,利用流形图谱与流形特征之间的关系,提出了一种基于图的半监督降维算法。算法以等式融合的方法推出了标记传播算法的另一种表达形式,用标记传播的结果作初始的数据映射,然后在图谱张成的线性空间中寻找最逼近初始映射的数据作为半监督降维的最后结果。

2 基于图的半监督降维算法

先给出流形降维半监督的描述,在非线形降维中,假定有一个 D 维的数据集 X ,而且这个数据集有内部维数 d (通常 $d \ll D$)。用数学语言描述就是, X 中的点位于或者靠近一个嵌套在 D 维空间的一个 d 维流形上。非线性降维就是要找出这个低维流形,以研究数据集主要的和本质的特征。但很多时候,已经知道数据集的一部分标记信息,希望利用上这些信息,使得降维朝着期望的方向进行,这就是半监督降维。半监督可以描述如下:已知高维样本数据 $\{x_i\}_{i=1}^n$,给定 l 个标记信息 $\{x_i, y_i\}_{i=1}^l$,需要找出样本的一个映射 $\{f_i\}_{i=1}^n$,使得映射值既逼近标记信息,又反映样本的流形结构信息。

本文的半监督降维算法分成如下两个步骤:1)用标记传播算法获取样本的初始映射;2)在图谱张成的线性空间中寻找最逼近初始映射的向量。这两个步骤都是必须的,如果跳过第一步直接做第二步会由于标记点的比例太小,而无法应用最小二乘法找出最好的线性表达系数。如果工作只做到第一步,则会由于标记传播算法使用能量衰减的图,导致标记点映射在非标记点的外围,使映射不平滑。所以这里采样两步相结合的算法,先用标记传播算法获得反映流形主要形状的映射,再用谱空间逼近获得最终的平滑映射。

2.1 获取初始映射值 f^0

初始映射值采用标记传播算法来获取。标记传播算法的计算公式可以通过迭代收敛获得,也可以在电路模型中利用欧姆定律获得,在这里通过构建既能反映流形特征又能反映标签逼近的两个等式约束,提出一种基于等式融合的推导方法,从而得到标记传播算法的计算公式。

1. 构建反映流形特征的等式约束

构造 k -邻域连接矩阵 W ,如果 j 在 i 的邻域内,则 $w_{ij} = \exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2})$,否则 $w_{ij} = 0$ 。其中 σ 是一个控制邻域影响的系数,取 $\sigma = \alpha \sqrt{D}$, D 是样本的维数, α 是与数据集有关的常数,一般默认为 1。计算 i 点的连接浓度信息 $d_i = \sum_k w_{ik}$,获取状态转移矩阵 P ,其中 $p_{ij} = \frac{w_{ij}}{d_i}$ 。要使得映射 f^0 符

合流形的特征,即要求

$$Pf^0 \approx f^0 \quad (1)$$

2. 构建反映标签逼近的等式约束

$$\text{记 } S = \begin{bmatrix} I_{l \times l} & 0 \\ 0 & 0_{n \times n} \end{bmatrix}, Y = \begin{bmatrix} y_{1 \times 1} \\ \vdots \\ 0_{n \times 1} \end{bmatrix}, \text{则约束可以写为}$$

$$Sf^0 = Y \quad (2)$$

其中, S 是一个选择矩阵,指明哪些点是标记点, Y 是对应的标记值。

3. 等式融合

式(2)中 S 是一个选择矩阵,选择标记点,同理 $(I-S)$ 也是一个选择矩阵,用来选择未标记点,在式(1)两边同时乘以 $(I-S)$,然后与式(2)相加得

$$(I-S)Pf^0 + Sf^0 = (I-S)f^0 + Y$$

整理得

$$f^0 = [I - (I-S)P]^{-1}Y \quad (3)$$

式(3)就是标记传播算法的另一种表现形式,利用式(3)获取初始映射值 f^0 。

2.2 用图谱空间逼近初始映射

先分析流形图谱与流形特征之间的关系。要使 f 能反映流形的主要特征,即要求 $Pf \rightarrow f$ 成立,也就是 $(I-P)f \rightarrow 0_{n \times 1}$,若 $M = I - P$,化为 $Mf \rightarrow 0_{n \times 1}$,而 Mf 趋近于 0 向量等效于要求 $\|Mf\|$ 最小,即要使 $(Mf)^T(Mf) \rightarrow 0$ 成立,那么展开得 $f^T M^T M f \rightarrow 0$,很明显 $\|f\|$ 越小越能满足式 $f^T M^T M f \rightarrow 0$ 的要求。为方便求解,加一个限制条件 $f^T f = 1$,那么问题变为有等式约束限制的优化问题。记 λ 为拉格朗日乘子系数,由拉格朗日乘子法得 $\frac{\partial (f^T M^T M f + \lambda (f^T f))}{\partial f} = 0$,求解后的 $M^T M f = -\lambda f$,即 f 是 $M^T M$ 对应的特征向量,所以 $M^T M$ 的特征向量组成的空间即图谱空间可以描述流形的主要特征。

对 $M^T M$ 做谱分解,取最小的 $m = \min(n * 0.2, l)$ 个特征值(去掉 0 特征值)对应的特征向量,记为 $E = [e_1 \ \dots \ e_m]$ 。

使用最小二乘法得到系数 $\alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix}$,使得 $E\alpha \rightarrow f^0$ 。

$$\alpha = (E^T E)^{-1} E^T f^0 \quad (4)$$

所求的映射是

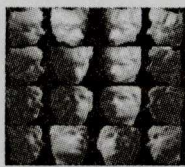
$$f = E\alpha \quad (5)$$

把式(3)和式(4)代入式(5)得到

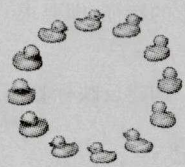
$$f = E(E^T E)^{-1} E^T [I - (I-S)P]^{-1} Y \quad (6)$$

3 实验和讨论

本文实验采用自然数据集:FACE 数据集和 COIL20 数据集,如图 1 所示(只显示了部分数据)。图 1(a)是 FACE 数据集,由 698 张 64×64 的头部方位图片构成;图 1(b)是 COIL20 数据集中的鸭子对象旋转 360° 得到的数据集,共包含 72 张 64×64 的图片数据。将图片转换成 4096 维的向量作为样本数据,对于 FACE 数据集抽取部分样本的方位作为标记信息,对于 COIL20 数据则用旋转的方位数据作为标记信息。



(a) FACE数据集



(b) COIL20数据集

图1 实验用的自然数据集

图2显示了FACE数据集的无监督降维效果,其中(a)是理想降维效果图,(b)~(f)为各种算法的无监督降维效果图。空心圈代表标记点,在无监督中没有利用标记点信息来指导降维,只用它们来观察降维的效果。由于没有利用标记信息,无监督降维后的效果与理想的效果差距很大,不但各种降维后效果图都不够规则,而且标记点信息不如理想效果图那样能够区分不同类别。

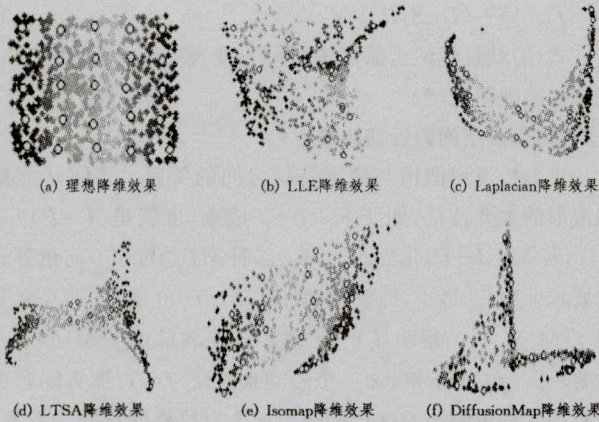


图2 FACE数据集的无监督降维

图3是FACE数据集上25个固定标记点的半监督降维效果图,其中(a)是理想降维效果,(b)是标记传播的方法的降维效果,(c)是谱空间逼近效果,(d)是本文算法的降维效果。(b)~(d)由于利用了标记数据信息指导数据降维,因此降维后效果显现的形状更接近理想效果图,都比无监督要好。图4为特定标记点的半监督降维图,从理想降维效果图中可以知道,降维后不但图形很规则,而且标记点能够很好地区分不同类别信息,其它3种效果图(b)~(d)也是基本上区别了类别信息,标记传播的方法和谱空间逼近的方法未标记点和标记点的映射没有平滑地衔接起来,尤其标记传播的方法导致很多标记点在外围,而本文算法中的标记点和未标记点的映射都呈比较均匀的分布。从降维后的图形和标记点的分类情况,可以看出使用本文的方法得到的结果更接近于(a)中的理想效果。

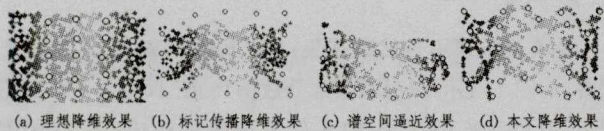


图3 FACE数据集特定标记点的半监督降维

下面用COIL20数据集做实验,图4(a)表明理想降维效果,图4(b)~(f)为各种无监督算法的降维效果。由于没有监督数据的引入,按照不同的无监督方法数据降维后与图4(a)理想降维效果相比,差距较大,尤其LLE、Laplacian、Diffu-

sionMap降维效果图均不能呈现圆周形状。

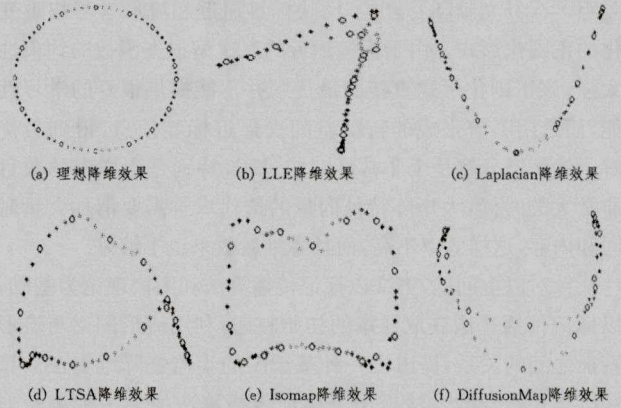


图4 COIL20数据集的无监督降维

图5是COIL20数据集特定标记点降维效果比较图,理想的降维效果是个椭圆形,标记点很好地融合未标记点。比较图(b)~(d),只有采用本文算法的图(d)形状更加接近理想效果图,而且标记点能平滑地融入流形中,图(c)虽然标记点比图(b)更好地融入流形,但是降维后图形和原始图差距更大,尤其是随机标记点的效果(见图8(c))。无论是特定标记点还是随机标记点,从标记点的位置和图形来分析,很明显本文算法映射的结果更平滑,更接近于理想效果。

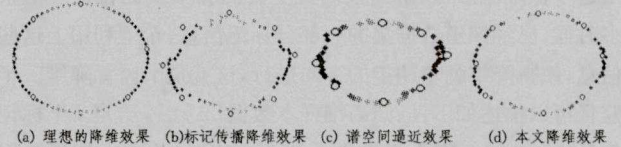


图5 COIL20数据集特定标记点的半监督降维

结束语 本文用等式融合的思路推出了标记传播算法的一种表达形式,并把标记传播算法和谐空间逼近的方法结合起来,提出了基于图的半监督降维算法。理论分析与实验结果都表明,此算法比无监督的降维算法、基于标记传播的降维算法、谱空间逼近算法更加容易取得好的降维效果,并且标记点信息能够很好地融合到流形中,不但使降维流形更加规则,而且能有效地区分不同类别信息。对高维复杂的流形数据做有监督的降维,提取有意义的特征,再对数据做分类、识别等是十分有意义的。

参考文献

- [1] van der M L J P, Postma E O, van den Herik H J. Dimensionality reduction; A comparative review [J]. Journal of Machine Learning Research, 2007(1)
- [2] Tenenbaum J B, de S V, Langford J C. A global geometric framework for nonlinear dimensionality reduction [J]. Science, 290, 2000; 2319-2323
- [3] Roweis S T, Saul L K. Nonlinear dimensionality reduction by locally linear embedding [J]. Science, 2000, 290; 2323-2326
- [4] Belkin M, Niyogi P. Laplacian eigenmaps for dimensionality reduction and data representation [J]. Neural Computation, 2003, 15(1); 1373-1396

(下转第296页)

了基于语义标签约束的特征点层次提取方法。不同于已有方法只是考虑极值理论进行特征点定义,算法对提取出来的外部特征点进行语义识别,从而能够通过语义标签来控制特征点的数量。此外,算法通过图割理论进行特征点优化,有效地提高了特征点提取的准确性;且通过实验分析表明,对于不同姿态的可变形物体,算法得到的特征点能够和物体的关节点相一致,符合人类的视觉感知结果。

但是,算法还存在需要改进的地方;本文只是采用测地路径和分割边界线的交点作为关节点位置,而没有考虑不同模型在关节点位置上的局部相似性。在后续工作中,可以考虑对不同模型分割边界线上的顶点集合进行聚类分析,提取出更为一致的关节点位置。另外一方面,目前算法只是考虑一些三维造型数据,而没有考虑到扫描得到的三维数据,由于这些数据存在局部遮挡等问题,使得特征点具有不完全性。因此,需要设计更为准确的语义标签识别算法,这将是后续研究的另一问题。

参 考 文 献

- [1] Chen D Y, Pei X T. On Visual Similarity Based 3D Model Retrieval[J]. Computer Graphics Forum, 2003, 22(3): 223-232
- [2] Karni Z, Gotsman C. Spectral compress of mesh geometry[C]// Proceedings of SIGGRAPH' 00. USA: ACM press, 2000, 279-286
- [3] Sagi K, George L, Ayellet T. Mesh Segmentation using feature point and core extraction[J]. The Visual Computer, 2005, 9: 649-659
- [4] Bronsteiny A, Bronsteiny M, Bustos B. SHREC 2010: robust feature detection and description benchmark[C]// Proceedings of the IEEE Conference on CVPR. New York: IEEE computer society press, 2010
- [5] Mikolajczyk K, Schmid C. Scale and Affine invariant interest point detectors[J]. International Journal of Computer Vision, 2004, 60(1): 63-86
- [6] Schmid C, Mohr R, Bauckhage C. Evaluation of interest point detectors[J]. International Journal of Computer Vision, 2000, 37(2): 151-172
- [7] Sipiran I, Bustos B. Harris 3D: A robust extension of the Harris operator for interest point detection on 3D meshes[J]. The Visual Computer, 2011, 27(11): 963-976
- [8] Castellani U, Cristani M, Fantoni S. Sparse points matching by combining 3D mesh saliency with statistical descriptors[J]. Computer Graphics Forum, 2008, 27(2): 643-652
- [9] Zaharescu A, Boyer E, Varanasi K. Surface feature detection and description with applications to mesh matching[C]// Proceedings of IEEE Conference on CVPR. New York: IEEE computer society press, 2009: 373-380
- [10] Zou G, Hua J, Dong M. Surface matching with salient keypoints in geodesic scale space[J]. Computer Animation and Virtual Worlds, 2008, 19(3): 399-410
- [11] Sun J, Ovsjanikov M, Guibas L. A concise and provably informative multi-scale signature based on heat diffusion[J]. Computer Graphics Forum, 2009, 28: 1383-1392
- [12] Mian A, Bennamoun M, Owens R. On the Repeatability and Quality of Keypoints for Local Feature-based 3D Object Retrieval from Cluttered Scenes[J]. International Journal of Computer Vision, 2009, 89(2): 348-361
- [13] 陈启华, 潘翔, 张三元. 语义知识约束的三维人体特征点检测和分割[J]. 计算机辅助设计与图形学报, 2011, 23(6): 1061-1068
- [14] Kristin P, Campbell C. Support Vector machines; hype or hallucination[C]// Proceedings of ACM SIGKDD Explorations Newsletter. New York: ACM Press, 2000, 2(2): 1-13
- [15] Kim V G, Yaron L, Thomas F. Blended Intrinsic Maps [J]. ACM Transaction on Graphics, 2011, 30(4)
- [16] Robert W, Sumner, Jovan P. Deformation Transfer for Triangle Meshes[J]. ACM Transactions on Graphics, 2004, 23(3): 399-405
- (上接第 282 页)
- [5] Brans M M. Charting a manifold[C]// Neural Information Processing Systems: Natural and Synthetic. Vancouver, Canada, 2000: 232-245
- [6] Zhang Zhen-yue, Zha Hong-yuan. Linear low-rank approximations and nonlinear dimensionality reduction [J]. Science in China Series A-Mathematics, 2005, 35(3): 273-285
- [7] Yan S C, Xu D, Zhang B, et al. Graph embedding and extensions: A general framework for dimensionality reduction [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2007, 29(1): 40-51
- [8] Zhu Xiao-jin, Ghahramani Z. Learning from labeled and unlabeled data with label propagation[R]. Technical Report 02-107, CMU-CALD. USA: Carnegie Mellon University, 2002
- [9] Zhu Xiao-jin, Lafferty J, Ghahramani Z. Semi-Supervised Learning: From Gaussian Fields to Gaussian Processes[R]. CMU. Technical Report, CMU-CS-03-175. USA: Carnegie Mellon University, 2003
- [10] Pothan, Alex, Fan C-J. Computing the Block Triangular Form of a Sparse Matrix [J]. ACM Transactions on Mathematical Software, 1990, 16(4): 303-324
- [11] 李岩波, 宋琼, 郭新辰. 基于流形距离的人工免疫半监督聚类算法[J]. 计算机科学, 2012, 39(11): 204-207
- [12] 刘志勇, 袁媛. 基于测地距离的半监督增强[J]. 计算机工程与应用, 2011, 47(21): 202-204
- [13] 任剑锋, 梁雪, 李淑红. 基于非线性流形学习和支持向量机的文本分类算法[J]. 计算机科学, 2012, 39(1): 261-263
- [14] 罗磊, 李跃华. 基于 LLE 的分类算法及其在被动毫米波目标识别中的应用[J]. 电子与信息学报, 2010, 32(6): 1306-1310
- [15] 王越, 王泉, 吕奇峰, 等. 基于初始聚类中心优化和维间加权的改进 k-means 算法[J]. 重庆理工大学学报: 自然科学版, 2013, 27(4): 77-80