

## 基于加权的不完备非负矩阵分解算法

杨志君, 叶东毅

(福州大学 数学与计算机科学学院, 福州 350108)

(flistorm@gmail.com)

**摘要:**非负矩阵分解(NMF)作为一种特征提取与数据降维的新方法,相较于一些传统算法,具有实现上的简便性,分解形式和分解结果上的可解释性等优点。但当样本矩阵不完备时,NMF无法对其进行直接分解。提出一种基于加权的不完备非负矩阵分解(NMFI)算法,该算法在处理不完备样本矩阵时,先采用随机修复的方法降低误差,再利用加权来控制各样本的权重,尽量削弱缺损数据对分解结果产生的干扰。此外,NMFI算法使用区域权重来进一步减少关键区域数据缺损对分解产生的影响。实验结果表明,NMFI算法能有效提取样本中残余数据的信息,减少缺损数据对分解结果的影响。

**关键词:**非负矩阵分解;不完备数据集;随机修复;加权;区域权重

**中图分类号:** TP391 **文献标志码:** A

## Weighted non-negative matrix factorization for incomplete dataset

YANG Zhi-jun, YE Dong-yi

(College of Mathematics and Computer Science, Fuzhou University, Fuzhou Fujian 350108, China)

**Abstract:** Nonnegative Matrix Factorization (NMF) is a new method for feature extraction and data dimension reduction. It has an advantage over traditional algorithms in the simple implementation and the interpretability of factorization form and factorization result. But NMF could not decompose the samples matrix when it is incomplete. However, when dealing with incomplete dataset, NMFI (Weighted Non-negative Matrix Factorization for Incomplete Dataset) made use of random repair to decrease the error and weighted method to control weights of the samples, which could weaken the disturbance of missing data as much as possible. In addition, NMFI used regional weight for further reducing the impact of missing data in critical region. The experimental results demonstrate that NMFI can effectively extract information from retained data and reduce the influence of missing data.

**Key words:** Nonnegative Matrix Factorization (NMF); incomplete dataset; random repair; weighting; regional weight

### 0 引言

非负矩阵分解(Nonnegative Matrix Factorization, NMF)<sup>[1]</sup>是机器学习领域中一种高效的数据降维方法,也是一种特征提取的新方法,与传统的矩阵分解方法不同,它的分解过程和结果都是非负的,这在心理学上符合人对整体的感知是由对组成整体的部分的感知构成的(纯加性的)。近年来非负矩阵分解方法在聚类分析、人脸识别、语音识别、语音分离<sup>[2-4]</sup>等领域取得了很好的应用。但NMF也存在一些缺陷,当样本矩阵由于各种原因而缺损时,NMF无法对其进行直接处理。另外,若其部分数据由于某种原因而产生较大偏差,则对分解结果的准确性会有一定程度的干扰。

针对以上问题,本文提出了一种基于加权的不完备非负矩阵分解算法NMFI。算法首先对不完备数据进行随机修复,使用其他样本中对应位置的未缺失数据来修复缺失数据,使修复之后的样本与原始样本误差尽量小。此外,考虑到样本的重要性各不相同,即无损或缺损程度较低的样本重要性应当比缺损程度高的样本高,算法采用加权的方法来控制各样本的权重,使缺失越大的样本权重越低,以避免缺损数据对分解结果产生较大干扰。最后,考虑了区域权重。所谓区域权重即在单个样本中某些数据的重要性明显低于其他数据,这

些数据缺失时对分解结果的影响比其他数据缺失时的影响要小,即对于该样本的权重而言,这些数据的贡献较小。

### 1 非负矩阵分解

给定一个 $n$ 行 $m$ 列的非负矩阵 $V$ ,将其近似地分解成两个非负矩阵之积,分别记为 $W$ 与 $H$ ,使得 $V \approx WH$ ,即:

$$V_{ij} \approx (WH)_{ij} = \sum_{a=1}^r W_{ia} H_{aj} \quad (1)$$

分解之后的矩阵 $W$ 与 $H$ 的维数分别为 $n \times r$ 和 $r \times m$ , $r$ 的选取与 $V$ 的秩有关, $r \ll \text{rank}(V)$ 。

如果用 $v$ 和 $h$ 分别表示矩阵 $V$ 与 $H$ 中对应的列向量,由式(1)可得 $V \approx Wh$ 。由此可见,原始矩阵 $V$ 中的列向量 $v$ 可以通过矩阵 $W$ 以一定的权重 $h$ 表示出来。因此,矩阵 $W$ 可以看作是基矩阵,其中的每一列代表基向量,而矩阵 $H$ 中的每一列代表与矩阵 $W$ 中基向量相对应的编码系数。为了找到一个矩阵分解 $W$ 与 $H$ ,使得 $\|V - WH\|^2$ 最小,且满足 $W \geq 0, H \geq 0$ ,可以使用如下的规则进行迭代计算<sup>[1]</sup>:

$$H_{aj} \leftarrow H_{aj} \frac{(W^T V)_{aj}}{(W^T W H)_{aj}} \quad (2)$$

$$W_{ia} \leftarrow W_{ia} \frac{(V H^T)_{ia}}{(W H H^T)_{ia}} \quad (3)$$

收稿日期:2009-10-09;修回日期:2010-01-19。 基金项目:国家自然科学基金资助项目(60805042)。

作者简介:杨志君(1985-),男,福建莆田人,硕士研究生,主要研究方向:非负矩阵分解; 叶东毅(1964-),男,福建南安人,教授,博士生导师,主要研究方向:计算智能、数据挖掘、神经网络。

在计算过程中,随机选取非负的初始矩阵  $W$  与  $H$ , 定义适当的迭代中止条件之后,按照式(2)和(3)进行迭代计算即可得到  $W$  和  $H$ 。

## 2 本文算法的基本思想

在实际应用中,经常由于某些主观或客观因素导致数据的缺失,如采集到的图像由于镜头存在污点而产生的缺损,道路监控中由于物体遮挡而产生的数据缺失等,对于这些数据 NMF 算法无法对其进行直接处理。鉴于此,本文提出一个基于加权的不完备非负矩阵分解算法 NMFI。NMFI 算法的目标是当样本不完备时,尽量提取样本中残余数据的信息,减少缺损数据对分解结果的影响。

### 2.1 算法的基本思想

在样本数据的采集和存储过程中,时常会出现数据缺失的现象,而 NMF 算法无法对这些不完备数据进行直接分解并分析,它要求数据矩阵必须是完备的。因此,首先要对不完备样本进行修复。一个简单的方法是采用单个数据值来修复缺失数据,再进行 NMF 分解。但由于样本的取值是多样性的,使用单一值来代替缺失数据会使修复后的样本与原始样本误差较大,进而影响分解结果。鉴于此,本文采用随机修复的方法对缺失数据进行修复,即使用其他样本中对应位置的未缺失数据来修复缺失数据,从而减少修复后的样本和原始样本之间的误差,避免缺损数据对分解结果产生较大干扰。

将不完备样本修复为完备样本后,若直接对其进行 NMF 分解,效果显然是不理想的。由于样本的缺失程度不尽相同,缺损较大的样本重要性应比无损或缺损较小的样本低,因此才能尽量减少缺损数据对分解结果的影响;反之,即加强较完备样本的影响。本文采用对各样本加权的方式,即控制各样本的权重来控制样本对分解结果的影响。

### 2.2 单一值修复与随机修复

为比较单一值修复和随机修复之间的差异,首先对使用单一值修复后的不完备样本矩阵进行 NMF 分解。文中所有实验均在 Windows Vista 环境下进行,CPU 为酷睿双核 2.20 GHz,内存 2 GB,Matlab 为 7.0 版。且 NMF 算法的迭代终止条件皆按如下定义。

设样本数为  $m$  时,算法第  $i$  次迭代各样本所产生的平均误差为:

$$T_i = \|V - W_i H_i\|^2 / m$$

其中  $W_i, H_i$  分别为第  $i$  次迭代分解产生的基矩阵和编码矩阵,则迭代中止条件为  $1 - (T_{i+1}/T_i) \leq 0.01\%$ 。

实验采用 Swimmer 数据集<sup>[5]</sup>,该数据集包含一系列黑白图像,每个图像有 4 个可移动的部分(四肢),每个部分可以旋转到 4 个不同的位置(关节)。且每个黑白图像的中心有一个包含 12 个像素的躯干和四个包含 6 个像素的可移动部分,每个可移动部分可位于 4 个不同位置,因此总共有 256 张图像包括所有的肢体位置组合。每张图像包含  $32 \times 32$  个像素。图 1 给出了该数据集的一个子集。

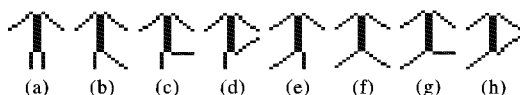


图 1 Swimmer 数据集

该数据集在文献[5]中是用来展示 NMF 算法发现局部特征(四肢)的能力。它遵守一些预定义的规则,如可分性和完全因子取样,即它是通过一种具有 NMF 样式的模型产生的。

实验时,取数据集中的前 20 张图像,保留前 2 张图像完整,对后 18 张图像采用随机去除元素的方法产生缺损图像。则使用单一值“1”修复缺损图像,并对其进行 NMF 分解后,分解产生的基图像如图 2。实验中基图像个数取为 12。

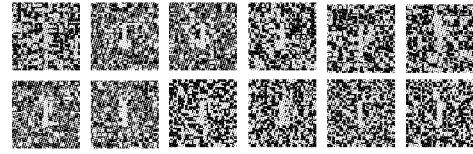


图 2 单一值修复时 NMF 分解的基图像

从图 2 可以看出,分解产生的基图像中较难分辨出四肢和躯干。这是由于使用单一值修复时,修复后的图像和原始图像之间存在较大误差,缺损数据会对分解结果产生较大影响。

NMFI 算法采用随机修复的方法对缺损数据进行修复,即使用其他样本中无缺损的数据来修复样本中同一位置的缺损数据。在实际应用中,进行 NMF 分解的各样本间通常具有一定程度的相似性,因此该方法可以有效减少修复后的图像和原始图像之间的误差,减少缺损数据对分解结果的影响。图 3 是对缺损图像进行随机修复后,使用 NMF 分解产生的基图像。

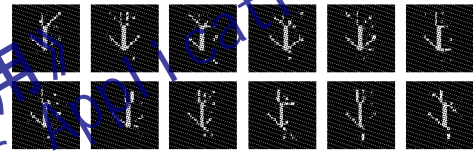


图 3 随机修复时 NMF 分解的基图像

从图 3 中可以看出,随机修复后分解产生的基图像四肢和躯干都比较明显。由于修复后的图像和原始图像之间误差较小,因此受到缺损数据的干扰较小。

### 2.3 样本加权

在不完备样本矩阵中,各样本的缺损程度不尽相同,缺损越大的样本其误差也越大,对分解结果的影响则应更小,即重要性更低。本文采用加权的方式来区分不同样本的重要性,重要的样本权重较高,反之则较低。

#### 2.3.1 相关公式及定义

假设  $V_n$  为需分解的样本矩阵, $n$  为样本数, $k \leq n$  为完备样本数,完备样本构成的样本矩阵为  $V_k$ ,  $(n-k)$  个不完备样本的权重分别为  $p_{k+1}, \dots, p_n$ , 且  $0 \leq p_i \leq 1, i = k+1, k+2, \dots, n$ 。  $W$  为分解产生的基矩阵,  $H_k$  为  $k$  个完备样本的编码矩阵,  $h_{k+1}, h_{k+2}, \dots, h_n$  为不完备样本基于基矩阵的编码向量,则算法的目标函数为:

$$F = \|V_k - WH_k\|^2 + \sum_{i=k+1}^n p_i \|v_i - Wh_i\|^2 \quad (4)$$

使用梯度下降法来取得  $F$  的局部最小值。

$$W_{ia} \leftarrow W_{ia} - \mu_w \frac{\partial F}{\partial W_{ia}} \quad (5)$$

$$H_{ia} \leftarrow H_{ia} - \mu_h \frac{\partial F}{\partial H_{ia}} \quad (6)$$

当步长  $\mu_h$  取式(7)的值时,式(5)转换为式(8)所示的乘法更新规则。其中  $w_i$  为  $W$  的第  $i$  列,  $h_t$  为  $H$  的第  $t$  列。

$$\mu_h = \frac{W_{ia}}{\sum_{i=1}^k w_i^T h_i H_{ia} + \sum_{i=k+1}^n p_i w_i^T h_i H_{ia}} \quad (7)$$

$$W_{ia} \leftarrow W_{ia} \frac{\sum_{t=1}^k V_{it} H_{at} + \sum_{t=k+1}^n p_t V_{it} H_{at}}{\sum_{t=1}^k w_t^T h_t H_{at} + \sum_{t=k+1}^n p_t w_t^T h_t H_{at}} \quad (8)$$

即:

$$W_{ia} \leftarrow W_{ia} \frac{(V H^T)_{ia}}{(W H H^T)_{ia}} \quad (9)$$

其中  $H' = [H_k \quad p_{k+1} h_{k+1} \quad p_{k+2} h_{k+2} \quad \cdots \quad p_n h_n]$ 。

当步长  $\mu_H$  取式(10)的值时,式(6)转换为式(11)所示的乘法更新规则,其中  $m$  为样本的维数。

$$\mu_H = \frac{H_{ia}}{\sum_{t=1}^m p_a w_t^T h_a W_{ia}} \quad (10)$$

$$H_{ia} \leftarrow H_{ia} \frac{(W^T v_a)_{ia}}{(W^T W h_a)_{ia}} \quad (11)$$

即:

$$H_{ia} \leftarrow H_{ia} \frac{(W^T V)_{ia}}{(W^T W H)_{ia}} \quad (12)$$

其中  $v_a$  为  $V$  的第  $a$  个列向量。

以上更新规则的收敛性证明请参阅文献[6]中所采用的设计辅助函数的方法,在此不作赘述。

为计算各样本的权重,作如下定义。

**定义1** 缺损度。样本中缺损的数据个数和样本总数数据个数的比值称为样本的缺损度。

**定义2** 高完整数。若样本  $x$  的缺损度为  $\theta$ , 样本矩阵中缺损度小于  $\theta$  的样本个数为  $\beta$ , 则称  $\beta$  为样本  $x$  的高完整数。

### 2.3.2 加权对分解结果的影响

从图3所示的基图像中可以看出,除四肢和躯干外,仍有由缺损数据引起的分解误差。这是由于在这个实验中不完备样本和完备样本的权重取值相同,即各样本的重要性没有差异,导致分解误差仍然较大。对样本的重要性产生影响的因素除了样本的缺损程度外,还与高完整数有关。若在一个样本矩阵中,完备样本的个数远大于不完备样本,显然,直接根据完备样本分解得到基图像和整个样本矩阵分解得到的基图像几乎是相同的。此时,不完备样本的权重可以为0。若绝大部分的样本缺损度均小于某个值  $\theta$ , 仅有少数样本高于  $\theta$ , 则高于  $\theta$  的少数样本权重可以为0。因此样本的权重和高完整数是密切相关的。

下面采用式(9)和(12)的乘法更新规则对图3中已修复的样本矩阵进行分解,各样本的权重取值为:

$$p = ((1 - \theta)/\beta)^2 \quad (13)$$

其中:  $\theta$  为样本缺损度;  $\beta$  为样本的高完整数。加权后的迭代中止条件均按如下定义。

设样本数为  $m$  时,算法第  $i$  次迭代各样本所产生的平均误差为:

$$T_i = (\|V_k - W H_k\|^2 + \sum_{i=k+1}^n p_i \|v_i - W h_i\|^2) / m$$

则迭代中止条件为  $1 - (T_{i+1}/T_i) \leq 0.01\%$ 。

分解产生的基图像如图4。从图4中可以看出因数据缺损而产生的分解误差明显减少,这是由于降低了缺损度较高的样本的权重,使得缺损数据对分解结果的影响进一步降低。

在某些情况下使用随机修复和使用单一值修复产生的样本矩阵,它们与原始矩阵的误差相差不甚明显。这时权重的影响尤为明显,合适的权重在降低缺损数据影响的同时,可以

尽量提取其余数据中的有用信息。利用下面的实验来解释这一点。

在文献[7]中, temporal image 数据集被用来验证 ONMF 算法在满秩分解条件下处理动态数据集的有效性。图5是该数据集的部分生成因子,每个生成因子对应于一张  $10 \times 10$  像素的图片,图片中包含一条横线或竖线。该数据集中的每张图片均是生成因子的某种线性组合,且产生的样本间没有明显的相似性,这也是实验采用该数据集的原因。

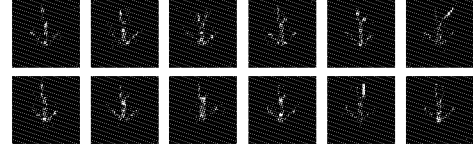


图4 加权分解的基图像

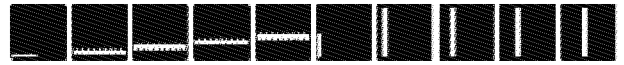


图5 temporal image 数据集的生成因子

从该数据集中取出 200 张图片,保留前 100 张图片无缺损,剩余的 100 张图片随机去除元素生成不完备样本。基图像个数取为 20,随机修复后对其进行 NMF 分解得到的基图像如图6。

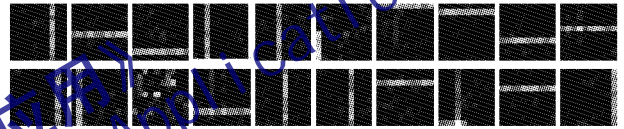


图6 NMF 分解产生的基图像

从图6中可以看出,缺损数据对分解产生的基图像有较大的影响,所有基图像均有除生成因子之外的分解误差。

若采用式(9)和(12)的乘法更新规则并对其采用加权的方式进行分解,各样本的权重取值为:

$$p = ((1 - \theta)/\beta)^2$$

则产生的基图像如图7。

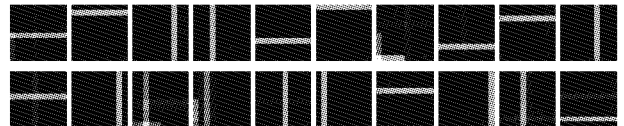


图7 采用加权方式后分解产生的基图像

从图7中可以看出,由于降低了不完备样本的权重,基图像中受缺损数据干扰而产生的分解误差已经明显降低。

### 2.4 区域性加权对分解结果的影响

所谓区域权重是指在一个样本中,某些区域对于分解产生的基图像具有较大的影响,而其余区域如背景则影响较小。如 Swimmer 数据集包含的样本中,除躯干和四肢围绕关节活动的区域外,其余均可视为背景。若有两个缺损度相同的样本,一个主要是关键区域数据缺损,而另一个则主要是非关键区域如背景数据缺损,则显然第二个样本的重要性要高于第一个样本。在类似的应用中,需要考虑区域性权重。

实验时,从 Swimmer 数据集中取 20 个样本,通过随机去除元素使所有样本的缺损度在 0.4 左右,其中前 10 个样本主要是背景区域缺损,后 10 个样本主要是关键区域缺损。另外,在单个样本中,背景区域元素个数为 688,关键区域元素个数为 366。背景区域与关键区域的重要程度之比为 1:9。

此时,若直接通过缺损元素个数来计算样本的缺损度,并按照式(13)得出各样本的权重,则随机修复后样本矩阵分解得到的基图像如图8。



从图8中可以看出,基图像中关键区域受到缺损数据的干扰较大,四肢也出现了不同程度的误差。这是由于实验时背景区域缺损的数据和关键区域缺损的数据对样本权重的影响已被默认为一致,导致关键区域缺损较大的样本对分解后的基图像产生较大的干扰。显然,应当相对降低关键区域缺损较大的样本的权重。

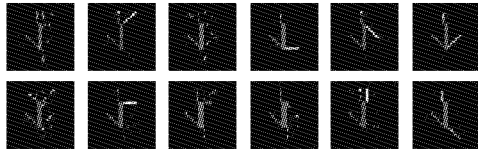


图8 直接计算缺损度后分解产生的基图像

接下来按区域权重来计算样本的缺损度,计算公式如下:

$$\theta = \frac{0.1}{688} \cdot x + \frac{0.9}{366} \cdot y \quad (14)$$

其中:  $x$  为背景区域缺损的元素个数;  $y$  为关键区域缺损的元素个数。

再根据式(13)的权重计算公式得出各样本权重,则样本矩阵随机修复后分解得到的基图像如图9。

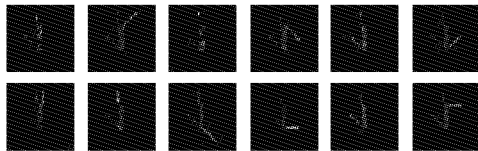


图9 根据区域权重分解产生的基图像

从图9中可以看出提取的信息有所减少,即四肢和躯干相对模糊,但基图像的关键区域受到缺损数据的干扰变小。这是由于关键区域重要性较背景区域高,后10个图像的权重较前10张图像的权重低,因此分解产生的基图像受关键区域数据缺损的干扰降低。

## 2.5 算法设计

根据前面的分析,本文提出的 NMFI 算法的具体步骤如下。

- 1) 扫描样本矩阵,根据区域重要性计算各样本的缺损度  $\theta$  及高完整数  $\beta$ 。
- 2) 计算各样本的权重  $p$ 。
- 3) 采用随机修复的方法对不完备样本进行修复。
- 4) 对基矩阵  $W$  和编码矩阵  $H$  随机初始化,并将  $H$  赋予  $H_{temp}$ 。
- 5) 将  $H_{temp}$  赋予  $H$ ,使用式(12)乘法更新规则更新  $H$ ,并将新的编码矩阵赋予  $H_{temp}$ 。
- 6) 使用式(9)乘法更新规则更新  $W$ ,其中  $H' = [H_k \ p_{k+1}h_{k+1} \ p_{k+2}h_{k+2} \ \cdots \ p_n h_n]$ ,即  $H'$  通过编码向量加权得到。
- 7) 符合迭代中止条件,算法结束,否则转5)。

## 3 Triangle 数据集的实验结果

在文献[11]中 Triangle 数据集被用于评价算法识别数据生成因子的能力,且各张图像均有关键区域和背景区域。该数据集的生成因子是三张  $32 \times 32$  的灰度图,每张图像包含一个大小不同的三角形,如图10。

图11是根据生成因子产生的该数据集的一部分图像。

实验时,从 Triangle 数据集中取出75张图像,每张图像为一个样本,即样本矩阵的一列。保留前25张图像的完整性,其余50张图像随机去除元素产生不完备样本,基图像个

数取为3。则采用单一值0对不完备样本矩阵进行修复时,NMF分解产生的基图像如图12。

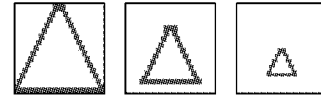


图10 Triangle 数据集生成因子

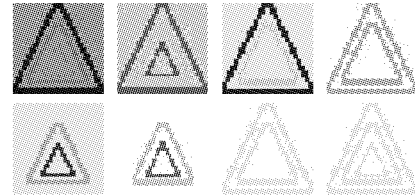
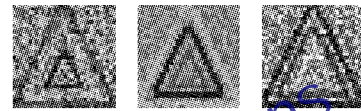


图11 Triangle 数据集部分图像



(a) 单一值修复时NMF分解产生的基图像



(b) 随机修复时NMF分解产生的基图像



(c) NMFI分解产生的基图像

图12 单一值修复时 NMF 分解产生的基图像

在 Triangle 数据集中,若采用单值修复,会使修复后的样本矩阵与原始图像相比有较大误差,缺损数据对分解结果干扰会比较明显,因此从图12(a)的分解结果中很难分辨出该数据集的生成因子。采用随机修复的方法对不完备样本进行修复后,NMF分解产生的基图像如图12(b)。从图12(b)中可以看出,随机修复使分解误差降低,但由于各样本的权重已默认为相同,缺损数据对分解结果的干扰仍然较大。

若设单一样本中三角形所在区域与背景区域的重要程度之比为9:1,则采用区域权重的方法按2.5节中的算法步骤,对不完备数据矩阵进行分解后产生的基图像如图12(c)。缺损度的计算公式如下:

$$\theta = \frac{0.1}{528} \cdot x + \frac{0.9}{496} \cdot y$$

其中:  $x$  为背景区域缺损的元素个数;  $y$  为关键区域缺损的元素个数。

从图12(c)中可以看出,由于削弱了关键区域缺损度较高样本的影响,分解产生的基图像误差进一步降低。

## 4 结语

本文在 NMF 算法的基础上提出了一种基于加权的不完备非负矩阵分解算法 NMFI。实验表明该算法能有效提取未缺损数据中的信息,避免缺损数据对分解结果产生较大的干扰。该算法中权重的取值至关重要,合适的权重不但能降低不完备数据的干扰,而且能最大限度提取残存数据中的信息。关于权重的取值方法有待我们今后进一步的研究。此外,对于图像恢复,本文算法的思想也提供了一个新的思路,为今后在这方面的研究提供一定的借鉴。

(下转第1286页)

### 3 实验及其结果分析

文献[6]采用了UCI数据集<sup>[9]</sup>中16个属性的Zoo数据和9个属性的Winsconsin breast cancer数据,本文为了在更复杂的数据上比较UpdateENBROD和ENBROD算法,采用UCI数据集中的Thyroid disease records数据。该数据集中有3772条记录,每个记录有30个属性值,被标记为negative(3540个记录,占93.8%)和sick(232个记录,占6.2%)。实验中,在标记为negative的数据中选分别随机选出100,200,...,600条数据,求出每个数据对象的 $S_{ROF}$ ,然后在标记为sick的数据中随机选出几条数据插入到前面的数据集中作为异常数据记录,分别采用UpdateENBROD和ENBROD求出每个数据对象的 $S_{ROF}$ ;然后随机删除一些数据对象,再次采用UpdateENBROD和ENBROD分别求出每个数据对象的 $S_{ROF}$ 。实验在CPU为Intel core 2 T5500 1.66 GHz,内存2 GB的PC机上进行,算法用C++实现。

表1给出了算法UpdateENBROD和ENBROD在数据记录数为600,  $k=80$ 发现所有异常记录时排序后的最大记录数。

表1 发现所有异常记录时的最大记录数

异常数据 的条数	发现所有异常数据记录时的最大记录数	
	UpdateENBROD	ENBROD
3	3	3
5	5	5
10	11	11

图3给出了当数据集中有数据对象更新时,算法UpdateENBROD和ENBROD在不同数据记录数下的运行时间曲线图。

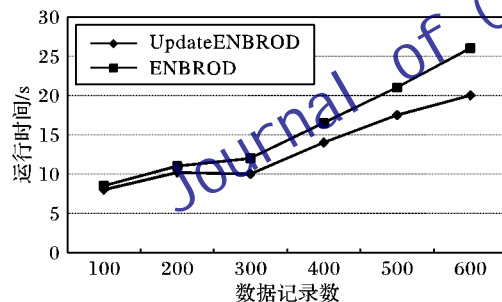


图3 UpdateENBROD和ENBROD的运行时间曲线图

由表1可以看出,当数据集里有3条和5条离群数据记录时,UpdateENBROD与ENBROD都能准确地发现异常的数据,当数据集中插入了10条异常数据,两个算法都在前11条里面发现了所有的异常数据,所以UpdateENBROD与ENBROD都能较高效地发现离群数据的记录,检测离群数据的能力完全一样,从而说明了算法UpdateENBROD的正确性。当数据集中的数据发生动态更新的时候,从图3可以看出,当数据集的记录数不大的时候,两个算法的运行时间相差不大,随着数据集记录数的增加,UpdateENBROD的运行时间小于ENBROD。

### 4 结语

在研究ENBROD算法的基础上,本文提出了一种动态数据环境下的离群点检测算法,当数据集中的数据存在一些更新操作的时候,不必重新计算所有数据点的ROF值,而只需要计算受影响点的ROF值。实验表明,该算法在动态数据环境下的时间性能优于ENBROD。

#### 参考文献:

- [1] HAWKINS D. Identification of outliers [M]. Berlin: Springer-Verlag, 1980.
- [2] 薛安荣,鞠时光,何伟华,等. 局部离群点挖掘算法研究[J]. 计算机学报, 2007, 30(8): 1455-1463.
- [3] BREUNIG M M, KRIEGER H P, NG T, et al. LOF: Identifying density-based local outliers [J]. ACM SIGMOD Record, 2000, 29(2): 93-104.
- [4] CHAWLA S, SUN P E L. SLOM: A new measure for local spatial outliers [J]. Knowledge and Information Systems, 2006, 9(4): 412-429.
- [5] BARBARÁ D, LI Y I, COUTO J. Coolcat: An entropy-based algorithm for categorical clustering [C]// Proceedings of the 11th International Conference on Information and Knowledge Management. New York: ACM Press, 2002: 582-589.
- [6] 于绍越,商琳. 基于信息熵的相对离群点的检测方法: ENBROD [J]. 南京大学学报: 自然科学版, 2008, 44(2): 212-218.
- [7] SHANNON C E. A mathematic theory of communication [J]. Bell System Technical Journal, 1948, XXVH(3): 379-423, 623-656.
- [8] 杨风召,朱扬勇,施伯乐. IncLOF: 动态环境下局部异常的增量挖掘算法[J]. 计算机研究与发展, 2004, 41(3): 477-484.
- [9] MURPHY P M, AHA D W. UCI repository of machine learning database [EB/OL]. (1994-06) [2009-10-06]. <http://www.ics.uci.edu/mllearn/MLRepository.html>.

(上接第1283页)

#### 参考文献:

- [1] LEE D D, SEUNG H S. Learning the parts of objects by nonnegative matrix factorization [J]. Nature, 1999, 401(6755): 788-791.
- [2] XU W E I, LIU X I N, GONG Y I H O N G. Document clustering based on non-negative matrix factorization [C]// Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval. New York: ACM Press, 2003: 267-273.
- [3] GUILLAMET D, VITRIA J. Non-negative matrix factorization for face recognition [C]// Proceedings of the 5th Catalanian Conference on AI: Topics in Artificial Intelligence, LNCS 2504. Berlin: Springer-Verlag, 2002: 336-344.
- [4] VIRTANEN T. Monaural sound source separation by non-negative matrix factorization with temporal continuity and sparseness criteria [J]. IEEE Transactions on Audio, Speech, and Language Process-
- ing, 2007, 15(3): 1066-1074.
- [5] DONOHO D, STODDEN V. When does nonnegative matrix factorization give a correct decomposition into parts? [C]// Proceedings of the 17th Annual Conference Neural Information Processing Systems. Cambridge, MA, USA: MIT Press, 2003: 1141-1148.
- [6] LEE D D, SEUNG H S. Algorithms for non-negative matrix factorization [C]// Advances in Neural Information Processing Systems. Cambridge, MA, USA: MIT Press, 2001, 13: 556-562.
- [7] CAO B I N, SHEN D O U, SUN J I A N -T A O, et al. Detect and track latent factors with online nonnegative matrix factorization [C]// Proceedings of the 20th International Joint Conference on Artificial Intelligence. San Francisco, CA, USA: Morgan Kaufmann Publishers, 2007: 2689-2694.
- [8] KLINGENBERG B, CURRY J, DOUGHERTY A. Non-negative matrix factorization: Ill-posedness and a geometric algorithm [J]. Pattern Recognition, 2009, 42(5): 918-928.