

学校编码: 10384
学号:

分类号_密级
UDC

廈門大學

碩 士 学 位 论 文

若干非负矩阵分解与半非负矩阵分解的
算法及应用

Several Algorithms and Applications of NMF and Semi NMF

指导教师姓名: 林鹭

专业名称: 计算数学

论文提交日期: 2016 年 4 月

论文答辩时间: 2016 年 5 月

学位授予日期: 2016 年 6 月

答辩委员会主席:

评 阅 人:

2016 年 6 月

厦门大学博硕士论文摘要库

厦门大学学位论文原创性声明

本人呈交的学位论文是本人在导师指导下,独立完成的研究成果。本人在论文写作中参考其他个人或集体已经发表的研究成果,均在文中以适当方式明确标明,并符合法律规范和《厦门大学研究生学术活动规范(试行)》。

另外,该学位论文为()课题(组)的研究成果,获得()课题(组)经费或实验室的资助,在()实验室完成。(请在以上括号内填写课题或课题组负责人或实验室名称,未有此项声明内容的,可以不作特别声明。)

声明人(签名):

年 月 日

厦门大学学位论文著作权使用声明

本人同意厦门大学根据《中华人民共和国学位条例暂行实施办法》等规定保留和使用此学位论文，并向主管部门或其指定机构送交学位论文（包括纸质版和电子版），允许学位论文进入厦门大学图书馆及其数据库被查阅、借阅。本人同意厦门大学将学位论文加入全国博士、硕士学位论文共建单位数据库进行检索，将学位论文的标题和摘要汇编出版，采用影印、缩印或者其它方式合理复制学位论文。

本学位论文属于：

1. 经厦门大学保密委员会审查核定的保密学位论文，
于 年 月 日解密，解密后适用上述授权。

2. 不保密，适用上述授权。

（请在以上相应括号内打“√”或填上相应内容。保密学位论文应是已经厦门大学保密委员会审定过的学位论文，未经厦门大学保密委员会审定的学位论文均为公开学位论文。此声明栏不填写的，默认为公开学位论文，均适用上述授权。）

声明人（签名）：

年 月 日

摘要

非负矩阵分解(Nonnegative Matrix Factorization, NMF)是近年来提出的一种新的大规模数据降维的方法;非负矩阵分解应用很广,如人脸识别,医学基因检测,文本聚类分析,模式识别以及盲源信号分离等。非负矩阵分解问题实质上是将非负矩阵近似为非负基矩阵 W 与非负系数矩阵 H 的乘积。这样数据矩阵 V 的列向量等于基矩阵 W 列向量的非负线性组合,这种基于基向量非负线性组合的表示方法反映了极为直观的“局部构成整体”的语义解释。将“局部构成整体”的思想用于一般的数据矩阵上,得到所谓的半非负矩阵分解(Semi NMF),该分解保留了 W 与 H 的可解释性。

本文第一个主要工作是:从线性互补问题出发,基于不动点方程投影梯度提出了三个非负矩阵分解算法。文中先把欧氏距离函数转换为若干个非负最小二乘问题,再利用 KKT 条件,将非负最小二乘问题转换为一个线性互补问题,然后基于线性互补问题提出不动点方程的梯度投影算法。分别依据最速下降法和最小梯度法确定搜索步长。文中证明了这两个非负矩阵分解算法的收敛性,并通过 ORL 人脸数据库进行数值实验,结果表明这三个算法在逼近误差上均优于 Lee 和 Seung 的乘性迭代算法,而在人脸识别的准确率与乘性迭代算法的结果相当。

本文的第二个主要工作是:对半非负矩阵分解,提出了基于罚函数的 Semi-Class 半非负矩阵分解算法。此外,结合最速下降法与最小梯度法搜索步长,进一步提出了基于不动点方程的两个半非负矩阵分解算法。文中证明了这三种算法的收敛性,并将这三种半非负矩阵分解算法应用于医疗检测报告中,同时将最大元准则和最近子空间准则分别作为分类的标准。数值实验表明,本文提出的三个半非负矩阵分解算法在这两个分类标准下的医学诊断准确率均优于 Chris Ding 的 Semi-NMF 算法。

关键词: 非负矩阵分解; 半非负矩阵分解; 不动点投影梯度算法

厦门大学博硕士学位论文摘要库

Abstract

Nonnegative matrix factorization (NMF) is a method proposed in recent year to reduce dimensions of massive data. NMF plays an important role in many real applications, such as face recognition, medical gene detection, text clustering, pattern recognition and blind source separation, etc. The essence of NMF is to find the nonnegative base matrix W and the nonnegative coefficient matrix H so that the nonnegative matrix W and H satisfies $V \approx WH$. This decompose form implies that the columns of V are the nonnegative linear combination of the columns of base matrix W . The representation based on nonnegative linear combination of base vector reflects the intuitive notions of “combining parts to form a whole”. Applying the idea of “combining parts to form a whole” to general matrix V , We obtain the so-called Semi Nonnegative matrix factorization, which also preserves the interpretability of W and H .

One major work of this paper is: From linear complementary problem, three projected gradient NMF algorithms based on fixed point equation are proposed. Firstly, the Euclidean distance function is equivalently converted to several nonnegative least square problems. With KKT condition, nonnegative least squares problems are equivalently converted to a linear complementary problem. Based on the linear complementary problem, a projected gradient NMF algorithm about fixed point equation is proposed. In the algorithm, the search step lengths are selected from the steepest descent method and the minimum gradient method respectively. The convergences of these two algorithms are proved. Numerical experiments on ORL face database show that the three algorithms on the approximation error are superior to multiplicative update algorithm by Lee and Seung. For the face recognition rate, these three algorithms nearly have the same performance than MU algorithm.

Another major work of this paper is: Base on penalty function, a Semi-Class algorithm is proposed. In addition, with the steepest descent and minimum gradient searching step lengths, two projected gradient Semi NMF algorithms on fixed point

equation are proposed. The convergences of these three algorithms are proved. All these three algorithms are applied on analysis medical experiments reports. Maximum rule and Nearest-Subspace method are used respectively as the classification benchmark, and different classification are obtained. Numerical experiments show that the proposed algorithms have a better performance than the Semi-NMF algorithm by Chris Ding on accuracy rate under the classification criteria.

Key words: NMF; Semi NMF; Fixed point Project Gradient algorithm

厦门大学博硕士学位论文摘要库

目录

摘要.....	I
目录.....	V
第一章 引言	1
1.1 非负矩阵分解与半非负矩阵分解研究背景	1
1.2 非负矩阵分解的发展历程	2
1.3 研究现状	2
1.4 本文的创新之处与使用符号定义	4
1.4.1 本文的组织结构.....	4
1.4.2 本文的创新之处.....	4
1.4.3 本文的使用符号定义	5
第二章 NMF 与 Semi NMF 算法介绍.....	7
2.1 预备知识	7
2.1.1 约束最优化问题.....	7
2.1.2 可行域.....	7
2.1.3 约束最优化问题的局部解与全局解.....	7
2.1.4 凸规划问题.....	8
2.1.5 约束最优化问题解的一阶必要条件.....	8
2.1.6 稀疏度定义.....	9
2.2 非负矩阵分解算法与半非负矩阵分解算法介绍	9
2.2.1 乘性迭代算法.....	9
2.2.2 半非负矩阵分解算法.....	9
第三章 非负矩阵分解及其应用	11
3.1 线性互补问题与不动点方程的等价性	11
3.2 基于不动点方程的 NMF 算法以及收敛性讨论	13
3.3 基于不动点方法的 SD 与 MG 算法	16
3.4 数值实验	20

3.4.1 实验数据	20
3.4.2 分类方法	21
3.4.3 KKT 剩余量	21
3.4.4 实验结果与分析	22
第四章 半非负矩阵分解及其应用	27
4.1 半非负矩阵的问题背景	27
4.2 基于罚函数法的半非负矩阵分解算法	27
4.3 基于不动点方程梯度下降的半非负矩阵分解算法	32
4.4 数值实验	34
4.4.1 实验数据	34
4.4.2 分类准则	35
4.4.3 生化全套检验报告数值试验	37
结 论	43
参考文献	45

Contents

Abstract.....	I
Contents	V
Chapter 1 Introduction.....	1
1.1 Backgroud of NMF and Semi NMF	1
1.2 Development progress of NMF	2
1.3 Present situation of research	2
1.4 Framework, innovation and symbol	4
1.4.1 Framework	4
1.4.2 Innovation	4
1.4.3 Symbol	5
Chapter 2 Summary of NMF and Semi NMF algorithm	7
2.1 Preliminary	7
2.1.1 Constrained optimization problems	7
2.1.2 Feasible region	7
2.1.3 Local and global solution of constrained optimization problems	7
2.1.4 Convex programming problem	8
2.1.5 KKT condition of constrained optimization problems	8
2.1.6 Definition of Sparsity.....	9
2.2 NMF and Semi NMF algorithms	9
2.2.1 MU algorithm.....	9
2.2.2 Semi-NMF algorithm.....	9
Chapter 3 Algorithm and application of NMF	10
3.1 Equivalence of LCP and fixed point equation	12
3.2 Convergence of algorithm based on fixed point equation	13
3.3 SD and MG algorithms based on fixed point equation.....	16
3.4 Numerical Experiments.....	21
3.4.1 Experimental data	21

3.4.2 Classificaiton method.....	21
3.4.3 KKT Residual	22
3.4.4 Result and analysis of experiment	22
Chapter 4 Algorithm and application of Semi NMF	27
4.1 Background of Semi NMF problem	27
4.2 Semi NMF algorithm with penalty function.....	27
4.3 Semi NMF algorithm based on SD and MG method	32
4.4 Numerical experiment	34
4.4.1 Experimental data	34
4.4.2 Classificaiton method.....	35
4.4.3 Numerical experiment of medical diagnostic report.....	37
Conclusions	43
Reference.....	45

第一章 引言

1.1 非负矩阵分解与半非负矩阵分解研究背景

随着社会进入大数据时代，需要处理和分析的大规模数据与日俱增，如卫星实时传输的大量图片数据，数据库中储存的大量文本数据，网络的视频、音乐、图像以及新闻等，这些信息常用多变量组成的向量数据表示，构成一组高维数据。高维数据提供了客观事物的详细信息，能够很好的保存不同事物的特征。然而，数据维数的增多会造成数据的冗余以及处理的困难，希望揭示隐藏在这些数据背后的客观规律成为越来越迫切的要求。因此，必须寻求一种处理高维数据的有效方法，以便从庞大的数据中挖掘出潜在的本质信息。

数据降维是解决这类问题的常用方法。主成分分析(PCA)方法与奇异值分解(SVD)方法均属于经典的数据降维方法。但对于非负数据，这两个方法得到的特征矩阵包含负元素，所以特征矩阵的列向量通常不具有可解释性。因此，探讨数据矩阵的压缩降维，使得降维后的特征矩阵能揭示数据向量之间的内在关系，且分解有很好的表示意义成为当前的一个研究热点，有重要的理论意义和实用价值。

非负矩阵分解(Nonnegative Matrix Factorization, NMF)是近年来提出的一种对大规模非负数据降维的方法。非负矩阵分解已经广泛地应用在人脸识别^[1]，图像分析^[2]，文本潜在分析^[3]，盲源信号分类^[4]，生物基因探测^[5]，以及数据挖掘^[6]等领域。具体描述如下：给定一个非负矩阵 $V \in R_+^{m \times n}$ 和一个正整数 r ，满足 $r \ll \min(m, n)$ ，求两个非负矩阵 $W \in R_+^{m \times r}$ 以及 $H \in R_+^{r \times n}$ ，使得

$$\begin{aligned} & \min_{W, H} f(W, H) \\ & \text{subject to } W \geq 0, H \geq 0 \end{aligned}$$

其中 $f(W, H)$ 是关于 W 和 H 的目标函数， W 称为基矩阵，而 H 称为系数矩阵。可以看出非负矩阵分解是一种数据降维方法。 $V \approx WH$ 按列表示为 $V_j \approx \sum_{k=1}^r W_k H_{kj}$ ， $j=1, 2, \dots, n$ ，即： V_j 由 W_1, W_2, \dots, W_r 非负线性组合得到，这样的表示方法反映了的“局部表示整体”的语义解释。 H_{kj} 的大小，代表着 W_k 对 V_j 贡献的大小。如果 H_{kj} 远大于 H_j 的其他分量，说明 V_j 主要是由 W_k 所构成， W_k 就可看做是 V_j 的主要

特征列；如果 H_{kj} 很小或者接近0，则代表着 W_k 起的作用不大。

那么，对一般的数据矩阵可否设计降维方法，也使得降维后的特征矩阵能揭示数据向量之间的内在关系，且有很好的表示意义呢？半非负矩阵的分解算法^[7] (Semi NMF)就是这样一种方法，它将数据矩阵 V 分解为特征矩阵 W 与非负矩阵 H 的乘积，使得 V_j 由 W_1, W_2, \dots, W_r 非负线性组合得到，这样的表示方法同样反映了极为直观的“局部表示整体”的语义解释。

1.2 非负矩阵分解的发展历程

1994年，Paatero和Tapper^[8]对环境数据进行分解，建立了以下优化模型：

$$\min_{W, H \geq 0} \|Q(V - WH)\|_F^2$$

他们针对上述模型，提出了带非负约束的交替方向最小二乘迭代算法(ALS)。但他们并没有从理论上证明算法的收敛性。Lee和Seung于1999年首次提出非负矩阵分解的概念^[9]，他们通过人脸图片数值实验，发现分解得到的特征列(W 的列向量)的图像形似人脸的“部件”。同时，由于系数矩阵的元素非负，这样人脸就由不同的“部件”叠加而成，这反映了极为直观的“局部表示整体”的语义解释。2001年，他们又提出交替方向的乘性迭代算法(MU)^[10]，证明了算法的收敛性。

1.3 研究现状

经过十几年的发展，NMF的研究日益深入，学者们相继提出各类NMF算法。根据对NMF问题的约束条件是否仅限于非负性，将算法大致分为基本NMF(Basic NMF, BNMF)与改进NMF算法(Improved NMF, INMF)^[11]。解的稀疏性可以降低数据的冗余度与存储量，而且更加凸显数据的潜在局部特征。要求NMF算法分解结果具有稀疏性是常见的除非负外的其他约束条件，常通过添加罚项到目标函数中实现。Heiler等人基于经典的凸二次规划方法，得到一类稀疏且收敛的NMF算法^[12]。Hoyer从罚函数的角度出发提出带稀疏约束条件的NMF算法^[13]，他将欧氏距离($\frac{1}{2}\|V - WH\|_F^2$)加上由1范数定义的稀疏性罚项一起作为目标函数，对 W 采取梯度投影的方法，而对 H 采取类似最大期望(expectation Maximization, EM)算法^[6]，构造了非负稀疏编码算法。Liu等在Hoyer的基础上提出了稀疏NMF算法^[14]，其思路与Hoyer类似，但区别在于：(1) Liu等选择的

目标函数为 K-L 散度和 1 范数定义的稀疏性罚函数组合而成；(2) 他们对 W 和 H 都采用了类似 EM 算法的方式优化，得到的迭代格式与 MU 算法类似。Heiler 等对 Hoyer 和 Liu 的工作又做了进一步发展^[6]，他们以欧氏距离加上由 1 范数定义的稀疏性罚项为目标函数，把一个稀疏性区间对应的可行域表示为不同二阶锥的集合运算，并将该集合的限制条件加入 NMF 问题中，在二阶锥规划的框架下构造了两个具有稀疏性的算法；数值实验结果表明，一个算法效率较高但出现震荡，另一个效率较低但能保证收敛性。由于过度稀疏的数据不具备可解释性，如何很好的控制数据的稀疏性与数据的解释性之间的矛盾就成为一个新的问题。Pascual 和 Montano 等人在上述的工作基础上提出了非平滑 NMF(Non-Smooth NMF, NSNMF)算法^[15]，他们在 NMF 问题中增加了用来控制平滑度的矩阵 S ，如果 S 非常平滑，则算法得到的 W 和 H 都是稀疏的。实验结果表明，NSNMF 算法很好地平衡了解析结果的稀疏性与数据的表述之间的矛盾。Kim 和 Park 提出基于非负最小二乘问题的分解算法^[16]，他们将 NMF 问题等价转化为两个非负最小二乘问题，并在目标函数加上了罚函数项，利用激活集算法，得到新的 INMF 算法。Stan Z. Li 等提出了局部非负矩阵分解^[17]，文章加入的约束条件能够使基之间的冗余度最小化、基的总体“活跃”程度最大化，从而使算法得到较为稀疏的解。C. J. Lin 提出的基于投影梯度的非负矩阵分解算法^[18]，将 W 与 H 非负的限制条件换为带边界的限制条件，运用梯度下降投影方法，得到改进的 NMF 算法。Li 和 Zhang 提出的基于秩 1 分解的算法^[19]，将目标函数写成等价的按元素的形式，从而把矩阵优化问题转化为若干个二次函数的最优化问题，通过二次函数的性质，得到了新的 BNMF 算法。陈卫刚等利用可行方向(Feasible Direction, FD)方法和模拟退火(Simulated Annealing, SA)算法结合，以欧氏距离为目标函数，构造出 FD-SA-NMF 算法^[20]，数值试验表明该方法比 MU 算法收敛更快，且具有更好的稀疏性，但算法的时间复杂度较高。

近些年来，NMF 算法在数据分类上有大量应用，得到了许多分类效果较好的 NMF 算法。Wang 等提出了基于 Fisher 非负矩阵分解算法(FNMF)^[21]，该算法对样本数据进行分类，然后以欧氏距离加上类内与类间散度罚函数项作为目标函数，构造了带分类信息的 INMF 算法，算法的目的是使得样本的类内散度最小，类间散度最大，从而能够保证得到的解具有良好的分类效果。Xue 等^[22]也以欧

氏距离加上类内与类间散度罚函数项作为目标函数,但用梯度下降方法对算法进行推导,得到新的分类 INMF 算法。

关于半非负矩阵分解,文献不多见。2008年 Chris Ding^[7]提出了一种全新的半非负矩阵分解算法,证明了该算法的收敛性。之后 Nicolas Gillis^[23]中从理论上证明了分解系数 r 对半非负矩阵分解问题的影响;并且证明了当 $r=1$ 时,半非负矩阵分解问题是一个 NP-hard 问题。Nicolas Gillis 在文中还提到,若对系数矩阵 H 再加上正交性的限制,此时半非负矩阵分解问题等价于 K-means 模型,对应的半非负矩阵分解算法可看做是一类硬聚类方法。因此,仅对 H 做非负限制的半非负矩阵分解算法可看做是一类软聚类方法。半非负矩阵分解鲜有成果发表,处于起步阶段,如何寻求更高效且实用的半非负矩阵分解算法,尚有很多工作可做。

1.4 本文的创新之处与使用符号定义

1.4.1 本文的组织结构

本文由“引言”,“预备知识与符号说明”,“算法介绍”,“若干非负矩阵分解算法”,“若干半非负矩阵分解算法”,“结论”六部分组成。

第一章介绍本文的研究背景和意义,简单概述非负矩阵分解与半非负矩阵分解的研究成果,给出了本文的创新之处,及使用符号说明。

第二章介绍相关预备知识,简要叙述经典的非负矩阵分解与半非负矩阵分解算法。

第三章为本文的第一部分核心内容,提出了基于不动点梯度下降的三个非负矩阵分解算法,证明了算法的收敛性,通过 ORL 人脸图库进行数值实验。

第四章为本文的第二部分核心内容,提出了基于罚函数与基于不动点梯度下降的三个半非负矩阵分解算法,证明了算法的收敛性,最后将该方法应用于医疗检测报告上。

第五章是全文的结论,同时对未来的研究方向做了展望。

1.4.2 本文的创新之处

关于非负矩阵分解与半非负矩阵分解算法,多数文献着重点在于提出新算法,但算法的收敛性证明不常见,且较少分析分类效果。

(1) 从 NMF 问题所对应的线性互补问题出发,提出了三个基于不动点梯度投影的 NMF 算法,分别为 PMG、FP-SD 与 FP-MG 算法,并证明了三个算法的

Degree papers are in the “[Xiamen University Electronic Theses and Dissertations Database](#)”.

Fulltexts are available in the following ways:

1. If your library is a CALIS member libraries, please log on <http://etd.calis.edu.cn/> and submit requests online, or consult the interlibrary loan department in your library.
2. For users of non-CALIS member libraries, please mail to etd@xmu.edu.cn for delivery details.