

矩阵补全算法研究进展

史加荣¹ 郑秀云¹ 周水生²

(西安建筑科技大学理学院 西安 710055)¹ (西安电子科技大学数学与统计学院 西安 710071)²

摘要 作为压缩感知理论的重要发展,矩阵补全与恢复已成为信号与图像处理的一种新的强有力的工具。综述了矩阵补全算法的最新研究进展。首先分析了核范数最小化模型的几种主要的矩阵补全算法,并对这些算法的迭代过程及原理进行了详细的阐述。其次讨论了矩阵补全的低秩矩阵分解模型,并列出了近年来出现的求解此模型的新算法。然后补充了上述两种模型的衍生版本,指出了相应的求解方法。在数值实验中,对文中所讨论的主要矩阵补全算法的性能进行了比较。最后给出了矩阵补全算法的未来研究方向及重点。

关键词 矩阵补全,低秩,核范数最小化,低秩矩阵分解,压缩感知,低秩矩阵恢复

中图分类号 TP391 **文献标识码** A

Research Progress in Matrix Completion Algorithms

SHI Jia-rong¹ ZHENG Xiu-yun¹ ZHOU Shui-sheng²

(School of Science, Xi'an University of Architecture and Technology, Xi'an 710055, China)¹

(School of Mathematics and Statistics, Xidian University, Xi'an 710071, China)²

Abstract As an important development of compressed sensing theory, matrix completion and recovery has been a new and remarkable technique for signal and image processing. This paper made a survey on the latest research progress in matrix completion algorithms. Firstly, it analyzed several main algorithms to nuclear norm minimization model, and elaborated their iterative procedure and principle. Secondly, it discussed low-rank matrix factorization model of matrix completion and listed the corresponding new algorithms emerged in recent years. Then it complemented other versions derived from the above two models and pointed out the solving methods. In numerical experiments, performance comparisons were made on the main algorithms to matrix completion. Finally, it gave future research direction and focus for matrix completion algorithms.

Keywords Matrix completion, Low-rank, Nuclear norm minimization, Low-rank matrix factorization, Compressed sensing, Low-rank matrix recovery

1 引言

一般来说,不能根据信号的部分采样元素来恢复所有元素。但信号在一组基下是稀疏的且满足一定条件时,压缩感知理论(Compressed Sensing/Compressed Sampling, CS)证实了可以通过求解 l_1 最小化问题来精确地恢复所有元素^[1,2]。当信号用矩阵形式表示时,同样一般不能根据它的部分元素来恢复所有丢失元素。Candès 和陶哲轩等人证明了当矩阵的奇异值具有稀疏性(即矩阵是低秩的)且采样数目满足一定条件时,大多数矩阵可以通过求解核范数最小化问题来精确地恢复所有元素^[3]。由矩阵的部分元素来恢复所有元素这一问题也称为矩阵补全(Matrix Completion, MC)。若将部分采样元素这一约束推广到一般的线性约束函数,则矩阵补全就称为低秩矩阵恢复(Low-Rank Matrix Recovery, LRMR)^[4,5]。

矩阵补全广泛地应用在计算机视觉^[6-9]、机器学习^[10,11]、推荐系统^[12]和系统辨识^[13,14]等诸多科学与工程领域中。作为信号与图像处理技术的一个强大的新兴分支,矩阵补全已成为继压缩感知之后的又一种重要的信号获取工具。在数学形式上,矩阵补全可描述为一个仿射秩最小化问题,但此问题是 NP 难的。近年来涌现出了许多求解秩最小化问题的启发式方法,这些方法主要分为两类:一类是将秩函数凸松弛到矩阵核范数,建立核范数优化模型;另一类是事先给定矩阵的秩,建立低秩分解模型。

与向量的 l_1 范数类似,矩阵的核范数也是一个连续的、不可微的凸函数。梯度或次梯度下降法是求解核范数优化模型的主要方法。矩阵核范数最小化模型等价于半定规划问题^[15,16],求解半定规划常用的方法是基于二阶梯度信息的内点算法^[17],但此方法计算复杂度非常高。因此,求解压缩感知的众多基于一阶梯度的方法被纷纷应用到矩阵核范数最小

到稿日期:2013-06-02 返修日期:2013-09-17 本文受国家自然科学基金(61179040),陕西省教育厅专项科研计划项目(2013JK0587, 2013JK0588, 2010JK642)资助。

史加荣(1979—),男,博士,副教授,CCF 会员,主要研究方向为机器学习与模式识别, E-mail: shijiarong@xauat.edu.cn; 郑秀云(1982—),女,博士,讲师,主要研究方向为最优化方法与应用; 周水生(1972—),男,博士,教授,博士生导师,主要研究方向为机器学习与模式识别。

化问题中。低秩矩阵分解是矩阵补全的另一种重要的模型，它先将数据矩阵近似分解为两个低秩矩阵之积，再通过求解非凸的低秩逼近问题来恢复丢失元素。求解这类模型的经典方法有 Wiberg 算法、阻尼牛顿算法、Levenberg-Marquardt 算法和幂迭代法等^[6,7]。由于低秩矩阵分解模型是 Grassmann 或 Riemann 流形上的优化问题，因此也可以使用几何子空间方法求解。

本文对矩阵补全算法的最新研究进展进行了综述。在第 2 节，引入了矩阵补全算法所需的矩阵算子和核范数等方面的基础知识；第 3 节分析了求解核范数最小化问题的主要算法；第 4 节讨论了求解低秩矩阵分解模型的最新方法；第 5 节补充了前两种模型的推广版本；在第 6 节，通过数值实验来比较矩阵补全的主要算法的性能；最后对矩阵补全算法进行总结，并给出了进一步的研究方向和重点。

2 基本知识

先给出矩阵分解的一种重要且常用的形式：奇异值分解 (Singular Value Decomposition, SVD)。

定理 1 对于秩为 r 的矩阵 $X \in \mathbb{R}^{m \times n}$ ，它的奇异值分解为： $X = U\Sigma V^T$ ，其中 $U \in \mathbb{R}^{m \times r}$ 和 $V \in \mathbb{R}^{n \times r}$ 满足 $U^T U = V^T V = I$ ，(I 为 r 阶单位矩阵)， $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$ 且其对角线元素满足 $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ 。

称 σ_i 为 X 的第 i 大奇异值， $i = 1, 2, \dots, r$ 。矩阵 U 有 mr 个元素，又由于 U 的各个列向量具有规范正交性，因此它的自由度为 $mr - (1 + 2 + \dots + r)$ 。类似地， V 的自由度为 $nr - (1 + 2 + \dots + r)$ 。故 X 的自由度为

$$dr(X) = (mr - (1 + 2 + \dots + r)) + (nr - (1 + 2 + \dots + r)) + r = (m + n - r)r \quad (1)$$

定义 1 矩阵 $X, Y \in \mathbb{R}^{m \times n}$ ，它们的内积为 $\langle X, Y \rangle = \sum_{i=1}^m \sum_{j=1}^n X_{ij} Y_{ij}$ ，其中 X_{ij} 表示矩阵 X 的第 i 行第 j 列元素。

定义 2 设 $\sigma_1, \sigma_2, \dots, \sigma_r$ 为矩阵 X 的所有非零奇异值，则它的 Schatten p -范数 ($p > 0$) 为 $\|X\|_p = (\sum_{i=1}^r \sigma_i^p)^{1/p}$ 。

特殊地，当 $p = 1$ 时， $\|X\|_{s_1}$ 为 X 的核范数 (也称为迹范数、Ky Fan 范数)，可用 $\|X\|_*$ 表示；当 $p = 2$ 时， $\|X\|_{s_2}$ 为 X 的 Frobenious 范数，可用 $\|X\|_F$ 表示；当 $p \rightarrow +\infty$ 时， $\|X\|_{s_\infty} = \sigma_1$ 为 X 的算子范数 (或谱范数)，可用 $\|X\|$ 表示；规定 $\|X\|_{s_0} = r$ ，可用 $\text{rank}(X)$ 表示。矩阵 X 的 Frobenious 范数满足： $\|X\|_F = \sqrt{\langle X, X \rangle} = \sqrt{\text{trace}(X^T X)} = \sqrt{\sum_{i=1}^r \sigma_i^2}$ ，这里 $\text{trace}(\cdot)$ 表示矩阵的迹算子。矩阵 X 的核范数、Frobenious 范数和算子范数满足下列不等式^[15]：

$$\|X\| \leq \|X\|_F \leq \|X\|_* \leq \sqrt{r} \|X\|_F \leq r \|X\| \quad (2)$$

定义 3 在内积空间中， X 的 Schatten p -范数的对偶范数为

$$\|X\|_{s_p} = \max_Y \{\langle X, Y \rangle : \|Y\|_{s_p} \leq 1\} \quad (3)$$

易知 Frobenious 范数的对偶范数仍为 Frobenious 范数，算子范数的对偶范数为核范数。

定义 4 对于多元实函数 $f(x)$ ， $x \in \mathcal{C}$ ，它的凸包络为满足

$$g(x) \leq f(x), \forall x \in \mathcal{C} \quad (4)$$

的最大凸函数 g 。

定理 2^[16] 秩函数 $\text{rank}(X)$ 在集合 $\mathcal{C} = \{X \in \mathbb{R}^{m \times n} : \|X\| \leq 1\}$ 上的凸包络函数为 $\phi(X) = \|X\|_*$ 。

凸包络函数 $\phi(X)$ 是连续的凸函数，但不光滑。若非零矩阵 X 的算子范数 $\|X\| = k$ ，则由定理 2 可知 $\text{rank}(X)$ 的凸包络函数为 $\|X\|_*/k$ 。

定理 3^[18] 函数 $\phi(X) = \|X\|_*$ 的次梯度为

$$\partial \|X\|_* = \{UV^T + M : \|M\| \leq 1, U^T M = 0, M V = 0\} \quad (5)$$

其中， $U \in \mathbb{R}^{m \times r}$ 、 $V \in \mathbb{R}^{n \times r}$ 分别为 X 的左奇异矩阵和右奇异矩阵。

根据定理 3，可以证明下列结论。

定理 4^[18] 设秩为 r 的 $m \times n$ 维实矩阵 Q 的奇异值分解为： $Q = U\Sigma V^T = U \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r) V^T$ ，则核范数最小化问题 $\min_X \|X - Q\|_*/2 + \epsilon \|X\|_*$ 有闭形式解 $X^* = \mathcal{Q}_\epsilon(Q) = U \mathcal{L}_\epsilon(\Sigma) V^T$ ，其中 $\epsilon \geq 0$ ， $\mathcal{L}_\epsilon(\Sigma) = \text{diag}(\max(\sigma_1 - \epsilon, 0), \dots, \max(\sigma_r - \epsilon, 0))$ 。

3 核范数优化模型的矩阵补全算法

3.1 问题描述

对于低秩矩阵 $M \in \mathbb{R}^{m \times n}$ ，希望根据它的部分观测元素来恢复所有元素，此恢复过程可用下列仿射秩最小化问题来描述：

$$\min_X \text{rank}(X), \text{ s. t. } X_{ij} = M_{ij}, (i, j) \in \Omega \quad (6)$$

其中， $X \in \mathbb{R}^{m \times n}$ ，指标集 $\Omega \subset \{1, 2, \dots, m\} \times \{1, 2, \dots, n\}$ 。记观测或采样元素数目为 p ，则 $p = |\Omega|$ ，这里 $|\cdot|$ 表示集合的势。在上述优化问题中，约束条件也可以表示为 $\mathcal{P}_\Omega(X) = \mathcal{P}_\Omega(M)$ 或 $\mathcal{A}(X) = b$ ，其中 \mathcal{P}_Ω 为投影算子，其定义为

$$\mathcal{P}_\Omega(X_{ij}) = \begin{cases} X_{ij}, & \text{若 } (i, j) \in \Omega \\ 0, & \text{否则} \end{cases} \quad (7)$$

$\mathcal{A}: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$ 为线性算子， $b \in \mathbb{R}^p$ 。

矩阵的秩函数是非连续、非凸的，故直接求解秩最小化问题比较困难。根据定理 2，可将秩函数凸松弛到核范数，于是得到下列凸优化问题

$$\min_X \|X\|_*, \text{ s. t. } X_{ij} = M_{ij}, (i, j) \in \Omega \quad (8)$$

下面给出求解上述核范数最小化问题的主要算法。

3.2 半定规划法

当矩阵 X 给定时，由定义 3 知 X 的核范数等于下列优化问题的最优值

$$\max_Y \langle X, Y \rangle, \text{ s. t. } \|Y\| \leq 1 \quad (9)$$

此优化问题也等价于

$$\begin{aligned} & \max_Y \text{trace}(X^T Y) \\ & \text{ s. t. } \begin{pmatrix} I_m & Y \\ Y^T & I_n \end{pmatrix} \geq 0 \end{aligned} \quad (10)$$

其中，“ $A \geq 0$ ”表示 A 为半正定矩阵。上述半定规划问题的对偶规划为

$$\begin{aligned} & \min_{W_1, W_2} \frac{1}{2} (\text{trace}(W_1) + \text{trace}(W_2)) \\ & \text{ s. t. } \begin{pmatrix} W_1 & X \\ X^T & W_2 \end{pmatrix} \geq 0 \end{aligned} \quad (11)$$

因此，核范数优化问题(8)等价于下列半定规划

$$\begin{aligned} & \min_{x, W_1, W_2} \frac{1}{2} (\text{trace}(W_1) + \text{trace}(W_2)) \\ & \text{ s. t. } \begin{pmatrix} W_1 & X \\ X^T & W_2 \end{pmatrix} \geq 0, X_{ij} = M_{ij}, (i, j) \in \Omega \end{aligned} \quad (12)$$

求解半定规划常用的方法是内点算法,使用此方法求解上述最优化问题的计算复杂度为 $O(p(m+n)^3 + p^2(m+n)^2 + p^3)$ 。现有的求解半定规划的 MATLAB 软件包主要有 CVX、SDPT3 和 SeDuMi。下面以 CVX 为例,给出求解半定规划(12)的 MATLAB 代码。

```
%输入低秩矩阵 M、采样行指标向量 I 和列指标向量 J
[m,n]=size(M);J=J+m;
cvx_setup
cvx_begin sdp
variables WX(m+n,m+n)
minimize(trace(WX(1:m,1:m))+trace(WX(m+1:m+n,m+1:m+n)))
for i=1:length(I),WX(I(i),J(i))=M(I(i),J(i)-m);end
WX>=0
cvx_end
MR=full(WX(1:m,m+1:end));%MR 为 M 的恢复矩阵
```

由于求解半定规划的内点算法具有非常高的复杂度,因此当 m 或 n 或 p 比较大时,半定规划法需要较长的运行时间,甚至不可行。例如:取 $p=0.3mn$,当 $m=n=100$ 时,运行时间大约 1 分钟;当 $m=n=120$ 时,运行时间大约 5 分钟;当 $m=n \geq 200$ 时, MATLAB 会溢出。为此,人们提出了许多基于一阶梯度的算法。

3.3 奇异值阈值算法

Cai Jian-feng 和 Candès 等人提出了求解核范数最小化问题的奇异值阈值算法 (Singular Value Thresholding, SVT)^[18]。此算法先将最优化问题(8)正则化,即有

$$\min_{\mathbf{X}} \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X}\|_F^2, \text{ s. t. } \mathcal{P}_\Omega(\mathbf{X}) = \mathcal{P}_\Omega(\mathbf{M}) \quad (13)$$

其中, $\tau > 0$ 。当参数 $\tau \rightarrow +\infty$ 时,上述最优化问题的最优解收敛到(8)的最优解。构造最优化问题(13)的拉格朗日函数

$$L(\mathbf{X}, \mathbf{Y}) = \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X}\|_F^2 + \langle \mathbf{Y}, \mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}) \rangle \quad (14)$$

其中,拉格朗日乘子 $\mathbf{Y} \in \mathbb{R}^{m \times n}$ 。如果 $(\mathbf{X}^*, \mathbf{Y}^*)$ 为优化问题(13)的原-对偶问题的最优解,则有

$$\sup_{\mathbf{Y}} \inf_{\mathbf{X}} L(\mathbf{X}, \mathbf{Y}) = \inf_{\mathbf{X}} \sup_{\mathbf{Y}} L(\mathbf{X}, \mathbf{Y}) = L(\mathbf{X}^*, \mathbf{Y}^*) \quad (15)$$

SVT 算法使用交替迭代方法求解优化问题(13)或(15)。初始化 $\mathbf{Y}^0 = 0$, 当 \mathbf{Y}^{k-1} 给定时,

$$\begin{aligned} \mathbf{X}^k &= \arg \min_{\mathbf{X}} L(\mathbf{X}, \mathbf{Y}^{k-1}) \\ &= \arg \min_{\mathbf{X}} \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X}\|_F^2 - \langle \mathbf{Y}^{k-1}, \mathcal{P}_\Omega(\mathbf{X}) \rangle \\ &= \arg \min_{\mathbf{X}} \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X}\|_F^2 - \langle \mathbf{X}, \mathcal{P}_\Omega(\mathbf{Y}^{k-1}) \rangle \\ &= \mathcal{D}_\tau(\mathcal{P}_\Omega(\mathbf{Y}^{k-1})) \end{aligned} \quad (16)$$

当 \mathbf{X}^k 给定时,记 $g(\mathbf{Y}) = L(\mathbf{X}^k, \mathbf{Y})$, 则 $\partial g(\mathbf{Y}) = \partial_{\mathbf{Y}} L(\mathbf{X}^k, \mathbf{Y}) = \mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}^k)$ 。使用梯度下降法(Uzawa 算法)来更新 \mathbf{Y} :

$$\mathbf{Y}^k = \mathbf{Y}^{k-1} + \delta_k \partial g(\mathbf{Y}^{k-1}) = \mathbf{Y}^{k-1} + \delta_k \mathcal{P}_\Omega(\mathbf{M} - \mathbf{X}^k) \quad (17)$$

其中, $\delta_k > 0$ 为步长大小, $k \geq 1$ 。由更新公式(17)及 \mathbf{Y} 的初始化可知迭代公式(16)还等价于 $\mathbf{X}^k = \mathcal{D}_\tau(\mathbf{Y}^{k-1})$ 。

上述迭代算法收敛速度往往比较慢,文献[19, 20]分别给出了 SVT 算法的加速版本。

3.4 不动点延续算法

最优化问题(8)的另一种正则化版本为

$$\min_{\mathbf{X}} \mu \|\mathbf{X}\|_* + \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{X}) - \mathcal{P}_\Omega(\mathbf{M})\|_F^2 \quad (18)$$

其中, $\mu > 0$ 。上述优化问题的目标函数关于 \mathbf{X} 的次梯度为 $\mu \partial \|\mathbf{X}\|_* + \mathcal{P}_\Omega(\mathbf{X}) - \mathcal{P}_\Omega(\mathbf{M})$ 。若 \mathbf{X}^* 为最优化问题(18)的最优解,则它满足

$$\begin{aligned} 0 &\in \mu \tau \partial \|\mathbf{X}^*\|_* + \tau(\mathcal{P}_\Omega(\mathbf{X}^*) - \mathcal{P}_\Omega(\mathbf{M})) \\ &= \mu \tau \partial \|\mathbf{X}^*\|_* + \mathbf{X}^* - (\mathbf{X}^* - \tau(\mathcal{P}_\Omega(\mathbf{X}^*) - \mathcal{P}_\Omega(\mathbf{M}))) \end{aligned} \quad (19)$$

记 $\mathbf{Y}^* = \mathbf{X}^* - \tau(\mathcal{P}_\Omega(\mathbf{X}^*) - \mathcal{P}_\Omega(\mathbf{M}))$, 则 $0 \in \mu \tau \partial \|\mathbf{X}^*\|_* + \mathbf{X}^* - \mathbf{Y}^*$, 即 \mathbf{X}^* 是核范数最小化问题

$$\min_{\mathbf{X}} \mu \tau \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X} - \mathbf{Y}^*\|_F^2 \quad (20)$$

的最优解。综上分析,不动点延续算法(Fixed Point Continuation, FPC)^[21]的迭代过程如下

$$\begin{cases} \mathbf{Y}^k = \mathbf{X}^{k-1} - \tau(\mathcal{P}_\Omega(\mathbf{X}^{k-1}) - \mathcal{P}_\Omega(\mathbf{M})) \\ \mathbf{X}^k = \mathcal{D}_{\mu \tau}(\mathbf{Y}^k) \\ \mu_{k+1} = \max(\mu_k \eta_\mu, \bar{\mu}) \end{cases} \quad (21)$$

其中, $\bar{\mu} > 0$ 为参数 μ 取值的下限, $0 < \eta_\mu < 1$ 为常数。文献[22]还建议将 FPC 和 Bregman 迭代技术结合在一起。

3.5 加速近端梯度算法

$$\text{记 } f(\mathbf{X}) = \frac{1}{2} \|\mathcal{P}_\Omega(\mathbf{X}) - \mathcal{P}_\Omega(\mathbf{M})\|_F^2, g(\mathbf{X}) = \mu \|\mathbf{X}\|_*,$$

则最小化问题(18)的目标函数为 $F(\mathbf{X}) = f(\mathbf{X}) + g(\mathbf{X})$ 。显然 $f(\mathbf{X})$ 为凸光滑函数, $g(\mathbf{X})$ 为连续、不光滑的凸函数。对于任意给定的 $\mathbf{Y} \in \mathbb{R}^{m \times n}$, 考虑 $F(\mathbf{X})$ 在 \mathbf{Y} 处的部分二次逼近

$$\begin{aligned} F(\mathbf{X}) &\approx Q(\mathbf{X}, \mathbf{Y}) \\ &= f(\mathbf{Y}) + \langle \mathbf{X} - \mathbf{Y}, \partial f(\mathbf{Y}) \rangle + \frac{1}{2t} \|\mathbf{X} - \mathbf{Y}\|_F^2 + g(\mathbf{X}) \end{aligned} \quad (22)$$

根据上述逼近公式, Toh 等人提出了加速近端梯度算法 (Accelerated Proximal Gradient, APG)^[23]。APG 算法交替更新 \mathbf{X}, \mathbf{Y} 和 t , 其迭代过程如下

$$\begin{aligned} \mathbf{X}^k &= \arg \min_{\mathbf{X}} Q(\mathbf{X}, \mathbf{Y}^{k-1}) \\ &= \arg \min_{\mathbf{X}} \frac{1}{2} \|\mathbf{X} - (\mathbf{Y}^{k-1} - t_k \partial f(\mathbf{Y}^{k-1}))\|_F^2 + \mu t_k \|\mathbf{X}\|_* \\ &= \mathcal{D}_{\mu t_k}(\mathbf{Y}^{k-1} - t_k(\mathcal{P}_\Omega(\mathbf{Y}^{k-1}) - \mathcal{P}_\Omega(\mathbf{M}))) \end{aligned} \quad (23)$$

$$\mathbf{Y}^k = \mathbf{X}^k + \frac{t_{k-1} - 1}{t_k} (\mathbf{X}^k - \mathbf{X}^{k-1}) \quad (24)$$

$$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2} \quad (25)$$

在上述迭代过程中,文献[23]还提出并入类似线性搜索的加速算法。

3.6 增广拉格朗日乘子法

通过引入 $m \times n$ 维实矩阵变量 \mathbf{E} , 最小化问题(8)可重新表示为

$$\min_{\mathbf{X}, \mathbf{E}} \|\mathbf{X}\|_*, \text{ s. t. } \mathbf{X} + \mathbf{E} = \mathbf{M}, \mathcal{P}_\Omega(\mathbf{E}) = 0 \quad (26)$$

其中, \mathbf{M} 的丢失元素设置为 0。上述最优化问题的部分增广拉格朗日函数为

$$\begin{aligned} L(\mathbf{X}, \mathbf{E}, \mathbf{Y}, \mu) &= \|\mathbf{X}\|_* + \langle \mathbf{Y}, \mathbf{M} - \mathbf{X} - \mathbf{E} \rangle + \frac{\mu}{2} \|\mathbf{M} - \mathbf{X} - \mathbf{E}\|_F^2 \end{aligned} \quad (27)$$

其中,正则化参数 $\mu > 0$ 。初始化 $\mathbf{E}^0 = \mathbf{Y}^0 = 0$, $\mu_0 > 0$, $\rho > 1$, $\epsilon > 0$, 则不精确增广拉格朗日乘子法 (Inexact Augmented Lagrange Multipliers, IALM)^[24] 交替更新 $\mathbf{X}, \mathbf{E}, \mathbf{Y}, \mu$, 其迭代公式为

$$\begin{aligned} \mathbf{X}^k &= \arg \min_{\mathbf{X}} L(\mathbf{X}, \mathbf{E}^{k-1}, \mathbf{Y}^{k-1}, \mu_{k-1}) \\ &= \arg \min_{\mathbf{X}} \frac{1}{\mu_{k-1}} \|\mathbf{X}\|_* + \frac{1}{2} \|\mathbf{X} - (\mathbf{M} - \mathbf{E}^{k-1} + \mathbf{Y}^{k-1} / \mu_{k-1})\|_F^2 \\ &= \mathcal{D}_{1/\mu_{k-1}}(\mathbf{M} - \mathbf{E}^{k-1} + \mathbf{Y}^{k-1} / \mu_{k-1}) \end{aligned} \quad (28)$$

$$\begin{aligned} \mathbf{E}^k &= \arg \min_{\mathcal{P}_{\bar{\Omega}}(\mathbf{E})=0} L(\mathbf{X}^k, \mathbf{E}, \mathbf{Y}^{k-1}, \mu_{k-1}) \\ &= \arg \min_{\mathcal{P}_{\bar{\Omega}}(\mathbf{E})=0} \|\mathbf{E} - (\mathbf{M} - \mathbf{X}^k + \mathbf{Y}^{k-1} / \mu_{k-1})\|_F^2 \\ &= \mathcal{P}_{\bar{\Omega}}(\mathbf{M} - \mathbf{X}^k + \mathbf{Y}^{k-1} / \mu_{k-1}) \end{aligned} \quad (29)$$

$$\mathbf{Y}^k = \mathbf{Y}^{k-1} + \mu_{k-1}(\mathbf{M} - \mathbf{X}^k - \mathbf{E}^k) \quad (30)$$

$$\mu_k = \begin{cases} \rho \mu_{k-1}, & \text{若 } \mu_{k-1} \|\mathbf{E}^k - \mathbf{E}^{k-1}\|_F / \|\mathbf{M}\|_F < \epsilon \\ \mu_{k-1}, & \text{否则} \end{cases} \quad (31)$$

其中, $\bar{\Omega}$ 表示 Ω 的补集。在更新 \mathbf{X} 和 \mathbf{E} 的过程中, 若将式(28)和式(29)循环迭代, 直到满足终止条件再执行式(30)和式(31), 称这种方法为精确增广拉格朗日乘子法(Exact Augmented Lagrange Multipliers, EALM)。EALM 方法较好地求解了 $\min_{\mathbf{X}, \mathbf{E}} L(\mathbf{X}, \mathbf{E}, \mathbf{Y}^{k-1}, \mu_{k-1})$ 的近似解, 但对于大的 μ_{k-1} , 其收敛速度往往非常慢。实验结果证实 IALM 比 EALM 运行速度快, 也称 IALM 为交替方向方法(Alternating Direction Method, ADM)^[25,26]。

4 低秩矩阵分解模型的矩阵补全算法

4.1 问题描述

对于核范数最小化模型, 用矩阵核范数来度量秩的大小。在实际问题中, 若事先知道秩的取值, 则可以用低秩矩阵分解模型来恢复丢失元素。这种方法还可以避免矩阵的奇异值分解。对于给定的含丢失元素的低秩或近似低秩矩阵 \mathbf{M} , 对它作如下分解: $\mathbf{M} \approx \mathbf{X}\mathbf{Y}$, 其中 $\mathbf{X} \in \mathbb{R}^{m \times r}$, $\mathbf{Y} \in \mathbb{R}^{r \times n}$, r 为 \mathbf{M} 的事先规定或根据某种方法确定的秩。

为了恢复丢失元素, 建立下列最小化模型

$$\min_{\mathbf{X}, \mathbf{Y}} \|\mathcal{P}_{\bar{\Omega}}(\mathbf{M} - \mathbf{X}\mathbf{Y})\|_F^2, \text{ s. t. } \mathcal{P}_{\bar{\Omega}}(\mathbf{M}) = \mathcal{P}_{\bar{\Omega}}(\mathbf{X}\mathbf{Y}) \quad (32)$$

下面给出近年来求解上述低秩矩阵分解模型的几种最新的方法。

4.2 最优空间算法

维数为 r 的 \mathbb{R}^n 的所有子空间的集合为 Grassmann 流形。Grassmann 流形也是紧的 Riemann 流形, 其测地线可以计算。Grassmann 矩阵流形中的一个点为 $m \times r$ 维正交矩阵 \mathbf{A} 的等价类: $[\mathbf{A}] = \{\mathbf{A}\mathbf{Q}; \mathbf{Q} \text{ 为任意的 } r \text{ 阶正交方阵}\}$ 。基于 Grassmann 流形, 文献[27]提出了求解最优化问题(32)的最优空间算法(OptSpace), 此算法由以下 4 步组成。

第 1 步 剪切。记 $\mathbf{M}^{\Omega} = \mathcal{P}_{\bar{\Omega}}(\mathbf{M})$, 若 \mathbf{M}^{Ω} 的某个列向量的非零元素的个数大于 $2p/n$, 则在此列向量中随机保留 $2p/n$ 个非零元素, 其余非零元素值变为 0; 若某行向量的非零元素的个数大于 $2p/m$, 则在此列向量中随机保留 $2p/m$ 个非零元素, 其余非零元素值改为 0。剪切后的矩阵记为 $\tilde{\mathbf{M}}^{\Omega}$ 。

第 2 步 估计 \mathbf{M} 的秩。对 $\tilde{\mathbf{M}}^{\Omega}$ 进行奇异值分解: $\tilde{\mathbf{M}}^{\Omega} = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$, 其中 k 为 $\tilde{\mathbf{M}}^{\Omega}$ 的秩, 奇异值满足 $\sigma_1 \geq \dots \geq \sigma_k > 0$ 。根据 $\tilde{\mathbf{M}}^{\Omega}$ 的奇异值估计 \mathbf{M} 的秩:

$$r = \arg \min_i \{(\sigma_{i+1} + \sigma_i \sqrt{mn} \sqrt{i/p}) / \sigma_i; i=1, 2, \dots, k-1\}$$

第 3 步 秩 r 投影。 $\tilde{\mathbf{M}}^{\Omega}$ 的秩 r 投影为 $\mathcal{P}_{\Omega}(\tilde{\mathbf{M}}^{\Omega}) = \mathbf{X}^0 \mathbf{S}^0 (\mathbf{Y}^0)^T$, 其中 $\mathbf{X}^0 = \sqrt{m}(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r)$, $\mathbf{Y}^0 = \sqrt{n}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r)$, $\mathbf{S}^0 = (\sqrt{mn}/p) \text{diag}(\sigma_1, \dots, \sigma_r)$ 。

第 4 步 Grassmann 流形上的梯度下降法。定义差异性函数: $F(\mathbf{X}, \mathbf{Y}) = \min_{\mathbf{S} \in \mathbb{R}^{r \times r}} \mathcal{F}(\mathbf{X}, \mathbf{Y}, \mathbf{S})$, 其中 $\mathcal{F}(\mathbf{X}, \mathbf{Y}, \mathbf{S}) = \frac{1}{2} \|\mathcal{P}_{\bar{\Omega}}(\mathbf{M}^{\Omega} - \mathbf{X}\mathbf{S}\mathbf{Y}^T)\|_F^2 + \frac{\lambda}{2} \|\mathcal{P}_{\bar{\Omega}}(\mathbf{X}\mathbf{S}\mathbf{Y}^T)\|_F^2$, \mathbf{X}, \mathbf{Y} 满足 $\mathbf{X}^T \mathbf{X} = m\mathbf{I}_r$, $\mathbf{Y}^T \mathbf{Y} = n\mathbf{I}_r$ 。为简便起见, 考虑 $\lambda=0$ 。使用交替和梯度下降法求解优化问题 $\min_{\mathbf{X}, \mathbf{Y}, \mathbf{S}} \mathcal{F}(\mathbf{X}, \mathbf{Y}, \mathbf{S})$, 迭代过程如下

$$\begin{cases} \mathbf{S}^k = \arg \min_{\mathbf{S}} \mathcal{F}(\mathbf{X}^{k-1}, \mathbf{Y}^{k-1}, \mathbf{S}) \\ t_k = \arg \min_{t \geq 0} F(\mathbf{X}^{k-1} - t \partial_{\mathbf{X}} F(\mathbf{X}^{k-1}, \mathbf{Y}^{k-1}), \mathbf{Y}^{k-1} - t \partial_{\mathbf{Y}} F(\mathbf{X}^{k-1}, \mathbf{Y}^{k-1})) \\ \mathbf{X}^k = \sqrt{m} \text{orth}(\mathbf{X}^{k-1} - t_k \partial_{\mathbf{X}} F(\mathbf{X}^{k-1}, \mathbf{Y}^{k-1})) \\ \mathbf{Y}^k = \sqrt{n} \text{orth}(\mathbf{Y}^{k-1} - t_k \partial_{\mathbf{Y}} F(\mathbf{X}^{k-1}, \mathbf{Y}^{k-1})) \end{cases} \quad (33)$$

其中, $\partial_{\mathbf{X}} F(\mathbf{X}, \mathbf{Y}) = \mathcal{P}_{\bar{\Omega}}(\mathbf{X}\mathbf{S}\mathbf{Y}^T - \mathbf{M}^{\Omega})\mathbf{Y}\mathbf{S}^T$, $\partial_{\mathbf{Y}} F(\mathbf{X}, \mathbf{Y}) = \mathcal{P}_{\bar{\Omega}}(\mathbf{X}\mathbf{S}\mathbf{Y}^T - \mathbf{M}^{\Omega})\mathbf{X}\mathbf{S}$, $\text{orth}(\cdot)$ 表示矩阵列正交化算子。当迭代公式(33)达到终止条件时, \mathbf{M} 补全后的结果为 $\mathbf{X}^k \mathbf{S}^k (\mathbf{Y}^k)^T$ 。

4.3 Grassmann 秩 1 更新子空间估计算法

由于噪声的存在, 最优化模型(32)中的约束 $\mathcal{P}_{\bar{\Omega}}(\mathbf{M}) = \mathcal{P}_{\bar{\Omega}}(\mathbf{X}\mathbf{Y})$ 不一定成立。为此, 考虑下列优化模型

$$\min_{\mathbf{X}, \mathbf{Y}} \|\mathcal{P}_{\bar{\Omega}}(\mathbf{M} - \mathbf{X}\mathbf{Y})\|_F^2 \quad (34)$$

将矩阵 \mathbf{M} 和 \mathbf{Y} 按列向量表示: $\mathbf{M} = (\mathbf{M}_{(1)}, \mathbf{M}_{(2)}, \dots, \mathbf{M}_{(n)})$, $\mathbf{Y} = (\mathbf{Y}_{(1)}, \mathbf{Y}_{(2)}, \dots, \mathbf{Y}_{(n)})$ 。记 $\mathcal{H}_{\bar{\Omega}}$ 为向量、矩阵选择算子, 即 $\mathcal{H}_{\bar{\Omega}}(\mathbf{M}_{(i)})$ 按照 $\bar{\Omega}$ 来选取 $\mathbf{M}_{(i)}$ 的非丢失元素, $\mathcal{H}_{\bar{\Omega}}(\mathbf{X})$ 按照 $\mathbf{M}_{(i)}$ 的非丢失元素来选取 \mathbf{X} 的相应的行向量, $i=1, 2, \dots, n$ 。于是最优化问题(34)等价于

$$\min_{\mathbf{X}, \mathbf{Y}} \sum_{i=1}^n \|\mathcal{H}_{\bar{\Omega}}(\mathbf{M}_{(i)}) - \mathcal{H}_{\bar{\Omega}}(\mathbf{X})\mathbf{Y}_{(i)}\|^2 \quad (35)$$

其中, $\|\cdot\|$ 表示向量的 l_2 范数。文献[13]提出了求解上述最优化问题的 Grassmann 秩 1 更新子空间估计算法(Grassmannian Rank-One Update Subspace Estimation, GROUSE)。在 GROUSE 中, 更新矩阵 \mathbf{X} 的算法如下。

For $i=1:n$

$$\mathbf{Y}_{(i)}^* = \arg \min_{\mathbf{Y}_{(i)}} \|\mathcal{H}_{\bar{\Omega}}(\mathbf{M}_{(i)}) - \mathcal{H}_{\bar{\Omega}}(\mathbf{X})\mathbf{Y}_{(i)}\|^2 = (\mathcal{H}_{\bar{\Omega}}(\mathbf{X}))^{\dagger} \mathcal{H}_{\bar{\Omega}}(\mathbf{M}_{(i)});$$

预测 $\mathbf{M}_{(i)}$ 的补全向量 $\mathbf{p}_i = \mathbf{X}\mathbf{Y}_{(i)}^*$;

计算残差向量 $\mathbf{r}_i = \mathcal{H}_{\bar{\Omega}}(\mathbf{M}_{(i)}) - \mathcal{H}_{\bar{\Omega}}(\mathbf{X})\mathbf{Y}_{(i)}^*$;

根据 \mathbf{r}_i 和 $\mathcal{H}_{\bar{\Omega}}(\cdot)$ 得到 m 维向量 \mathbf{r}_i , 即 \mathbf{r}_i 对应的丢失元素取值为 0;

更新矩阵 \mathbf{X} : $\mathbf{X} \leftarrow \mathbf{X} + [(\cos(\|\mathbf{r}_i\| / \|\mathbf{p}_i\| \eta) - 1)\mathbf{p}_i / \|\mathbf{p}_i\| + \sin$

$(\|\mathbf{r}_i\| / \|\mathbf{p}_i\| \eta)\mathbf{r}_i / \|\mathbf{r}_i\|](\mathbf{Y}_{(i)}^*)^T / \|\mathbf{Y}_{(i)}^*\|$ 。

End

在上述算法中, “ \dagger ” 表示矩阵的伪逆, η 为步长大小。为了得到更优的解, 可将上述算法重复多遍。最终矩阵 \mathbf{Y} 的求解公式为

$$\mathbf{Y}^* = ((\mathcal{H}_{\bar{\Omega}}(\mathbf{X}))^{\dagger} \mathcal{H}_{\bar{\Omega}}(\mathbf{M}_{(1)}), \dots, (\mathcal{H}_{\bar{\Omega}}(\mathbf{X}))^{\dagger} \mathcal{H}_{\bar{\Omega}}(\mathbf{M}_{(n)})) \quad (36)$$

矩阵补全的其它的 Riemann 或 Grassmann 流形子空间方法可参见文献[28,29]。

4.4 奇异值投影算法

记 $\mathbf{Z} = \mathbf{X}\mathbf{Y}$, 则最小化问题(32)的鲁棒形式为

$$\min_{\mathbf{Z}} \psi(\mathbf{Z}) = \frac{1}{2} \|\mathcal{P}_{\bar{\Omega}}(\mathbf{Z}) - \mathcal{P}_{\bar{\Omega}}(\mathbf{M})\|_F^2 \quad (37)$$

s. t. $\mathbf{Z} \in \mathcal{C}(r) = \{\mathbf{U} \in \mathbb{R}^{m \times n}; \text{rank}(\mathbf{U}) \leq r\}$

由于集合 $\mathcal{C}(r)$ 是非凸、非连续的, 因此直接求解上述最优化问题是非常困难的。

利用启发式方法,文献[30]提出了求解最小化问题(37)的奇异值投影算法(Singular Value Projection, SVP)。在SVP中,先不考虑低秩约束,直接使用梯度下降法更新矩阵 Z 。 Z 的更新公式为

$$\hat{Z}^k = Z^{k-1} - \eta_k \partial \psi(Z^{k-1}) = Z^{k-1} - \eta_k (\mathcal{P}_\Omega(Z^{k-1}) - \mathcal{P}_\Omega(M)) \quad (38)$$

步长大小 η_k 可按如下方式选取: $\eta_k = mn/(p + \delta p)$,其中 $0 < \delta < 1/3$ 且取值依赖于 p 。再将 \hat{Z}^k 投影到 $\mathcal{C}(r)$ 上,即

$$Z^k = \arg \min_Z \{ \|Z - \hat{Z}^k\|_F; Z \in \mathcal{C}(r) \} \quad (39)$$

对 Z^k 进行奇异值分解,由其前 r 个最大的奇异值及对应的奇异向量可构成矩阵 Z^k 。

4.5 低秩矩阵拟合算法

通过引入辅助变量矩阵 Z ,最小化问题(32)等价于

$$\min_{X,Y,Z} \frac{1}{2} \|XY - Z\|_F^2, \text{ s. t. } \mathcal{P}_\Omega(Z) = \mathcal{P}_\Omega(M) \quad (40)$$

此优化问题的拉格朗日函数为

$$L(X,Y,Z,\Lambda) = \frac{1}{2} \|XY - Z\|_F^2 - \langle \Lambda, \mathcal{P}_\Omega(Z - M) \rangle \quad (41)$$

对 $L(X,Y,Z,\Lambda)$ 关于各个块变量求梯度,得到一阶最优性条件

$$\begin{cases} (XY - Z)Y^T = 0 \\ X^T(XY - Z) = 0 \\ \mathcal{P}_\Omega(Z - XY) = 0 \\ \mathcal{P}_\Omega(Z - M) = 0 \\ \mathcal{P}_\Omega(Z - XY) = \Lambda \end{cases} \quad (42)$$

使用非线性逐次超松弛迭代法求解上述方程系统,交替迭代公式如下

$$\begin{cases} X_+ \leftarrow ZY^T(YY^T)^{-1} \\ X_+(\omega) \leftarrow \omega X_+ + (1-\omega)X \\ Y_+ \leftarrow (X_+(\omega)^T X_+(\omega))^{-1} X_+(\omega)^T Z \\ Y_+(\omega) \leftarrow \omega Y_+ + (1-\omega)Y \\ Z_+(\omega) \leftarrow X_+(\omega)Y_+(\omega) + \mathcal{P}_\Omega(M - X_+(\omega)Y_+(\omega)) \end{cases} \quad (43)$$

其中,参数 $\omega \geq 1$ 。当 $\omega = 1$ 时,上述迭代方法即为 Gauss-Seidel 方法。迭代公式(43)称为求解矩阵补全的低秩矩阵拟合算法(Low-rank matrix Fitting, LmaFit)^[31]。

5 其它优化模型

由核范数最小化模型和低秩矩阵分解模型衍生出许多求解矩阵补全的其它模型,本节对部分模型及求解方法作简要介绍。

5.1 剪切核范数正则化模型

在核范数最小化模型中,同时最小化所有奇异值,这在实际问题中可能不能更好地逼近矩阵的秩。文献[32]提出了剪切核范数正则化模型(Truncated Nuclear Norm Regularization, TNNR),它在一定程度上避免了核范数优化模型的上述缺陷。

定义5^[32] 矩阵 $X \in \mathbb{R}^{m \times n}$ 的奇异值按大小排列为 $\sigma_1, \sigma_2, \dots, \sigma_{\min(m,n)}$,则它的 r -剪切核范数为

$$\|X\|_{-r} = \sum_{i=r+1}^{\min(m,n)} \sigma_i \quad (44)$$

在核范数最小化模型(8)中,将矩阵的核范数替换成矩阵的剪切核范数,便得如下的 TNNR 模型

$$\min_X \|X\|_{-r}, \text{ s. t. } \mathcal{P}_\Omega(X) = \mathcal{P}_\Omega(M) \quad (45)$$

上述最优化问题也等价于

$$\min_X (\|X\|_* - \max_{AA^T=BB^T=I} \text{trace}(AXB^T)), \text{ s. t. } \mathcal{P}_\Omega(X) = \mathcal{P}_\Omega(M) \quad (46)$$

其中, $A \in \mathbb{R}^{r \times m}, B \in \mathbb{R}^{r \times n}$ 。初始化 $X^0 = \mathcal{P}_\Omega(M)$,上述最优化问题可转化为交替求解两个子优化问题

$$\max_{A,B} \text{trace}(AX^{k-1}B^T), \text{ s. t. } AA^T = I_r, BB^T = I_r \quad (47)$$

$$\min_X \|X\|_* - \langle X, B_k^T A_k \rangle, \text{ s. t. } \mathcal{P}_\Omega(X) = \mathcal{P}_\Omega(M) \quad (48)$$

最优化问题(47)的最优解记为 $A^k = (u_1, u_2, \dots, u_r)^T$, $B^k = (v_1, v_2, \dots, v_r)^T$,其中单位向量 u_i, v_i 分别为 X^{k-1} 的第 i 大奇异值对应的左、右奇异向量。对于最小化问题(48),文献[32]提出了两种求解方法,即 ADM 和基于线性搜索的 APG。

5.2 Schatten p -范数优化模型

在核范数优化模型(8)中,将矩阵的 Schatten 1-范数推广到 Schatten p -范数($0 < p \leq 2$),则有最小化模型

$$\min_X \|X\|_{sp}, \text{ s. t. } \mathcal{P}_\Omega(X) = \mathcal{P}_\Omega(M) \quad (49)$$

上述最优化问题的目标函数也可用迹函数来表示,即 $\|X\|_{sp} = (\text{trace}((X^T X)^{p/2}))^{1/p}$ 。文献[33]提出了求解上述最小化模型的拉格朗日乘子法,并证明了所提算法的收敛性。

5.3 最小秩逼近原子分解模型

令 Γ 为 $\mathbb{C}^{m \times n}$ 中所有秩为1、Frobenious 范数也为1的矩阵构成的集合。对 Γ 的元素进行精炼,即任意两个相异的矩阵不共线,精炼后的集合称为原子集,用 O 表示。 $\forall Z \in \mathbb{C}^{m \times n}$,它的原子分解形式为 $Z = \sum_i \alpha_i \phi_i, \phi_i \in O$ 。易知 $\text{rank}(Z) = \min_{\Phi} \{ |\Phi|; \Phi \subset O, Z \in \text{span}(\Phi) \}$ 。

秩最小化问题(6)对应的鲁棒优化问题为

$$\min_X \| \mathcal{A}(X) - b \|, \text{ s. t. } \text{rank}(X) \leq r \quad (50)$$

文献[34]提出了使用最小秩逼近原子分解模型(Atomic Decomposition for Minimum Rank Approximation, ADMiRA)求解上述最优化问题。对于初始化的 X^0 和 Φ^0 ,ADMiRA 方法交替更新 Φ 和 X :

$$\begin{cases} \tilde{\Phi}^k = \arg \max_{\Phi \subset O} \{ \| \mathcal{P}_\Phi \mathcal{A}^* (\mathcal{A}(X^{k-1}) - b) \|_F; |\Phi| \leq 2r \} \cup \Phi^{k-1} \\ \tilde{X}^k = \arg \min_X \{ \| \mathcal{A}(X) - b \|; X \in \text{span}(\tilde{\Phi}^k) \} \\ \Phi^k = \arg \min_{\Phi \subset O} \{ \| \mathcal{P}_\Phi \tilde{X}^k \|_F; |\Phi| \leq r \} \\ X^k = \mathcal{P}_{\Phi^k} \tilde{X}^k \end{cases} \quad (51)$$

其中, \mathcal{A}^* 为 \mathcal{A} 的伴随算子, \mathcal{P}_Φ 表示 Φ 张成子空间上的正交投影。

5.4 最大范数约束的优化模型

对于秩为 r 的矩阵 $Z \in \mathbb{R}^{m \times n}$,补充它的几种范数: $\|Z\|_\infty = \max_{i,j} |Z_{ij}|, \|Z\|_{2,\infty} = \max_i \|Z_{(i)}\|, \|Z\|_{\max} = \min \{ \|U\|_{2,\infty}, \|V\|_{2,\infty}; Z = UV^T, U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r} \}$,其中 $Z_{(i)}$ 为 Z 的第 i 列构成的向量。对于给定的 $R \geq \alpha > 0$,文献[35]提出的最大范数约束优化模型为

$$\min_Z \| \mathcal{P}_\Omega(M - Z) \|_F^2, \text{ s. t. } \|Z\|_\infty \leq \alpha, \|Z\|_{\max} \leq R \quad (52)$$

上述优化问题可以转化为半定规划模型

$$\begin{aligned} & \min_{W_1, W_2} \| \mathcal{P}_\Omega(M - Z) \|_F^2 \\ & \text{ s. t. } \begin{pmatrix} W_1 & Z \\ Z^T & W_2 \end{pmatrix} \geq 0 \\ & \text{diag}(W_1) \leq R, \text{diag}(W_2) \leq R, \|Z\|_\infty \leq \alpha \end{aligned} \quad (53)$$

或分解形式模型

$$\begin{aligned} \min_X & \| \mathcal{P}_0(\mathbf{M}-\mathbf{XY}) \|_F^2 \\ \text{s. t. } & \max\{ \| \mathbf{X} \|_{2,\infty}, \| \mathbf{Y} \|_{2,\infty} \} \leq R \\ & \max_{i,j} |(XY)_{ij}| \leq \alpha \end{aligned} \quad (54)$$

可使用一阶梯度下降法求解最小化问题(54)。

5.5 低秩逼近模型

若矩阵 \mathbf{X} 的秩为 r , 则它可以分解为: $\mathbf{X}=\mathbf{LR}^T$, 其中 $\mathbf{L} \in \mathbb{R}^{m \times r}, \mathbf{R} \in \mathbb{R}^{n \times r}$. 此时, \mathbf{X} 的核范数满足下列性质^[10]

$$\| \mathbf{X} \|_* = \min_{\mathbf{X}=\mathbf{LR}^T} \| \mathbf{L} \|_F \| \mathbf{R} \|_F = \frac{1}{2} \min_{\mathbf{X}=\mathbf{LR}^T} (\| \mathbf{L} \|_F^2 + \| \mathbf{R} \|_F^2) \quad (55)$$

最优化问题(8)对应的低秩逼近模型为

$$\min_{\mathbf{L}, \mathbf{R}} \| \mathbf{L} \|_F^2 + \| \mathbf{R} \|_F^2, \text{ s. t. } \mathcal{A}(\mathbf{LR}^T) = \mathbf{b} \quad (56)$$

如果 r 比式(8)的最优解矩阵的秩充分大, 则最优化问题(8)与式(56)等价。文献[15]提出了求解式(56)的增广拉格朗日乘法, 文献[36]对式(56)进行正则化, 并使用梯度下降法求解。

5.6 非负矩阵分解模型

非负矩阵分解是处理非负低秩矩阵的一种流行的方法, 它将非负数据矩阵 \mathbf{M} 分解为非负低秩矩阵 \mathbf{X} 与 \mathbf{Y} 之积。基于非负矩阵分解的矩阵补全描述如下

$$\begin{aligned} \min_{\mathbf{X}, \mathbf{Y}} & \| \mathcal{P}_0(\mathbf{M}-\mathbf{XY}) \|_F^2 \\ \text{s. t. } & X_{ik} \geq 0, Y_{kj} \geq 0, i=1, \dots, m, \\ & k=1, \dots, r, j=1, \dots, n \end{aligned} \quad (57)$$

其中, $\mathbf{X} \in \mathbb{R}^{m \times r}, \mathbf{Y} \in \mathbb{R}^{r \times n}, r \leq \min(m, n)$. 对于上述最优化问题, 文献[37]先将它转化为交替求解两个非负最小二乘问题, 再对每个非负最小二乘问题使用内点梯度法求解; 文献[38]提出使用 ADM 方法求解。

此外, 矩阵补全的基于其它模型的算法还有: 鲁棒低秩矩阵分解框架的交替方向法^[39]、正则化算法^[40]、迭代重新加权算法^[41]和原-对偶算法^[42]等。

6 实验分析

先在人工数据集上比较矩阵补全的主要算法的性能, 再将部分算法应用到图像恢复中。

6.1 人工数据集

随机生成秩为 r 的 $m \times n$ 维矩阵 \mathbf{M} , 生成方式如下: $\mathbf{M}=\mathbf{M}_L \mathbf{M}_R$, 其中 $\mathbf{M}_L \in \mathbb{R}^{m \times r}, \mathbf{M}_R \in \mathbb{R}^{r \times n}, r \leq \min(m, n)$. 矩阵 \mathbf{M}_L 、 \mathbf{M}_R 的各元素之间相互独立且服从均值为 0、标准差为 1 的正态分布。按照均匀分布对 \mathbf{M} 进行随机采样, 采样指标集合为 Ω , 采样数目记为 p . 根据核范数优化模型和低秩矩阵分解模型来补全 \mathbf{M} 的未采样元素, 补全后的矩阵记为 $\tilde{\mathbf{M}}$.

对于各种矩阵补全算法(GROUSE 除外), 设终止条件为 $\| \mathcal{P}_0(\mathbf{M}^k) - \mathcal{P}_0(\mathbf{M}) \|_F / \| \mathcal{P}_0(\mathbf{M}) \|_F \leq \epsilon$, 其中 \mathbf{M}^k 为 \mathbf{M} 的第 k 次迭代时的补全结果, ϵ 为容许误差。使用相对误差 $\| \tilde{\mathbf{M}} - \mathbf{M} \|_F / \| \mathbf{M} \|_F$ 来度量 $\tilde{\mathbf{M}}$ 的补全性能。在每种算法中, 对于给定的秩 r 和采样数目 p , 将矩阵补全算法重复 10 次, 最后报告运行时间、迭代次数和相对误差的平均值。在实验中, 取 $m=n=1000, \epsilon=10^{-4}$, 则 \mathbf{M} 的自由度为 $dr=r(2000-r)$. 考虑 $r=10, 50, 100$ 这 3 种情形, 实验结果如表 1 所列。

表 1 各种矩阵补全算法的性能比较

算法	秩	p/dr	p/n^2	运行时间(s)	迭代次数	相对误差
SVT	10	6	0.12	85.0	117	1.66e-4
	50	4	0.39	559.3	113.1	1.62e-4
	100	3	0.57	1359.1	128.7	1.70e-4
FPC	10	6	0.12	18.9	91.6	4.69e-4
	50	4	0.39	83.7	70.1	3.05e-5
	100	3	0.57	132.3	79	2.25e-5
APG	10	6	0.12	10.2	37.7	2.71e-4
	50	4	0.39	28.2	51	3.77e-2
	100	3	0.57	360.7	49	9.60e-3
IALM	10	6	0.12	8.1	45	1.40e-4
	50	4	0.39	30.6	25.6	1.44e-4
	100	3	0.57	89.9	28	1.53e-4
OptSpace	10	6	0.12	66.2	24.6	4.94e-5
	50	4	0.39	6804.5	34	2.00e-5
GROUSE	10	6	0.12	4.2	—	2.05e-6
	50	4	0.39	39.7	—	1.57e-8
	100	3	0.57	137.1	—	4.56e-5
SVP	10	6	0.12	6.4	26.1	1.56e-4
	50	4	0.39	28.4	29	7.07e-5
	100	3	0.57	75.8	33	6.25e-5
LmaFit	10	6	0.12	1.8	28	1.54e-4
	50	4	0.39	4.6	21	1.43e-4
	100	3	0.57	8.6	21	1.58e-4

对于低秩矩阵分解模型的 OptSpace 方法, 当 $r=50$, $p/dr=4$ 时, 由于运行时间非常长, 因此只报告了一次实验的结果; 当 $r=100, p/dr=3$ 时, 运行时间更长, 没有报告实验结果。从表 1 中可以看出: 对于核范数优化模型的 4 种算法, APG 的相对误差最大, 其它 3 种方法都达到了 10^{-4} 或更低的数量级, 而 IALM 运行时间最短, 且迭代次数较少。对于低秩矩阵分解模型, 4 种算法都达到了较好的相对误差, 而 LmaFit 的运行时间最短, 迭代次数也较低。IALM 与 LmaFit 相比, 后者的运算速度更快, 但需要事先确定矩阵的秩, 而如何确定不完全矩阵的秩是一个公开问题。

6.2 图像恢复

通过 6.1 节的分析可知, IALM 和 LmaFit 分别代表两类模型的最有效的算法。本节将这两种算法应用到图像恢复中。“熊猫”灰度图像的分辨率为 194×259 , 采样数目记为 p , 则采样概率为 $sp=p/(194 \times 259)$. 考虑 4 种采样概率, 即 $sp=0.1, 0.3, 0.5, 0.7$, 分别使用 IALM 和 LamFit 来补全丢失元素。用峰值信噪比 (psnr) 来度量算法的恢复性能, 其中在 LamFit 算法中, 考虑 psnr 最大时秩 r 的取值情形。实验结果如图 1 所示。

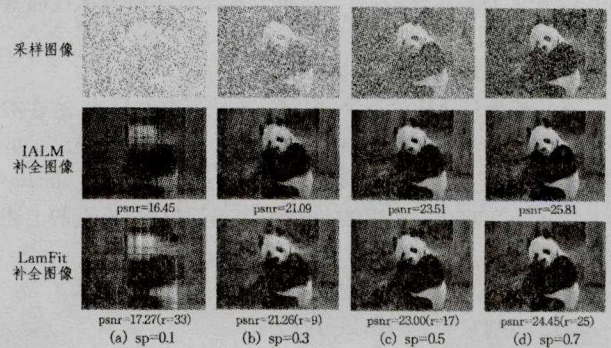


图 1 不同采样概率下 IALM 和 LamFit 的恢复性能比较

在图 1 的采样图像中, 未采样元素或丢失元素用白像素值(即 255)来表示。从图 1 可以看出: 当 $sp=0.1$ 和 0.3 时, LamFit 的恢复性能优于 IALM; 当 $sp=0.5$ 和 0.7 时, IALM

的恢复性能优于 LamFit。对于 $sp=0.1$, 两种方法均未能较好地恢复熊猫的轮廓, 这可能是由于图像的秩取值较大且未充分采样所致。此外, 在 LamFit 方法中需要事先给定矩阵的秩, 而秩的大小直接影响着恢复性能。

结束语 本文报告了矩阵补全算法的最新研究进展。将求解矩阵补全的主要算法按模型归纳为核范数优化模型和低秩矩阵分解模型两类。对于这两类模型, 总结了各种主流的算法, 并通过数值实验比较了它们的性能。本文也补充了这两类模型的衍生版本及相应的求解方法。

对于大规模问题, 设计可扩展的快速算法是当前矩阵补全算法的研究核心, 这涉及到部分奇异值分解的快速算法和稀疏矩阵的存储与运算。此外, 在今后的矩阵补全算法研究中, 下面几个方向值得关注: (a) 建立适合不同问题背景的矩阵补全优化模型; (b) 将矩阵补全推广到高阶张量上, 即研究张量补全和非负张量补全算法^[43-45]; (c) Grassmann 流形子空间方法求解矩阵补全问题; (d) 如何根据矩阵的部分元素来确定它的秩^[46]。

参 考 文 献

- [1] Donoho D. Compressed sensing [J]. *IEEE Transactions on Information Theory*, 2006, 52(4): 1289-1306
- [2] Candès E J, Michael W. An introduction to compressive sampling [J]. *IEEE Signal Processing Magazine*, 2008, 25(2): 21-30
- [3] Candès E J, Tao T. The power of convex relaxation: Near-optimal matrix completion [J]. *IEEE Transactions on Information Theory*, 2010, 56(5): 2053-2080
- [4] Candès E J, Plan Y. Matrix completion with noise [J]. *Proceedings of the IEEE*, 2010, 98(6): 925-936
- [5] 史加荣, 郑秀云, 魏宗田, 等. 低秩矩阵恢复算法综述[J]. *计算机应用研究*, 2013, 30(6): 1601-1605
- [6] Chen P. Optimization algorithms on subspaces: revisiting missing data problem in low-rank matrix [J]. *International Journal of Computer Vision*, 2008, 80(1): 125-142
- [7] Vidal R, Tron R, Hartley R. Multiframe motion segmentation with missing data using PowerFactorization and GPCA [J]. *International Journal of Computer Vision*, 2008, 79(1): 85-105
- [8] Ji Hui, Liu Chao-qiang, Shen Zuo-wei, et al. Robust video denoising using low rank matrix completion [C]//*Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2010: 1791-1798
- [9] Jing Guo-dong, Shi Yun-hui, Yin Bao-cai. Image super-resolution reconstruction based on sparse representation and low-rank matrix completion [J]. *Journal of Information & Computational Science*, 2012, 9(13): 3859-3866
- [10] Rennie D M, Srebro N. Fast maximum margin matrix factorization for collaborative prediction [C]//*Proceedings of International Conference on Machine Learning (ICML)*. 2005: 713-719
- [11] Cabral R S, Torre F D, Costeira J P, et al. Matrix completion for multi-label image classification [C]//*Proceedings of Neural Information Processing Systems (NIPS)*. 2011: 190-198
- [12] Keshavan R H, Oh S. OptSpace: A gradient descent algorithm on the Grassman manifold for matrix completion [EB/OL]. <http://arxiv.org/pdf/0910.5260v2.pdf>, 2009
- [13] Balzano L, Nowak R, Recht B. Online identification and tracking of subspaces from highly incomplete information [C]//*Proceedings of Annual Allerton Conference on Communication, Control, and Computing*. 2010: 704-711
- [14] Liu Zhang, Hansson A, Vandenberghe L. Nuclear norm system identification with missing inputs and outputs [EB/OL]. <http://www.ee.ucla.edu/~vandenbe/publications/subspace.pdf>, 2012
- [15] Recht B, Fazel M, Parrilo P A. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization [J]. *SIAM Review*, 2010, 52(3): 471-501
- [16] Fazel M. Matrix rank minimization with applications [D]. Stanford University, 2002
- [17] Liu Zhang, Vandenberghe L. Interior-point method for nuclear norm approximation with application to system identification [J]. *SIAM Journal on Matrix Analysis and Applications*, 2009, 31(3): 1235-1256
- [18] Cai Jian-feng, Candès E J, Shen Zuo-wei. A singular value thresholding algorithm for matrix completion [J]. *SIAM Journal on Optimization*, 2010, 20(4): 1956-1982
- [19] Hu Yao, Zhang De-bing, Liu Jun, et al. Accelerated singular value thresholding for matrix completion [C]//*Proceedings of ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2012: 298-306
- [20] Cai Jian-feng, Osher S. Fast singular value thresholding without singular value decomposition [J]. *Methods and Applications of Analysis*, 2013, 20(4): 335-352
- [21] Donald Goldfarb D, Ma Shi-qian. Convergence of fixed-point continuation algorithms for matrix rank minimization [J]. *Foundations of Computational Mathematics*, 2011, 11(2): 183-210
- [22] Ma Shi-qian, Goldfarb D, Chen Li-feng. Fixed point and Bregman iterative methods for matrix rank minimization [J]. *Mathematical Programming*, 2011, 128(1/2): 321-353
- [23] Toh K C, Yun S. An accelerated proximal gradient algorithm for nuclear norm regularized linear least squares problems [J]. *Pacific Journal of Optimization*, 2010, 6: 615-640
- [24] Lin Zhou-chen, Chen Min-ming, Ma Yi, et al. The augmented Lagrange multiplier method for exact recovery of corrupted low-rank matrices [EB/OL]. <http://arxiv.org/pdf/1009.5055>, 2010
- [25] Chen Cai-hua, He Bing-sheng, Yuan Xiao-ming. Matrix completion via an alternating direction method [J]. *IMA Journal of Numerical Analysis*, 2012, 32(1): 227-245
- [26] Yang Jun-feng, Yuan Xiao-ming. Linearized augmented Lagrangian and alternating direction methods for nuclear norm minimization [J]. *Mathematics of Computation*, 2013, 82(281): 301-329
- [27] Keshavan R H, Montanari A, Sewoong Oh. Matrix completion from a few entries [J]. *IEEE Transactions on Information Theory*, 2010, 56(6): 2980-2998
- [28] Dai Wei, Kerman E, Milenkovic O. A geometric approach to low-rank matrix completion [J]. *IEEE Transactions on Information Theory*, 2012, 58(1): 237-247
- [29] Boutilier N, Absil P A. RTRMC: A Riemannian trust-region method for low-rank matrix completion [C]//*Proceedings of Neural Information Processing Systems (NIPS)*. 2011
- [30] Jain P, Meka R, Dhillon I. Guaranteed rank minimization via singular value projection [EB/OL]. <http://arxiv.org/abs/0909.5457>, 2009
- [31] Wen Zai-wen, Yin Wo-tao, Zhang Yin. Solving a low-rank factorization model for matrix completion by a nonlinear successive

- over-relaxation algorithm[J]. *Mathematical Programming Computation*, 2012, 4(4): 333-361
- [32] Zhang De-bing, Hu Yao, Ye Jie-ping, et al. Matrix completion by truncated nuclear norm regularization[C]//*Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2012; 2192-2199
- [33] Nie Fei-ping, Huang H, Ding C. Low-rank matrix recovery via efficient Schatten p-norm minimization [C] // *Proceedings of AAAI Conference on Artificial Intelligence*. 2012; 655-661
- [34] Lee K, Bresler Y. ADMiRA: Atomic decomposition for minimum rank approximation[J]. *IEEE Transactions on Information Theory*, 2010, 56(9): 4402-4416
- [35] Cai T T, Zhou Wen-xin. Matrix completion via max-norm constrained optimization [EB/OL]. <http://arxiv.org/pdf/1303.0341v1.pdf>, 2013
- [36] Recht B, Christopher Re C. Parallel stochastic gradient algorithms for large-scale matrix completion [EB/OL]. <http://pages.cs.wisc.edu/~brecht/papers/11.Rec.Re.IPGM.pdf>, 2013
- [37] 史加荣, 焦李成, 尚凡华. 不完全非负矩阵分解的加速算法[J]. *电子学报*, 2011, 39(2): 291-295
- [38] Xu Yang-yang, Yin Wo-tao, Wen Zai-wen, et al. An alternating direction algorithm for matrix completion with nonnegative factors [J]. *Frontiers of Mathematics in China*, 2012, 7(2): 365-384
- [39] 尚凡华. 基于低秩结构学习数据表示[D]. 西安: 西安电子科技大学, 2012
- [40] Mazumder R, Hastie T, Tibshirani R. Spectral regularization algorithms for learning large incomplete matrices[J]. *Journal of Machine Learning Research*, 2010, 11: 2287-2322
- [41] Mohan K, Fazel M. Iterative reweighted algorithms for matrix rank minimization [J]. *Journal of Machine Learning Research*, 2012, 13: 3441-3473
- [42] Xin Yu, Tommi Jaakkola T. Primal-Dual methods for sparse constrained matrix completion[J]. *Journal of Machine Learning Research-Proceedings Track*, 2012, 22: 1323-1331
- [43] Liu Ji, Musialski P, Wonka P, et al. Tensor completion for estimating missing values in visual data[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2013, 35(1): 208-220
- [44] Signoretto M, Van de Plas, R, De Moor B, et al. Tensor versus matrix completion; a comparison with application to spectral data[J]. *IEEE Signal Processing Letters*, 2011, 18(7): 403-406
- [45] 史加荣, 焦李成, 尚凡华. 张量补全算法及其在人脸识别中的应用[J]. *模式识别与人工智能*, 2011, 24(2): 255-261
- [46] Julià C, Sappa A D, Lumbleras F, et al. Rank estimation in missing data matrix problems[J]. *Journal of Mathematical Imaging and Vision*, 2011, 39: 140-160

(上接第 12 页)

- [25] Bilenko M, Mooney R, Cohen W, et al. Adaptive name matching in information integration[J]. *IEEE Intelligent Systems*, 2003, 18(5): 16-23
- [26] Gravano L, Ipeirotis P G, Koudas N, et al. Text joins in an RD-BMS for web data integration[C]//*Proceedings of the 12th international conference on World Wide Web*. ACM, 2003; 90-101
- [27] Gill L. OX-LINK: The Oxford Medical Record Linkage System [C]//*Proc. Int'l Record Linkage Workshop and Exposition*. 1997; 15-33
- [28] 刁兴春, 谭明超, 曹建军. 一种融合多种编辑距离的字符串相似度计算方法[J]. *计算机应用研究*, 2010, 27(12): 4523-4525
- [29] Wang J, Li G, Yu J X, et al. Entity matching: how similar is similar[J]. *Proceedings of the VLDB Endowment*, 2011, 4(10): 622-633
- [30] Herzog T N, Scheuren F J, Winkler W E. *Data quality and record linkage techniques*[M]. Springer, 2007
- [31] Verykios V S, Moustakides G V, Elfeky M G. A Bayesian decision model for cost optimal record matching [J]. *The VLDB Journal*, 2003, 12(1): 28-40
- [32] Naumann F, Herschel M. An introduction to duplicate detection [J]. *Synthesis Lectures on Data Management*, 2010, 2(1): 1-87
- [33] Fan W, Jia X, Li J, et al. Reasoning about record matching rules [J]. *Proceedings of the VLDB Endowment*, 2009, 2(1): 407-418
- [34] Cochinwala M, Kurien V, Lalk G, et al. Efficient data reconciliation[J]. *Information Sciences*, 2001, 137(1): 1-15
- [35] Christen P. Automatic record linkage using seeded nearest neighbour and support vector machine classification[C]//*Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2008; 151-159
- [36] Christen P. Automatic training example selection for scalable unsupervised record linkage[M]//*Advances in Knowledge Discovery and Data Mining*. Springer Berlin Heidelberg, 2008; 511-518
- [37] Arasu A, Götz M, Kaushik R. On active learning of record matching packages[C]//*Proceedings of the 2010 International Conference on Management of Data*. ACM, 2010; 783-794
- [38] Whang S E, Garcia-Molina H. Entity resolution with evolving rules[J]. *Proceedings of the VLDB Endowment*, 2010, 3(1/2): 1326-1337
- [39] Monge A E. Matching algorithms within a duplicate detection system[J]. *IEEE Data Engineering Bulletin*, 2000, 23(4): 14-20
- [40] Hassanzadeh O, Miller R J. Creating probabilistic databases from duplicated data[J]. *The VLDB Journal—The International Journal on Very Large Data Bases*, 2009, 18(5): 1141-1166
- [41] Hernández M A, Stolfo S J. The merge/purge problem for large databases[C]//*ACM SIGMOD Record*. ACM, 1995, 24(2): 127-138
- [42] Chaudhuri S, Ganti V, Motwani R. Robust identification of fuzzy duplicates[C]//*Proceedings of 21st International Conference on Data Engineering 2005*. IEEE, 2005; 865-876
- [43] Kalashnikov D V, Mehrotra S. Domain-independent data cleaning via analysis of entity-relationship graph[J]. *ACM Transactions on Database Systems (TODS)*, 2006, 31(2): 716-767
- [44] Dong X, Halevy A, Madhavan J. Reference reconciliation in complex information spaces [C] // *Proceedings of the 2005 ACM SIGMOD International Conference on Management of Data*. ACM, 2005; 85-96
- [45] Bhattacharya I, Getoor L. Collective entity resolution in relational data[J]. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 2007, 1(1): 5
- [46] Rastogi V, Dalvi N, Garofalakis M. Large-scale collective entity matching[J]. *Proceedings of the VLDB Endowment*, 2011, 4(4): 208-218
- [47] Fan W. Dependencies revisited for improving data quality[C]//*Proceedings of the twenty-seventh ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*. ACM, 2008; 159-170