



兰州大学

本科毕业论文

论文题目（中文）基于稀疏图码的离群值去除

论文题目（英文）Outliers Removal Based on
Sparse-Graph Codes

学生姓名 庄启源

指导教师 李朋

学 院 数学与统计学院

专 业 数学（基础理论班）

年 级 2020 级

兰州大学教务处

诚信责任书

本人郑重声明：本人所呈交的毕业论文（设计），是在导师的指导下独立进行研究所取得的成果。毕业论文（设计）中凡引用他人已经发表或未发表的成果、数据、观点等，均已明确注明出处。除文中已经注明引用的内容外，不包含任何其他个人、集体已经发表或未发表的论文。

本声明的法律责任由本人承担。

论文作者签名： 庄启源 日期： 2024年5月22日

关于毕业论文（设计）使用授权的声明

本人在导师指导下所完成的论文及相关的职务作品，知识产权归属兰州大学。本人完全了解兰州大学有关保存、使用毕业论文（设计）的规定，同意学校保存或向国家有关部门或机构递交论文的纸质版和电子版，允许论文被查阅和借阅；本人授权兰州大学可以将本毕业论文（设计）的全部或部分内容编入有关数据库进行检索，可以采用任何复制手段保存和汇编本毕业论文（设计）。本人离校后发表、使用毕业论文（设计）或与该毕业论文（设计）直接相关的学术论文或成果时，第一署名单位仍然为兰州大学。

本毕业论文（设计）研究内容：

可以公开

不宜公开，已在学位办公室办理保密申请，解密后适用本授权书。

（请在以上选项内选择其中一项打“√”）

论文作者签名： 庄启源

导师签名： 李明

日期： 2024年5月22日

日期： 2024年5月22日

基于稀疏图码的离群值去除

中文摘要

随着压缩感知理论研究的推进以及在工业界的大量应用，众多解码方法被提出以更精确快速地重构稀疏度为 K 的高维稀疏信号 $\mathbf{x} \in \mathbb{R}^N$ 。为了更好地解决现实会遇到的信号恢复问题，我们研究在无噪声情形下当观测信号存在少量离群值时稀疏信号的精确恢复。本文提出了一种名为 QTP 的离群值去除算法，能够有效消除离群值在信号恢复中的巨大偏差影响。在恢复过程中，我们应用稀疏图码相关理论，借由 DFT 矩阵和左正则二部图给出精心设计的测量矩阵。由此实现的重构算法仅需对观测值进行几次简单的解码迭代即可精确恢复稀疏系数。实验展示了在无噪声情况下，我们的框架能够高效且相当精确地恢复存在离群值的观测信号。本文的所有源代码和数据公开在https://github.com/waterEand/thesis_lzu。

关键词：压缩感知；离群值；稀疏信号恢复；稀疏图码

Outliers Removal Based on Sparse-Graph Codes

Abstract

With the advancement of research in compressed perception theory and its numerous applications in industry, numerous decoding methods have been proposed to fast reconstruct a high-dimensional K -sparse signal $\mathbf{x} \in \mathbb{R}^N$ accurately. In order to solve a signal recovery problem that may be encountered in reality more comprehensively, we study the exact recovery of sparse signals from the observations containing a few outliers without noise. In this work, an outlier removal algorithm named QTP is proposed to effectively eliminate the influence of large deviations of outliers in signal recovery. Through recovery, the measurement matrix is carefully designed through sparse-graph codes combined with the DFT matrix and the left regular bipartite graph. A reconstruction algorithm is also implemented, by which the sparse coefficients can be recovered in a few iterations by performing simple error decoding over the observations. It is shown in experiments that in the absence of noise, our framework is able to recover observed signals in the presence of outliers efficiently and quite accurately. All source codes and data of our work are available at https://github.com/waterEand/thesis_lzu.

Keywords: Compressed sensing; Outliers; Sparse signal recovery; Sparse-graph codes

目 录

中文摘要	I
英文摘要	II
第一章 绪 论	1
1.1 问题介绍	1
1.1.1 压缩感知问题	1
1.1.2 次线性时间的支撑集恢复算法	1
1.1.3 基于稀疏图码的分治法	2
1.2 研究目的与主要贡献	3
1.3 论文结构与记号	3
1.3.1 本文结构	3
1.3.2 相关记号	4
第二章 稀疏图编码	5
第三章 基于带离群值观测的稀疏图解码	8
3.1 节点类型检测	8
3.2 离群值去除	9
3.3 信号恢复	10
第四章 数值实验结果	11
4.1 信号恢复效率	11
4.1.1 不同测量次数下的数值效果	11
4.1.2 不同信号稀疏度下的数值效果	12
4.2 算法参数大小对数值效果的影响	12
4.3 与其他信号恢复算法的比较	13
第五章 总结与展望	15
参考文献	16

附录	18
A.1 主算法代码	18
A.2 对比算法代码	21
A.3 信号初始化代码	23
致谢	24

第一章 絮 论

1.1 问题介绍

1.1.1 压缩感知问题

压缩感知 (Compressed Sensing, CS) ^[1] 是一种利用信号稀疏性的信号采样理论。在有噪声情形，这一问题可被表达为：对一个长度为 N 的稀疏信号 \mathbf{x} ，使用一个随机测量矩阵采样后得到一个长度为 M 的观测向量，再由该观测值估计初始信号。信号的压缩观测过程可表示为：

$$\mathbf{y} = \mathbf{Ax} + \mathbf{w} \quad (1)$$

这一过程被称为编码。其中 \mathbf{A} 是一个维数为 $M \times N$ 的已知测量矩阵（观测矩阵）， \mathbf{w} 是一个可能的噪声向量， \mathbf{y} 是观测值。

通常情况下，若 \mathbf{x} 没有特殊的结构或额外的信息，我们无法通过远少于该信号维度的测量（采样）次数来恢复 \mathbf{x} 。但是，如果信号是稀疏的或可压缩的，即 \mathbf{x} 向量中只有 K 个非零元素 (\mathbf{x} 的稀疏度为 K) 且 $K \ll N$ ，那么我们就可以由极少数采样精确地恢复原信号。这一过程也被称为解码。这种通过低维观测值重建高维信号的压缩感知方法被应用在各个领域中，如医疗成像 ^[2]、天体物理学成像 ^[3]、语音与图像处理 ^[4] 等。

1.1.2 次线性时间的支撑集恢复算法

在压缩感知问题中，一个好的测量矩阵设计和高效的解码算法是极其重要的。这两个目标可以被具体描述为：

1. 要求能够保证信号精确恢复所需的观测次数 M (即测量矩阵 \mathbf{A} 的行数) 尽可能小；
2. 构建合适的测量矩阵 \mathbf{A} 使得恢复信号所需的时间复杂度尽可能小。

为同时达到这两个目的的最优，研究人员利用信号的稀疏性，提出了大量测量矩阵设计思路和信号重构算法来从低维观测值中恢复信号。其中，大部分相关文献提出的稀疏信号恢复模型都主要基于 ℓ_2/ℓ_1 范数或 ℓ_1/ℓ_1 范数的近似误差指标，而基于支撑集恢复 ^[5] 的工作就相对较少。对于在不同观测情况下精确恢复信号支撑集的充要条件，学者们采取了最优解码器 ^[6,7]、 ℓ_1 极小 ^[8]、贪心算法 ^[9] 等策略进行研究。

Wainwright ^[7] 提出，在高斯噪声下，当测量矩阵的元素服从独立同分布 (i.i.d.) 的高斯分布时，信号支撑集恢复的充要条件是 $O(K \log(N/K))$ 次测量 (即 $M = O(K \log(N/K))$)。Bakshi、Jaggi ^[10] 等人提出的 SHO-FA 算法使有噪声情形下的信号恢复仅需 $O(K)$ 次测量且编码解码时间复杂度分别为 $O(N)$ 和 $O(K)$ ，这一结果达到了信息论中的阶次最优。Li、

Yin^[11] 等人建立了一种新的压缩感知框架，能够在无噪声情况下，用 $2K$ 次测量精确地恢复任意稀疏度为 K 的信号，解码时间复杂度达到 $O(K)$ ；在有噪声情形，用 $O(K \log(N/K))$ 次测量能达到同等效果，解码时间复杂度为 $O(K \log(N/K))$ 。并且文献^[11] 提出的模型在 K 为次线性时，即存在 $\delta \in (0, 1)$ 使得 $K = O(N^\delta)$ ，测量次数和计算时间都与信号维数 N 成次线性关系，实现了概率保证下两者同时的阶次最优。

本文中我们将应用基于稀疏图码的算法^[11, 12] 恢复离群值去除后的信号，使得总体测量次数和恢复时间都达到次线性的效果。

1.1.3 基于稀疏图码的分治法

我们通过通信系统中的稀疏图码来看待压缩感知问题，并用简单的“分治法”（分而治之）^[11] 来处理稀疏信号。如图1.1的 (a) 所示，我们用不同的颜色表示稀疏向量 \mathbf{x} 中的不同元素，其中红色、绿色、蓝色方格表示非零元素，白色方格表示零元素。在测量矩阵中，灰色方格表示随机生成的元素。相应生成的观测向量中每个元素则是红、绿、蓝颜色的混合。

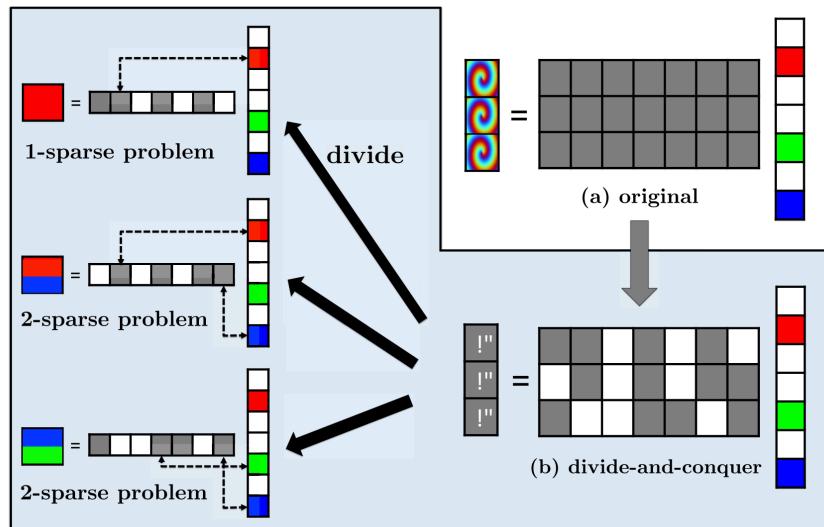


图 1.1 “分治法”的概念抽象图。子图 (a) 描绘一个稀疏度为 3 的恢复问题，其中测量矩阵着灰色，表示元素是随机生成的。由于不同颜色成分（红色、绿色、蓝色）随机混合，导致得到的观测结果呈现混合色。在子图 (b) 中，我们通过在每行中放置三个零元素来稀疏化测量矩阵，即图中的白色方格。由此产生的测量矩阵将稀疏度为 3 的恢复问题分解为多个稀疏度小于 3 的子问题。若子问题的观测值是单一颜色，我们可以立即恢复这一元素。在这一图中，第一个测量中的红色元素可以被立即恢复，之后再从第二个红蓝元素中把红色剥离，恢复蓝色，以此类推。

在图1.1的 (b) 中，我们用零元素（图中的空白格）来稀疏化测量矩阵的每一行。这样设计的测量矩阵就导致了观测到的向量 \mathbf{y} 中包含的元素一部分是单一颜色，一部分是几种颜色的混合。我们的设计理念就是把稀疏信号中的非零元分散到多个单一颜色的测量中（例如第一个测量中的红色），然后从把这些单一颜色从混合颜色中剥离出来（例如第二个测量

中的红蓝混合和第三个测量中的蓝绿混合), 如此迭代逐步解码其他颜色的未知元素, 这也就是“分治”的过程。稀疏图码在这里的作用本质上是将一般的稀疏信号重构问题分解为多个可以轻松解决的子问题, 最后再把这些子问题的结果融合以恢复所有非零元素。这样的设计兼具了稀疏图编码在测量成本(容量逼近)和计算复杂性(基于快速剥离的解码)方面的特性, 从而使稀疏测量矩阵达到了低测量次数和低计算成本的效果。

1.2 研究目的与主要贡献

考虑到压缩感知常被应用在图像采样的任务中, 且采样过程可能会出现探头损坏, 导致得到的观测值会出现离群值(Outliers)^[13]的情况, 本文研究在无噪声情形下, 观测值中存在部分离群值的稀疏信号的解码恢复。该观测过程可表示为:

$$\mathbf{y} = \mathbf{Ax} + \mathbf{f} \quad (2)$$

其中 \mathbf{f} 是一个维数为 M 的离群向量, \mathbf{f} 中的非零元素(离群值)个数 T 远小于 M , 满足 $|\text{supp}(\mathbf{f})| = T = \eta M$ ($0 \leq \eta < \frac{1}{2}$), 且这些离群值的分布方差显著大于 \mathbf{x} 中的非零元素。我们的目标是尽可能去除 \mathbf{y} 的非零元素中的所有离群值, 并着眼于精确地恢复稀疏度为 K 长度为 N 的信号的支撑集(稀疏体系)以及相对应的稀疏系数。这一问题可以通过以下模型求解:

$$\min_{\|\mathbf{x}\|_0 \leq K} \|\mathbf{Ax} - \mathbf{y}\|_1. \quad (*)$$

由于 $\mathbf{f} = \mathbf{y} - \mathbf{Ax}$ 是稀疏的, 即 $\|\mathbf{y} - \mathbf{Ax}\|_0 \ll M$, 我们将 ℓ_0 范数松弛为 ℓ_1 范数。

支撑集恢复结果的一个经典误差评价指标是信号支撑集无法被精确恢复的概率, 即 $\Pr(\text{supp}(\hat{\mathbf{x}}) \neq \text{supp}(\mathbf{x}))$, 其中 $\text{supp}(\mathbf{x}) := \{k : x[k] \neq 0, 0 \leq k \leq N-1\}$ 。也就是说, 对于任意给定的稀疏度为 K 长度为 N 的信号 \mathbf{x} , 我们设计一个测量矩阵 \mathbf{A} , 并计算出一个估计向量 $\hat{\mathbf{x}}$, 并且希望这一个估计的支撑集与原信号 \mathbf{x} 的支撑集完全匹配的概率接近于 1。同时, 我们也要恢复 \mathbf{x} 中的稀疏系数, 也就是要恢复信号中非零元的精确值。兼顾以上两个要求, 我们以信号被精确恢复的概率为评价指标, 具体定义将在第四章给出。

在文章中, 我们将一种基于分位数截断的离群值去除方法^[14]与基于稀疏图码的剥离解码算法结合, 提出一种全新的能够有效恢复存在离群值的信号的算法, 并将其命名为 QTP(Quantile-Truncated Peeling) 算法。我们模拟的数值实验表明, 同等情况下, QTP 算法的恢复效率远高于传统的 PLAD 算法^[15, 16], 恢复误差也极小, 且更具稳定性。

1.3 论文结构与记号

1.3.1 本文结构

本文主要讨论当观测值中存在离群值情况下, 用稀疏图码的思想来重构稀疏信号。文章共五章, 主要结构如下:

第一章介绍了问题背景，研究意义，本文的主要贡献及相关记号设定。

第二章介绍了稀疏图编码过程，主要内容是测量矩阵的设计，包括一些基础定义如节点类型、克罗内克积等，为第三章基于稀疏图解码算法做铺垫。

第三章详细阐述了 QTP 算法的细节及具体应用。我们的算法可分为两阶段，首先由基于分位数的算法1去除观测值中的离群值，再根据1.1.3小节的思想把稀疏度为 K 的信号恢复问题转化为求解多个稀疏度为 1 维数为 N 的信号恢复问题，也就是剥离解码器（算法2）。

第四章是对 QTP 算法的数值模拟以及与传统的 PLAD 算法的对比。在这一章我们探究了 QTP 算法的性质，也证明了其在时间和精度上的有效性。

第五章对本文内容作了总结，并提出未来可能推广到的领域和方向。

1.3.2 相关记号

我们对本文中相关记号做以下规定：粗体大写字母（如 \mathbf{A} ）表示矩阵，粗体小写字母（如 \mathbf{x} ）表示向量，大写斜体字母（如 N ）表示矩阵或向量的维数，小写斜体字母（如 x ）表示向量或矩阵中的元素。花体大写字母（如 \mathcal{A} ）表示集合，特别的，花体大写字母 \mathcal{G} 表示二部图。 $|\mathcal{A}|$ 表示集合 \mathcal{A} 中的元素个数， $|y|$ 表示 y 的模（或绝对值）。

对于任意正整数 n ， $[n]$ 表示集合 $\{0, 1, 2, \dots, n-1\}$ 。一个元素类型为复数的向量 $\mathbf{x} \in \mathbb{C}^N$ 可表示为 $\mathbf{x} = [x[0], x[1], \dots, x[N-1]]^T$ ；其中向量 \mathbf{x} 中的第 i 个元素用 $x[i]$ 或 x_i 表示。对于任意矩阵 \mathbf{A} ，我们用 \mathbf{A}_i 表示其第 i 列（向量）。对任意向量 $\mathbf{x} \in \mathbb{R}^N$ ， ℓ_p 范数 ($p \geq 1$) 定义为 $\|\mathbf{x}\|_p = (\sum_{i=1}^N |x_i|^p)^{\frac{1}{p}}$ 。

第二章 稀疏图编码

为了对稀疏信号 \mathbf{x} 进行编码（压缩），即 $\mathbf{y} = \mathbf{Ax} + \mathbf{f}$ 的过程，我们根据1.1.3小节在本章中借助简单例子介绍测量矩阵 \mathbf{A} 的具体设计方法。

考虑一个长度 $N = 16$, 稀疏度 $K = 4$ 的稀疏信号 \mathbf{x} 。其中包含有非零元素 $x[1] = 1$, $x[4] = 4$, $x[8] = 2$, $x[13] = 7$ 。接下来我们构建一个左节点数为 16, 右节点数为 9 的二部图。如图2.1所示，这个图有以下性质：

- 每个左节点 $x[k]$ 都连接至少一个右节点。
- 每个右节点 y_r 的值是与其相连接的的左节点值的和。

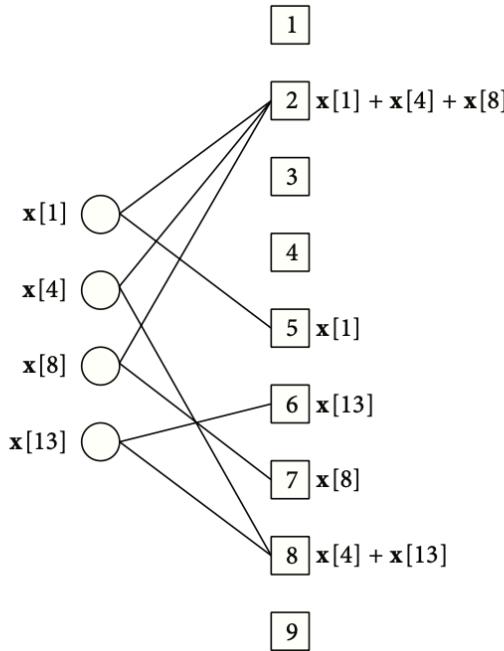


图 2.1 稀疏测量二部图

这样就形成了对于长度为 16 的稀疏信号 \mathbf{x} （左节点）的 9 次测量：

$$y_1 = y_3 = y_4 = y_9 = 0,$$

$$y_2 = x[1] + x[4] + x[8],$$

$$y_5 = x[1],$$

$$y_6 = x[13],$$

$$y_7 = x[8],$$

$$y_8 = x[4] + x[13].$$

根据稀疏测量二部图的连通性，我们将与右节点分为以下三种类型：

1. 零节点：若一个右节点不包含任何非零元素，则称其为零节点，如图2.1的 y_1 。
2. 单节点：若一个右节点仅包含一个非零元素，则称其为单节点，如图2.1的 y_5 。
3. 多节点：若一个右节点包含一个以上的非零元素，则称其为多节点，如图2.1的 y_2 。

其实单节点就是章节1.1.3描述的单一颜色的元素，多节点就是混合颜色。我们从单节点入手进行逐步剥离解码，具体算法将在第3.3章介绍。

为了判断观测向量中的元素（右节点）属于哪种节点，我们给每个左节点不同的2维向量加权，于是每个右节点可被表示为几个二维向量之和。为此我们设计一个列数与信号 \mathbf{x} 长度相同的检测矩阵 \mathbf{S} ：

$$\mathbf{S} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & \dots & 1 \\ 1 & W & W^2 & W^3 & W^4 & \dots & W^{15} \end{bmatrix}, \quad (3)$$

其中 $W = e^{i\frac{2\pi}{N}}$ 是一个 N 次单位根，这里的例子中 $N = 16$ 。 \mathbf{S} 是一个 16×16 的 DFT 矩阵的前两行，通过 \mathbf{S} 检测节点的具体算法将在章节3.3介绍。

除此之外，我们还需要定义一个编码矩阵 \mathbf{H} 来表示二部图左右节点的连接， \mathbf{H} 在这个例子中是一个 9×16 的由 0 和 1 组成的邻接矩阵。

为了将 \mathbf{H} 和 \mathbf{S} 结合在一起形成完整的测量矩阵 \mathbf{A} ，下面我们介绍行张量算子 \boxtimes (row-tensor operator) 的定义。一般情况下，检测矩阵 $\mathbf{S} = [\mathbf{s}_0, \dots, \mathbf{s}_{N-1}] \in \mathbb{C}^{M_2 \times N}$ ，邻接矩阵 $\mathbf{H} = [\mathbf{h}_0, \dots, \mathbf{h}_{N-1}] \in \mathbb{C}^{M_1 \times N}$ 。行张量运算 $\mathbf{H} \boxtimes \mathbf{S}$ 本质上就是将矩阵 \mathbf{H} 的每一行通过与矩阵 \mathbf{S} 中每个相应的列进行逐个元素的增广。我们得到的行张量积是一个维数为 $M_1 M_2 \times N$ 的矩阵：

$$\mathbf{H} \boxtimes \mathbf{S} = [\mathbf{h}_0 \otimes \mathbf{s}_0, \dots, \mathbf{h}_{N-1} \otimes \mathbf{s}_{N-1}],$$

其中 \otimes 是一个标准的克罗内克积 (Kronecker product)。下面我们举一个简单例子来具象化这一过程。我们设定

$$\mathbf{H} = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 1 & 0 & 1 \end{bmatrix}, \quad (4)$$

$$\mathbf{S} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & W & W^2 & W^3 & W^4 & W^5 & W^6 \end{bmatrix}, \quad (5)$$

这里 $W = e^{i\frac{2\pi}{7}}$ 。那么我们就可以计算出行张量积：

$$\mathbf{H} \boxtimes \mathbf{S} = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & W & 0 & W^3 & 0 & W^5 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 1 \\ 0 & W & 0 & W^3 & 0 & 0 & W^6 \\ 1 & 0 & 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 0 & W^3 & W^4 & 0 & W^6 \end{bmatrix}. \quad (6)$$

在本例中， \mathbf{H} 有三行，分别与 $\mathbf{H} \boxtimes \mathbf{S}$ 的前两行，中间两行，后两行对应。

用这种方式，我们可以生成一个测量矩阵 $\mathbf{A} = \mathbf{H} \boxtimes \mathbf{S} \in \mathbb{C}^{M_1 M_2 \times N}$ ，这样得到观测值 $\mathbf{y} \in \mathbb{C}^{M_1 M_2}$ 。

定义 2.1 (测量矩阵) 设 $M = RP$ 且 $P, R \in \mathbb{N}^*$ 。给定一个 $R \times N$ 的编码矩阵 \mathbf{H} 和一个 $P \times N$ 的检测矩阵 \mathbf{S} ，那么 $M \times N$ 的测量矩阵 \mathbf{A} 就可被定义为：

$$\mathbf{A} = \mathbf{H} \boxtimes \mathbf{S}, \quad (7)$$

其中 \boxtimes 是行张量算子，并且编码矩阵和检测矩阵满足下述条件：

- 编码矩阵 $\mathbf{H} = [H_{r,n}]_{R \times N}$ 是一个二部图 \mathcal{G} 的邻接矩阵。图 \mathcal{G} 包含 N 个左节点 $V_1 := [N]$ 和 R 个右节点 $V_2 := [R]$ ，边的集合为 $\mathcal{E} := V_1 \times V_2$ 。
- 检测矩阵 $\mathbf{S} := [\mathbf{s}_0, \dots, \mathbf{s}_{N-1}]$ 是一个 $N \times N$ 的 DFT 矩阵的前 P 行。

第三章 基于带离群值观测的稀疏图解码

本章将介绍基于带离群值观测的信号恢复策略。我们考虑模型 (*), 并设计具体求解算法。

3.1 节点类型检测

我们继续借用图2.1的例子来模拟算法。由公式3可知, 我们可以通过 \mathbf{S} 把每个右节点都记为一个二维向量 $\mathbf{y}_r = [y_r[0], y_r[1]]^T$, 我们称每个这样的向量为一个测量对 (measurement bin)。在上述例子中, 右节点 1, 2, 5 的测量就可以被记为:

$$\begin{aligned}\mathbf{y}_1 &= \mathbf{0}, \\ \mathbf{y}_2 &= x[1] \times \begin{bmatrix} 1 \\ W \end{bmatrix} + x[4] \times \begin{bmatrix} 1 \\ W^4 \end{bmatrix} + x[8] \times \begin{bmatrix} 1 \\ W^8 \end{bmatrix}, \\ \mathbf{y}_5 &= x[1] \times \begin{bmatrix} 1 \\ W \end{bmatrix}.\end{aligned}$$

于是通过对测量对的检测, 我们能够有效地确定右节点是零节点、单节点还是多节点。我们由右节点 1, 2, 5 的例子介绍节点检测的方法:

1. 零节点对: 参考右节点 1。若测量对是一个全 0 向量, 即 $\mathbf{y}_r = \mathbf{0}$, 则该节点为零节点。
2. 单节点对: 参考右节点 5。这个测量对满足比值检测:

$$\hat{k} = \frac{\angle y_5[1]/y_5[0]}{2\pi/16} = 1,$$

其中 \hat{k} 为整数, 则该节点就是单节点。于是

$$\hat{x}[\hat{k}] = y_5[0].$$

我们也就从中顺利恢复了信号 \mathbf{x} 中的非零元素 x_1 。另一种检测方法是, 这样的单节点对的两个元素的模一定相等, 即 $|y_3[0]| = |y_3[1]|$ 。因为 W 是单位根, $|W^r| = 1$ 。

3. 多节点对: 参考右节点 2。对这个测量对进行比值检测:

$$\hat{k} = \frac{\angle y_2[1]/y_2[0]}{2\pi/16} = 4.85995.$$

并且两个元素的模并不相等, $|y_2[1]| \neq |y_2[0]|$ 。在这种比值检测的结果 \hat{k} 不为非零整数或两个元素不同模的情况下, 该节点为多节点。

综上，要判断任意非零节点对 y_r 的类型，我们只需计算

$$\hat{k} = \frac{\angle y_r[1]/y_r[0]}{2\pi/N}, \quad (8)$$

当 \hat{k} 为整数，该节点是单节点；反之则为多节点。

3.2 离群值去除

离群值 (outliers)，也被称为异常值，是指与其他观察结果显著不同的数据点。存在离群值情况下，压缩观测过程可被描述为 $\mathbf{y} = \mathbf{Ax} + \mathbf{f}$ ，其中 $\mathbf{f} = [f_1, f_2, \dots, f_M]^T$ ，且 $|\text{supp}(\mathbf{f})| = T = \eta M$ ($0 \leq \eta < \frac{1}{2}$)。我们设初始稀疏度为 K 的稀疏信号 \mathbf{x} 中的非零元素服从均值为 0，标准差为 1 的高斯分布，即 $x_i \sim N(0, 1), i \in [N]$ 。稀疏度为 T 的离群向量 \mathbf{f} 中的非零元素满足 $f_j \sim N(0, 100), j \in [M]$ 。

从中我们可以显然得到： \mathbf{f} 中大部分非零元素的绝对值远大于 \mathbf{x} 中元素的绝对值。为了尽可能剔除这些离群值，又考虑到 \mathbf{y} 的稀疏性，我们可以尝试取 \mathbf{y} 的非零元素绝对值的某个分位数 θ_p (p 是一个分位数算子，常取 0.5、0.75)，并且设定一个倍率 α 得到一个分界数 $\alpha \cdot \theta_p$ ，抹去 \mathbf{y} 的非零元素中所有绝对值大于 $\alpha \cdot \theta_p$ 的数的集合。这里“抹去”的定义是将这些非零元素设为 0，也就是这些元素不参与信号重构的过程。

用数学语言描述，我们把 \mathbf{y} 的有效非零元素限制在一个指标集 $S = \{s : |y_s| < \alpha \theta_p(\{|y_j| : j \in \text{supp}(\mathbf{y})\})\}$ 上，其他指标的对应元素全部设为 0，即截断方法^[14]。例如，当 $\alpha = 1, p = \frac{1}{2}$ 时，我们就抹去所有绝对值大于非零元素绝对值中位数的数。这样就理论上去除了全部绝对值过大的离群值，接下来对被限制的 \mathbf{y} 应用信号重构算法即可。

由于离群向量 \mathbf{f} 的稀疏度 T 是未知的，我们无法简单地抹去绝对值最大的前 T 个元素，因此采取上述设置分位数的方式。我们又知道 \mathbf{y} 和 \mathbf{f} 都是稀疏的，所以 \mathbf{y} 和 \mathbf{f} 的支撑集大概率不会存在元素位置重合的情况，不会抹去需要参与信号恢复的非离群值。因此，这样直接设为 0 的方式是合理且可取的。算法伪代码如下：

算法 1 离群值去除算法

输入：观测值 $\mathbf{y} \in \mathbb{C}^{2R}$ (测量对形式)，常数 α 和 p 。

- 1: 首先计算得到 \mathbf{y} 的非零元素模的分位数 $q = \theta_p(\{|y_j[0]| : y_j[0] \neq 0, y_j \in \mathbf{y}\})$;
 - 2: **for** $r = 1$ to R **do**
 - 3: **if** $|y_r[0]| > \alpha \cdot q$ **then**
 - 4: 将 y_r 抹去为 $\mathbf{0}$ 向量;
 - 5: **end if**
 - 6: **end for**
-

3.3 信号恢复

本节将聚焦于离群值已去除的信号恢复。经过上述铺垫，我们按照以下步骤解码即可（参考1.1.3小节）：

- 1) 选取二部图中所有使得右节点度数为 1 的边（找到所有单节点）
- 2) 去除（剥离）所有这些边以及与之对应的左右节点对。
- 3) 去除（剥离）步骤2)中剥离的左节点的其他未被剥离的边。
- 4) 找到步骤3)中去除的左节点连接的所有右节点，将这些左节点的值从右节点的值中减去。

伪代码如下：

算法 2 剥离解码器

输入：观测值 $\mathbf{y} \in \mathbb{C}^{2R}$ （测量对形式），待恢复信号 $\hat{\mathbf{x}} = \mathbf{0}$ 。

```

1: for  $i = 1, 2, \dots$  直到  $\mathbf{y}^{(i)} = \mathbf{0}$  或找不到单节点 do
2:   for  $r = 1$  to  $R$  do
3:     由公式 8 计算  $\hat{k}$ ，判断  $\mathbf{y}_r^{(i)}$  是否为单节点对；
4:     if  $\mathbf{y}_r^{(i)}$  是单节点 then
5:        $\hat{\mathbf{x}}[\hat{k}] = \mathbf{y}_r^{(i)}[0]$  ;
6:       for  $r' = 1$  to  $R$  do
7:         找到连接左节点  $\hat{k}$  的右节点  $r'$ ;
8:         剥离过程  $\mathbf{y}_{r'}^{(i+1)} = \mathbf{y}_{r'}^{(i)} - \hat{\mathbf{x}}[\hat{k}] \mathbf{s}_{\hat{k}}$ ，这里  $\mathbf{s}_{\hat{k}}$  是检测矩阵  $\mathbf{S}$  的第  $\hat{k}$  列;
9:       end for
10:      else
11:        继续下一个节点对  $\mathbf{y}_{r+1}^{(i)}$ 。
12:      end if
13:    end for
14:  end for

```

完整的恢复过程是两阶段的，先用算法1去除离群值，再用算法2恢复信号。在离群值去除过程中，我们把分位数乘上一定倍率作为分割点；在信号恢复中，我们应用稀疏图码中的剥离解码方法。因此我们把这一算法命名为 QTP（Quantile-Truncated Peeling）算法。

第四章 数值实验结果

本章我们对无噪声条件下存在离群值的信号恢复问题进行数值模拟，从实验角度出发论证理论内容。在本章的实验中，我们通过上一章描述的 QTP 算法对不同情况的离群值问题进行求解。先验证算法在不同情形下的高效性，之后比较离群值去除算法参数不同时的恢复效果，最后与同条件下其他算法做对比来证明我们算法的优越性。

我们将实验的参数统一设置为： $N = 200$ ；初始稀疏信号 $\mathbf{x} \in \mathbb{R}^N$ 且其中非零元素 $x_i \sim N(0, 1)$ ；编码矩阵 $\mathbf{H} \in \{0, 1\}^{M \times N}$ ，且 \mathbf{H} 为每列有且仅有三个 1 的矩阵，即正则度为 3 的左正则二部图邻接矩阵；检测矩阵 \mathbf{S} 为 $N \times N$ 的 DFT 矩阵前两行；离群向量 \mathbf{f} 中的非零元素 $f_j \sim N(0, 100)$ 。本文中数值实验均由 Python 3.11.3 实现，在 macOS Sonoma 14.2.1 操作系统下，用一块 Apple M2 Pro (CPU) 芯片完成运行。

为了评估信号评估恢复的效果，我们定义评价指标精确恢复率为多次重复实验中恢复结果误差 $err < thr_err$ 的次数与总实验次数的比值。这里

$$err = \frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_2}{\|\mathbf{x}\|_2}, \quad (9)$$

其中 \mathbf{x} 是初始信号， $\hat{\mathbf{x}}$ 为恢复结果； thr_err 是界定一次信号恢复是否成功的阈值。在实验的大多数情况下我们取 $thr_err = 10^{-7}$ ；而由于我们将对比的传统算法（PLAD 算法）无法达到这样的准确度，我们在对比实验中加入 $thr_err = 10^{-2}$ 的实验。当 $err < thr_err$ ，我们称信号被精确恢复一次。

4.1 信号恢复效率

在本节中我们用信号精确恢复率和恢复用时两个指标来评估 QTP 算法的高效性。

4.1.1 不同测量次数下的数值效果

在第一个实验中，测量值数目 M 作为自变量。我们设定 \mathbf{x} 的稀疏度为 2，离群值比例 $\eta = 0.01$ ，离群值去除算法中 $\alpha = 2, p = \frac{1}{2}$ ，重复实验次数为 200 次。结果如下：

表 4.1 不同测量次数下的恢复效率

测量值数目 M	10	20	40	80	100	120
精确恢复率	70.0%	84.5%	90.5%	92.5%	95.5%	95.0%
平均用时 (毫秒)	48.55	79.25	89.70	95.40	115.25	126.50

由表4.1可知仅需 100 次的测量就可 95% 精确地恢复存在一个离群值的稀疏度为 2 的信号 $\mathbf{x} \in \mathbb{R}^{200}$, 且在 500 次随机实验中每次恢复所需时间仅 0.115 秒左右。

4.1.2 不同信号稀疏度下的数值效果

在第二个实验中, \mathbf{x} 的信号稀疏度 K 作为自变量。我们设定 $M = 100$; 离群值去除算法中 $\alpha = 2, p = \frac{3}{4}$; 重复实验次数为 500 次; 其他与上一个实验相同。结果如下:

表 4.2 不同信号稀疏度下的恢复效率

稀疏度 K	1	2	3	4	5
精确恢复率	98.0%	95.4%	92.2%	91.8%	88.8%
平均用时 (毫秒)	124.0	129.0	129.6	125.4	130.4

在表4.2我们可以看到信号稀疏度越小, 恢复越精确; 平均每次信号恢复时间为 0.124 ~ 0.130 秒, 总体上稀疏度 K 越大恢复所需时间越长。

4.2 算法参数大小对数值效果的影响

本节实验中我们探究 QTP 算法的参数大小 (包括 α 和 p) 对恢复效果的影响。

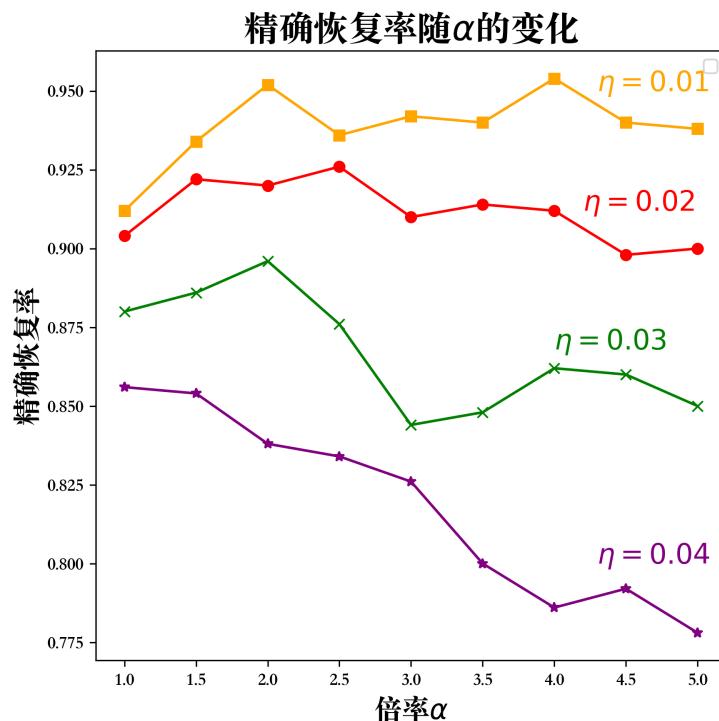


图 4.1 不同离群值比例下的恢复效果随算法参数变化展示。图中横坐标表示不同的倍率 α , 纵坐标表示精确恢复率, 四条曲线分别表示离群值比例取 1%, 2%, 3%, 4% 的情形。

这里我们选取 \mathbf{x} 的稀疏度 $K = 2$, 分位数 θ_p 中的分位数算子 $p = 0.5$, 重复实验次数为 500 次。这里只讨论倍率 α 的影响, 是因为由章节3.2可知, 离群值去除的分界数是 $\alpha\theta_p$, 因此模拟的数值实验中 α 的变化和 p 的变化是等价的。我们通过离群值比例 η 不同的实验, 对精确恢复率随倍率 α 的变化进行展示。结果如图4.1所示。

在离群向量 \mathbf{f} 中稀疏度 $T = |\text{supp}(\mathbf{f})| = \eta M$ 。显然, 随着离群值比例 η 的增大, 信号恢复的难度也越来越大。直觉上, 离群值比例越高, 分位数 θ_p 就越大, 倍率 α 就要更小(或 p 更小)来保证去除所有离群值。由图可知, 随离群值比例增大, 精确恢复率下降明显。我们需要调小倍率 α (或调小 p) 来保证较好的恢复效果。当离群值比例超过 0.03 时, 精确恢复概率无法达到期望。

因此, 当我们在解决离群值极少的任务时, 可以选择较大的倍率 α 或分位数算子 p ; 而在大多数情况下, 选择小的倍率 α 是更加稳定的。

4.3 与其他信号恢复算法的比较

为了说明 QTP 算法的存在意义与优越性, 我们将算法与其他传统算法对比, 比较同等情况下两者的精确恢复率和恢复所需时间。我们选择对比算法为 PLAD 算法 (Penalized Least Absolute Deviation), 该算法模型为:

$$\min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{Ax}\|_1 + \lambda \|\mathbf{x}\|_1. \quad (10)$$

这是信号恢复领域非常经典且能够解决离群值问题的算法。其他算法如 LASSO、DS 等算法都无法去除偏差较大的离群值, 故不做比较。

在本节实验中, 我们设置 \mathbf{x} 的稀疏度 $K = 3$, $M = 100$, 离群值比例取 0.02 和 0.03, 编码矩阵 H 和重复实验次数与上一节相同。QTP 算法的参数 α, p 根据上一节实验取对应情况下精确恢复率最高的值。PLAD 算法中测量矩阵 \mathbf{A} 考虑高斯矩阵和二元稀疏矩阵用于比较。对于 PLAD 算法, 我们在每次实验对不同 λ 的值进行微调, 以期得到最好效果。经检验, 对于高斯矩阵情况, 取 $\lambda \in (0.006, 0.01)$; 对二元稀疏矩阵, 取 $\lambda \in (0.0008, 0.0016)$ 。考虑到 PLAD 精度较低问题, 本节中两组实验分别设定 $thr_err = 10^{-7}$ 和 $thr_err = 10^{-2}$ 。实验结果如下 (表中 XX.X% 是精确恢复率):

表 4.3 QTP 算法与 PLAD 算法的对比

thr_err	离群值比例 η	$\text{PLAD}_G^{\text{FA}}$	$\text{PLAD}_G^{\text{SL}}$	$\text{PLAD}_B^{\text{FA}}$	$\text{PLAD}_B^{\text{SL}}$	QTP
10^{-7}	0.02	0%	0%	0%	0%	92.6%
10^{-2}	0.02	80.8%	91.6%	0.2%	87.2%	92.8%
	0.03	79.0%	90.2%	0%	79.4%	89.6%
平均用时 (毫秒)		175.4	584.7	608.2	2318.6	129.0

上表中, $\text{PLAD}_G^{\text{FA}}$ 表示测量矩阵为高斯矩阵 \mathbf{A} , 每次恢复迭代次数 I 为 500 次; $\text{PLAD}_G^{\text{SL}}$ 表示测量矩阵为高斯矩阵, 而 $I = 2000$ 。 $\text{PLAD}_B^{\text{FA}}$ 表示测量矩阵为二元稀疏矩阵, 且与 QTP 算法中的 \mathbf{H} 一致, 迭代次数 $I = 2000$; $\text{PLAD}_B^{\text{SL}}$ 表示测量矩阵同为 \mathbf{H} , 而 $I = 10000$ 。

实验结果显示, 对于每一次信号恢复, QTP 算法的精确程度远高于 PLAD 算法, 达到 $err < 10^{-7}$ 。而在 $thr_err = 10^{-2}$ 下, 相同恢复时间的 PLAD 算法精确恢复率大大低于 QTP 算法。要使 PLAD 算法的精确恢复率达到本文提出的 QTP 算法, 平均恢复所需用时远高于 (4 倍以上) QTP 算法, 因此我们的 QTP 算法相较 PLAD 更能胜任实时 (快速) 的信号恢复任务。

另外, 在对 PLAD 算法中的参数 λ 进行微调的过程中, 我们发现 PLAD 的恢复效果受参数 λ 取值的影响极严重。 λ 取值的轻微波动 (小于 10^{-4}) 都可能导致精确恢复率下降 10% 以上。并且, 一旦离群值比例轻微变化, 相同的 λ 取值极可能导致精确恢复下降到低于 50%。相对地, 由4.2节可知, QTP 算法的恢复效果受参数取值影响并不严重。因此, 我们的 QTP 算法相较 PLAD 在实际应用中将更加稳定。

第五章 总结与展望

本文研究了存在少量离群值时稀疏信号的恢复。我们先介绍了用二部图生成邻接矩阵，再将一种基于分位数截断的离群值去除方法与前人提出的基于稀疏图码的次线性时间解码算法结合，通过先去除再恢复的两阶段法达到精确恢复信号的目的。我们将这一完整的计算过程命名为 **QTP** 算法。最后的数值实验论证了算法的有效性，其中与 **PLAD** 算法的比较展示了 **QTP** 算法的性能：恢复更加精确，仅需极短的 CPU 时间，且不会因算法参数变化产生极大波动。

在展望方面，我们可以进一步研究拓展到稀疏低秩矩阵恢复的情形，与文献 [17] 改进的算法结合，也可能有较好的效果。另外，文中没有讨论编码矩阵 **H** 的具体设计策略，我们选择的是左正则二部图生成的二元稀疏矩阵，而其他满足 **RIP** 条件的矩阵可能会有更好的效果。

参考文献

- [1] Donoho D. Compressed sensing[J]. IEEE Transactions on Information Theory, 2006, 52(4):1289–1306.
- [2] Lustig M, Donoho D, Pauly J M. Sparse mri: The application of compressed sensing for rapid mr imaging[J]. Magnetic resonance in medicine, 2007, 58(6):1182—1195.
- [3] Wiaux Y, Jacques L, Puy G, et al. Compressed sensing imaging techniques for radio interferometry[J]. Monthly Notices of the Royal Astronomical Society, 2009, 395(3):1733–1742.
- [4] Elad M. Sparse and redundant representations: from theory to applications in signal and image processing[M]. Springer Science & Business Media, 2010.
- [5] Gilbert A, Indyk P. Sparse recovery using sparse matrices[J]. Proceedings of the IEEE, 2010, 98(6):937–947.
- [6] Gastpar M, Bresler Y. On the necessary density for spectrum-blind nonuniform sampling subject to quantization[C]. In 2000 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.00CH37100), volume 1, 348–351 vol.1, 2000.
- [7] Wainwright M J. Information-theoretic limits on sparsity recovery in the high-dimensional and noisy setting[Z], 2007.
- [8] Candès E J, Plan Y. Near-ideal model selection by 1 minimization[J]. The Annals of Statistics, 2009, 37(5A):2145 – 2177.
- [9] Cai T T, Wang L. Orthogonal matching pursuit for sparse signal recovery with noise[J]. IEEE Transactions on Information Theory, 2011, 57(7):4680–4688.
- [10] Bakshi M, Jaggi S, Cai S, et al. Sho-fa: Robust compressive sensing with order-optimal complexity, measurements, and bits[J]. IEEE Transactions on Information Theory, 2016, 62(12):7419–7444.
- [11] Li X, Yin D, Pawar S, et al. Sub-linear time support recovery for compressed sensing using sparse-graph codes[J]. IEEE Transactions on Information Theory, 2019, 65(10):6580–6619.
- [12] 周华乔, 徐义晗, 孙一凡, 等. 基于稀疏图码的物联网邻居节点发现 [J]. 计算机应用研究, 2022, 39(6):1829–1833.
- [13] Laska J N, Davenport M A, Baraniuk R G. Exact signal recovery from sparsely corrupted measurements through the pursuit of justice[C]. In 2009 Conference Record of the Forty-Third Asilomar Conference on Signals, Systems and Computers, 1556–1560, 2009.
- [14] Li Y, Chi Y, Zhang H, et al. Non-convex low-rank matrix recovery with arbitrary outliers via median-truncated gradient descent[J]. Information and Inference: A Journal of the IMA, 2019, 9(2):289–325.
- [15] Yang J, Zhang Y. Alternating direction algorithms for ℓ_1 -problems in compressive sensing[J]. SIAM Journal on Scientific Computing, 2011, 33(1):250–278.
- [16] Wang L. The ℓ_1 penalized lad estimator for high dimensional linear regression[J]. Journal of Multivariate Analysis, 2013, 120:135–151.

- [17] Liu X, Venkataraman R. Sketching sparse low-rank matrices with near-optimal sample- and time-complexity using message passing[J]. IEEE Transactions on Information Theory, 2023, 69(9):6071–6097.

附录

A.1 主算法代码

QTP 算法文件。

```
1 # peeling_decoder_vector.py
2
3 import random
4 import time
5 import numpy as np
6 import scipy as sp
7 import cmath
8 from signal_vectors import signal_vec, outliers
9 from scipy.sparse import csr_matrix
10 from numpy import linalg as LA
11
12 # parameters to be modified
13 n = 200 # N
14 R = 100 # M
15 k = 2 # sparsity of x
16 density_H = 3
17 num_outliers = 1
18 med_times = 2
19 percent = 50 # p
20 var_outliers = 100 # variance of outliers
21 thr_err = 1e-7
22 num_epoch = 500
23
24 def input_vector(n, k):
25     x = signal_vec(n, k)
26     return np.squeeze(x)
27
28 def H_measure(row, col, k):
29     res = np.zeros((row, col))
30     for i in range(col):
31         rvs = sp.stats.bernoulli(1).rvs
32         S = sp.sparse.random(row, 1, density=k/row, data_rvs=rvs)
```

```
33         res[:, i] = S.toarray().reshape((row, ))
34     return res
35
36 def S_matrix(n, W):
37     S = np.zeros((2, n), dtype=complex)
38     for i in range(n):
39         S[:, i] = np.array([[1], [W**i]], dtype=complex).reshape(2, )
40     return S
41
42 def kron_product(a, b):
43     res = np.empty(0)
44     for elem in a:
45         res = np.append(res, np.dot(elem, b.T))
46     return res.T
47
48 def B_matrix(H, S, n):
49     B = np.empty((2 * H.shape[0], H.shape[1]), dtype=complex)
50     for i in range(n):
51         B[:, i] = kron_product(H[:, i], S[:, i])
52     return B
53
54 def isSingleton(y, j, n):
55     if (y[2*j].real == 0):
56         return 0
57     l_hat = cmath.phase(y[2*j+1]/y[2*j]) / (2*np.pi/n)
58     if (abs(round(l_hat) - l_hat) < 0.001):
59         return l_hat
60     return 0
61
62 def add_outliers(y_hat, y, num_outliers, med_times, var, R, W):
63     for i in range(num_outliers):
64         rand_row = random.randint(0, R-1)
65         out_val = np.random.normal(0, var_outliers)
66         y[rand_row] += out_val
67         y_hat[rand_row*2] += out_val
68         y_hat[rand_row*2+1] += out_val * W
69         y_supp = np.transpose(np.nonzero(y_hat))
70         y_supp_elem = []
71         for index in y_supp:
72             if index%2 == 0:
```

```

73             y_supp_elem.append(abs(y_hat[index][0]))
74             y_supp_elem = np.array(y_supp_elem)
75             med = np.percentile(y_supp_elem, percent)
76             for index in y_supp:
77                 if index%2 == 0:
78                     if abs(y_hat[index]) > med_times*med:
79                         y_hat[index] = 0
80                         y_hat[index+1] = 0.
81             return y_hat
82
83     def peeling_decoder():
84         total_err = 0.
85         num_success = 0
86         start_time = time.time()
87
88         for epoch in range(num_epoch):
89             x = input_vector(n, k)
90             W = cmath.exp(2 * np.pi * cmath.sqrt(-1) / n)
91             H = H_measure(R, n, density_H) # (R, n)
92             S = S_matrix(n, W)
93             B = B_matrix(H, S, n)
94             y = np.dot(H, x)
95             y = np.squeeze(y)
96             y_hat = np.dot(B, x) # y_hat = np.zeros(R*2, dtype=complex)
97             y_hat = np.squeeze(y_hat)
98             y_hat = add_outliers(y_hat, y, num_outliers, med_times,
99                                   var_outliers, R=R, W=W)
100            x_hat = np.zeros(n)
101
102            err = 1.
103            while (err > thr_err):
104                flag = False
105                for r in range(R):
106                    if isSingleton(y_hat, r, n) != 0:
107                        flag = True
108                        l_hat = round(isSingleton(y_hat, r, n))
109                        x_hat[l_hat] = y_hat[2*r].real
110                        y_hat[2*r] = 0
111                        y_hat[2*r+1] = 0
112                        H[r, l_hat] = 0

```

```

112         for j in range(R):
113             if (H[j, l_hat] > 0 and j != r):
114                 r1 = j
115                 H[j, l_hat] = 0
116                 y_hat[r1*2] = y_hat[r1*2] - x_hat[l_hat]
117                 y_hat[r1*2+1] = y_hat[r1*2+1] - x_hat[
118                     l_hat] * np.power(W, l_hat)
119             break
120         else:
121             continue
122         if (np.all(y_hat.real < 1e-5)):
123             break
124         elif(not flag):
125             print("There is no single-ton but still multi-ton in
126                   y_hat!\n")
127             break
128
129         err = LA.norm(x_hat - x) / LA.norm(x)
130         if err < thr_err:
131             num_success += 1
132
133         if (epoch+1) % 100 == 0:
134             print("time of succ in {} epochs: {}".format(epoch+1),
135                   num_success)
136
137         end_time = time.time()
138         print("The average time cost per-epoch is: ",
139               (end_time -
140                start_time)/num_epoch)

141     if __name__ == '__main__':
142         peeling_decoder()

```

A.2 对比算法代码

用于对比的 ℓ_1 -PLAD 算法文件。

```

1  # PLAD.py
2
3  import numpy as np
4  import time
5  from signal_vectors import signal_vec, outliers

```

```

6
7     n = 200
8     m = 100
9     lam = 0.006
10    k = 2
11    thr_err = 1e-2 # 1e-7
12    num_epoch = 500
13
14    def A_Gauss(m, n, mean, var):
15        Matrix = np.random.normal(mean, var**0.5, size=[m, n])
16        Matrix = np.mat(Matrix)
17        return Matrix
18
19    def SoftThreshold(b, lambd):
20        xx = np.maximum(np.abs(b) - lambd, 0)
21        return np.multiply(np.sign(b), xx)
22
23    def PLAD(lambd, iterN, m, n, A, x0, e):
24        b = A * x0 + e
25        x = np.zeros((n, 1), dtype=complex)
26        z = np.zeros((m, 1), dtype=complex)
27        AA = np.dot(A.T, A)
28        a, v = np.linalg.eig(AA)
29        L = np.max(a)
30        t = 1 / L ** 2
31        for k in range(iterN):
32            alfa = 1/L
33            x = SoftThreshold(x - alfa * A.T * z, alfa * lambd)
34            z = z + t * (A * x - b)
35            for i in range(m):
36                if z[i, 0] > t:
37                    z[i, 0] = t
38                if z[i, 0] < -t:
39                    z[i, 0] = -t
40        return x, x0, np.linalg.norm((x-x0), ord=2) / np.linalg.norm(x0,
41                           ord=2), max(np.abs(A.T*e))
42
43    if __name__ == '__main__':
44        num_succ_PLAD = 0
45        time_sum = 0.

```

```

45
46     for epoch in range(num_epoch):
47         A = A_Gauss(m, n, 0, 1 / m)
48         x0 = signal_vec(n, k)
49         e = outliers(m, 3, 100)
50
51         start_time = time.time()
52
53         _, _, E_PLAD, _ = PLAD(lam, 2000, m, n, A, x0, e)
54
55         end_time = time.time()
56         time_sum += (end_time - start_time)
57
58         if E_PLAD < thr_err:
59             num_succ_PLAD += 1
60
61         if (epoch + 1) % 100 == 0:
62             print("time of succ in {} epochs: {}".format(epoch+1),
63                   num_succ_PLAD)
64
65     print("The average time cost per-epoch is: ", time_sum/num_epoch)

```

A.3 信号初始化代码

信号及离群值的随机生成初始化文件。

```

1 # signal_vectors.py
2
3 import numpy as np
4 import scipy as sp
5
6 # k: signal vectors are k-sparse
7 def signal_vec(n, k, var=1):
8     rvs = sp.stats.norm(loc=0, scale=var).rvs
9     S = sp.sparse.random(n, 1, density=k/n, data_rvs=rvs)
10    return S.toarray()
11
12 def outliers(m, k, var=100):
13     rvs = sp.stats.norm(loc=0, scale=var).rvs
14     S = sp.sparse.random(m, 1, density=k / m, data_rvs=rvs)
15     return S.toarray()

```

致 谢

五月，没有课，没有考试。前几天去至公楼领了学士服，才意识到，我们真的要毕业了。四年前的我大概会很羡慕现在的自己：不需要在清晨起床，没有一个接一个的作业 **ddl**，毕业去向已定，拥有令人艳羡的友情、亲情、爱情。是的，我可以在 11 点起床慵懒地冲一杯咖啡，可以在疯狂星期四点一份无需外送费的企业专送肯德基，可以听着网易云日推做毕设，可以在晚上和女友蹦蹦跳跳地回寝室。我对生活、对这一切没有哪怕一点点的不满；如果这样的生活能够延续更多年，那将是莫大的恩赐。

2023 年 10 月，我带着懵懂与迷茫找到论文导师。后面的日子里，他带领我读论文，上讨论班，在我遇到瓶颈、代码卡壳时和我一起面对困难，为我提供可行的建议，给予我鼓励。再读整篇论文，半年前的我一定想不到最后能做出这样虽有些许瑕疵但令自己满意的工作。对此我衷心地感谢我的导师——李朋副教授。正是李老师的无私引领和帮助，让我接触到此前从未见识过的科研工作，我的科研素养和动手能力都大大提高。我也非常感谢同课题组的师兄师姐同学们，感谢龚瑞师姐对我的论文提供详尽的修改和建议，感谢温坤锴、高奇斌同学在代码方面提供的帮助，感谢王程正、郑翔雨师兄对我的鼓励。另外，感谢我的室友们，芮国玺、滕达、祝康乐同学，在我面对同龄人压力时带来欢乐和鼓舞。他们所有人让我在兰大的最后半年得以在一种舒适的科研氛围中结束。

在兰州的这些年，自然也会思念家乡、亲人。有人说，过得不好，才会想家。我觉得未必吧。哪怕朋友再多，事情再忙，哪怕周围的一切占据了心的领地，总有些瞬间能够让我们想起，家。譬如家乡的一则新闻，不经意与别人谈及儿时，父母的一通来电。思念，未必是脆弱的，也未必夹杂着失落与不安。大抵是不经意间觉得，回家挺好的。四年里，我的父亲庄世恩先生、母亲胡维娜女士，给予了我极大的情绪和经济上的支持。感谢他们，因为他们，我能够在学校里心无旁骛地学习、玩耍，能够在家里放下防备、生活幸福安逸。感谢其他亲人，他们对我的陪伴、支持、褒奖都是我不断前行的动力。

很多人会在评价本科生涯时制定一些标准：去了几个地方玩，交了几个朋友，拿了多少 **GPA**，获了多少奖状，毕业去向如何，发了多少文章，做了几段实习，成长了多少……一言蔽之，获得成就的数量。总有人告诉我们，绩点有多重要，一定要做科研，人脉就是一切，趁着年轻走天下。我感谢兰州大学，感谢一切告诉我这些规则的人——网友、老师、同学、家人、朋友。在这些制度的驱使下我学到了太多，专业知识，为人处事方法，都将让我受益终生。

然而，我们也因为这些别人拥有而自己没做到的焦虑、懊悔。最让我们深感遗憾的是，这类量化的数据确实是在校园与社会“横行霸道”的资本，似乎其他柔性的特点与经历没有资格代表我们。我感谢自己，几年里没有在这样的重压下成为自己憎恶的精致利己主义

者。感谢你，庄启源同学，你带领我发现无数美好，让我保持着对自己、对生活的热爱，对世界的无限期许。

这篇论文的结尾，也伴随着兰大四年生活的尾声。特别感谢我的女友，黄蔚瑄女士。从 2023 年 3 月的相识，她陪伴我走过春夏秋冬。我们共同旅行，誓要与彼此看尽万水千山；一起学习，期望在南京读研也能同频。她给我的生活带来光和温暖。与她同行，我不再是日夜与焦虑、痛苦为伴，而慢慢成长为一个乐观、爱笑、负责的青年。论文撰写期间，我们互相鼓励，驱散阴霾。我坚信，我们能一起克服所有困苦，创造我们的世界，属于我们的独一无二的当下和未来。

最后，感谢缘分，感谢命运。能够在兰大完成本科学业，我无比自豪；遇到无数善良、无私的人，我深感荣幸。我期待能将这份善良传递；对如此幸运的我来说，助人不是情分，已是本分。

庄启源
2024 年 5 月于兰州大学 30 号学生公寓

毕业论文（设计）成绩表

导师评语

该论文研究了基于稀疏矩阵测量下的带有离群值的稀疏信号重构。由于测量矩阵是稀疏的，因而存储空间很小。对于离群值，可以使用中位数算子去除大的离群值；当去除离群值之后，可以使用逐层剥离解码重构稀疏信号。使用这个两阶段方法去除离群值并且重构稀疏信号，其性能明显优于经典的惩罚最小绝对偏差(PLAD)方法。该论文设计了新的方法，高效的解决了带离群值的稀疏信号重构问题，是一篇优秀的学位论文。

建议成绩 优

指导教师（签字） 李明

答辩委员会意见

答辩通过。

答辩委员会负责人（签字） 董小强

成绩 优秀

学院（盖章）

2024年5月23日



兰州大学本科生优秀毕业论文电子版使用授权书

《基于稀疏圆码的离群值去除》是本人在兰州大学的本科毕业论文，现已通过答辩。本人作为此论文的著作权人，同意向兰州大学图书馆提交该论文的电子版和印刷本各一份。

根据《中华人民共和国著作权法》的规定，本人授权兰州大学图书馆对该论文电子版享有以下权力：（同意者画√）

1、同意提交全文。可以在

公开（半年） 延时公开（1年， 2年）

期限之后，由图书馆在校园网上提供全文浏览。

2、不同意提交电子版论文。（选此项者，须由作者本人出具不能公开的证明，导师签字，院系所加盖公章。否则，电子版论文正常提交。）

图书馆承诺：

1、不对论文从事收集、保存、发布以外的其他活动；

2、未经著作权人同意，不得从事营利性活动。

院系：数学与统计学院

学号：320200942581

作者（授权人）签名：庄启源

时间：2024年6月3日

被授权人：兰州大学图书馆