Randomly Projected Convex Clustering Model: Motivation, Realization, and Cluster Recovery Guarantees

Ziwen Wang

ZWWANG@MATH.CUHK.EDU.HK

Department of Mathematics The Chinese University of Hong Kong Hong Kong

Yancheng Yuan*

YANCHENG.YUAN@POLYU.EDU.HK

Department of Applied Mathematics The Hong Kong Polytechnic University Hong Kong

Jiaming Ma

22051002R@CONNECT.POLYU.HK

Department of Applied Mathematics The Hong Kong Polytechnic University Hong Kong

Tieyong Zeng

ZENG@MATH.CUHK.EDU.HK

Department of Mathematics
The Chinese University of Hong Kong
Hong Kong

Defeng Sun

DEFENG.SUN@POLYU.EDU.HK

Department of Applied Mathematics The Hong Kong Polytechnic University Hong Kong

Abstract

In this paper, we propose a randomly projected convex clustering model for clustering a collection of n high dimensional data points in \mathbb{R}^d with K hidden clusters. Compared to the convex clustering model for clustering original data with dimension d, we prove that, under some mild conditions, the perfect recovery of the cluster membership assignments of the convex clustering model, if exists, can be preserved by the randomly projected convex clustering model with embedding dimension $m = O(\epsilon^{-2}\log(n))$, where $0 < \epsilon < 1$ is some given parameter. We further prove that the embedding dimension can be improved to be $O(\epsilon^{-2}\log(K))$, which is independent of the number of data points. Extensive numerical experiment results will be presented in this paper to demonstrate the robustness and superior performance of the randomly projected convex clustering model. The numerical results presented in this paper also demonstrate that the randomly projected convex clustering model can outperform the randomly projected K-means model in practice.

Keywords: convex clustering, Johnson-Lindenstrauss lemma, unsupervised learning.

1. Introduction

Clustering is a fundamental and important problem in data science. Among many others, K-means is arguably the most popular model. It has been widely known that K-means may

^{*.} Corresponding author.

suffer from the nonconvexity of the model and is very sensitive to the initialization. More critically, K-means requires the number of clusters as a prior, which is not practical in many applications. Recently, researchers have proposed the convex clustering model, which aims to overcome the aforementioned challenges (Pelckmans et al., 2005; Hocking et al., 2011; Lindsten et al., 2011).

Given a collection of n data points with d features $A = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\} \subseteq \mathbb{R}^d$, the general weighted convex clustering model (CCM) solves the following convex optimization problem

$$\min_{x_1,...,x_n \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{a}_i\|^2 + \gamma \sum_{i < j} w_{ij} \|\mathbf{x}_i - \mathbf{x}_j\|_q,$$
 (CCM)

where $w_{ij} = w_{ji} \ge 0$ are given weights depending on the input data $A, \gamma > 0$ is a tuning parameter which controls the strength of the fusion penalty, and $\|\cdot\|_q$ is the vector q-norm $(q \ge 1)$. In this paper, we focus on the convex clustering model with q = 2. We denote $\|\cdot\|$ as the vector 2-norm. One choice of the weights is setting $w_{ij} = 1$ for all $1 \le i < j \le n$, and the resulting model is usually called the convex clustering model with uniform weights. In practice, the following k-nearest neighbors-based weights are popular due to their robustness and computational efficiency:

$$w_{ij} = \begin{cases} \exp(-\phi \|\mathbf{a}_i - \mathbf{a}_j\|^2), & \text{if } (i,j) \in \mathcal{E}(k), \\ 0, & \text{otherwise,} \end{cases}$$
 (1)

here, $\mathcal{E}(k) := \{(i, j) \mid \text{ if } \mathbf{a}_i \text{ (or } \mathbf{a}_j) \text{ is in } \mathbf{a}_j\text{'s (or } \mathbf{a}_i\text{'s) k-nearest neighbors}, 1 \le i \ne j \le n \}.$

Extensive investigation has been conducted for the convex clustering model in recent years and impressive progress has been achieved from the perspectives of both the recovery properties and efficient numerical algorithms. From the theoretical understanding perspective, some deterministic and statistical cluster recovery guarantees have been established (Zhu et al., 2014; Tan and Witten, 2015; Panahi et al., 2017; Radchenko and Mukherjee, 2017; Chiquet et al., 2017; Chi and Steinerberger, 2019; Sun et al., 2021; Chi et al., 2020; Jiang et al., 2020; Dunlap and Mourrat, 2022). More specifically, under some mild conditions, there exists a nonempty interval of the tuning parameter γ such that the convex clustering model can perfectly recover the cluster membership of the data (Panahi et al., 2017; Sun et al., 2021). From the perspective of optimization algorithms, impressive progress has been achieved in solving the convex clustering model with a large number of data points but with moderate feature dimensions (say with $d \leq 100$ in (CCM)). Along this direction, Chi and Lange (2015) adopted the alternating direction method of multipliers (ADMM) and proposed an alternating minimization algorithm (AMA). Later, Yuan et al. (2018) designed a semismooth Newton based augmented Lagrangian (SSNAL) method that can solve the convex clustering model efficiently with high accuracy. More recently, by taking the advantage of the structured sparsity of the convex clustering model, Yuan et al. (2022) proposed dimension reduction techniques (in the sense of the number of data points) called adaptive sieving (AS) and enhanced adaptive sieving (EAS), which further accelerate SSNAL (and other algorithms). Consequently, the existing algorithms can be scalable with respect to the number of data points. However, it is still very challenging to solve the convex clustering model when the dimension of the data features is high (i.e. d is large in (CCM)).

In this paper, we will design a dimension reduction technique for overcoming the computational challenges of the convex clustering model for clustering high dimensional data. Our

approach is inspired by the Johnson-Lindenstrauss (JL) lemma (Johnson and Lindenstrauss, 1984) and the fact that the recovery guarantees of the convex clustering model mainly depend on the pair-wise distances among the data points and centroids. In particular, we will propose a randomly projected (weighted) convex clustering model which clusters the data with a much smaller dimension obtained by applying a random projection mapping to the input data. Among other advantages, we want to mention that random projection is a computationally efficient approach to obtaining the embedding data. Importantly, we will prove that the randomly projected convex clustering model will preserve the recovery guarantees of the original convex clustering model. In other words, if there exists a nonempty interval of the parameter γ such that the convex clustering model (CCM) perfectly recovers the cluster memberships of the input data, so will be the randomly projected model in high probability. This is a very interesting and inspiring result since we can obtain the clustering results of the original high dimensional data by solving a more tractable randomly projected convex clustering model with much smaller dimensions. Moreover, we will establish the cluster recovery guarantees for the randomly projected convex clustering model where the embedding dimension can be independent of the number of data points. Extensive numerical experiment results will be presented in this paper to justify the theoretical guarantees and to demonstrate the superior performance and robustness of the proposed model. To further demonstrate the superior performance of the randomly projected convex clustering model, we also compare its performance to the randomly projected K-means model (Cohen et al., 2015; Makarychev et al., 2022).

We summarize the main contributions of this paper as follows:

- 1. We propose a randomly projected convex clustering model which is much more computationally tractable than the convex clustering model (CCM).
- 2. We establish the recovery guarantees of the randomly projected convex clustering model under mild conditions. We further prove that the embedding dimension can be independent of the number of data points.
- 3. We conduct extensive numerical experiments to justify the established theoretical guarantees and demonstrate the superior performance of the proposed randomly projected convex clustering model.

The rest of the paper is organized as follows: In Section 2, we introduce some concepts and notation and then review some necessary preliminary results of the recovery guarantees of the convex clustering model and the JL lemma. In Section 3, we will propose a randomly projected convex clustering model and prove its theoretical recovery guarantees. We will then present the numerical results in Section 4. We will conclude the paper and include some discussion of future research directions in Section 5.

2. Preliminaries

In this section, we first introduce some commonly used notation and then introduce some results about the convex clustering model and the Johnson-Lindenstrauss lemma.

2.1 Problem Settings

In this paper, we focus on the following problem setting.

General problem setting: Cluster a collection of n given data points $A = \{\mathbf{a}_1, \dots, \mathbf{a}_n\} \subseteq \mathbb{R}^d$ with a hidden clustering partition $\mathcal{V} = \{V_1, V_2, \dots, V_K\}$.

We define some notation in Table 1, which will be commonly used later in this paper.

Table 1: Some commonly used notation. In this table, we assume by default that $1 \le \alpha \ne \beta \le K$.

Notation	Definition
I_{α}	$\{i \mid \mathbf{a}_i \in V_{\alpha}\}$
n_{lpha}	cardinality of I_{α}
[m] for a given integer $m > 0$	$[m] := \{1, 2, \dots, m\}$
$\mathbf{a}^{(\alpha)}$	$\frac{1}{n_{\alpha}}\sum_{i\in I_{\alpha}}\mathbf{a}_{i}$
$\mathbf{a}^{(0)}$	$\frac{\frac{1}{n}\sum_{i=1}^{n}\mathbf{a}_{i}}{n}$
$w^{(lpha,eta)}$	$\sum_{i \in I_{\alpha}} \sum_{j \in I_{\beta}} w_{ij}$
$ar{w}^{(eta)}$	$\frac{1}{n_{\beta}} \sum_{1 \le l \le K, l \ne \beta} w^{(\beta, l)}$
$w_i^{(eta)} (i \in [n])$	$\sum_{j\in I_{\beta}} w_{ij}$
$\mu_{ij}^{(\alpha)} (i, j \in I_{\alpha})$	$\sum_{\beta=1,\beta\neq\alpha}^{K} \left w_i^{(\beta)} - w_j^{(\beta)} \right $
$C(n,k) \ (1 \le k \le n)$	$\frac{n!}{k!(n-k)!}$

Following the settings in (Sun et al., 2021), we assume the following assumptions hold throughout this paper.

Assumption 1 In the general problem setting, the mean vector $\mathbf{a}^{(0)}$ and the centroids $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(K)}$ are all distinct.

Assumption 2 The specified weights w_{ij} in the model (CCM) satisfy

$$w_{ij} > 0$$
 and $n_{\alpha}w_{ij} > \mu_{ij}^{(\alpha)}, \quad \forall i, j \in I_{\alpha}, 1 \le \alpha \le K.$ (2)

A quick comment is that Assumption 2 holds automatically for uniform weights. The next definition will be useful for the discussion of the convex clustering model.

Definition 1 We say that a map $\psi : \mathbb{R}^d \to \mathbb{R}^{\bar{d}}$ perfectly recovers \mathcal{V} on the data A if $\psi(\mathbf{a}_i) = \psi(\mathbf{a}_j)$ is equivalent to \mathbf{a}_i and \mathbf{a}_j belonging to the same V_{α} for some $1 \leq \alpha \leq K$. We call a partition $\mathcal{W} = \{W_1, \ldots, W_L\}$ of A a coarsening of \mathcal{V} if there exists a partition $\{\alpha_1, \ldots, \alpha_L\}$ of [K] such that $W_l = \bigcup_{i \in \alpha_l} V_i$ for all $1 \leq l \leq L$. We call \mathcal{W} a non-trivial coarsening of \mathcal{V} if L > 1.

2.2 Recovery guarantees for convex clustering model (CCM)

In this section, we review the recovery guarantees of the weighted convex clustering model.

Theorem 2 ((Sun et al., 2021, Theorem 5)) In the general problem setting, denote the optimal solution of the convex clustering model (CCM) with some given parameter $\gamma \geq 0$ by $\{\mathbf{x}_i^*(\gamma)\}_{i=1}^n$ and define the map $\phi_{\gamma}(\mathbf{a}_i) = \mathbf{x}_i^*(\gamma)$ for $i = 1, \ldots, n$. Define

$$\gamma_{\min} := \max_{1 \leq \alpha \leq K} \max_{i,j \in I_{\alpha}} \left\{ \frac{\|\mathbf{a}_{i} - \mathbf{a}_{j}\|}{n_{\alpha} w_{ij} - \mu_{ij}^{(\alpha)}} \right\}, \ \gamma_{\max} := \min_{1 \leq \alpha < \beta \leq K} \left\{ \frac{\|\mathbf{a}^{(\alpha)} - \mathbf{a}^{(\beta)}\|}{\bar{w}^{(\alpha)} + \bar{w}^{(\beta)}} \right\},$$

$$\gamma_{\max 2} := \max_{1 \leq \alpha \leq K} \frac{\|\bar{\mathbf{a}} - \mathbf{a}^{(\alpha)}\|}{\bar{w}^{(\alpha)}}, \ r := \frac{\gamma_{\max}}{\gamma_{\min}}, \ r_{2} := \frac{\gamma_{\max} 2}{\gamma_{\min}}.$$

$$(3)$$

Under Assumption 1 and Assumption 2, we have

- 1. If r > 1 and $\gamma \in [\gamma_{\min}, \gamma_{\max})$, then the map ϕ_{γ} perfectly recovers \mathcal{V} .
- 2. If $r_2 > 1$ and $\gamma \in [\gamma_{\min}, \gamma_{\max 2})$, then the map ϕ_{γ} recovers a non-trivial coarsening of \mathcal{V} .

2.3 Johnson-Lindenstrauss Lemma and the Random Projection

In this section, we introduce the Johnson-Lindenstrauss (JL) lemma, which is a key tool for this paper. Consider a collection of high-dimensional data points $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$, the JL lemma shows the existence of a mapping $f: X \to \mathbb{R}^m$ such that for all points $\mathbf{x}_i \neq \mathbf{x}_j \in X$, $\|\mathbf{x}_i - \mathbf{x}_j\|^2$ are approximately maintained in a m dimensional space within a distortion tolerance $\epsilon \in (0,1)$. More surprisingly, the required embedding dimension $m = O(\epsilon^{-2} \log(n))$ is independent of d.

Theorem 3 (JL lemma (Johnson and Lindenstrauss, 1984, Lemma 1)) For any given collection of n data points $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$ and any $\epsilon \in (0,1)$, there exists an ϵ - isometry embedding $f : \mathbb{R}^d \to \mathbb{R}^m$ with $m = O\left(\min\{d, \epsilon^{-2} \log n\}\right)$. In other words, $\forall \mathbf{x}_i, \mathbf{x}_j \in X$,

$$(1 - \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|^2 \le \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|^2 \le (1 + \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|^2.$$
 (4)

We also call a mapping f satisfies (4) an ϵ -JL Transform (or ϵ -JLT in short) on X. The mapping f can be found in randomized polynomial time (Dasgupta and Gupta, 2003). Moreover, if the mapping f must be linear, then $m = \Omega\left(\min\{d, \epsilon^{-2} \log n\}\right)$ is optimal (Larsen and Nelson, 2016). The following Distributional Johnson-Lindenstrauss (DJL) lemma is useful.

Theorem 4 (DJL lemma) For any $\epsilon \in (0,1)$, $\delta \in (0,1/2)$ and integer d > 1, there exists a distribution $\mathcal{D}_{\epsilon,\delta}$ over matrices $\Pi \in \mathbb{R}^{m \times d}$ for $m = O\left(\epsilon^{-2}\log(1/\delta)\right)$ such that for any $z \in \mathbb{R}^d$ with ||z|| = 1,

$$\mathbb{P}_{\Pi \sim \mathcal{D}_{\epsilon, \delta}} \left[\left| \|\Pi z\|^2 - 1 \right| > \epsilon \right] < \delta. \tag{5}$$

We call a distribution $\mathcal{D}_{\epsilon,\delta}$ that satisfies (5) a DJL distribution.

For later convenience, we include the following proposition, which is a direct consequence of Theorem 4 and the union bound in the probability theory.

Proposition 5 (Random projection for multiple sets) Assume that there are l nonempty collections of data points X_1, \ldots, X_l in \mathbb{R}^d with $|X_j| = n_j$ $(1 \leq j \leq l)$. Denote $X = \bigcup_{j=1}^l X_j$. Given any $0 < \epsilon < 1$ and $0 < \delta < \frac{1}{\sum_{j=1}^l n_j}$, and let $D_{\epsilon,\delta}$ be a DJL distribution over $\mathbb{R}^{m \times d}$ with $m = O(\epsilon^{-2} \log(1/\delta))$. We have

$$\mathbb{P}_{\Pi \sim \mathcal{D}_{\epsilon, \delta}} \left[(1 - \epsilon) \|\mathbf{x}\|^2 \le \|\Pi \mathbf{x}\|^2 \le (1 + \epsilon) \|\mathbf{x}\|^2, \ \forall \mathbf{x} \in X \right] \ge 1 - (\sum_{j=1}^l n_j) \delta > 0.$$
 (6)

Thus, there exists a matrix $\Pi \in \mathbb{R}^{m \times d}$ such that

$$(1 - \epsilon) \|\mathbf{x}\|^2 \le \|\Pi\mathbf{x}\|^2 \le (1 + \epsilon) \|\mathbf{x}\|^2, \ \forall \mathbf{x} \in X.$$

3. A Randomly Projected Convex Clustering Model

The convex clustering model (CCM) has promising recovery guarantees, but solving the model can be computationally challenging, especially when the feature dimension d is high. In this section, we will propose a randomly projected convex clustering model of (CCM) with much smaller feature dimensions. We will prove that the recovery guarantees will be preserved with a high probability for the random projected convex clustering model. More specifically, for the given collection of data points $A = \{\mathbf{a}_1, \dots, \mathbf{a}_n\} \subseteq \mathbb{R}^d$ considered in the general problem setting and a given $\epsilon \in (0,1)$, we will construct an ϵ -isometry mapping $f: \mathbb{R}^d \to \mathbb{R}^m$ with $m = O\left(\min\left\{d, \log(n)/\epsilon^2\right\}\right)$, where m can be much smaller than d. We solve the following projected convex clustering model

$$\min_{\hat{X} \in \mathbb{R}^{m \times n}} \frac{1}{2} \sum_{i=1}^{n} \|\hat{\mathbf{x}}_i - f(\mathbf{a}_i)\|^2 + \gamma \sum_{i < j} w_{ij} \|\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_j\|.$$
 (RPCCM)

In this paper, we will choose f as a random projection matrix motivated by the DJL lemma and call the corresponding model (RPCCM) a randomly projected convex clustering model.

3.1 An ϵ -isometry Mapping for the Convex Clustering Model

A key observation is that the recovery guarantees of the convex clustering (e.g. Theorem 2) mainly depend on the distances between data points within the same cluster and the distance between the centroids of different clusters. Thus, the recovery guarantees of the convex clustering model (CCM) can be inherited by the model (RPCCM) if we can construct an ϵ -isometry mapping for some small enough $\epsilon > 0$ for the data points A and the corresponding centroids. The next proposition shows the existence of a desired ϵ -isometry mapping for the convex clustering model.

Proposition 6 For the general problem setting, define $X_{\alpha} := \{\mathbf{a}_i - \mathbf{a}_j \mid i, j \in I_{\alpha}, i < j\}$ $(1 \leq \alpha \leq K)$, and $X_c := \{\mathbf{a}^{(\alpha)} - \mathbf{a}^{(\beta)} \mid 0 \leq \alpha < \beta \leq K\}$. Denote $N_1 = \sum_{\alpha=1}^K |X_{\alpha}| = \sum_{\alpha=1}^K C(n_{\alpha}, 2) < C(n, 2)$, and $N_2 = |X_c| = C(K+1, 2)$. For any $0 < \epsilon < 1$, let $\delta = \frac{1}{(N_1 + N_2)^p}$, where p > 1, and let $D_{\epsilon, \delta}$ be a DJL distribution over $\mathbb{R}^{m \times d}$, where $m = O(\epsilon^{-2} \log(1/\delta)) = O(p\epsilon^{-2} \log(N_1 + N_2))$. Then for any $\Pi \in \mathbb{R}^{m \times d}$ randomly drawn from $D_{\epsilon, \delta}$, with probability at least $1 - \frac{1}{(N_1 + N_2)^{p-1}}$ that

$$(1 - \epsilon) \|\mathbf{a}_i - \mathbf{a}_j\|^2 \le \|\Pi(\mathbf{a}_i - \mathbf{a}_j)\|^2 \le (1 + \epsilon) \|\mathbf{a}_i - \mathbf{a}_j\|^2, \ \mathbf{a}_i, \mathbf{a}_j \in V_\alpha, 0 \le \alpha \le K,$$
 (7a)

$$(1 - \epsilon) \|\mathbf{a}^{(\alpha)} - \mathbf{a}^{(\beta)}\|^2 \le \|\Pi(\mathbf{a}^{(\alpha)} - \mathbf{a}^{(\beta)})\|^2 \le (1 + \epsilon) \|\mathbf{a}^{(\alpha)} - \mathbf{a}^{(\beta)}\|, \ 1 \le \alpha \ne \beta \le K.$$
 (7b)

If $N_2 \leq n/2$, then it is enough to take $\delta = \frac{2}{n^{p+1}}$ and $m = O(\epsilon^{-2}(p+1)\log(n))$, and the inequalities (7) hold with probability at least $1 - \frac{1}{n^{p-1}}$.

The above proposition can be proved as a consequence of Proposition 5. In practice, only the pair-wise distance between the input data points can be checked after a projection

matrix Π is randomly sampled (which covers the condition (7a)). But an insight is K should be much much smaller than n (which is the reason for us to do clustering). The next corollary shows that the condition (7b) can be satisfied in much higher probability if (7a) holds.

Proposition 7 Let $\Pi \in \mathbb{R}^{m \times d}$ be a projection matrix sampled from a DJL distribution $D_{\epsilon,\delta}$ with $m = O(\epsilon^{-2} \log(1/\delta))$ and $\delta = \frac{1}{(N_1 + N_2)^p}$. Let E_1 be the event that Π satisfies (7a) and E_2 be the event that Π satisfies (7b), respectively. Then, the conditional probability $\mathbb{P}[E_2 \mid E_1]$ satisfies

$$\mathbb{P}\left[E_2 \mid E_1\right] \ge 1 - \frac{N_2}{(N_1 + N_2)^p - N_1}.$$
 (8)

If we further assume that $N_2 \leq n/2$ and $\delta = \frac{2}{n^{p+1}}$, then

$$\mathbb{P}\left[E_2 \mid E_1\right] \ge 1 - \frac{1}{n^p - n + 1}.\tag{9}$$

Proof Direct calculation gives that

$$\mathbb{P}[E_2 \mid E_1] = 1 - \mathbb{P}[E_2^c \mid E_1]
= 1 - \frac{\mathbb{P}[E_1 \cap E_2^c]}{\mathbb{P}[E_1]}
\ge 1 - \frac{\mathbb{P}[E_2^c]}{\mathbb{P}[E_1]}
\ge 1 - \frac{N_2 \delta}{1 - N_1 \delta}
= 1 - \frac{N_2}{(N_1 + N_2)^p - N_1}.$$

The inequality (9) can be proved similarly.

Remark 8 The DJL distribution plays a role in the construction of the ϵ -isometry mapping. Indeed, a vast amount of variants of the DJL lemma have been explored by designing the structure of DJL distributions, including the subgaussians (Indyk and Motwani, 1998; Achlioptas, 2003; Matoušek, 2008), the Fast JL Transform (Ailon and Chazelle, 2009; Ailon and Liberty, 2009, 2013), and the Sparse JL Transform (Dasgupta et al., 2010; Kane and Nelson, 2010, 2014; Cohen et al., 2018). The particular choice of the DJL distribution is beyond the concern of this paper. In this paper, we will follow (Matoušek, 2008) and consider the linear random projection matrix $\Pi = \frac{1}{\sqrt{m}}R \in \mathbb{R}^{m \times d}$, where R_{ij} are independent random variables with zero mean and a uniform subgaussian tail.

3.2 Cluster Recovery Guarantees of the Randomly Projected Convex Clustering Model for the General Problem Setting

Next, we will establish the cluster recovery guarantees of the randomly projected convex clustering model (RPCCM) for the general problem setting. For later convenience, we introduce some useful notation.

Definition 9 In the general problem setting, we consider the randomly projected convex clustering model (RPCCM) with some specified weights $w_{ij} = w_{ji} \ge 0$ ($1 \le i \ne j \le n$) and a randomly sampled projection matrix $\Pi \in \mathbb{R}^{m \times d}$ (for some $m \ge 1$). Without explicitly mentioning the dependence on Π , we define

$$\hat{\gamma}_{\min} := \max_{1 \leq \alpha \leq K} \max_{i,j \in I_{\alpha}} \left\{ \frac{\|\Pi(\mathbf{a}_{i} - \mathbf{a}_{j})\|}{n_{\alpha}w_{ij} - \mu_{ij}^{(\alpha)}} \right\}, \quad \hat{\gamma}_{\max} := \min_{1 \leq \alpha < \beta \leq K} \left\{ \frac{\|\Pi(\mathbf{a}^{(\alpha)} - \mathbf{a}^{(\beta)})\|}{\bar{w}^{(\alpha)} + \bar{w}^{(\beta)}} \right\},$$

$$\hat{\gamma}_{\max 2} := \max_{1 \leq \alpha \leq K} \frac{\|\Pi(\mathbf{a}^{(0)} - \mathbf{a}^{(\alpha)})\|}{\bar{w}^{(\alpha)}}.$$
(10)

The next theorem shows that the recovery properties of the original convex clustering model can be preserved by the randomly projected convex clustering model in high probability. For convenience, we assume the following assumption holds for the rest of this paper.

Assumption 3 The inequality n > K(K+1) holds, where n is the number of data points and K is the number of hidden clusters.

The above assumption is mild and it is consistent with the purpose of the clustering task.

Theorem 10 Consider the general problem setting and the models (CCM) and (RPCCM) with the same specified weights w_{ij} . For any $0 < \epsilon < 1$, let $\delta = \frac{2}{n^p}$ with some p > 2, and let $D_{\epsilon,\delta}$ be a DJL distribution over $\mathbb{R}^{m \times d}$ with $m = O(p\epsilon^{-2}\log(n))$. Here and below in this theorem, the notation $O(\cdot)$ depends on the same absolute constant. Without loss of generality, we assume that m < d (or equivalently, we can assume $\sqrt{\frac{O(p\log(n))}{d}} < 1$ and $\epsilon \in (\sqrt{\frac{O(p\log(n))}{d}}, 1)$). Define

$$\epsilon_{\min} = \sqrt{\frac{O(p\log(n))}{d}}, \ \epsilon_{\sup} = \frac{r^2 - 1}{r^2 + 1}, \ and \ \epsilon_{\sup 2} = \frac{r_2^2 - 1}{r_2^2 + 1},$$
(11)

where r and r_2 are the constants defined in (3). Let $\Pi \in \mathbb{R}^{m \times d}$ be a random projection matrix drawn from $D_{\epsilon,\delta}$. Denote the optimal solution of the model (RPCCM) with Π and $\gamma \geq 0$ by $\{\hat{\mathbf{x}}_i^*(\gamma)\}_{i=1}^n$ and define the map $\hat{\phi}_{\gamma}(\mathbf{a}_i) = \hat{\mathbf{x}}_i^*$. Then, we have

- 1. If $r > \sqrt{\frac{1+\epsilon_{\min}}{1-\epsilon_{\min}}}$, then $\epsilon_{\min} < \epsilon_{\sup}$. For any $\epsilon \in [\epsilon_{\min}, \epsilon_{\sup})$, and $\hat{\gamma} \in [\sqrt{1+\epsilon}\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max})$, with probability over $1 \frac{1}{n^{p-2}}$, the map $\hat{\phi}_{\hat{\gamma}}$ perfectly recovers \mathcal{V} .
- 2. If $r_2 > \sqrt{\frac{1+\epsilon_{\min}}{1-\epsilon_{\min}}}$, then $\epsilon_{\min} < \epsilon_{\sup 2}$. For any $\epsilon \in [\epsilon_{\min}, \epsilon_{\sup 2})$, and $\hat{\gamma} \in [\sqrt{1+\epsilon}\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max 2})$, with probability over $1 \frac{1}{n^{p-2}}$, the map $\hat{\phi}_{\hat{\gamma}}$ recovers a non-trivial coarsening of \mathcal{V} .

Proof It directly follows Proposition 6 that, with probability over $1 - \frac{1}{n^{p-2}}$, the following statements hold:

(i) The centroids $\{\Pi \mathbf{a}^{(0)}, \dots, \Pi \mathbf{a}^{(K)}\}$ of the embedded data are distinct.

(ii) The parameters $\hat{\gamma}_{\min}$, $\hat{\gamma}_{\max}$, and $\hat{\gamma}_{\max 2}$ defined in (10) satisfy the following inequalities:

$$\begin{split} \sqrt{1-\epsilon}\gamma_{\min} & \leq \hat{\gamma}_{\min} \leq \sqrt{1+\epsilon}\gamma_{\min}, \quad \sqrt{1-\epsilon}\gamma_{\max} \leq \hat{\gamma}_{\max} \leq \sqrt{1+\epsilon}\gamma_{\max}, \\ \sqrt{1-\epsilon}\gamma_{\max 2} & \leq \hat{\gamma}_{\max 2} \leq \sqrt{1+\epsilon}\gamma_{\max 2}. \end{split}$$

The above implies that

$$\begin{bmatrix}
\sqrt{1+\epsilon}\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max} &) \subseteq [\hat{\gamma}_{\min}, \hat{\gamma}_{\max}), \\
\sqrt{1+\epsilon}\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max} &) \subseteq [\hat{\gamma}_{\min}, \hat{\gamma}_{\max} &).
\end{bmatrix}$$
(12)

Now, we prove the first part of the theorem. We claim here that it is sufficient to show: if $r > \sqrt{\frac{1+\epsilon_{\min}}{1-\epsilon_{\min}}}$, then $\epsilon_{\min} < \epsilon_{\sup}$, and for any $\epsilon \in [\epsilon_{\min}, \epsilon_{\sup})$, $[\sqrt{1+\epsilon}\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max})$ is nonempty. In fact, if $[\sqrt{1+\epsilon}\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max})$ is nonempty, then by the first inclusion of (12), $[\hat{\gamma}_{\min}, \hat{\gamma}_{\max})$ is nonempty. Applying Theorem 2 to the embbedded data ΠA implies that for any $\hat{\gamma} \in [\hat{\gamma}_{\min}, \hat{\gamma}_{\max})$, the map $\hat{\phi}_{\hat{\gamma}}$ perfectly recovers \mathcal{V} .

On the one hand, we have

$$r > \sqrt{\frac{1+\epsilon_{\min}}{1-\epsilon_{\min}}} \implies (1-\epsilon_{\min})r^2 > 1+\epsilon_{\min}$$

 $\implies (r^2-1) > \epsilon_{\min}(r^2+1)$
 $\implies \epsilon_{\min} < \epsilon_{\sup}.$

This implies that the interval $[\epsilon_{\min}, \epsilon_{\sup})$ is nonempty. On the other hand, we have

$$\epsilon < \epsilon_{\text{sup}} \implies \epsilon < \frac{r^2 - 1}{r^2 + 1}
\implies \frac{1 + \epsilon}{1 - \epsilon} < r^2
\implies \frac{1 + \epsilon}{1 - \epsilon} < \frac{\gamma_{\text{max}}^2}{\gamma_{\text{min}}^2}
\implies \sqrt{1 + \epsilon} \gamma_{\text{min}} < \sqrt{1 - \epsilon} \gamma_{\text{max}}.$$

Thus, we have proved the first part of the theorem. The second part of the theorem can be proved in a similar way.

We can obtain an ϵ -isometry mapping in randomized polynomial time (Dasgupta and Gupta, 2003) and we can check the ϵ -isometry of the mapping on the data A (but not for the centroids). The following corollary is useful. The proof of the corollary follows directly from Theorem 10 and Proposition 7.

Corollary 11 Let $\Pi \in \mathbb{R}^{m \times d}$ be a random projection matrix drawn from $D_{\epsilon,\delta}$ as in Theorem 10. If we further assume that Π satisfies (7a), then, the statements of Theorem 10 hold with probability at least $1 - \frac{1}{n^{p-1}-n+1}$ under the same assumptions.

Remark 12 We want to make some remarks on the obtained recovery guarantees of the model (RPCCM).

1. For convenience, we assumed $K(K+1) \leq n$. But we can also easily obtain recovery guarantees regarding N_1 and N_2 as defined in Proposition 6.

- 2. The embedding dimension is $m = O(p\epsilon^{-2}\log n)$, which only depends on ϵ , n, and p, but it is independent of the data dimension d. Also, the embedding dimension grows very slowly with respect to n (in $O(\log(n))$).
- 3. We derives the lower bound ϵ_{\min} and the upper bound ϵ_{\sup} of ϵ for perfect recovery of the model (RPCCM). In particular, the lower bound ϵ_{\min} can be very small for high dimensional data. The upper bound ϵ_{\sup} depends on the ratio of γ_{\max} and γ_{\min} , and it is independent of the scale of the data.
- 4. We want to mention that, the weights used in the model (CCM) and the model (RPCCM) are the same.
- 5. The dimension reduction based on the JL lemma has been also investigated for the K-means model (Cohen et al., 2015). However, for the K-means model, only the cost (the optimal objective function value of the K-means model) can be preserved up to a tolerance ε > 0. Here, we proved that the perfect recovery guarantee of the convex clustering model can be inherited. A comparison of the empirical performance between the randomly projected K-means model and the randomly projected convex clustering model can be found later in the numerical experiments.

The embedding dimension m in Theorem 10 depends on n of the order $O(\log(n))$. Next, we will further improves it from $O(\log(n))$ to $O(\log(K))$. The key insights come from the estimate of the spectral norm of the random matrices. The following lemma is useful, which is a direct consequence of Theorem 4.6.1 and Lemma 3.4.2 in (Vershynin, 2018).

Lemma 13 (Two-sided bound on sub-gaussian matrices) Let $\Pi = \frac{1}{\sqrt{m}}R \in \mathbb{R}^{m \times d}$ $(m \leq d)$, where R_{ij} are independent random variables with $\mathbb{E}[R_{ij}] = 0$, $\operatorname{Var}[R_{ij}] = 1$ and the subgaussian norm $\kappa := \|R_{ij}\|_{\psi_2} := \inf\{s > 0 : \mathbb{E}[\exp(R_{ij}^2/s^2)] \leq 2\}$. Let $s_1(\Pi)$ be the largest singular value of Π , and let $s_m(\Pi)$ be the smallest singular value of Π . Then for any $t \geq 0$, we have

$$\underline{S}(m,d,t) \le s_m(\Pi) \le s_1(\Pi) \le \bar{S}(m,d,t) \tag{13}$$

with probability at least $1-2\exp\left(-t^2\right)$. Here, $C_{\kappa}^2>0$ is a constant that only depends on κ , and $\bar{S}(m,d,t)=\frac{\sqrt{d}+C_{\kappa}^2t}{\sqrt{m}}+C_{\kappa}^2$, and $\underline{S}(m,d,t)=\frac{\sqrt{d}-C_{\kappa}^2t}{\sqrt{m}}-C_{\kappa}^2$.

The next theorem shows that the embedding dimension can be independent of the number of data points n.

Theorem 14 Consider the general problem setting and the models (CCM) and (RPCCM) with the same specified weights w_{ij} . For any $0 < \epsilon < 1$, let $\Pi = \frac{1}{\sqrt{m}}R \in \mathbb{R}^{m \times d}$ with $m = O(p\epsilon^{-2}\log(K))$, where R_{ij} are independent random variables with $\mathbb{E}\left[R_{ij}\right] = 0$, $\operatorname{Var}\left[R_{ij}\right] = 1$, and with the subgaussian norm $\kappa = \|R_{i,j}\|_{\psi_2}$. Here and below in this theorem, the notation $O(\cdot)$ depends on the same absolute constant. Without loss of generality, we assume that m < d (or equivalently, we can assume that $\sqrt{\frac{O(p\log(K))}{d}} < 1$ and $\epsilon \in (\sqrt{\frac{O(p\log(K))}{d}}, 1)$). Let $\Pi \in \mathbb{R}^{m \times d}$ be a random projection matrix drawn from $D_{\epsilon,\delta}$. By Lemma 13, there exists

some constant $C_{\kappa}^2 > 0$ that only depends on κ such that the inequality (13) holds for any $t \geq 0$ with probability over $1 - 2 \exp(-t^2)$. Define

$$C_{0} = \frac{\sqrt{O(p \log K)}}{\sqrt{d} + C_{\kappa}^{2} t}, \ \tilde{\epsilon}_{\sup} = r C_{0} \sqrt{\frac{r^{2} C_{0}^{2}}{4} + C_{\kappa}^{2} C_{0} + 1} - C_{\kappa}^{2} C_{0} - \frac{r^{2} C_{0}^{2}}{2},$$

$$\tilde{\epsilon}_{\sup 2} = r_{2} C_{0} \sqrt{\frac{r_{2}^{2} C_{0}^{2}}{4} + C_{\kappa}^{2} C_{0} + 1} - C_{\kappa}^{2} C_{0} - \frac{r_{2}^{2} C_{0}^{2}}{2}, \ \tilde{\epsilon}_{\min} = \sqrt{\frac{O(p \log K)}{d}}.$$
(14)

Denote the optimal solution of the model (RPCCM) with Π and $\gamma \geq 0$ by $\{\hat{\mathbf{x}}_i^*(\gamma)\}_{i=1}^n$ and define the map $\hat{\phi}_{\gamma}(\mathbf{a}_i) = \hat{\mathbf{x}}_i^*$. Then, we have:

- 1. If $r > \frac{1 + C_{\kappa}^2 + \frac{C_{\kappa}^2 t}{\sqrt{d}}}{\sqrt{1 \tilde{\epsilon}_{\min}}}$, then $\tilde{\epsilon}_{\min} < \tilde{\epsilon}_{\sup}$. For any $\epsilon \in [\tilde{\epsilon}_{\min}, \tilde{\epsilon}_{\sup})$, and $\hat{\gamma} \in [\bar{S}(m, d, t)\gamma_{\min}, \sqrt{1 \epsilon}\gamma_{\max})$, with probability over $1 \frac{1}{K^{p-2}} 2\exp(-t^2)$, the map $\hat{\phi}_{\gamma}$ perfectly recovers \mathcal{V} .
- 2. If $r_2 > \frac{1 + C_{\kappa}^2 + \frac{C_{\kappa}^2 t}{\sqrt{1 \tilde{\epsilon}_{\min}}}}{\sqrt{1 \tilde{\epsilon}_{\min}}}$, then $\tilde{\epsilon}_{\min} < \tilde{\epsilon}_{\sup 2}$. For any $\epsilon \in [\tilde{\epsilon}_{\min}, \tilde{\epsilon}_{\sup 2})$, and $\hat{\gamma} \in [\bar{S}(m, d, t)\gamma_{\min}, \sqrt{1 \epsilon}\gamma_{\max 2})$, with probability over $1 \frac{1}{K^{p-2}} 2\exp(-t^2)$, the map $\hat{\phi}_{\gamma}$ recovers a non-trivial coarsening of \mathcal{V} .

Proof With probability over $1 - \frac{1}{K^{p-2}} - 2\exp(-t^2)$, we have

- (i) The centroids $\{\Pi \mathbf{a}^{(0)}, \dots, \Pi \mathbf{a}^{(K)}\}$ of the embedded data are distinct.
- (ii) The parameters $\hat{\gamma}_{\min}$, $\hat{\gamma}_{\max}$, and $\hat{\gamma}_{\max 2}$ defined in (10) satisfy the following inequalities:

$$\hat{\gamma}_{\min} \le s_1(\Pi)\gamma_{\min} \le \bar{S}(m,d,t)\gamma_{\min}, \quad \sqrt{1-\epsilon}\gamma_{\max} \le \hat{\gamma}_{\max}, \quad \sqrt{1-\epsilon}\gamma_{\max} \ge \hat{\gamma}_{\max}.$$

The above implies that

$$\begin{bmatrix} \bar{S}(m,d,t)\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max}) \subseteq [\hat{\gamma}_{\min}, \hat{\gamma}_{\max}), \\ \bar{S}(m,d,t)\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max}2) \subseteq [\hat{\gamma}_{\min}, \hat{\gamma}_{\max}2). \end{bmatrix}$$
(15)

Now, we prove the first part of the theorem. We claim here that it is sufficient to show: if $r > \frac{1+C_{\kappa}^2+\frac{C_{\kappa}^2t}{\sqrt{d}}}{\sqrt{1-\tilde{\epsilon}_{\min}}}$, then $\tilde{\epsilon}_{\min} < \tilde{\epsilon}_{\sup}$, and for any $\epsilon \in [\tilde{\epsilon}_{\min}, \tilde{\epsilon}_{\sup})$, the interval $[\bar{S}(m,d,t)\gamma_{\min},\sqrt{1-\epsilon}\gamma_{\max})$ is nonempty. In fact, if $[\bar{S}(m,d,t)\gamma_{\min},\sqrt{1-\epsilon}\gamma_{\max})$ is nonempty, then by the first inclusion of (15), $[\hat{\gamma}_{\min},\hat{\gamma}_{\max})$ is nonempty. Applying Theorem 2 to the embbedded data ΠA implies that for any $\hat{\gamma} \in [\hat{\gamma}_{\min},\hat{\gamma}_{\max})$, the map $\hat{\phi}_{\hat{\gamma}}$ perfectly recovers \mathcal{V} .

On the one hand, by definition of $\tilde{\epsilon}_{\min}$ and C_0 , we have $\frac{1}{\sqrt{d}} = \frac{\tilde{\epsilon}_{\min}}{\sqrt{O(p \log(K))}}$, and $C_0^{-1} = \tilde{\epsilon}_{\min}^{-1} + \frac{C_{\kappa}^2 t}{\sqrt{O(P \log(k))}}$. As a result,

$$\begin{split} r > \frac{C_{\kappa}^2 + 1 + \frac{C_{\kappa}^2 t}{\sqrt{d}}}{\sqrt{1 - \tilde{\epsilon}_{\min}}} & \implies r > \frac{C_{\kappa}^2 + \tilde{\epsilon}_{\min} \tilde{\epsilon}_{\min}^{-1} + \frac{\tilde{\epsilon}_{\min} C_{\kappa}^2 t}{\sqrt{O(P \log(k))}}}{\sqrt{1 - \tilde{\epsilon}_{\min}}} \\ & \implies r > \frac{C_{\kappa}^2 + \tilde{\epsilon}_{\min} \left(\tilde{\epsilon}_{\min}^{-1} + \frac{C_{\kappa}^2 t}{\sqrt{O(P \log(k))}}\right)}{\sqrt{1 - \tilde{\epsilon}_{\min}}} \\ & \implies r > \frac{C_{\kappa}^2 + \tilde{\epsilon}_{\min} C_0^{-1}}{\sqrt{1 - \tilde{\epsilon}_{\min}}} \\ & \implies c_{\kappa}^2 + \tilde{\epsilon}_{\min} C_0^{-1} < \sqrt{1 - \tilde{\epsilon}_{\min}} r \\ & \implies (C_0^{-1} \tilde{\epsilon}_{\min} + C_{\kappa}^2)^2 < (1 - \tilde{\epsilon}_{\min}) r^2 \\ & \implies \tilde{\epsilon}_{\min}^2 + C_0 \left(2C_{\kappa}^2 + r^2 C_0\right) \tilde{\epsilon}_{\min} + C_0^2 \left(-r^2 + C_{\kappa}^4\right) < 0. \end{split}$$

In other words, $\tilde{\epsilon}_{\min}$ satisfies the following inequality

$$x^{2} + C_{0} \left(2C_{\kappa}^{2} + r^{2}C_{0}\right)x + C_{0}^{2} \left(-r^{2} + C_{\kappa}^{4}\right) < 0.$$

It is not difficult to check the solutions to the above inequality is $x \in (x_1, x_2)$, where

$$\begin{aligned} x_1 &= -rC_0\sqrt{\frac{r^2C_0^2}{4} + C_\kappa^2C_0 + 1} - C_\kappa^2C_0 - \frac{r^2C_0^2}{2} < 0, \\ x_2 &= rC_0\sqrt{\frac{r^2C_0^2}{4} + C_\kappa^2C_0 + 1} - C_\kappa^2C_0 - \frac{r^2C_0^2}{2} \in (0, 1). \end{aligned}$$

One may realize that $x_2 = \tilde{\epsilon}_{\sup}$. This implies that $\tilde{\epsilon}_{\min} < \tilde{\epsilon}_{\sup}$. On the other hand, we have

$$\begin{split} \epsilon \in \left[\tilde{\epsilon}_{\min}, \tilde{\epsilon}_{\sup}\right) \subseteq \left(x_{1}, \tilde{\epsilon}_{\sup}\right) &\implies \epsilon^{2} + C_{0}\left(2C_{\kappa}^{2} + r^{2}C_{0}\right)\epsilon + C_{0}^{2}\left(-r^{2} + C_{\kappa}^{4}\right) < 0 \\ &\implies \left(C_{0}^{-1}\epsilon + C_{\kappa}^{2}\right)^{2} < (1 - \epsilon)r^{2} \\ &\implies \left(C_{0}^{-1}\epsilon + C_{\kappa}^{2}\right) < \sqrt{1 - \epsilon r} \\ &\implies \left(C_{0}^{-1}\epsilon + C_{\kappa}^{2}\right)\gamma_{\min} < \sqrt{1 - \epsilon}\gamma_{\max} \\ &\implies \left(\frac{\sqrt{d} + C_{\kappa}^{2}t}{\sqrt{O(p\log(K)}}\epsilon + C_{\kappa}^{2}\right)\gamma_{\min} < \sqrt{1 - \epsilon}\gamma_{\max} \\ &\implies \left(\frac{\sqrt{d} + C_{\kappa}^{2}t}{\sqrt{m}} + C_{\kappa}^{2}\right)\gamma_{\min} < \sqrt{1 - \epsilon}\gamma_{\max} \\ &\implies \bar{S}(m, d, t)\gamma_{\min} < \sqrt{1 - \epsilon}\gamma_{\max}. \end{split}$$

Thus, we have proved the first part of the theorem. The second part of the theorem can be proved in a similar way.

Remark 15 Here, we want to make some remarks on the obtained results.

- 1. The embedding dimension in Theorem 14 is independent of the number of data points n, which is important for clustering an extremely large number of data points.
- 2. The results of this theorem and Theorem 10 further demonstrate that the ratio $r = \gamma_{\text{max}}/\gamma_{\text{min}}$ is a data scale-invariant measure to characterize the difficulty of clustering a given collection of data. Since the embedding dimension of the JL lemma depends on $O(\epsilon^{-2})$ and $\epsilon \in (0,1)$, the value $\tilde{\epsilon}_{\min}$ (and ϵ_{\min}) can be interpreted as the lowest possible dimension reduction ratio obtained by the JL lemma. Since the JL lemma is optimal if the ϵ -isometry mapping is linear, thus, the condition $r > \frac{1+C_{\kappa}^2+\frac{C_{\kappa}^2t}{\sqrt{d}}}{\sqrt{1-\tilde{\epsilon}_{\min}}}$ in Theorem 14 (and $r > \sqrt{\frac{1+\epsilon_{\min}}{1-\epsilon_{\min}}}$ in Theorem 10) shows that the dimension reduction results obtained in this paper are intrinsically depending on the difficulties of clustering the input data.
- 3. For the K-means model, Cohen et al. (2015) proved that the cost can be preserved up to a $(9+\epsilon)$ approximation bounds if the embedding dimension $m = O(\epsilon^{-2} \log K)$. This bound has been improved to $(1+\epsilon)$ if the embedding dimension is $m = O(\epsilon^{-2} \log(K/\epsilon))$ (Makarychev et al., 2022). However, it is still unknown whether the randomly projected K-means model can preserve the cluster membership assignments or not.

4. Numerical Experiments

In this section, we present extensive numerical experiment results to show the practical performance of our model (RPCCM). We first consider high-dimensional data randomly generated from a mixture of spherical Gaussians $\mathcal{N}(\boldsymbol{\mu}_k, \sigma_k^2 I_d)$ with K distinct means $\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_K \in \mathbb{R}^d$.

In the realization of dimension reduction, by default, we randomly sample $\Pi \in \mathbb{R}^{m \times d}$ through

$$\Pi = \frac{1}{\sqrt{m}} G \in \mathbb{R}^{m \times d},\tag{16}$$

where G_{ij} are sampled from i.i.d. standard normal distribution, and $m = ceil(C\epsilon^{-2}\log(n))$. Here, C > 0 is a constant, and $\epsilon \in (0,1)$ is the distortion parameter, which will be specified in the experiments.

In this section, we will set the weights of the convex clustering model as follows:

$$w_{ij} = \begin{cases} \exp(-\phi \|\mathbf{a}_i - \mathbf{a}_j\|^2) & \text{if } (i, j) \in \mathcal{E}, \\ 0 & \text{otherwise}, \end{cases}$$
(17)

where $\mathcal{E} := \{(i, j) \mid \text{if } \mathbf{a}_i \text{ (or } \mathbf{a}_j) \text{ is in } \mathbf{a}_j\text{'s (or } \mathbf{a}_i\text{'s) k-nearest neighbors}, 1 \leq i \neq j \leq n\}.$ We set $\phi = \frac{1}{d}$ by default to rescale the weights and k will be specified in the experiments.

We adopt the semismooth Newton based augmented Lagrangian method (SSNAL) (Sun et al., 2021), which is a state-of-the-art algorithm for solving models (CCM) and (RPCCM). We adopt the duality gap as the stopping criterion (see (Yuan et al., 2022) for details) with a tolerance $\epsilon_{\text{tol}} = 10^{-6}$.

We organize our numerical experiment results as follows: In Section 4.1, we first justify the quality of the random projection matrix Π for preserving the pairwise distances for the data points and centroids. After that, we verify the recovery guarantees of the model (RPCCM). We further compare the cluster recovery performance of the model (RPCCM) to the randomly projected K-means model (RP K-means). In Section 4.2, we will numerically demonstrate that the embedding dimension can be $O(\epsilon^{-2} \log(K))$. In Section 4.3, we test the robustness of the model (RPCCM) with different problem scales and embedding dimensions. Lastly, we test the performance of the model on real data in Section 4.4.

4.1 Numerical Verification for the Randomly Projected Convex Clustering Model with $m = O(\epsilon^{-2}\log(n))$

In this section, we verify the theoretical performance of the model (RPCCM) by conducting numerical experiments on one simulated balanced Gaussian data $A \in \mathbb{R}^{2000 \times 1000}$. Data A is generated from a mixture of K = 20 spherical Gaussians $\mathcal{N}(\mathbf{e}_k, 0.005I_{2000})$ with equal probability $w_k = \frac{1}{20}$, for all $k = 1, \ldots, 20$. Here, $\mathbf{e}_k \in \mathbb{R}^{2000}$ is the k-th column of the identity matrix I_{2000} . Note that we know the true cluster assignments of the simulated data. Let $X_A = \{\mathbf{a}_i - \mathbf{a}_j \mid 1 \le i < j \le n\}$, $X_\alpha = \{\mathbf{a}_i - \mathbf{a}_j \mid i, j \in I_\alpha, i \ne j\}$, $\alpha = 1, \ldots, K$, and $X_{\mathcal{C}(A)} = \{\mathbf{a}^{(\alpha)} - \mathbf{a}^{(\beta)} \mid 1 \le \alpha < \beta \le K\}$. Let $X_{\mathcal{V}} = \bigcup_{\alpha=1}^{20} X_\alpha$. The size of X_A is $\mathcal{C}(1000, 2) = 499500$, the size of $X_{\mathcal{V}}$ is $\sum_{\alpha=1}^{20} \mathcal{C}(n_\alpha, 2) = 24926$, and the size of $X_{\mathcal{C}(A)}$ is $\mathcal{C}(20, 2) = 210$. The visualization of this data set is in Figure 1a. For all the visualizations of the high-dimensional data points in this paper, we adopt the t-SNE (van der Maaten

and Hinton, 2008) to project them to \mathbb{R}^3 . Motivated by the assumptions of the recovery guarantees, we will set the weights w_{ij} as (17) with a graph

$$\mathcal{E}_{A} := \bigcup_{i=1}^{1000} \{(i,j) \mid \text{ if } \mathbf{a}_{i} \text{ (or } \mathbf{a}_{j}) \text{ is in } \mathbf{a}_{j}\text{'s (or } \mathbf{a}_{i}\text{'s) } 20\text{-nearest neighbors}, 1 \leq i \neq j \leq 1000\}$$

$$\bigcup_{\alpha=1}^{20} \{(i,j) \mid i,j \in I_{\alpha}, i \neq j\}.$$

$$(18)$$

4.1.1 Quality of the Random Projection Matrix

We will verify the robustness of Π for pair-wise distance preservation. For this purpose, we will generate the projection matrices Π following (16) with $m = \text{ceil}(9\epsilon^{-2}\log(1000))$ and $\epsilon \in \{0.2, 0.4, 0.6, 0.8, 0.95\}$. In other words, we will test the random projection matrices with $m \in \{1555, 389, 173, 98, 69\}$. We first randomly generate a projection matrix and visualize the embedded data for each m in Figure 1b, 1c, 1d, 1e, and 1f, respectively. From

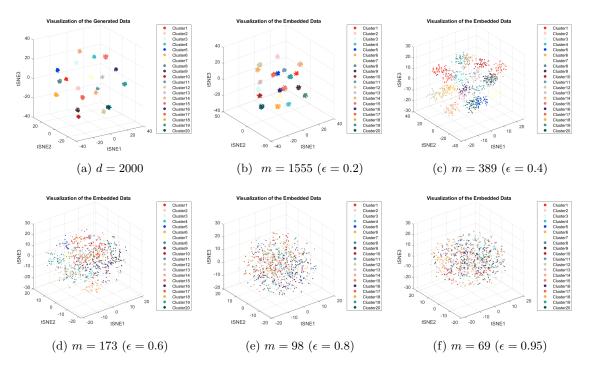


Figure 1: Visualization for A and five embedded data.

the figures, one may see that as the distortion parameter ϵ increases (in other words, m decreases), different clusters in the embedded data become less separate, which is intuitive. Moreover, we can observe that the random projection matrix can preserve the pairwise distances structure of the input data A very well if we set a relatively small distortion parameter ϵ . To further demonstrate the robustness, we will randomly generate 1000 independent samples of the random projection matrix Π for every parameter setting, and test the successful probability of the squared-norm preservation of the points in the sets X_A , X_V , and $X_{\mathcal{C}(A)}$ within the desired distortion ϵ . The results are summarized in Table 2. The results demonstrate the robustness of the random projection matrices for pair-wise

distance preservation. On the one hand, the square-norm can be preserved for almost all points (with a percentage over 99.999%). On the other hand, the success rate for a random projection matrix to preserve the square-norm for all the points in $X_{\mathcal{V}}$ and the centroids $X_{\mathcal{C}}$ are very high.

Table 2: The numerical performance of the random projection matrix Π for preserving the squared norm of the points in X_A , $X_{\mathcal{V}}$, and $X_{\mathcal{C}(A)}$ within the desired distortion. In the table, p_{X_A} , $p_{X_{\mathcal{V}}}$, and $p_{X_{\mathcal{C}(A)}}$ are the successful probability for preserving the squared norm of all the points. $X_A\%$, $X_{\mathcal{V}}\%$, and $X_{\mathcal{C}(A)}\%$ are the average percentage of the points whose squared norm are preserved within the desired distortion.

Dimension (distortion)	p_{X_A}	$X_A\%$	$p_{X_{\mathcal{V}}}$	$X_{\mathcal{V}}\%$	$p_{X_{\mathcal{C}(A)}}$	$X_{\mathcal{C}(A)}\%$
$m = 1555 (\epsilon = 0.2)$	950/1000	99.999%	1000/1000	99.999%	1000/1000	100%
$m = 389 (\epsilon = 0.4)$	855/1000	99.999%	993/1000	99.999%	1000/1000	100%
$m = 173 (\epsilon = 0.6)$	705/1000	99.999%	982/1000	99.999%	1000/1000	100%
$m = 98 (\epsilon = 0.8)$	501/1000	99.999%	951/1000	99.999%	1000/1000	100%
$m = 69 (\epsilon = 0.95)$	248/1000	99.999%	907/1000	99.999%	1000/1000	100%

4.1.2 Verification of the Recovery Guarantees of the Randomly Projected Convex Clustering Model

Next, we will verify the recovery guarantees of the model (RPCCM) established in Theorem 10. Since the effectiveness of the random projection matrix Π for pair-wise distance preservation has already been verified, now, we will randomly sample a projection matrix Π for each m in the experiments described below. We first compute the upper bound γ_{max} and the lower bound γ_{min} of γ defined by (3) and their ratio $r = \frac{\gamma_{\text{max}}}{\gamma_{\text{min}}}$ on the original data A. The values are

$$\gamma_{\min} = 0.1620, \quad \gamma_{\max} = 1.2474, \quad r = 7.6985,$$

which imply that the model (CCM) with our designed weights w_{ij} can perfectly recover the true cluster membership of A for any $\gamma \in [0.1620, 1.2474)$. The large ratio r implies the feasibility of the model (RPCCM) with some suitable $\epsilon \in (0, 1)$ under the same weights w_{ij} . We then estimate the values ϵ_{\min} and ϵ_{\sup} defined in Theorem 10, which are

$$\epsilon_{\min} = 0.1763, \quad \epsilon_{\sup} = 0.9668.$$

The results in Theorem 10 imply that for $0.1763 \le \epsilon < 0.9668$, and $\gamma \in [\sqrt{1+\epsilon}\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max})$, the model (RPCCM) with $m = O(9\epsilon^{-2}\log(1000))$ can perfectly recover the true cluster membership of A with high probability. Here, we take $m = ceil(9\epsilon^{-2}\log(1000))$.

Since the estimated valid interval of distortions is $(\epsilon_{\min}, \epsilon_{\sup}) = (0.1763, 0.9668)$, we choose $\epsilon \in \{0.2, 0.4, 0.6, 0.8, 0.95\}$ for verification. The corresponding embedding dimensions are $m \in \{1555, 389, 173, 98, 69\}$. To verify the recovery guarantees of the models (CCM) and (RPCCM), we will generate a clustering path of the model (CCM) on the original data A and a clustering path of the model (RPCCM) on the embedded data for each m. In particular, we will generate all clustering paths with $\gamma \in [10:-0.1:0.1]$. We will compute the number of clusters K, the rand index, and the adjusted rand index against γ

on the clustering paths. The results are shown in Figure 2. To better verify the recovery guarantees, we compute the estimated range $\left[\sqrt{1+\epsilon}\gamma_{\min},\sqrt{1-\epsilon}\gamma_{\max}\right]$ in Theorem 10 for perfect recovery for different m in Table 3.

Table 3: Estimated ranges of γ for perfect recovery guarantees of RPCCM. The range $\left[\sqrt{1+\epsilon}\gamma_{\min},\sqrt{1-\epsilon}\gamma_{\max}\right]$ is estimated using Theorem 10 by the model (RPCCM) and the range $\left[\hat{\gamma}_{\min},\hat{\gamma}_{\max}\right]$ defined by (10) is implicitly estimated using Theorem 2 by (CCM).

Dimension (distortion)	$\left[\sqrt{1+\epsilon}\gamma_{\min},\sqrt{1-\epsilon}\gamma_{\max}\right)$	$[\hat{\gamma}_{\min}, \hat{\gamma}_{\max})$
$m = 1555 (\epsilon = 0.2)$	[0.1775, 1.1157]	[0.1631,1.2334)
$m = 389 (\epsilon = 0.4)$	[0.1917, 0.9669)	[0.1699,1.2680)
$m = 173 (\epsilon = 0.6)$	[0.2049, 0.7889)	[0.1610,1.1783)
$m = 98 (\epsilon = 0.8)$	[0.2174, 0.5578)	[0.1618,1.2101)
$m = 69 (\epsilon = 0.95)$	[0.2263, 0.2789)	[0.1707, 1.2443)

From the results in Figure 2 and Table 3, we can see that the model (RPCCM) indeed performs perfect cluster recovery when γ is chosen in the interval $\left[\sqrt{1+\epsilon}\gamma_{\min},\sqrt{1-\epsilon}\gamma_{\max}\right]$.

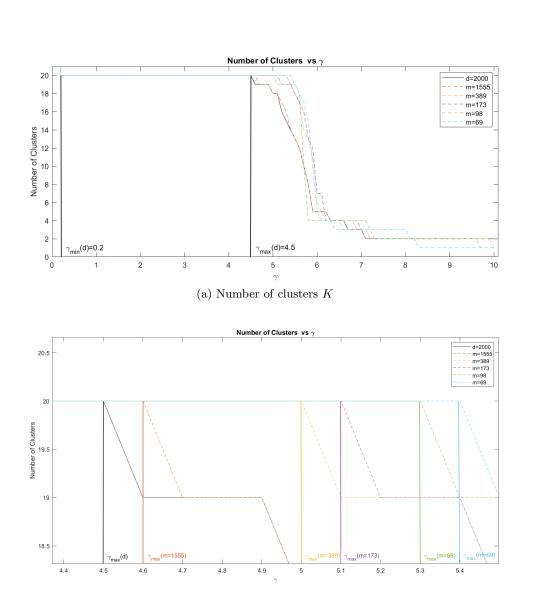
In a word, the recovery guarantees of the convex clustering model (CCM) on the original data A can be preserved by the model (RPCCM) with a much smaller dimension and the performance of the model (RPCCM) is attractive in practice.

Remark 16 We want to add a remark here on the empirical performance of the model (RPCCM). As shown in Table 3, the interval $[\hat{\gamma}_{\min}, \hat{\gamma}_{\max})$ of the model (RPCCM) for the perfect recovery can be larger. The empirical performance can be robust with respect to the embedding dimension. This can be demonstrated by the results in Figure 2.

4.1.3 Comparison between the Randomly Projected Convex Clustering Model and the Randomly Projected K-means Model

To further demonstrate the superior performance of the model (RPCCM), we compare the clustering performance between the model (RPCCM) and the RP K-means on the data A. Since we know the true number of clusters is K=20, we compare the clustering quality of the two models for $K \in \{16, 17, 18, 19, 20\}$. More specifically, we will compare their performance in terms of the rand index and the adjusted rand index against different numbers of clusters. For the implementation of K-means and RP K-means in this paper, we use the "kmeans" package from Matlab with parameters 'MaxIter'=10000 and 'Replicates'=30. We summarize the results in Table 4 and Table 5.

From the results in Table 4 and Table 5, we can see that the performance of the model (RPCCM) is better and more robust than RP K-means, even when the number of clusters is not correctly classified ($K \in \{16, 17, 18, 19\}$). Neither K-means nor RP K-means can perform a perfect recovery based on our experiments, and the recovery performance of RP K-means becomes less reliable as m decreases. As a comparison, the recovery results of the model (CCM) are robustly inherited by the model (RPCCM), and the model (RPCCM) with all five m could perform perfectly recovery on A with some suitable γ on the path. The embedding dimension is as low as m=69, which can greatly reduce the computational cost.





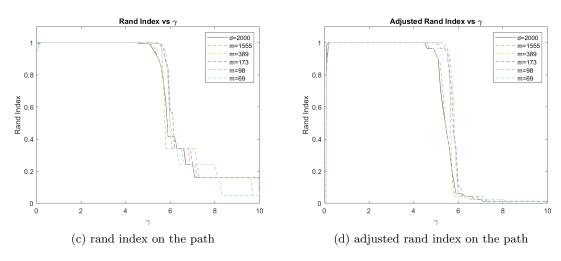


Figure 2: Clustering performance on the clustering path.

Table 4: The rand index value against the number of clusters $(K = \{16, 17, 18, 19, 20\})$ for CCM and K-means, and RPCCM and RP K-means with each m on data A. For CCM and RPCCM, if K is identified by some γ (maybe not unique) on the clustering path, we will pick an adjusted rand index value as the record. Otherwise, if there is no γ on the clustering path such that K is identified, we will denote it by '/'.

Model	K = 16	K = 17	K = 18	K = 19	K = 20
CCM $(d = 2000)$	0.9637	/	0.9929	0.9965	1.0000
RPCCM $(m = 1555)$	0.9637	0.9857	0.9929	0.9965	1.0000
RPCCM $(m = 389)$	/	0.9786	0.9929	0.9965	1.0000
RPCCM $(m = 173)$	0.9637	0.9786	0.9929	0.9965	1.0000
RPCCM $(m = 98)$	/	/	0.9893	0.9965	1.0000
RPCCM $(m = 69)$	/	0.9786	0.9929	0.9965	1.0000
K-means $(d = 2000)$	0.9695	0.9702	0.9684	0.9733	0.9851
RP K-means $(m = 1555)$	0.9619	0.9808	0.9804	0.9754	0.9836
RP K-means ($m = 389$)	0.9620	0.9644	0.9824	0.9778	0.9811
RP K-means ($m = 173$)	0.9367	0.9401	0.9344	0.9458	0.9473
RP K-means $(m = 98)$	0.9045	0.8992	0.9031	0.9122	0.9107
RP K-means $(m = 69)$	0.8971	0.8995	0.9019	0.9029	0.9040

Table 5: The adjusted rand index value against the number of clusters $(K = \{16, 17, 18, 19, 20\})$ for CCM and K-means, and RPCCM and RP K-means with each m on data A. For CCM and RPCCM, if K is identified by some γ (maybe not unique) on the clustering path, we will pick an adjusted rand index value as the record. Otherwise, if there is no γ on the clustering path such that K is identified, we will denote it by γ .

Model	K = 16	K = 17	K = 18	K = 19	K = 20
CCM $(d = 2000)$	0.7154	/	0.9299	0.9645	1.0000
RPCCM $(m = 1555)$	0.7154	0.8670	0.9299	0.9645	1.0000
RPCCM $(m = 389)$	/	0.8125	0.9299	0.9645	1.0000
RPCCM $(m = 173)$	0.7154	0.8125	0.9299	0.9645	1.0000
RPCCM $(m = 98)$	/	/	0.8975	0.9645	1.0000
RPCCM $(m = 69)$	/	0.8125	0.9299	0.9645	1.0000
K-means $(d = 2000)$	0.7493	0.7525	0.7355	0.7674	0.8578
RP K-means $(m = 1555)$	0.6989	0.8284	0.8237	0.7801	0.8426
RP K-means ($m = 389$)	0.6791	0.6949	0.8367	0.7901	0.8164
RP K-means $(m = 173)$	0.4669	0.4525	0.4053	0.4654	0.4964
RP K-means $(m = 98)$	0.1284	0.0858	0.1007	0.1353	0.1293
RP K-means $(m = 69)$	0.0585	0.0807	0.0494	0.0645	0.0759

4.2 Numerical Verification for the Randomly Projected Convex Clustering Model with $m = O(\epsilon^{-2} \log(K))$

In this section, we will further verify the recovery guarantees established in Theorem 14 for the model (RPCCM). In other words, we want to numerically verify that the embedding dimension m of the model (RPCCM) can be further improved from $O(\epsilon^{-2}\log(n))$ to $O(\epsilon^{-2}\log(K))$. For simplicity, we choose $m = ceil(10\epsilon^{-2}\log(n))$ and $\tilde{m} = ceil(10\epsilon^{-2}\log(K))$, respectively. Here, $\epsilon \in (0,1)$ is some given distortion.

We will conduct experiments on a collection of data points $A' := \{\mathbf{a}_1', \dots, \mathbf{a}_{10000}'\} \subseteq \mathbb{R}^{100}$, where each \mathbf{a}_i' is randomly sampled from a balanced Gaussian mixture. In particular, we set K = 10, $\boldsymbol{\mu}_k = \mathbf{e}_k$, $\sigma_k^2 = 0.1$, and $w_k = \frac{1}{10}$, for $k = 1, \dots, 10$ for the Gaussian mixture. Let $X_\alpha' = \{\mathbf{a}_i' - \mathbf{a}_j' \mid i, j \in I_\alpha, i \neq j\}$, $\alpha = 1, \dots, 10$, and $X_{\mathcal{C}(A)}' = \{\mathbf{a}'^{(\alpha)} - \{\mathbf{a}'^{(\beta)} \mid 1 \leq \alpha < \beta \leq 10\}$, and denote $X_\mathcal{V}' = \bigcup_{\alpha=1}^{10} X_\alpha'$. Similarly, inspired by the assumptions in Theorem 14, we will set the weights w_{ij} as (17) with a graph

$$\mathcal{E}_{A'} := \bigcup_{i=1}^{10000} \{(i,j) \mid \text{ if } \mathbf{a}'_i \text{ (or } \mathbf{a}'_j) \text{ is in } \mathbf{a}_j \text{'s (or } \mathbf{a}'_i \text{'s) 10-nearest neighbors, } 1 \leq i \neq j \leq 10000 \}$$

$$\bigcup_{\alpha=1}^{10} \{(i,j) \mid i,j \in I_{\alpha}, i \neq j \}.$$
(10)

First, we compute the values γ_{max} and γ_{min} defined by (3) and their ratio $r = \frac{\gamma_{\text{max}}}{\gamma_{\text{min}}}$ on the original data A'. The values are

$$\gamma_{\min} = 0.0093, \quad \gamma_{\max} = 0.0887, \quad r = 9.5397,$$

which implies that the model (CCM) with above weights w_{ij} can perfectly recover the true cluster membership of A' for any $\gamma \in [0.0093, 0.0887)$. The large ratio r implies the feasibility of the model (RPCCM) with some suitable $\epsilon \in (0, 1)$ under the same weights w_{ij} .

Next, we will calculate the theoretically valid embedding dimensions for both cases, respectively. In order to achieve this goal, we will calculate the values ϵ_{\min} and ϵ_{\sup} defined in Theorem 10 and the values $\tilde{\epsilon}_{\min}$ and $\tilde{\epsilon}_{\sup}$ defined in Theorem 14, respectively.

If we take the embedding dimension as $m = O(\epsilon^{-2}\log(n)) = ceil(10\epsilon^{-2}\log(10000))$. The values ϵ_{\min} and ϵ_{\sup} defined in Theorem 10 on the data A' are $\epsilon_{\min} = 0.9597$ and $\epsilon_{\max} = 0.9782$. This implies that for $\epsilon \in [0.9597, 0.9782)$ and $\gamma \in [\sqrt{1+\epsilon}\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max})$, the model (RPCCM) with the corresponding embedding dimension m can perform the perfect clustering recovery on A' with high probability. The lowest possible dimension reduction ratio ϵ_{\min} is very close to 1, which implies that we can hardly obtain a sufficient dimension reduction effect by Theorem (10). In fact, the lowest possible embedding dimension guaranteed by Theorem (10) is $m = ceil(10\epsilon_{\sup}^{-2}\log(10000)) = 97$. We will choose a valid distortion $\epsilon = 0.975 \in [0.9597, 0.9782)$. and test with m = 97. We will compute the theoretically estimated interval $[\sqrt{1+\epsilon}\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max})$ in Theorem 10 for perfect recovery with $\epsilon = 0.975$. Then, we will randomly sample a Π and test whether the model (RPCCM) could perform the perfect clustering recovery for γ in the estimated interval. Results are listed in Table 6.

Now, we move on to consider taking $\tilde{m} = O(\epsilon^{-2} \log(K)) = ceil(10\epsilon^{-2} \log(10))$. For a random matrix $\Pi \in \mathbb{R}^{\tilde{m} \times d}$ defined as (16), it follows from Theorem II.13 in (Davidson and Szarek, 2001) and Theorem 2.6 in (Rudelson and Vershynin, 2010) that, the two-side

bounds $\bar{S}(\tilde{m}, d, t)$ and $\underline{S}(\tilde{m}, d, t)$ in (13) are

$$\bar{S}(\tilde{m}, d, t) = \frac{\sqrt{100} + t}{\sqrt{\tilde{m}}} + 1 = \frac{10 + t}{\sqrt{\tilde{m}}} + 1, \quad \underline{S}(\tilde{m}, d, t) = \frac{\sqrt{100} - t}{\sqrt{\tilde{m}}} - 1 = \frac{10 - t}{\sqrt{\tilde{m}}} - 1.$$

By setting t = 2, with a probability over $1 - 2\exp(-2^2) = 0.9634$, we have

$$s_1(\Pi) \le \bar{S}(\tilde{m}, 100, 2) = \frac{12}{\sqrt{\tilde{m}}} + 1,$$

 $s_{\tilde{m}}(\Pi) \ge \underline{S}(\tilde{m}, 100, 2) = \frac{8}{\sqrt{\tilde{m}}} - 1,$

and the values $\tilde{\epsilon}_{\min}$ and $\tilde{\epsilon}_{\sup}$ defined in Theorem 14 are then estimated to be $\tilde{\epsilon}_{\min} = 0.4799$ and $\tilde{\epsilon}_{\max} = 0.8863$. This implies that for any $\epsilon \in [0.4799, 0.8863)$, and $\gamma \in [\bar{S}(\tilde{m}, 100, 2)\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max})$, the model (RPCCM) with embedding dimension \tilde{m} can perform the perfect clustering recovery of the data A' with high probability. We choose $\epsilon \in \{0.70, 0.85\}$ in the valid interval [0.4799, 0.8863). In other words, we will test with $\tilde{m} \in \{47, 32\}$. For each \tilde{m} , we will first randomly sample 1000 independent Π , and then test the successful probability $p_{X'_{\mathcal{C}(A)}}$ of the squared-norm preservation of the points in the set $X'_{\mathcal{C}(A)}$ within the desired distortion, as well as the successful probability p_S of the two-side bounds of extreme singulars of the random projection matrices. We will then compute the estimated range $[\bar{S}(\tilde{m}, 100, 2)\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max})$ in Theorem 14 for perfect recovery. Finally, we will randomly sample a random projection matrix Π for each \tilde{m} and test test whether the model (RPCCM) could do perfect recovery with $\gamma \in [\bar{S}(\tilde{m}, 100, 2)\gamma_{\min}, \sqrt{1-\epsilon}\gamma_{\max})$. Results are listed in Table 7.

Table 6: The numerical performance of the model (RPCCM) with embedding dimension $m = O(\epsilon^{-2} \log(n))$.

Dimension (distortion)	$\sqrt{1+\epsilon}\gamma_{\min},\sqrt{1-\epsilon}\gamma_{\max}$	Perfect recovery
$m = 97 (\epsilon = 0.975)$	[0.0131, 0.0140)	✓

Table 7: The numerical performance of the model (RPCCM) with embedding dimension $\tilde{m} = O(\epsilon^{-2} \log(K))$.

Dimension (distortion)	$p_{X'_{\mathcal{C}(A)}}$	p_S	$\left[\tilde{S}(\tilde{m}, 100, 2) \gamma_{\min}, \sqrt{1 - \epsilon} \gamma_{\max} \right)$	Perfect recovery
$\tilde{m} = 47 (\epsilon = 0.70)$	921/1000	1000/1000	[0.0256, 0.0486)	✓
$\tilde{m} = 32 (\epsilon = 0.85)$	915/1000	1000/1000	[0.0290, 0.0344)	✓

From the results in Table 6 and Table 7, we may observe that, under the settings in this section, we can only reduce the original dimension d=100 to m=97 theoretically if we take $m=O(\epsilon^{-2}\log(n))$. In contrast, if we take $m=O(\epsilon^{-2}\log(K))$, we can reduce the data dimension from d=100 to m=32. The above experiments demonstrate that the embedding dimension of the model (RPCCM) can be further improved from $O(\epsilon^{-2}\log(K))$ to $O(\epsilon^{-2}\log(K))$.

4.3 Robustness of the Randomly Projected Convex Clustering Model under Practical Settings

In this section, we will focus on further demonstrating the robustness of the model (RPCCM) under practical settings. We will demonstrate from two perspectives: The robustness of different problem scales and embedding dimensions. First of all, we will conduct some analysis on the practical settings for (RPCCM), in terms of weights w_{ij} and the embedding dimension m. In terms of experiments, we will first exploit the potential of the model (RPCCM) by choosing lower embedding dimensions on data A. Then, we will test on six more simulated balanced Gaussian data with different dimension d, size n, and ground-truth cluster number K. We will also provide numerical experiments on some unbalanced Gaussian data. The datasets are described in details later.

4.3.1 Practical Settings of the Randomly Projected Convex Clustering Model

Recall the settings we use in the numerical verification of the model (RPCCM) on data A: 1. For weights w_{ij} , we choose the Gaussian kernel weights (17) with a well-designed graph (18). 2. For the embedding dimension m, we set $m = ceil(9\epsilon^{-2}\log(n))$, where $\epsilon \in (0,1)$ is some desired distortion. These settings guarantee the conditions in the recovery guarantee of the model (RPCCM): (1) $w_{ij} > 0$ and $n_{\alpha}w_{ij} > \mu_{ij}^{(\alpha)}$ for all $i, j \in I_{\alpha}, \alpha \in [K]$. (2) With high probability, a random projection matrix Π could preserve the squared norm for all the points in $X_{\mathcal{V}}$ and the centroids $X_{\mathcal{C}}$ within the desired distortion ϵ .

In practical implementations of the model (RPCCM), there are two challenges: First, we have no idea about the true cluster assignments of data. Second, computational efficiency should be taken into consideration. To overcome these challenges, we explore some robust and efficient practical settings. For the weights w_{ij} , since the Gaussian kernel weights (17) with a k-nearest neighbors graph has already demonstrated its robustness in the past literature (Chi and Lange, 2015; Yuan et al., 2018; Sun et al., 2021), we simply choose the weights with a 10-nearest neighbors graph by default. We will focus more on testing the robustness regarding the embedding dimension m.

Although the mentioned two conditions for recovery guarantees might no longer hold in practical settings, our experimental results show that the practical performance of the model (RPCCM) could still be robust. This motivates us to explore tighter and more general recovery guarantees of the model (RPCCM) in further work.

4.3.2 Robustness of the Randomly Projected Convex Clustering Model with Lower Embedding Dimensions

We will test the robustness of the model (RPCCM) on A regarding m. We choose the same desired distortions $\epsilon \in \{0.2, 0.4, 0.6, 0.8, 0.95\}$ as in the previous section but set $m = ceil(\epsilon^{-2}\log(n))$. In other words, the corresponding embedding dimensions are $m \in \{173, 44, 20, 11, 8\}$, which are much lower than the previous setting with $m = ceil(9\epsilon^{-2}\log(n))$. For each m, we first randomly sample ten random projections Π . Then, we compute the averaged percentage of the squared norm of points that are successfully jointly preserved within the desired distortion ϵ in X_A , X_V , and $X_{\mathcal{C}(A)}$. The results are listed in Table 8. We can observe from the results that over 93% of points on average could still be preserved

jointly within the desired distortion ϵ . Next, we test the practical clustering performance of the model (RPCCM) regarding all the ten randomly sampled projection matrices Π on a clustering path generated by $\gamma = [10:-0.2:2]$. The results are summarized in Table 8. From Table 8, we can see that for each m, the model (RPCCM) can perform perfect recovery robustly for all the ten randomly sampled projection matrices. These results show that the practical performance of the randomly projected convex clustering model is very robust.

Table 8: Averaged percentage of points in X_A , X_V , and $X_{\mathcal{C}(A)}$ that the square-norm of these points can be jointly preserved by one random Π within the desired distortion, and the recovery results.

Dimension (distortion)	$X_A\%$	$X_{\mathcal{V}}\%$	$X_{\mathcal{C}(A)}\%$	Perfect recovery
$m = 173 (\epsilon = 0.2)$	93.70%	93.72%	93.12%	10/10
$m = 44 (\epsilon = 0.4)$	94.32%	94.44%	94.79%	10/10
$m = 20 (\epsilon = 0.6)$	94.78%	94.85%	95.38%	10/10
$m = 11 (\epsilon = 0.8)$	95.04%	95.09%	95.07%	10/10
$m = 8 (\epsilon = 0.95)$	95.16%	95.26%	95.17%	10/10

4.3.3 Robustness of the Randomly Projected Convex Clustering Model with Different Problem Scales

We will test the robustness of the model (RPCCM) with different problem scales. We first test on balanced Gaussian data of different scales. In particular, we choose the scale $(d, n, K) \in \{(10^2, 10^3, 10), (10^3, 10^3, 10), (10^4, 10^3, 10), (10^3, 10^3, 2), (10^3, 10^3, 50), (10^4, 10^4, 50)\}$, and we set $\mu_k = \mathbf{e}_k$, $\sigma_k^2 = 0.005$, and $w_k = \frac{1}{K}$, for $k = 1, \ldots, K$. The above six data sets are visualized in Figure 3.

For each data, we randomly sample ten random projection matrices $\Pi \in \mathbb{R}^{m \times d}$ for every m = 10, 20, 50. The clustering performance of the model (RPCCM) along $\gamma = [10: -0.2: 2]$ is summarized in Table 9. The results show that the model (RPCCM) is robust to the scale of the data in practice.

Table 9: Clustering performance of RPCCM with m = 10, 20, 50 along $\gamma = [10: -0.2: 2]$ on six balanced Gaussian data.

Dimension	Perfect recovery
m = 50	60/60
m = 20	60/60
m = 10	60/60

We also test on an unbalanced Gaussian data generated from 20 spherical Gaussians $\mathcal{N}(\mathbf{e}_k, 0.005I_d)$ for all $k = 1, \dots, 20$, containing 7700 samples in total. In particular, there are 2000 samples for each of the first three clusters, and there are 100 samples for each of the rest 17 clusters. Again, for each m = 10, 20, 50, we randomly sample ten projection

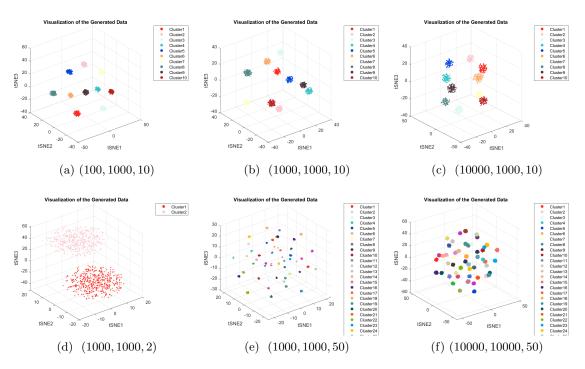


Figure 3: Visualization for six balanced Gaussian data of scale (d, n, k).

matrices $\Pi \in \mathbb{R}^{m \times d}$. We then compare the clustering performance of the model (RPCCM) and the RP 20-means model. We generate the clustering path of the model (RPCCM) with $\gamma = [10:-0.2:2]$. The results are summarized in Table 10, which demonstrate the effectiveness and robustness of the model (RPCCM).

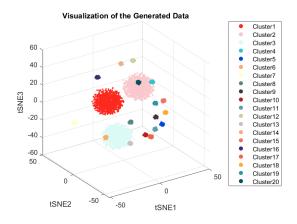


Figure 4: Visualization for the unbalanced Gaussian data.

Table 10: Clustering performance of RPCCM and RP 20-means on the unbalanced Gaussian data (in terms of averaged rand index and adjusted rand index).

Clustering model	Perfect recovery	rand index	adjusted rand index
RPCCM $(m = 50)$	10/10	1.0000	1.0000
RPCCM $(m=20)$	10/10	1.0000	1.0000
RPCCM $(m = 10)$	10/10	1.0000	1.0000
RP 20-means $(m = 50)$	0/10	0.8211	0.2343
RP 20-means $(m=20)$	0/10	0.7899	0.1061
RP 20-means $(m=10)$	0/10	0.7771	0.0503

4.4 Practical Performance of the Randomly Projected Convex Clustering Model on the Real Data

In this section, we will test the practical performance of the model (RPCCM) on the lung cancer data (Lee et al., 2010). The lung cancer data contains the microarray gene expressions of 12625 genes for 56 subjects belonging to one of four disease subgroups: Normal subjects (Normal), pulmonary carcinoid tumors (Carcinoid), colon metastases (Colon), and small cell carcinoma (Small Cell). In the models (CCM) and (RPCCM), we will compute the weights w_{ij} following (17) with a 5-nearest neighbors graph. For the embedding dimension of the model (RPCCM), we will set $m \in \{10, 20, 100, 500\}$. For each m, we will randomly sample a random projection matrix $\Pi \in \mathbb{R}^{m \times d}$ following (16). We will then test the practical performance of the models (CCM) and (RPCCM) by generating a clustering path with $\gamma \in [1:1:35] \cup [36:20:556]$. We visualize the clustering paths in Figure 5.

From the visualizations, we can observe that the convex clustering model (CCM) performs well on this real data set, where only one data point from the Carcinoid cluster is clustered wrongly. A possible reason is that this wrongly clustered data point is closer to the SmallCell cluster. Moreover, the superior performance of the convex clustering model can be properly preserved by the model (RPCCM), even for a very low embedding dimension. More detailed numbers can be found in Table 11.

We also compare the clustering performance of the model (RPCCM) with the RP K-means model. For the sake of fairness, we will test with the true number of clusters K=4. The results are summarized in Table11. The results show that the performance of the RP K-means model becomes less reliable as m decreases, while the model (RPCCM) is robust.

5. Conclusion and Future Works

In this paper, we proposed a randomly projected convex clustering model for clustering high dimensional data. We proved that, under some mild conditions, the perfect recovery of the cluster membership assignments of the convex clustering model on the original data, if exists, can be preserved by the randomly projected convex clustering model with a much smaller embedding dimension. In particular, we proved that the embedding dimension can be $m = O(\epsilon^{-2} \log(n))$, where n is the number of data points and $0 < \epsilon < 1$ is some given tolerance. We further proved that the embedding dimension can be $m = O(\epsilon^{-2} \log K)$, where K is the number of hidden clusters, which is independent of the number of data points.

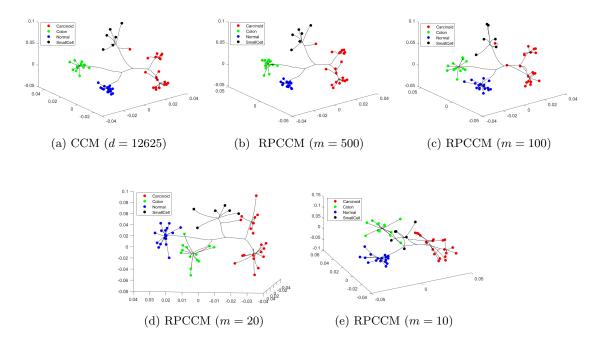


Figure 5: Visualization of the clustering paths.

Table 11: Clustering performance of RPCCM and RP 4-means on the lung cancer data. Here, accuracy means the ratio of correctly clustered data points, and γ^* is some value of γ corresponding to the best clustering results of RPCCM.

Clustering model	accuracy	rand index	adjusted rand index	γ^*
RPCCM $(d = 12625)$	55/56	0.9838	0.9586	76
RPCCM $(m = 500)$	55/56	0.9838	0.9586	76
RPCCM $(m = 100)$	55/56	0.9838	0.9586	76
RPCCM $(m=20)$	55/56	0.9838	0.9586	76
RPCCM $(m = 10)$	55/56	0.9838	0.9586	96
4-means $(d = 12625)$	55/56	0.9838	0.9586	/
RP 4-means $(m = 500)$	55/56	0.9838	0.9586	/
RP 4-means $(m = 100)$	54/56	0.9701	0.9245	/
RP 4-means $(m=20)$	48/56	0.9000	0.7421	/
RP 4-means $(m=10)$	43/56	0.8753	0.6795	/

Extensive numerical experiment results were presented in this paper to demonstrate the robustness and superior performance of the randomly projected convex clustering model. The numerical results presented in this paper also demonstrated that the randomly projected convex clustering model can outperform the randomly projected K-means model in practice.

It is worthwhile pointing out that the practical performance of the convex clustering model and the randomly projected convex clustering model depends on the quality of the input data features. We regard it as a future research direction to investigate a new technique that can do dimension reduction and feature representation learning simultaneously.

Acknowledgments and Disclosure of Funding

The research of Yancheng Yuan is supported in part by The Hong Kong Polytechnic University under grant P0038284. The research of Defeng Sun is supported in part by the Hong Kong Research Grant Council under grant 15304721.

References

- D. Achlioptas. Database-friendly random projections: Johnson-Lindenstrauss with binary coins. *Journal of Computer and System Sciences*, 66(4):671–687, 2003.
- N. Ailon and B. Chazelle. The fast Johnson-Lindenstrauss transform and approximate nearest neighbors. SIAM Journal on Computing, 39(1):302–322, 2009.
- N. Ailon and E. Liberty. Fast dimension reduction using Rademacher series on dual BCH codes. *Discrete & Computational Geometry*, 42(4):615–630, 2009.
- N. Ailon and E. Liberty. An almost optimal unrestricted fast Johnson-Lindenstrauss transform. *ACM Transactions on Algorithms*, 9(3):21, 2013.
- E. C. Chi and K. Lange. Splitting methods for convex clustering. Journal of Computational and Graphical Statistics, 24(4):994–1013, 2015.
- E. C. Chi and S. Steinerberger. Recovering trees with convex clustering. SIAM Journal on Mathematics of Data Science, 1(3):383–407, 2019.
- E. C. Chi, B. R. Gaines, W. W. Sun, H. Zhou, and J. Yang. Provable convex co-clustering of tensors. *Journal of Machine Learning Research*, 21 (214):1–58, 2020.
- J. Chiquet, P. Gutierrez, and G. Rigaill. Fast tree inference with weighted fusion penalties. Journal of Computational and Graphical Statistics, 26(1):205–216, 2017.
- M. B. Cohen, S. Elder, C. Musco, C. Musco, and M. Persu. Dimensionality reduction for k-means clustering and low rank approximation. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 163–172, 2015.
- M. B. Cohen, T. Jayram, and J. Nelson. Simple analyses of the sparse Johnson-Lindenstrauss transform. In 1st Symposium on Simplicity in Algorithms, pages 15:1–15:9, 2018.

- A. Dasgupta, R. Kumar, and T. Sarlós. A sparse Johnson-Lindenstrauss transform. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pages 341–350, 2010.
- S. Dasgupta and A. Gupta. An elementary proof of a theorem of Johnson and Lindenstrauss. Random Structures & Algorithms, 22(1):60–65, 2003.
- K. R. Davidson and S. J. Szarek. Local operator theory, random matrices and banach spaces. *Handbook of the geometry of Banach spaces*, 1:317–366, 2001.
- A. Dunlap and J.-C. Mourrat. Local versions of sum-of-norms clustering. SIAM Journal on Mathematics of Data Science, 4(4):1250–1271, 2022.
- T. D. Hocking, A. Joulin, F. Bach, and J.-P. Vert. Clusterpath an algorithm for clustering using convex fusion penalties. In *International Conference on Machine Learning*, pages 745–752, 2011.
- P. Indyk and R. Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, pages 604–613, 1998.
- T. Jiang, S. Vavasis, and C. W. Zhai. Recovery of a mixture of Gaussians by sum-of-norms clustering. *Journal of Machine Learning Research*, 21(225):1–16, 2020.
- W. B. Johnson and J. Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. In *Contemporary Mathematics*, volume 26, pages 189–206. American Mathematical Society, 1984.
- D. M. Kane and J. Nelson. A derandomized sparse Johnson-Lindenstrauss transform. arXiv preprint arXiv:1006.3585, 2010.
- D. M. Kane and J. Nelson. Sparser Johnson-Lindenstrauss transforms. *Journal of the ACM (JACM)*, 61(1):1–23, 2014.
- K. G. Larsen and J. Nelson. The Johnson-Lindenstrauss lemma is optimal for linear dimensionality reduction. In 43rd International Colloquium on Automata, Languages, and Programming, pages 82:1–82:11, 2016.
- M. Lee, H. Shen, J. Z. Huang, and J. S. Marron. Biclustering via sparse singular value decomposition. *Biometrics*, 66(4):1087–1095, 2010.
- F. Lindsten, H. Ohlsson, and L. Ljung. Clustering using sum-of-norms regularization: With application to particle filter output computation. In 2011 IEEE Statistical Signal Processing Workshop, pages 201–204, 2011.
- K. Makarychev, Y. Makarychev, and I. Razenshteyn. Performance of Johnson–Lindenstrauss transform for k-means and k-medians clustering. SIAM Journal on Computing, 0(0):STOC19–269–STOC19–297, 2022.
- J. Matoušek. On variants of the Johnson-Lindenstrauss lemma. Random Structures & Algorithms, 33(2):142–156, 2008.

- A. Panahi, D. Dubhashi, F. D. Johansson, and C. Bhattacharyya. Clustering by sum of norms: Stochastic incremental algorithm, convergence and cluster recovery. In *Interna*tional Conference on Machine Learning, pages 2769–2777, 2017.
- K. Pelckmans, J. De Brabanter, J. A. Suykens, and B. De Moor. Convex clustering shrinkage. In *PASCAL Workshop on Statistics and Optimization of Clustering Workshop*, 2005.
- P. Radchenko and G. Mukherjee. Convex clustering via l_1 fusion penalization. Journal of the Royal Statistical Society. Series B (Statistical Methodology), 79(5):1527–1546, 2017.
- M. Rudelson and R. Vershynin. Non-asymptotic theory of random matrices: Extreme singular values. In Proceedings of the International Congress of Mathematicians 2010 (ICM 2010) (In 4 Volumes) Vol. I: Plenary Lectures and Ceremonies Vols. II–IV: Invited Lectures, pages 1576–1602, 2010.
- D. F. Sun, K.-C. Toh, and Y. C. Yuan. Convex clustering: Model, theoretical guarantee and efficient algorithm. *Journal of Machine Learning Research*, 22(9):1–32, 2021.
- K. M. Tan and D. Witten. Statistical properties of convex clustering. *Electronic Journal of Statistics*, 9(2):2324–2347, 2015.
- L. van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(11):2579–2605, 2008.
- R. Vershynin. *High-dimensional probability: An Introduction with applications in data science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018.
- Y. C. Yuan, D. F. Sun, and K.-C. Toh. An efficient semismooth Newton based algorithm for convex clustering. In *International Conference on Machine Learning*, pages 5718–5726, 2018.
- Y. C. Yuan, T.-H. Chang, D. F. Sun, and K.-C. Toh. A dimension reduction technique for large-scale structured sparse optimization problems with application to convex clustering. SIAM Journal on Optimization, 32(3):2294–2318, 2022.
- C. Zhu, H. Xu, C. Leng, and S. Yan. Convex optimization procedure for clustering: Theoretical revisit. Advances in Neural Information Processing Systems, 27:1619–1627, 2014.