

On deterministic conditions for subspace clustering under missing data

Wenqi Wang, Shuchin Aeron and Vaneeet Aggarwal

Abstract—In this paper we present deterministic analysis of sufficient conditions for sparse subspace clustering under missing data, when data is assumed to come from a Union of Subspaces (UoS) model. In this context we consider two cases, namely Case I when all the points are sampled at the same co-ordinates, and Case II when points are sampled at different locations. We show that results for Case I directly follow from several existing results in the literature, while results for Case II are not as straightforward and we provide a set of dual conditions under which, perfect clustering holds true. We provide extensive set of simulation results for clustering as well as completion of data under missing entries, under the UoS model. Our experimental results indicate that in contrast to the full data case, accurate clustering does not imply accurate subspace identification and completion, indicating the natural order of relative hardness of these problems.

I. INTRODUCTION

In this paper we consider the problem of data clustering under the union of subspaces (UOS) model [1], [2], when each data vector is sampled in an element-wise manner. This is referred to as the case of *missing data*. In other words we are looking to harvest a union of subspaces structure from the data, when the data is missing. Such a problem has been recently considered in a number of papers [3], [4], [5], [6]. This setting has implications to data completion under the union of subspaces model in contrast to the single subspace model that has been prevalent in the matrix completion literature. In contrast to statistical analysis in [3], [4], [5], this paper uses a variant of the sparse subspace clustering (SSC) algorithm [2] to give sufficient deterministic conditions for accurate subspace clustering under missing data. In contrast to [6], which does not provide any specific conditions for success of SSC under missing data, in this paper we provide implications of the deterministic conditions for several specific cases of sampling. Further through extensive simulations we demonstrate for the first time that accurate clustering under missing data does not imply accurate subspace clustering and completion thereby indicating the natural order of hardness of these problems under missing data.

II. PROBLEM SET-UP

We are given a set of data points collected as columns of a matrix \mathbf{X} , i.e. $\mathbf{X}_i, i = 1, 2, \dots, N$ such that, $\mathbf{X}_i \in \bigcup_{\ell=1}^L \mathbb{S}^{(\ell)}$, where $\mathbb{S}^{(\ell)}$ is a subspace of dimension d_ℓ in \mathbb{R}^d , for $\ell = 1, 2, \dots, L$. Further, each data point \mathbf{X}_i is sampled at Ω_i co-ordinates (randomly or deterministically) and the problem is to identify the subspaces $\mathbb{S}^{(\ell)}$ under missing data. In order to derive meaningful performance guarantees for the proposed

algorithm, we consider the following generative model for the data. Let $\mathbf{X}^{(\ell)}$ denote the set of vectors in \mathbf{X} which belong to subspace ℓ . Let

$$\mathbf{X}^{(\ell)} = \mathbf{U}^{(\ell)} \mathbf{A}^{(\ell)},$$

where the N_ℓ columns of $\mathbf{A}^{(\ell)}$ are drawn from the unit sphere $\mathcal{S}^{d_\ell-1}$ and $\mathbf{U}^{(\ell)}$ is a matrix with orthonormal columns, whose columns span the subspace, $\mathbb{S}^{(\ell)}$. Under missing data, point $\mathbf{X}_i^{(\ell)}$ is sampled at locations $\Omega_i^{(\ell)}$ and

$$\mathbf{X}_{\Omega_i}^{(\ell)} = \mathbf{I}_{\Omega_i^{(\ell)}} \mathbf{U}^{(\ell)} \mathbf{a}_i^{(\ell)} \quad (1)$$

where $\mathbf{I}_{\Omega_i^{(\ell)}}$ is a diagonal matrix with $\mathbf{I}_{\Omega_i^{(\ell)}}(k, k) = 1$ iff $k \in \Omega_i^{(\ell)}$. It is essentially a zero filled $\mathbf{X}_i^{(\ell)}$.

For the sake of exposition, in the following we will assume that $d_\ell = d$ for all $\ell = 1, 2, \dots, L$.

Given this set-up the problem is to accurately cluster the data points such that within each cluster the data points belong to the same (original) subspace.

A. The Algorithm: SSC-LP

Our algorithm as presented below is a minor variation of the SSC-EWZF algorithm in [6] in that instead of solving a Lasso problem we solve a linear program (LP) to estimate the coefficient matrix. The steps of the algorithm are as follows.

- 1) For each i solve for

$$\arg \min \|\mathbf{c}_i\|_1 : \mathbf{I}_{\Omega_i} \mathbf{X}_{\Omega_i} = \mathbf{I}_{\Omega_i} \mathbf{X}_{-i, \Omega} \mathbf{c}_i$$

where \mathbf{X}_{Ω_i} denotes the data point \mathbf{X}_i with zeros filling at non-sampled locations and $\mathbf{X}_{-i, \Omega}$ denotes the zero-filled data points except the i data point.

- 2) Collect the \mathbf{c}_i into a matrix \mathbf{C} and apply spectral clustering [7] to $\mathbf{A} = |\mathbf{C}| + |\mathbf{C}|^\top$

III. ANALYSIS OF THE ALGORITHM

A. Case I

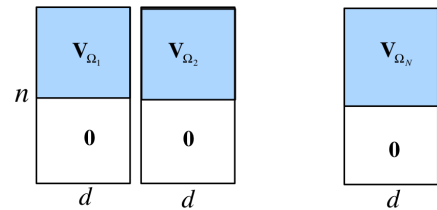


Fig. 1: Schematic depicting Case 1. All $\Omega_i = \Omega$.

Let $\Omega_i = \Omega$ for all i . In this case, all the points are sampled at the same locations. Let point i belong to subspace ℓ .

Let $\mathbf{V}_\Omega^{(\ell)} = \mathbf{I}_\Omega \mathbf{U}^{(\ell)} \in \mathbb{R}^{|\Omega| \times d}$ where \mathbf{I}_Ω denotes the projection onto the (non-zero) co-ordinates in Ω . The following theorem summarizes the conditions under which SSC-EWZF algorithm results in correct clustering.

Theorem III.1. *Let $\Omega_i^{(\ell)} = \Omega$ for all i, ℓ and $|\Omega| \geq d$. Then SSC-LP leads to correct clustering if for all $i \in [N_\ell]$, $k \neq \ell$, the following holds,*

$$\left| \frac{\lambda_i^{(\ell)\top} (\mathbf{V}_\Omega^{(\ell)})^\dagger \mathbf{V}_\Omega^{(k)} \mathbf{a}_j^{(k)}}{\|\lambda_i^{(\ell)}\|_2} \right| < r_{in}(\mathcal{P}(\mathbf{A}_{-i}^{(\ell)})), \quad (2)$$

where,

- 1) $\lambda_i^{(\ell)} \in \arg \max_{\lambda} \langle \mathbf{a}_i^{(\ell)}, \lambda \rangle : \|\mathbf{A}_{-i}^{(\ell)\top} \lambda\|_\infty \leq 1$
- 2) $\mathcal{P}(\mathbf{A}_{-i}^{(\ell)}) = \{\mathbf{v} : \mathbf{v} = \mathbf{A}_{-i}^{(\ell)} \mathbf{b} : \|\mathbf{b}\|_1 = 1\}$, and $r_{in}(\mathcal{P}(\mathbf{A}_{-i}^{(\ell)}))$ denotes the in-radius of the centrosymmetric body.

Proof: The proof directly follows from the proof of [Theorem 6 in [8]]. ■

Implications of the condition in Theorem III.1:

- 1) We note that the RHS in Equation (2) is the same as the RHS for the full observation case, see [1]. Therefore, the performance degradation in clustering comes from increase in the LHS on an average under missing data.
- 2) The increase in RHS depends on how badly conditioned $\mathbf{V}_\Omega^{(\ell)}$ is. If some of the singular values of $\mathbf{V}_\Omega^{(\ell)}$ are very small, it will make the LHS in Equation 2 large. To further understand this, let us consider a *semi-random* model where in each subspace the data points are generated by choosing $\mathbf{a}_i^{(\ell)}$ uniformly randomly on a unit sphere [1]. In this case it is easy to show that $\frac{\lambda_i^{(\ell)\top}}{\|\lambda_i^{(\ell)}\|_2}$ are also uniformly distributed on the unit sphere [8] and the *expected value* of the LHS becomes $\frac{\|(\mathbf{V}_\Omega^{(\ell)})^\dagger \mathbf{V}_\Omega^{(k)}\|_F}{d}$. This is essentially the (unnormalized) co-ordinate restricted coherence between subspaces ℓ and k .
- 3) If the subspaces are sufficiently incoherent with the standard basis or satisfy an RIP like property for large enough $|\Omega|$, then the condition number of $\mathbf{V}_\Omega^{(\ell)}$ is controlled and one can expect to obtain similar performance as the full observation case.
- 4) Note that *while in this setting one can ensure perfect clustering, one cannot ensure either perfect subspace identification or completion*. This is because in this case the deterministic necessary conditions for identification and completion [9] are not satisfied. This indicates that clustering is an easier problem compared to subspace identification and completion under missing data.

Here we would like to mention that the notion of in-radius is related to the notion of **permeance** [10] of data points in a given subspace, quantifying how well the data is distributed inside each subspace. In-radius can be thought of as a worst-case permeance that doesn't scale with the number of data points, while permeance scales with the number of data points and is more of an averaged criteria. Perhaps this is the reason that the primal-dual analysis of SSC under full observation is

not able to support the empirical evidence that as the number of points per subspace increases the clustering error goes down dramatically. For subspace clustering such the effect of the number of data points was shown more explicitly in a recent paper [11]. A connection between these two quantities, namely the in-radius and permeance for subspace clustering under missing data will be undertaken in a future work.

We will analyze the more general case next. Analysis of the general case is hard and to gain *insights* we break it up into several results. For this we introduce the following notation.

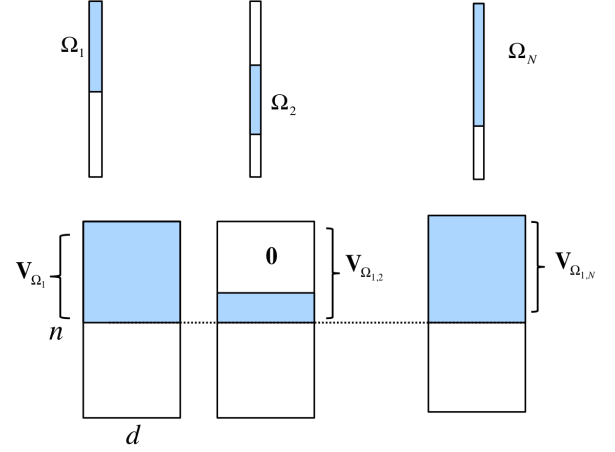


Fig. 2: Schematic depicting the case when all data is sampled at different locations and the notation for the basis $\mathbf{V}_{\Omega_{i,j}}$ zero-filled and restricted to Ω_i . Here we show for the case when $i = 1$.

Notation:

- 1) Let $\mathbf{V}_{\Omega_i^{(\ell)}} = \mathbf{I}_{\Omega_i^{(\ell)}} \mathbf{U}^{(\ell)}$. For other data points, \mathbf{X}_j if \mathbf{X}_j belongs to subspace $k \in \{1, 2, \dots, L\}$, $\mathbf{V}_{\Omega_{i,j}^{(k)}} = \mathbf{I}_{\Omega_i^{(\ell)}} \mathbf{U}^{(k)}$ if $\Omega_i^{(\ell)} \subseteq \Omega_j^{(k)}$, else $\mathbf{V}_{\Omega_{i,j}^{(k)}} = \mathbf{I}_{\Omega_i^{(\ell)} \cap \Omega_j^{(k)}} \mathbf{U}^{(k)}$, where $\mathbf{I}_{\Omega_i^{(\ell)} \cap \Omega_j^{(k)}}$ is an $|\Omega_i^{(\ell)}| \times |\Omega_j^{(k)}|$ diagonal matrix with diagonal entries 1 at locations in $\mathbf{I}_{\Omega_i^{(\ell)} \cap \Omega_j^{(k)}}$ and 0 otherwise. This is depicted in Figure, 2.
- 2) Let $\mathbf{V}_{\Omega_i^{(\ell)}} = \mathbf{Q}_i^{(\ell)} \Sigma_i^{(\ell)} \mathbf{R}_i^{(\ell)\top}$ and $\mathbf{V}_{\Omega_{i,j}^{(\ell)}} = \mathbf{Q}_{i,j}^{(\ell)} \Sigma_{i,j}^{(\ell)} \mathbf{R}_{i,j}^{(\ell)\top}$ be the SVDs of the truncated basis matrices. We will refer to $\mathbf{Q}_i^{(\ell)}, \mathbf{Q}_{i,j}^{(\ell)}$ as corresponding to *co-ordinate restricted subspaces*.
- 3) Let

$$\tilde{\mathbf{a}}_j^{(\ell)} = \mathbf{Q}_i^{(\ell)\top} \mathbf{V}_{\Omega_{i,j}^{(\ell)}} \mathbf{a}_j^{(\ell)} \quad \forall j \quad (3)$$

$$= \mathbf{Q}_i^{(\ell)\top} \mathbf{Q}_j^{(\ell)} \underbrace{\Sigma_{i,j}^{(\ell)} \mathbf{R}_{i,j}^{(\ell)\top}}_{\tilde{\mathbf{a}}_{i,j}^{(\ell)}} \mathbf{a}_j^{(\ell)} \quad (4)$$

and let $\tilde{\mathbf{A}}_{-i}^{(\ell)}$ be the $r \times (N_\ell - 1)$ matrix with columns as $\tilde{\mathbf{a}}_j^{(\ell)}, j \neq i$.

B. Case II

First we consider the case when $|\Omega_i| = d$, i.e. we want to understand, if a given point is sampled at only d locations,

when can SSC-LP correctly assign it to the cluster of points from the same subspace. Without loss of generality let the point come from the subspace ℓ .

Assumption: For all possible supports $\Omega_i^{(\ell)}$, we assume that $\mathbf{Q}_i^{(\ell)} \in \mathbb{R}^{d \times d}$ is invertible. This assumption will be satisfied if the columns of $\mathbf{U}^{(\ell)}$ are sufficiently incoherent with the standard basis. Under this assumption, we have the following theorem.

Theorem III.2. *Let $|\Omega_i^{(\ell)}| = d$. Then SSC-LP correctly assigns it to the data points coming from subspace ℓ , if for all $k \neq \ell$ and j , the following holds,*

$$\left| \frac{\lambda_i^{(\ell)\top} \mathbf{Q}_i^{(\ell)\top} \mathbf{V}_{\Omega_{i,j}^{(k)}} \mathbf{a}_j^{(k)}}{\|\lambda_i^{(\ell)}\|_2} \right| < r_{in}(\mathcal{P}(\tilde{\mathbf{A}}_{-i}^{(\ell)})), \quad (5)$$

where,

- 1) $\lambda_i^{(\ell)} \in \arg \max_{\lambda} \langle \tilde{\mathbf{a}}_i^{(\ell)}, \lambda \rangle : \|\tilde{\mathbf{A}}_{-i}^{(\ell)\top} \lambda\|_{\infty} \leq 1$
- 2) $\mathcal{P}(\tilde{\mathbf{A}}_{-i}^{(\ell)}) = \{\mathbf{v} : \mathbf{v} = \tilde{\mathbf{A}}_{-i}^{(\ell)} \mathbf{b} : \|\mathbf{b}\|_1 = 1\}$, and $r_{in}(\mathcal{P}(\tilde{\mathbf{A}}_{-i}^{(\ell)}))$ denotes the in-radius of the centrosymmetric body.

Proof: The proof follows from slight modifications of the arguments of the proof of Theorem III.1, which in turn follows from the proof of [Theorem 6 in [8]]. ■

Implications of condition in Theorem III.2:

- The in-radius $r(\mathcal{P}^\circ(\tilde{\mathbf{A}}_{-i}^{(\ell)}))$ now depends on the sampling pattern within the subspace and it is different for different set of points. Specifically note that it depends on the *intrinsic subspace coherence* between the sampled points within the same subspace. In general it can be very small if the number of points $\mathbf{X}_j^{(\ell)}$ (within the same subspace) that have good amount of overlap with $\mathbf{X}_i^{(\ell)}$ is small.
- The subspace coherence also depends on the *relative* sampling pattern of points in other subspaces with respect to the point under consideration. This is reflected in the coherence term $\mathbf{Q}_i^{(\ell)\top} \mathbf{V}_{\Omega_{i,j}^{(k)}}$. Unlike Case I, for the semi-random model the distribution of the normalized dual direction is not uniform and we cannot comment on the average case performance. However, a worst-case condition can be obtained here from Equation 5 which says that if

$$\|\mathbf{Q}_i^{(\ell)\top} \mathbf{V}_{\Omega_{i,j}^{(k)}}\|_2 \leq r_{in}(\mathcal{P}^\circ(\tilde{\mathbf{A}}_{-i}^{(\ell)})) \quad (6)$$

then SSC-LP leads to correct clustering. This ofcourse requires that the co-ordinate projected subspaces be sufficiently incoherent as well as *disjoint*.

C. Case III

With these **insights** let us now proceed to handle the more general case when $|\Omega_i^{(\ell)}| > d$. In order to do that let us first provide a geometric result using Lemma VI.1 (see Appendix) and then find a useful characterization of the dual certificate.

Again let $\mathbf{V}_{\Omega_i^{(\ell)}} = \mathbf{Q}_i^{(\ell)} \Sigma_i^{(\ell)} \mathbf{R}_i^{(\ell)\top}$ be the SVD of $\mathbf{V}_{\Omega_i^{(\ell)}}$. Let $\mathbf{X}_{\Omega_i^{(\ell)},j}$ denote the zero filled $\mathbf{I}_{\Omega_j} \mathbf{X}_j$ restricted to the indices

in $\Omega_i^{(\ell)}$ for all $j \in [N]$. Let $\mathbf{X}_{-i,\Omega_i^{(\ell)}}$ denote the matrix of these points except point i .

Then SSC-LP solves $\arg \min \|\mathbf{c}_i\|_1 : \mathbf{X}_{i,\Omega_i^{(\ell)}} = \mathbf{X}_{-i,\Omega_i^{(\ell)}} \mathbf{c}_i$. Let $\mathbf{c}_i^{(\ell)}$ and $\nu_i^{(\ell)}$ be solutions to,

$$\arg \min \|\mathbf{b}\|_1 : \mathbf{X}_{i,\Omega_i^{(\ell)}}^{(\ell)} = \mathbf{X}_{-i,\Omega_i^{(\ell)}}^{(\ell)} \mathbf{b}$$

$$\arg \max \langle \mathbf{X}_{i,\Omega_i^{(\ell)}}^{(\ell)}, \nu \rangle : \|\mathbf{X}_{-i,\Omega_i^{(\ell)}}^{(\ell)\top} \nu\|_{\infty} \leq 1$$

Assuming that *strong duality* holds, the vectors

$$\mathbf{c} = [\mathbf{0}, \dots, \mathbf{0}, \mathbf{c}_i^{(\ell)}, \dots, \mathbf{0}]^\top$$

and $\nu_i^{(\ell)}$ satisfy the conditions of Lemma VI.1 (See Appendix), hence guaranteeing correct subspace recovery if for all $k \neq \ell, j$, the following condition is satisfied.

$$|\nu_i^{(\ell)\top} \mathbf{X}_{\Omega_i^{(\ell)},j}^{(k)}| < 1 \quad (7)$$

As such this condition does not provide any further intuition or insights. In order to arrive at results that are similar in flavor to Theorems III.1 and III.2, let us further analyze this condition. We have the following **Assumption:** The matrix $\mathbf{Q}_i^{(\ell)}$ has full column rank. This assumption is satisfied if columns of $\mathbf{U}^{(\ell)}$ are sufficiently incoherent with respect to the co-ordinate basis.

Let $\mathbf{Q}_i^{(\ell)} \mathbf{Q}_i^{(\ell)\top}$ denote the projection matrix that projects a vector onto the subspace spanned by columns of $\mathbf{Q}_i^{(\ell)}$. $\mathbf{Q}_i^{(\ell)\top}$ and $(\mathbf{I} - \mathbf{Q}_i^{(\ell)} \mathbf{Q}_i^{(\ell)\top})$ denote the null space of $\mathbf{Q}_i^{(\ell)\top}$. Note that,

$$\nu_i^{(\ell)} = \mathbf{Q}_i^{(\ell)} \mathbf{Q}_i^{(\ell)\top} \nu_i^{(\ell)} + (\mathbf{I} - \mathbf{Q}_i^{(\ell)} \mathbf{Q}_i^{(\ell)\top}) \nu_i^{(\ell)}.$$

Then the condition for correct clustering becomes,

$$|(\mathbf{Q}_i^{(\ell)} \mathbf{Q}_i^{(\ell)\top} \nu_i^{(\ell)} + (\mathbf{I} - \mathbf{Q}_i^{(\ell)} \mathbf{Q}_i^{(\ell)\top}) \nu_i^{(\ell)})^\top \mathbf{X}_{\Omega_i^{(\ell)},j}^{(k)}| \leq 1 \quad (8)$$

By triangle inequality if for some $0 \leq \alpha \leq 1$,

$$|(\nu_i^{(\ell)})^\top \mathbf{Q} \mathbf{Q}^\top \mathbf{X}_{\Omega_i^{(\ell)},j}^{(k)}| < \alpha \quad (9)$$

and

$$|(\nu_i^{(\ell)})^\top (\mathbf{I} - \mathbf{Q} \mathbf{Q}^\top) \mathbf{X}_{\Omega_i^{(\ell)},j}^{(k)}| \leq (1 - \alpha) \quad (10)$$

then the point $\mathbf{X}_i^{(\ell)}$ is correctly clustered. Based on this we have the following theorem.

Theorem III.3. *If for any point \mathbf{X}_i belonging to subspace ℓ , the following conditions are satisfied,*

$$\|\mathbf{Q}_i^{(\ell)\top} \mathbf{V}_{\Omega_{i,j}^{(k)}}\|_2 < \alpha r_{in} \left(\mathcal{P}(\mathbf{X}_{-i,\Omega_i^{(\ell)}}^{(\ell)}) \right) \quad (11)$$

$$\|(\mathbf{I} - \mathbf{Q}_i^{(\ell)} \mathbf{Q}_i^{(\ell)\top}) \mathbf{V}_{\Omega_{i,j}^{(k)}}\|_2 \leq (1 - \alpha) r_{in} \left(\mathcal{P}(\mathbf{X}_{-i,\Omega_i^{(\ell)}}^{(\ell)}) \right) \quad (12)$$

for some $0 \leq \alpha \leq 1$, then SSC-LP succeeds. Here $\mathcal{P}(\mathbf{X}_{-i,\Omega_i^{(\ell)}}^{(\ell)}) = \{\mathbf{v} : \mathbf{v} = \mathbf{X}_{-i,\Omega_i^{(\ell)}}^{(\ell)} \mathbf{b} : \|\mathbf{b}\|_1 = 1\}$, and $r_{in}(\mathcal{P}(\mathbf{X}_{-i,\Omega_i^{(\ell)}}^{(\ell)}))$ denotes the in-radius of the centrosymmetric body

We note the following points.

- 1) Compared to Theorems III.1 and III.2, the conditions stated in Theorem III.3 are worst-case and slightly weaker in the sense that they require the co-ordinate projected subspaces to be disjoint.
- 2) Also note that for the most general case it is required that *co-ordinate restricted subspace coherence* be small BUT ALSO that the *projection error of the other set of points* onto the co-ordinate restricted subspace for the point under consideration be small. When $|\Omega_i^{(\ell)}|$ is large, Equation 11 is more likely to be satisfied compared to Equation 6.
- 3) In the most general case, while more observations (per vector) is required, in this case under correct clustering if the sampling patterns and the number of data points satisfy the necessary and sufficient conditions in [9], [12] then it is possible to also correctly identify the subspaces and also ensure completion

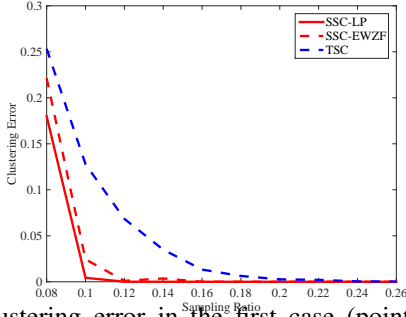


Fig. 3: Clustering error in the first case (points sampled at same pn co-ordinates) for varying p .

IV. NUMERICAL RESULTS

In this section, we will see the numerical performance of the proposed algorithm, SSC-LP, as compared to two other baseline algorithms for data clustering and completion under the UoS model with missing data. The data is generated by using L rank d subspaces, each composed of N_l vectors of dimension n . We assume $n = 50$, $L = 3$, $d = 3$, and $N_l = 150$ for the numerical results. Data in each subspace is generated by a multiplication of standard entry-wise Gaussian distributed $n \times d$ matrix and a standard entry-wise Gaussian distributed $d \times N_l$ matrix.

We compare our algorithm with two other algorithms. The first is SSC-EWZF [6], which solves the Lasso problem rather than the linear program in this paper. This algorithm is selected as it gives the lowest clustering error among the different algorithms considered in [6] (e.g., SSC-EC, SSC-CEC, ZF-SSC, MC-SSC). The performance of SSC-EWZF algorithm largely depends on the choice of λ , which is chosen to be

$$\lambda = \frac{\alpha}{\max_{i \neq j} |X_{\Omega_i}^\top X_{\Omega_j}|_{ij}}, \quad (13)$$

where α is a tuning parameter [6], set to be 7.34 for our experiment (selected by performing an optimized performance for different values of α). The second is TSC algorithm proposed in [13], which builds adjacency matrix by thresholding correlations. The thresholding parameter q is given as

$$q = \sqrt{N_l \log(N_l)}, \quad (14)$$

such that q is of an order smaller than N_l and larger than $\log N_l$ [13]. Since our proposed algorithm considers only observed entries when building adjacency matrix for each data, TSC algorithm could be a good comparison in respect to the way of building adjacency matrix.

In our simulations, we consider two cases. The first case illustrates Case 1 in this paper, where all the all the points are sampled at the same co-ordinates which are the first pn co-ordinates. The second case corresponds to Case 3 in this paper, where missing data in each point is randomly sampled at a rounded value of pn co-ordinates thus giving a sampling rate of approximately p and a missing rate of approximately $1 - p$. All the results are averaged over 100 runs for the choice of the data and the sampled elements.

The comparisons for data clustering and completion are performed using three metrics as explained further. The first metric is the clustering error. Clustering error is the ratio calculated by the number of wrongly classified data divided by the number of total data, same as defined in [6]. The clustering error for different sampling ratio p in the first case is shown in Fig. 3, where SSC-LP shows the best clustering accuracy for all sampling ratios from 8% to 26%. Furthermore, the plot indicates that SSC-LP clusters data perfectly with $p = 0.1$, equivalently at 5 observations out of 50 entries, while SSC-EWZF and TSC require 12 and 11 observations, respectively. Since the rank of each cluster is 3 and only 5 observations for each point are needed, we see that observing data at same co-ordinates need much less data for efficient clustering. We note that we cannot identify the subspace or complete the data with these observations further illustrating that clustering requires less number of observations than that required for data completion.

The clustering error for random sampling, the scenario described in the second case, with sampling ratio from 0.25 to 0.95 is shown in Fig 4(a), where we note that clustering error with our proposed algorithm, SSC-LP, is the minimum among the three algorithms. Furthermore, the plot shows that the sampling ratio at which the clustering error hits zero for the algorithms SSC-LP, SSC-EWZF and TSC are 0.38, 0.42, and 0.54, respectively. Thus, the proposed algorithm required least number of observations to efficiently cluster the data. We further note that the amount of data needed to cluster efficiently for random sampling (38%) is larger than that for observing data at the same co-ordinates (10%).

The second metric is the completion error. Let the recovered matrix using a clustering algorithm be the output of matrix completion using SVT method [14] on the subspaces found as a result of the subspace clustering and the true matrix be the ground truth of the matrix with missing data. Then the recovery difference is defined as the matrix difference between recovered matrix and true matrix. Thus the completion error is measured by ratio of the Frobenius Norm of the recovery difference to the Frobenius Norm of the true matrix. The completion error for different values of p can be seen in Fig 4(b), where we see completion error is positively correlated with clustering error and small percentage clustering error can result in large percentage completion error. Similar to the clustering error, SSC-LP has the lowest completion error

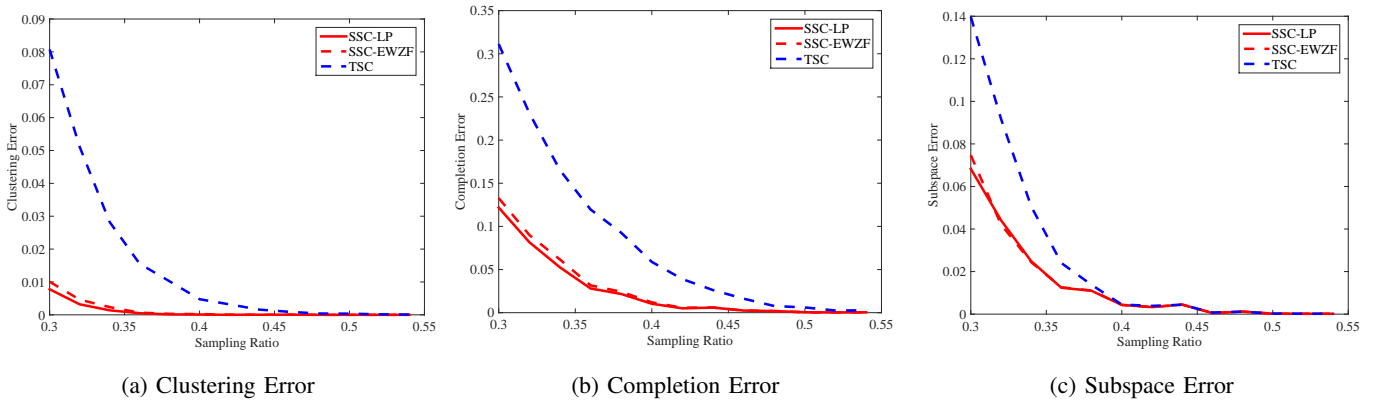


Fig. 4: Different error metrics for the case of random sampling.

among the three algorithms and the completion errors for SSC-LP, SSC-EWZF and TSC becomes zero at sampling ratio of around 0.50, 0.50, and 0.54 respectively, which are larger than the corresponding thresholds for the clustering errors.

The third metric is the subspace error. Subspace of a matrix with missing data can be recovered by finding the orthonormal basis of the recovered matrix. Projection difference is calculated by minus the projection of orthonormal basis of the recovered matrix from the orthonormal basis of the true matrix, and consequently, angle based subspace error is measured by arcsin of the normal of the projection difference, same as defined in [15]. Mathematically, subspace error is expressed as:

$$\theta = \arcsin(\|(\mathbf{B} - \mathbf{A}\mathbf{A}^\top\mathbf{B})\|_2)$$

where θ is the angle based subspace error, \mathbf{A} is the orthonormal basis of the completed matrix and \mathbf{B} is the orthonormal basis of the true matrix. With the simulation result in Fig 4(c), that subspace errors for all three algorithms start at 46% sampling ratio. For any sampling ratio lower than the threshold, subspace error for our proposed algorithm is the smallest among the three algorithms.

V. CONCLUSIONS

This paper proposes an algorithm for sparse subspace clustering under missing data, when data is assumed to come from a Union of Subspaces (UoS) model, using a linear programming method, SSC-LP. Deterministic analysis of sufficient conditions when this algorithm leads to correct clustering is presented. Extensive set of simulation results for clustering as well as completion of data under missing entries, under the UoS model are provided which demonstrate the effectiveness of the algorithm, and demonstrate that accurate clustering does not imply accurate subspace identification.

In future we will derive the performance bounds in terms of average cases analysis of the deterministic conditions. In particular we will determine how the in-radius changes under element-wise sampling and how does the un-normalized subspace coherence behaves under these models.

VI. APPENDIX

The following Lemma appears in [1].

Lemma VI.1. Let $\mathbf{A} \in \mathbb{R}^{n \times N}$ and $\mathbf{y} \in \mathbb{R}^n$ be given. If there exists \mathbf{c} obeying $\mathbf{y} = \mathbf{A}\mathbf{c}$ with support $S \subset T$ and a dual certificate $\boldsymbol{\nu}$ satisfying

$$\mathbf{A}_S^\top \boldsymbol{\nu} = \text{sign}(\mathbf{c}_S), \|\mathbf{A}_{T \cap S^c}^\top \boldsymbol{\nu}\|_\infty \leq 1, \|\mathbf{A}_{T^c}^\top \boldsymbol{\nu}\|_\infty < 1$$

then all optimal solutions \mathbf{c}^* to $\arg \min_{\mathbf{c}} \|\mathbf{c}\|_1 : \mathbf{A}\mathbf{y} = \mathbf{c}$ satisfy $\mathbf{c}_{T^c}^* = \mathbf{0}$.

REFERENCES

- [1] Mahdi Soltanolkotabi and Emmanuel J. Candes, "A geometric analysis of subspace clustering with outliers," *Ann. Statist.*, vol. 40, no. 4, pp. 2195–2238, 08 2012.
- [2] E Elhamifar and R Vidal, "Sparse Subspace Clustering: Algorithm, Theory, and Applications," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 35, no. 11, pp. 2765–2781, 2013.
- [3] L. Balzano, A. Szlam, B. Recht, and R. Nowak, "K-subspaces with missing data," in *Statistical Signal Processing Workshop (SSP), 2012 IEEE*, Aug 2012, pp. 612–615.
- [4] D. Pimentel, R. Nowak, and L. Balzano, "On the sample complexity of subspace clustering with missing data," in *Statistical Signal Processing (SSP), 2014 IEEE Workshop on*, June 2014, pp. 280–283.
- [5] Brian Eriksson, Laura Balzano, and Robert D. Nowak, "High-rank matrix completion and subspace clustering with missing data," *CoRR*, vol. abs/1112.5629, 2011.
- [6] Congyuan Yang, Daniel Robinson, and Rene Vidal, "Sparse subspace clustering with missing entries," in *Proceedings of The 32nd International Conference on Machine Learning*, 2015, pp. 2463?–2472.
- [7] Ulrike Luxburg, "A tutorial on spectral clustering," vol. 17, no. 4, pp. 395–416, 2007.
- [8] Reinhard Heckel, Michael Tschannen, and Helmut Bölcskei, "Dimensionality-reduced subspace clustering," *Journal of Machine Learning Research*, July 2015.
- [9] D. L. Pimentel-Alarcón, R. D. Nowak, and N. Boston, "Deterministic Conditions for Subspace Identifiability from Incomplete Sampling," *ArXiv e-prints*, Oct. 2014.
- [10] Gilad Lerman, Michael B. McCoy, Joel A. Tropp, and Teng Zhang, "Robust computation of linear models, or how to find a needle in a haystack," *CoRR*, vol. abs/1202.4044, 2012.
- [11] D. Park, C. Caramanis, and S. Sanghavi, "Greedy Subspace Clustering," *ArXiv e-prints*, Oct. 2014.
- [12] D. L. Pimentel-Alarcón, N. Boston, and R. D. Nowak, "A Characterization of Deterministic Sampling Patterns for Low-Rank Matrix Completion," *ArXiv e-prints*, Mar. 2015.
- [13] Reinhard Heckel and Helmut Bölcskei, "Robust subspace clustering via thresholding," *IEEE Transactions on Information Theory*, 2015.
- [14] Emmanuel J Candès and Benjamin Recht, "Exact matrix completion via convex optimization," *Foundations of Computational mathematics*, vol. 9, no. 6, pp. 717–772, 2009.
- [15] Ake Björck and Gene H Golub, "Numerical methods for computing angles between linear subspaces," *Mathematics of computation*, vol. 27, no. 123, pp. 579–594, 1973.